

Bounds for the normal approximation of the maximum likelihood estimator



Andreas Anastasiou

Supervisor: Professor Gesine Reinert

Department of Statistics

University of Oxford

This dissertation is submitted for the degree of

Doctor of Philosophy

To my parents, Kyriacos and Maria, and my siblings, Anna and Tassos.

Acknowledgements

The first and most important person to thank is my supervisor, Professor Gesine Reinert. I am deeply grateful to her for the guidance, encouragement and support she offered me throughout my time as her student.

I would like to thank Professor Christophe Ley for a fruitful collaboration which lead to a paper, of which the results can be found in Chapter 5 of the thesis.

I would also like to thank Dr Robert Gaunt for various insightful comments and suggestions.

Last, but not of least importance, I would like to thank my parents who encouraged me to continue with postgraduate studies.

The research is conducted with the support of a Teaching Assistantship Bursary from the Department of Statistics, University of Oxford, and the Engineering and Physical Sciences Research Council (EPSRC) grant EP/K503113/1.

Abstract

The asymptotic normality of the maximum likelihood estimator (MLE) under regularity conditions is a long established and famous result. This is a qualitative result and the assessment of such a normal approximation is our main interest. For this task we partly use Stein's method, which is a probabilistic technique that can be used to explicitly measure the distributional distance between two distributions. Since its first appearance in 1972, the method has been developed for various distributions; here we use the results related to Stein's method for normal approximation.

In this thesis, we derive explicit upper bounds on the distributional distance between the distribution of the MLE and the normal distribution. First, the focus is on independent and identically distributed random variables from both discrete and continuous single-parameter distributions with particular attention to exponential families. For discrete distributions, the case where the MLE can be on the boundary of the parameter space is treated through a perturbation approach, which allows us to obtain bounds on the distributional distance of interest. The bounds are of order $\frac{1}{\sqrt{n}}$, where n is the number of observations. Simulation-based results are given to illustrate the power of the bound. Furthermore, often the MLE can not be obtained analytically and optimisation methods (such as the Newton-Raphson algorithm) are used. Even in such cases, order $\frac{1}{\sqrt{n}}$ bounds are given for the distributional distance related to the MLE.

The case of multi-parameter distributions follows smoothly after the detailed discussion

related to a scalar parameter. Apart from extending our approach to a multi-parameter setting, we also cover the case of independent but not necessarily identically distributed (i.n.i.d.) random vectors with specific focus on the widely applicable linear regression models.

Going back to the single-parameter setting a different approach to get an upper bound on the distributional distance between the distribution of the MLE and the normal distribution, based on the Delta method, is also developed. The MLE for a Generalised Gamma distribution gives an illustration of the results obtained through this Delta method approach.

Finally, we relax the independence assumption and results for the case of locally dependent random variables are obtained. An example of correlated sums of normally distributed random variables illustrates the bounds. Again, results that do not require an analytic expression of the MLE to be known are given. We end this thesis with ideas currently in progress and further open research questions.

Table of contents

List of figures	ix
List of tables	x
1 Introduction	1
1.1 Motivation for the thesis	1
1.2 Stein's method as a useful tool	3
1.3 Outline of the thesis	5
2 Background	8
2.1 Central Limit theorem	8
2.2 Maximum likelihood estimation	10
2.3 Stein's method	15
3 Bounds on the distance to normal for the MLE	24
3.1 The bounded Wasserstein and the Kolmogorov distance	25
3.2 Bounds in terms of the bounded Wasserstein distance	27
3.3 Single-parameter exponential families	36
3.3.1 The exponential distribution	40
3.3.2 Empirical results	47
3.4 Discrete distributions: The boundary issue	51
3.4.1 The perturbation approach	52

3.4.2	Example: The Poisson distribution	60
3.5	Bounds on the mean squared error	68
4	Multi-parameter distributions	79
4.1	Non-identically distributed random vectors	81
4.1.1	A general bound	82
4.1.2	Linear regression	97
4.2	Identically distributed random vectors	102
4.2.1	The normal distribution	104
4.2.2	Bounds when the MLE is not known explicitly	114
5	An approach using the Delta method	133
5.1	New bounds on the distance to normal for the MLE	134
5.2	Comparison in exponential families	143
5.2.1	Bounds for single-parameter exponential families	143
5.2.2	Bounds for the exponential distribution	149
6	Locally dependent random variables	151
6.1	A general approach	152
6.2	Locally dependent normal random variables	160
6.3	An alternative bound	170
6.4	Results for not analytically known MLE	174
6.4.1	Bounded support	175
6.4.2	Bounded parameter space	179
7	Conclusion	183
7.1	Summary	183
7.2	Open problems for future exploration	186
7.2.1	Further applications	186
7.2.2	Beyond the classical framework of the MLE	189

Table of contents	viii
References	192

List of figures

3.1	The density function of $\sqrt{ni(\theta_0)}(\hat{\theta}_n(\mathbf{X}) - \theta_0)$ for various sample sizes. In this case X_1, X_2, \dots, X_n are i.i.d. random variables which follow the $\text{Exp}(1)$ distribution.	49
3.2	The density function of $\sqrt{ni(\theta_0)}(\hat{\theta}_n(\mathbf{X}) - \theta_0)$ for various sample sizes. In this case X_1, X_2, \dots, X_n are i.i.d. random variables which follow the $\text{Exp}(0.5)$ distribution.	50
3.3	The Q-Q plot, when $n = 1000$, and the bar plot of the error for different sample sizes for the exponential distribution example, $\text{Exp}(1)$	50
3.4	The Q-Q plot, when $n = 1000$, and the bar plot of the error for different sample sizes for the exponential distribution example, $\text{Exp}(0.5)$	51
3.5	The density function of $\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0)$ for various sample sizes. In this case X_1, X_2, \dots, X_n are i.i.d. random variables which follow the Poisson(2) distribution.	67
3.6	The Q-Q plot, when $n = 1000$, and the bar plot of the error for different sample sizes for the Poisson distribution with mean equal to 2.	68
6.1	Structure of a $2m$ -dependent sequence.	153

List of tables

3.1	Simulation results from the $\text{Exp}(1)$ distribution	47
3.2	Simulation results from the $\text{Exp}(0.5)$ distribution treated as a non-canonical exponential family	48
3.3	Simulation results from the $\text{Poisson}(2)$ distribution	66
3.4	Simulation results from the $\text{Beta}(1.5, 1)$ distribution	78

Chapter 1

Introduction

1.1 Motivation for the thesis

Many statistical procedures rely on asymptotic results as exact finite-sample results are often difficult to derive. Yet a sample size is never infinite. In recent years, there has been a growing interest in assessing the quality of such approximations. The Central Limit Theorem (CLT), which is considered as one of the best known asymptotic developments in the area of probability and statistics, states that the sum of a large number of independent random variables, all of which have the same mean and variance, or in many alternative cases, random variables with specific types of dependence, will tend to be approximately normally distributed regardless of the random variables' exact distribution. The result can be extended to multidimensional random vectors.

Using the CLT, the Weak Law of Large Numbers and Taylor expansions, one is able to show that the asymptotic distribution of the maximum likelihood estimator (MLE) is normal under sufficient regularity conditions. This is an asymptotic result and the quality of the approximation may depend on the sample size as well as on the distribution of the random variables. The assessment of the quality of the normal approximation

related to the distribution of the MLE has been the motivation behind the research presented in this thesis.

The asymptotic normality of the MLE was first discussed in Fisher (1925). It is a fundamental qualitative result. Our aim is to find upper bounds on the distributional distance between the distribution of the MLE and the normal distribution under various general frameworks. We cover a variety of different scenarios related to the dimensionality of the MLE and also the dependence structure of possibly heterogeneous random variables (or vectors). Complementing the qualitative asymptotic result with a quantitative statement obtained from our bounds can be helpful to assess whether using the limiting distribution of the MLE is an acceptable approximation or not. From the opposite point of view, the results presented in this thesis can save both money and time by giving a good indication on whether a larger sample size is indeed necessary, for a good approximation to hold. These quantitative results for the MLE which are developed and explained in this thesis have already influenced researchers; see for example Pinelis and Molzon (2016).

The wide applicability of the maximum likelihood estimation method adds to the importance of our results. Among others, an MLE is used in ordinary and generalised linear models, time series analysis and a large number of other situations related to hypothesis testing and confidence intervals. They appear in a broad category of different fields, such as econometrics, computational biology and data modelling in physics and psychology. In addition to the usefulness of the MLE, in many cases its existence is easily ensured. Also, apart from the attractive asymptotic distribution property, the MLE satisfies many other exact sample size properties such as the functional invariance; if $\hat{\theta}$ is the MLE for θ , then the MLE for $g(\theta)$ is $g(\hat{\theta})$, with $g(\theta)$ being any transformation of the parameter. The numerous and useful properties of the MLE justify the fact that it is considered to be the most powerful and thus widely used estimation method.

1.2 Stein's method as a useful tool

For our results, we partly employ Stein's method, which was first introduced in 1972 by Charles Stein. It is a quite robust probabilistic technique that can measure the error in the approximation of one known distribution by another in a wide range of metrics. For the scope of this thesis, Stein's method for normal approximation is applied to a random variable

$$W = \frac{1}{\sqrt{n}} \sum_{i=1}^n X_i,$$

where $X_i, i = 1, 2, \dots, n$ are mean zero random variables such that $\text{Var}(W) = 1$. We will obtain upper bounds on the quantity

$$|\mathbb{E}[h(W)] - \mathbb{E}[h(Z)]|, \quad (1.1)$$

where Z follows the standard normal distribution and h belongs to a class of test functions, H . Depending on H , we get results in different metrics. Throughout this thesis, we mainly restrict ourselves to smooth test functions. Our approach is related to the work from Berry and Esséen for the known traditional convergence in distribution (Kolmogorov distance) result, which is of the form

$$\sup_{z \in \mathbb{R}} |\mathbb{P}(W \leq z) - \mathbb{P}(Z \leq z)| \leq C \frac{\mathbb{E}|X_1|^3}{\sqrt{n}}. \quad (1.2)$$

Esséen's original bound on C was 7.59, but has been decreased over the years and the latest improvement is $C \leq 0.4748$ as shown by Shevtsova (2014) using Fourier techniques. To see the result in (1.2) from Stein's method perspective, just take as test functions the class of indicator functions of half-lines $\{\mathbb{1}_{[\cdot, \leq x]}, x \in \mathbb{R}\}$. By doing that, the expectations in (1.1) become cumulative distribution functions leading to the result in (1.2).

The Stein's method results are used here to give upper bounds on the distributional distance related to a sum of random variables. Since the MLE is not in general a sum, we find a way to represent the standardised MLE such that it contains a quantity which is a sum, plus a term which we can control. For this purpose we use second-order Taylor expansions of the first derivative of the log-likelihood function evaluated at the MLE, about the true value of the parameter. This approach gives as sum component the first derivative of the log-likelihood function, namely the score function, evaluated at the true value of the unknown parameter. We need to take the dependence structure of the random variables into account and use the appropriate results from Stein's method to bound the score function. The remainder terms get bounded using other techniques based mainly on Taylor expansions, conditional expectations and known probability inequalities, such as the Markov and Cauchy-Schwarz inequalities. This is the main approach followed in this thesis. However, another idea based on the widely known Delta method is also developed and used to give bounds on the required distributional distance between the exact distribution of the MLE and its limiting normal distribution in cases where the MLE can be expressed as a function of a sum of random variables. Stein's method and Taylor expansions are still our main tools.

It is reasonable to compare Stein's method over other known techniques that can give answers on distributional distance, such as moment generating or characteristic functions. Although more details are given later, perhaps the biggest advantage of Stein's method compared to these techniques is its appropriateness on various types of dependence structure. Another advantage of Stein's method is the flexibility to manipulate the distribution using the random quantities of which the quantity under interest is composed. In contrast, characteristic function methods are based on manipulating distributions through their transforms, leading to loss of probabilistic intuition. The immediate explicit upper bound on the distributional distance is another advantage (Reinert, 1998).

1.3 Outline of the thesis

In Chapter 2, some background information for the CLT, the MLE and Stein's method is provided. The notation used throughout the thesis is explained and assumptions under which existence and uniqueness of the MLE is secured are also given; these are split into two cases depending on whether the parameter space is open and connected, or compact. We then explain Fisher's original asymptotic result for the distribution of the MLE, under its most basic form of independent and identically distributed (i.i.d.) random variables. The sufficient regularity conditions under which this asymptotic result holds are also discussed and briefly explained.

The work presented in Chapter 3 is an already accepted paper (with Gesine Reinert) for *Bernoulli*. We begin the chapter presenting the bounded Wasserstein distance, which is the metric used for the single-parameter case. The relationship of this metric with the Kolmogorov distance, which leads to insights for confidence intervals, is given before proceeding to the main result of the chapter for the bounds on the distributional distance of interest. Special focus on the broad classification of exponential family distributions is given. Simulation-based results related to the exponential distribution illustrate the behaviour of the bounds through a graphical representation. We cover both the canonical and non-canonical parametrisation and compare these two situations. In some discrete distributions there is positive probability of the MLE being on the boundary of the parameter space, which leads to problems related to the applicability of Taylor expansions. We overcome this problem using a perturbation of both the parameter and the data, forcing them to be interior to the parameter space. The Poisson distribution with parameter space, Θ , the semi-closed interval $[0, \infty)$ serves as an illustration of these results. Under further assumptions, Chapter 3 gives a procedure to get an upper bound that does not contain any terms related to the MLE. The approach is based on finding upper bounds for the mean squared error. The new bound obtained is still applicable

in cases where an analytic expression of the MLE is not available. The results are applied to the Beta distribution with one of the two shape parameters being unknown. Simulation-based results are also given.

Chapter 4 treats the multi-parameter case for independent but not necessarily identically distributed (i.n.i.d.), high-dimensional random vectors. This is a single-author paper submitted to the *Journal of Multivariate Analysis*. We explain how an adaptation of the method developed in Chapter 3 for the case of a vector parameter is feasible. A multivariate Stein's method result by Reinert and Röllin (2009), multivariate Taylor expansions and conditional expectations lead to the final upper bound. We give a general upper bound but also a specific upper bound for linear regression models and we apply the specific bound to the simplest case of the straight-line regression with two unknown parameters (the intercept and the slope). The special case where the random vectors are identically distributed is also covered; simpler regularity conditions than those for the i.n.i.d. case are given in order for the normal approximation to hold. Following a similar approach as the one followed in the i.n.i.d. case we obtain the required, simpler upper bound. As an example, the normal distribution with both mean and variance unknown is treated. Furthermore, we give an upper bound when the vector MLE is not known explicitly. The approach is again based on firstly bounding the mean squared error of the vector MLE. The bound might seem slightly complicated and tedious to calculate. However, for the Beta distribution with both parameters being unknown, an expression for the bound is easily obtained.

Chapter 5 is based on joint work with Christophe Ley which has been submitted as a paper to the *Latin American Journal of Probability and Mathematical Statistics (ALEA)*. This chapter uses the Delta method to obtain an expression of the MLE as a sum of random variables. More precisely, we use a one-to-one differentiable function in such a way that when applied to the MLE it gives a sum of random variables, allowing

the use of Stein's method for this quantity. There will also be error terms and those will be bounded through alternative techniques, based again on our main tool; Taylor expansions. We give the bound for the situation of a scalar parameter; special focus is again given to exponential families. The results are then applied to the generalised gamma distribution, of which the exponential distribution is a special case. This example allows a comparison of the results in this chapter with the results obtained in Chapter 3 for the exponential distribution.

In Chapter 6, the independence assumption of the previous chapters is relaxed to a local dependence structure between the random variables. We obtain bounds for the distributional distance of interest under this dependence structure. In an example related to normally distributed locally dependent random variables the bound behaves very well. In addition, as in the previous chapters, under further assumptions we give bounds when the MLE is not explicitly known.

Chapter 7 includes a summary of the results presented in this thesis and some ideas for future work. We cover applications as well as further theoretical extension of the results obtained in this thesis.

Chapter 2

Background

This chapter gives the background for the work presented in this thesis. In Section 2.1, a brief description of the Central Limit Theorem is given as it is the main tool used to prove the asymptotic normality of the MLE. In Section 2.2, we explain the main asymptotic result for the distribution of the MLE, which then incites us to find a way to assess the quality of such an approximation. For this purpose, Stein's method for normal approximation is employed. An introduction to the main results of this effective probabilistic method is fundamental for the understanding of the following chapters and is given in Section 2.3.

2.1 Central Limit theorem

The Central Limit Theorem was first discussed in its initial form related to Bernoulli random variables in a notable article by Abraham De Moivre in 1733 and since then it has become a cornerstone in probability and statistical theory. In general words, it states that the suitably scaled sample mean of a large number of random variables follows approximately the normal distribution regardless of the distribution of the random variables. As long as the assumptions are satisfied, the CLT has a number of

alternative versions depending on the dependence structure of the random variables, their dimensionality as well as whether we have identically distributed random variables (vectors) or not. From now on, $N(\mu, \sigma^2)$ denotes the normal distribution with mean μ and variance σ^2 . Below, we give the CLT in its classical form of i.i.d. random variables, known also as Lindeberg-Lévy CLT.

Theorem 2.1. *Let X_1, X_2, \dots be i.i.d. random variables with mean $\mu \in \mathbb{R}$ and variance $0 < \sigma^2 < \infty$. For $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ and $K \sim N(0, \sigma^2)$,*

$$\sqrt{n}(\bar{X} - \mu) \xrightarrow[n \rightarrow \infty]{d} K, \quad (2.1)$$

where \xrightarrow{d} denotes convergence in distribution.

The case $\sigma^2 = 0$ could be included by interpreting $N(0, 0)$ as a point mass at zero. A flavour of the result in (2.1) could be given through Abraham De Moivre's original article. We take a fair coin and we flip it many times. The distribution of each trial is Bernoulli with mean the probability to get 'heads', which is $\frac{1}{2}$. We add 1 (starting from zero) whenever the outcome of the coin is 'heads'. Then the distribution of the resulting sum over the total number of trials (the average) is well-approximated by the normal distribution with mean equal to $\frac{1}{2}$ and variance equal to $\frac{1}{4n}$. It is notable that through the CLT the exact distribution of the suitably scaled sum of discrete random variables can be approximated by the normal distribution, which is a continuous probability distribution.

The CLT is an asymptotic result and provides only an approximation for finite sample sizes. The quality of this approximation may depend on the magnitude of the sample size as well as the distribution of the random variables of which we take the sum. A quantification of this qualitative result is given through the celebrated Berry-Esséen theorem, which just requires the third central moment of $X_i, i = 1, 2, \dots, n$ to exist and to be finite. An upper bound of order $\frac{1}{\sqrt{n}}$ on the Kolmogorov distance between the exact distribution of the sample mean and the limiting normal distribution can be obtained; see

(1.2) for the bound. Generalising this result, we can obtain upper bounds on a variety of distributional distances using Stein's method; see Section 2.3 for a review of the method.

As our main interest is the MLE, which is not in general a sum of random variables, the CLT can not be directly applied to get the asymptotic normality of the MLE. The next section introduces the maximum likelihood estimation method and gives sufficient regularity conditions under which the asymptotic normality of the MLE holds.

2.2 Maximum likelihood estimation

In Statistics, maximum likelihood estimation is a widely used method of estimating d parameters $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_d)$ found in a parametric statistical model. Having a fixed set of data and a defined statistical model, the maximum likelihood estimation gives a set of values $\hat{\boldsymbol{\theta}}_n(\mathbf{X})$ for the parameters of the model which maximises the likelihood function. Such a vector $\hat{\boldsymbol{\theta}}_n(\mathbf{X})$, whenever it exists, is called a maximum likelihood estimator (MLE).

To set the scene in the simplest scenario, suppose that $\mathbf{X} = (X_1, X_2, \dots, X_n)$ is a sample of n i.i.d. random variables with a joint density function $f(\mathbf{x}|\boldsymbol{\theta})$, where $\boldsymbol{\theta}$ is an unknown parameter. For $X_i = x_i, i = 1, 2, \dots, n$ being some observed values, we define the so-called likelihood function by $L(\boldsymbol{\theta}; \mathbf{x}) = f(\mathbf{x}|\boldsymbol{\theta})$, where $\mathbf{x} = (x_1, x_2, \dots, x_n)$. If the distribution of the random variables is discrete (instead of a density, we have a probability mass function), then the likelihood function is the probability of observing the given data \mathbf{x} . Let $\boldsymbol{\theta}_0 = (\theta_{0_1}, \theta_{0_2}, \dots, \theta_{0_d})$, be the true (still unknown) value of the parameter vector $\boldsymbol{\theta}$. The aim is to find an estimate of $\boldsymbol{\theta}_0$ from the data \mathbf{x} by maximising the likelihood. Rather than maximising the likelihood function, it is often more convenient to work with its logarithm, called the log-likelihood function and denoted by $l(\boldsymbol{\theta}; \mathbf{x})$. The logarithmic form makes calculations easier, particularly since $L(\boldsymbol{\theta}; \mathbf{x})$

will often (specifically in the case of independent data points) be a product of marginal probability density (mass) functions.

For many models the MLE exists and is unique; this is known as the ‘regular’ case. For a number of statistical models, however, uniqueness or even existence of the MLE is not secured. In the following proposition, we specify a number of sufficient conditions for both existence and uniqueness of the MLE, as found in Makelainen et al. (1981).

Proposition 2.1. *Let the likelihood function $L(\boldsymbol{\theta}; \mathbf{x})$ be a twice continuously differentiable function with $\boldsymbol{\theta}$ varying in a connected open subset $\Theta \subset \mathbb{R}^d$, called the parameter space. Suppose that:*

- (i) $\lim_{\boldsymbol{\theta} \rightarrow \partial\Theta} L(\boldsymbol{\theta}; \mathbf{x}) = 0$, where $\partial\Theta$ is the boundary of the parameter space Θ ,
- (ii) With $\frac{\partial}{\partial \theta_i}$ denoting now partial derivatives, the Hessian matrix

$$\mathbf{H}(\boldsymbol{\theta}; \mathbf{x}) = \left\{ \frac{\partial^2}{\partial \theta_i \partial \theta_j} L(\boldsymbol{\theta}; \mathbf{x}) \right\}_{i,j=1,\dots,d}$$

of second partial derivatives is negative definite at every point $\boldsymbol{\theta} \in \Theta$ for which the gradient vector

$$\nabla(L(\boldsymbol{\theta}; \mathbf{x})) = \left\{ \frac{\partial}{\partial \theta_i} L(\boldsymbol{\theta}; \mathbf{x}) \right\}_{i=1,\dots,d}$$

vanishes. Then

1. *there is a unique maximum likelihood estimate $\hat{\boldsymbol{\theta}}_n(\mathbf{x}) \in \Theta$, and*
2. *the likelihood function attains*
 - (a) *no other maxima in Θ ;*
 - (b) *no minima or other stationary points in Θ ;*
 - (c) *its infimum value 0 on the boundary $\partial\Theta$ and nowhere else.*

If only $\nabla(L(\boldsymbol{\theta}; \mathbf{x})) = \mathbf{0}$ is assumed, then there may be multiple solutions (stationary points) (Billingsley, 1961).

Remark 2.1. Condition (ii) can be translated in terms of the log-likelihood function $l(\boldsymbol{\theta}; \mathbf{x})$. Notice that

$$\nabla(l(\boldsymbol{\theta}; \mathbf{x})) = \frac{\nabla(L(\boldsymbol{\theta}; \mathbf{x}))}{L(\boldsymbol{\theta}; \mathbf{x})}.$$

Therefore, the vector of the first order partial derivatives of the likelihood function, vanishes whenever the respective vector related to the log-likelihood function vanishes, as long as $0 < L(\boldsymbol{\theta}; \mathbf{x}) < \infty$. In addition, the Hessian matrix for $l(\boldsymbol{\theta}; \mathbf{x})$ is

$$\mathbf{H}_l(\boldsymbol{\theta}; \mathbf{x}) = \frac{\mathbf{H}(\boldsymbol{\theta}; \mathbf{x})L(\boldsymbol{\theta}; \mathbf{x}) - [\nabla(L(\boldsymbol{\theta}; \mathbf{x}))]^2}{[L(\boldsymbol{\theta}; \mathbf{x})]^2}.$$

For $\nabla(L(\boldsymbol{\theta}; \mathbf{x})) = \mathbf{0}$, notice that $\mathbf{H}_l(\boldsymbol{\theta}; \mathbf{x}) = \frac{\mathbf{H}(\boldsymbol{\theta}; \mathbf{x})}{L(\boldsymbol{\theta}; \mathbf{x})}$ and thus $\mathbf{H}(\boldsymbol{\theta}; \mathbf{x})$ is negative definite whenever $\mathbf{H}_l(\boldsymbol{\theta}; \mathbf{x})$ is a negative definite matrix and $L(\boldsymbol{\theta}; \mathbf{x}) > 0$. Therefore, for $L(\boldsymbol{\theta}; \mathbf{x}) > 0$ we can easily express condition (ii) of Proposition 2.1 in terms of the log-likelihood function, which simplifies the calculations. It becomes

(ii)' With $\frac{\partial}{\partial \theta_i}$ denoting partial derivatives, the Hessian matrix of the log-likelihood function

$$\mathbf{H}_l(\boldsymbol{\theta}; \mathbf{x}) = \left\{ \frac{\partial^2}{\partial \theta_i \partial \theta_j} l(\boldsymbol{\theta}; \mathbf{x}) \right\}_{i,j=1,\dots,d}$$

of second partial derivatives is negative definite at every point $\boldsymbol{\theta} \in \Theta$ for which the gradient vector

$$\nabla(l(\boldsymbol{\theta}; \mathbf{x})) = \left\{ \frac{\partial}{\partial \theta_i} l(\boldsymbol{\theta}; \mathbf{x}) \right\}_{i=1,\dots,d}$$

vanishes.

Proposition 2.1 provides sufficient conditions for the existence and uniqueness of the MLE when the parameter space, Θ , is an open and connected subset of \mathbb{R}^d . However, in some statistical models, the parameter space is compact. We use Kiefer (2008) for

the following proposition related to the existence and uniqueness of the MLE when it is obtained from parametric statistical models with compact parameter space.

Proposition 2.2. *If the parameter space Θ is compact and if the likelihood function is continuous on Θ then there exists a MLE. In addition, if the parameter space Θ is also convex and the likelihood function is strictly concave in θ , then the MLE is unique.*

The first result, related to the existence of the MLE, arises as a direct application of the Weierstrass Extreme Value Theorem, which states that if a real valued function f is continuous in a closed and bounded set, then it attains a maximum and a minimum, each at least once.

For this chapter and Chapters 3, 5 and 6, unless otherwise stated, the parameter θ is assumed to be scalar. The vector-parameter case is covered in Chapter 4. In the single-parameter setting, we denote the derivatives of the log-likelihood function with respect to θ , by $l'(\theta; x), l''(\theta; x), \dots, l^{(j)}(\theta; x)$, for j any integer greater than 2. Unless otherwise stated, the parameter varies in an open interval (a, b) , where $a, b \in \mathbb{R} \cup \{-\infty, \infty\}$ and $a < b$. Following (Casella and Berger, 2002, p.516), we also make the following assumptions:

- (R1) the parameter is identifiable, meaning that if $\theta \neq \theta'$, then $\exists x : f(x|\theta) \neq f(x|\theta')$;
- (R2) the density $f(x|\theta)$ is three times differentiable with respect to θ , the third derivative is continuous in θ and $\int f(x|\theta) dx$ can be differentiated three times under the integral sign;
- (R3) for any $\theta_0 \in \Theta$ and for S denoting the support of $f(x|\theta)$, there exists a positive number ε_0 and a function $M(x)$ (both of which may depend on θ_0) such that

$$\left| \frac{d^3}{d\theta^3} \log f(x|\theta) \right| \leq M(x) \quad \forall x \in S, \quad \forall \theta \in \Theta \text{ such that } |\theta - \theta_0| < \varepsilon_0,$$

with $E_{\theta_0}[M(X)] < \infty$;

(R4) $i(\theta_0) \neq 0$, where $i(\theta)$ is the expected Fisher Information number for one random variable.

The requirement (R2) that $\int f(x|\theta) dx$ can be differentiated three times under the integral sign means that $\int_{\mathbb{R}} \frac{d^j}{d\theta^j} f(x|\theta) dx = \frac{d^j}{d\theta^j} \int_{\mathbb{R}} f(x|\theta) dx = 0, j \in \{1, 2, 3\}$. This condition ensures that if the expressions exist, then $E_{\theta}[l'(\theta; \mathbf{X})] = 0$ and $\text{Var}_{\theta}[l'(\theta; \mathbf{X})] = n i(\theta)$. The motivation of the work presented in this thesis are the results given in Theorem 2.2. The efficiency and asymptotic normality of the MLE have first been discussed in Fisher (1925). Here the i.i.d. case is presented; see Theorem 4.1 for the case of independent but not necessarily identically distributed random vectors. In addition, for dependent random variables the asymptotic result is given in Theorem 6.2. Here, $\mathbf{X} = (X_1, X_2, \dots, X_n)$ is the random vector and the MLE, $\hat{\theta}_n(\mathbf{X})$, is of course also random as a function of \mathbf{X} .

Theorem 2.2. *Let X_1, X_2, \dots, X_n be i.i.d. random variables with probability density (or mass) function $f(x_i|\theta)$, where θ is the scalar parameter. Assume that the MLE exists, it is unique and (R1)-(R4) are satisfied. Then for $Z \sim N(0, 1)$*

$$(a) \frac{1}{\sqrt{n}} l'(\theta_0; \mathbf{X}) \xrightarrow[n \rightarrow \infty]{d} \sqrt{i(\theta_0)} Z, \quad (b) \sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \xrightarrow[n \rightarrow \infty]{d} Z. \quad (2.2)$$

Thus, in the case of i.i.d. random variables, under the regularity conditions (R1)-(R4), the asymptotic distribution of the suitably scaled MLE is the standard normal distribution. However, this is an asymptotic result and we will never have an infinite sample size. We will assess the quality of this approximation using partly Stein's method.

2.3 Stein's method

Stein's method was first introduced by Charles Stein in Stein (1972), while the monograph Stein (1986) explains in detail the method and it is in our opinion the most notable contribution. Stein's method is a powerful technique to assess whether a random variable, W , has a distribution close to a target distribution. The distance between the distribution of W and the target one is measured using a Zolotarev-type distance. If F, G are two random variables with values in \mathbb{R} and H is a class of separating functions, then a Zolotarev-type distance between the laws of F and G , induced by H , is given by the quantity

$$d_H(F, G) = \{ |E[h(F)] - E[h(G)]| : h \in H \}. \quad (2.3)$$

For example, taking $H = H_K = \{ \mathbb{1}_{[\cdot, \leq x]}, x \in \mathbb{R} \}$ in (2.3), gives the Kolmogorov-Smirnov distance

$$d_K(F, G) = \sup_{z \in \mathbb{R}} |\mathbb{P}(F \leq z) - \mathbb{P}(G \leq z)|,$$

which is the maximum distance between two cumulative distribution functions. With

$$H_W = \{ h : \mathbb{R} \rightarrow \mathbb{R} : |h(x) - h(y)| \leq |x - y| \}, \quad (2.4)$$

the set of Lipschitz functions with constant equal to one, (2.3) gives the Wasserstein distance,

$$d_W(F, G) = \sup \{ |E[h(F)] - E[h(G)]| : h \in H_W \}. \quad (2.5)$$

Apart from just giving the general results, in the thesis we also focus on the Wasserstein as well as on the bounded Wasserstein (or Fortet-Mourier) distance, which is obtained for

$$H = H_{bW} = \{ h : \mathbb{R} \rightarrow \mathbb{R} : \|h\|_{\text{Lip}} + \|h\| \leq 1 \}, \quad (2.6)$$

where from now on $\|\cdot\|$ denotes the supremum norm ($\|\cdot\|_\infty$) and $\|h\|_{\text{Lip}} = \sup_{\substack{x \neq y \\ x, y \in \mathbb{R}}} \frac{|h(x) - h(y)|}{|x - y|}$.

The bounded Wasserstein distance is denoted from now on by

$$d_{bW}(F, G) = \sup \{ |E[h(F)] - E[h(G)]| : h \in H_{bW} \}, \quad (2.7)$$

with H_{bW} as in (2.6); see for example Nourdin and Peccati (2012). Using Rademacher's theorem, if $\|h\|_{\text{Lip}} \leq 1$, then h is differentiable almost everywhere, with h' denoting its derivative. Hence first-order Taylor expansion is possible.

An alternative to Stein's method are characteristic functions, which are quite powerful in examining if the distributions of two random variables are similar. Closeness of the characteristic functions implies closeness of the distributions. However, such methods lead to loss of probabilistic intuition as they rely on transformations of the quantity of interest. In addition, these methods are not easy to apply in the presence of a dependence structure; for limit theorems for sums of m -dependent random variables see Heinrich (1982). On the other hand, Stein's method is based on a characterisation of the target distribution and this makes it easier to manipulate the distribution through random quantities of which W is composed. Therefore, probabilistic intuition is not lost. In addition, Stein's method is quite powerful in the presence of dependence, providing explicit upper bounds through coupling constructions such as the exchangeable pair, the zero-bias and the size bias couplings.

Stein's method for normal approximation—the target distribution is the normal—is based on the well-known characterisation of the normal distribution that $K \sim N(0, \sigma^2)$ if and only if

$$E[Kt(K)] - \sigma^2 E[t'(K)] = 0,$$

for all absolutely continuous and differentiable functions t for which the above expectations exist. If the distribution of the random variable, W , is close to the $N(0, \sigma^2)$

distribution, then both $|\mathbb{E}[h(W)] - \mathbb{E}[h(K)]|$ and $|\mathbb{E}[Wt(W)] - \sigma^2 \mathbb{E}[t'(W)]|$ should be small, where h is a test function. These two expressions, more precisely the functions h and t , are now related through the so-called Stein equation for the $N(0, \sigma^2)$ distribution

$$\sigma^2 t'(w) - wt(w) = h(w) - \mathbb{E}[h(K)]. \quad (2.8)$$

In order to obtain (2.3) for a given function h , it suffices to find the solution t_h of (2.8), evaluate $\sigma^2 t'_h(w) - wt_h(w)$ at W and take the absolute value of the expectation.

Lemma 2.1. (Stein, 1972). *For $Z \sim N(0, 1)$, the unique bounded solution of the Stein equation in (2.8) for $\sigma^2 = 1$ is given by*

$$f_h(w) = -e^{\frac{w^2}{2}} \int_w^\infty e^{-\frac{x^2}{2}} (h(x) - \mathbb{E}[h(Z)]) dx. \quad (2.9)$$

Proof. See Chen et al. (2011). The idea is to solve a first order ordinary differential equation. Multiplying both sides of (2.8) for $\sigma^2 = 1$ by the integrating factor, which is $e^{-\frac{w^2}{2}}$, yields

$$\frac{d}{dw} \left(f(w) e^{-\frac{w^2}{2}} \right) = e^{-\frac{w^2}{2}} (h(w) - \mathbb{E}[h(Z)]).$$

The next step is to integrate both sides to get, for $c \in \mathbb{R}$, that

$$\begin{aligned} f_h(w) e^{-\frac{w^2}{2}} &= \int_{-\infty}^w e^{-\frac{x^2}{2}} (h(x) - \mathbb{E}[h(Z)]) dx + c \\ \Rightarrow f_h(w) &= e^{\frac{w^2}{2}} \int_{-\infty}^w e^{-\frac{x^2}{2}} (h(x) - \mathbb{E}[h(Z)]) dx + c e^{\frac{w^2}{2}}. \end{aligned}$$

Now, $|f_h(w)|$ is bounded if and only if $c = 0$. Thus, using also that

$$\int_{-\infty}^\infty e^{-\frac{x^2}{2}} (h(x) - \mathbb{E}[h(Z)]) dx = 0,$$

yields that the unique bounded solution is the one in (2.9). □

If $h : \mathbb{R} \rightarrow \mathbb{R}$, is absolutely continuous and bounded then the solution f_h in (2.9) satisfies

$$\|f_h''\| \leq 2\|h'\|, \quad (2.10)$$

see Chen et al. (2011), (2.13) in p.16. In the general case (σ^2 not necessarily equal to 1) the unique bounded solution, $t_h(w)$, of the Stein equation in (2.8) satisfies

$$t_h(w) = \frac{1}{\sigma} f_h\left(\frac{w}{\sigma}\right).$$

Thus, using (2.10),

$$\|t_h''\| \leq 2 \frac{\|h'\|}{\sigma^3}. \quad (2.11)$$

A general result lies at the centre of our work and is given in Lemma 2.2.

Lemma 2.2. *Let Y_1, Y_2, \dots, Y_n be i.i.d. random variables with $E(Y_i) = 0$, $\text{Var}(Y_i) = \sigma^2 > 0$ and $E|Y_i|^3 < \infty$. Let $W = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i$ and $K \sim N(0, \sigma^2)$. Then for any function $h : \mathbb{R} \rightarrow \mathbb{R}$, such that h is absolutely continuous and bounded*

$$|E[h(W)] - E[h(K)]| \leq \frac{\|h'\|}{\sqrt{n}} \left(2 + \frac{1}{\sigma^3} [E|Y_1|^3] \right). \quad (2.12)$$

Proof. The method of the proof is standard using the famous ‘leave one out’ approach and follows the steps in (Chen et al., 2011, p.5), where a sketch of the proof for $\sigma^2 = 1$ is given. We will show (2.12) by finding an upper bound on $|\sigma^2 E[t_h'(W)] - E[W t_h(W)]|$. Let $W_i = W - \frac{Y_i}{\sqrt{n}} = \frac{1}{\sqrt{n}} \sum_{j=1, j \neq i}^n Y_j$. A second-order Taylor expansion of $t_h(W)$ about W_i gives

$$\begin{aligned} E[W t_h(W)] &= \frac{1}{\sqrt{n}} E \left[\sum_{i=1}^n Y_i t_h(W) \right] \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n E \left[Y_i (t_h(W_i) + (W - W_i) t_h'(W_i) + \frac{1}{2} (W - W_i)^2 t_h''(W_i^*)) \right], \end{aligned}$$

where W^* is between W and W_i . Having $W - W_i = \frac{Y_i}{\sqrt{n}}$, with Y_i and W_i being independent yields

$$\begin{aligned} E[W t_h(W)] &= \frac{1}{\sqrt{n}} \sum_{i=1}^n E[Y_i] E[t_h(W_i)] + \frac{1}{n} \sum_{i=1}^n E[Y_i^2] E[t'_h(W_i)] + \frac{1}{2n^{\frac{3}{2}}} \sum_{i=1}^n E[Y_i^3 t''_h(W^*)] \\ &= \frac{\sigma^2}{n} \sum_{i=1}^n E[t'_h(W_i)] + \frac{1}{2n^{\frac{3}{2}}} \sum_{i=1}^n E[Y_i^3 t''_h(W^*)]. \end{aligned} \quad (2.13)$$

Using (2.13),

$$\begin{aligned} &|\sigma^2 E[t'_h(W)] - E[W t_h(W)]| \\ &= \left| \sigma^2 E[t'_h(W)] - \frac{\sigma^2}{n} \sum_{i=1}^n E[t'_h(W_i)] - \frac{1}{2n^{\frac{3}{2}}} \sum_{i=1}^n E[Y_i^3 t''_h(W^*)] \right| \\ &= \left| \frac{\sigma^2}{n} \sum_{i=1}^n E[t'_h(W) - t'_h(W_i)] - \frac{1}{2n^{\frac{3}{2}}} \sum_{i=1}^n E[Y_i^3 t''_h(W^*)] \right| \\ &\leq \frac{\sigma^2}{n^{\frac{3}{2}}} \|t''_h\| \sum_{i=1}^n E|Y_i| + \frac{1}{2n^{\frac{3}{2}}} \|t''_h\| \sum_{i=1}^n E|Y_i^3|, \end{aligned}$$

by first-order Taylor expansion of $t'_h(W)$ about W_i . Applying the Cauchy-Schwarz inequality to the first term and using (2.11) yields

$$\begin{aligned} |\sigma^2 E[t'_h(W)] - E[W t_h(W)]| &\leq \frac{\sigma^3}{\sqrt{n}} \|t''_h\| + \frac{1}{2\sqrt{n}} \|t''_h\| E|Y_1^3| \\ &\leq \frac{\|h'\|}{\sqrt{n}} \left(2 + \frac{1}{\sigma^3} [E|Y_1|^3] \right). \end{aligned}$$

□

Remark 2.2. Using $Y_i = \frac{d}{d\theta} \log f(X_i|\theta_0)$, notice that (2.12) can be used to assess the quality of (a) in (2.2).

Lemma 2.2 gives upper bounds in the case of independent random variables. However, sometimes independence is not secured and thus we need to use results related to dependent random variables. In Chapter 6 we give bounds when there is a local dependence structure between the random variables. Below, the notion of local dependence is

introduced before we give the Stein's method result that is used in the case of locally dependent random variables.

An m -dependent sequence of random variables $\{X_i, i \in \mathbb{N}\}$ is such that for each $i \in \mathbb{N}$ the sets of random variables $\{X_j, j \leq i\}$ and $\{X_j, j > i + m\}$ are independent. The Stein's method result for the case of locally dependent random variables is based on the local dependence condition (LD) which is explained below.

For a set of random variables $\{\xi_i, i = 1, 2, \dots, n\}$, for any $A \subset \{1, 2, \dots, n\}$ we define

$$A^c = \{i \in \{1, 2, \dots, n\} : i \notin A\}$$

$$\xi_A = \{\xi_i : i \in A\}$$

and we introduce the local dependence condition

(LD) For each $i \in \{1, 2, \dots, n\}$ there exist $A_i \subset B_i \subset \{1, 2, \dots, n\}$ such that ξ_i is independent of $\xi_{A_i^c}$ and ξ_{A_i} is independent of $\xi_{B_i^c}$.

Whenever this condition holds we set

$$\eta_i = \sum_{j \in A_i} \xi_j, \quad \tau_i = \sum_{j \in B_i} \xi_j. \quad (2.14)$$

Lemma 2.3 below gives an upper bound for the Wasserstein distance between the distribution of a sum of m -dependent random variables satisfying (LD). The random variables are assumed to have mean zero with the variance of their sum being equal to one. The proof of the lemma is beyond the scope of the thesis and can be found in (Chen et al., 2011, p.134).

Lemma 2.3. *Let $\{\xi_i, i = 1, 2, \dots, n\}$ be a set of random variables with mean zero and $\text{Var}(W) = 1$, where $W = \sum_{i=1}^n \xi_i$. If (LD) holds, then with η_i and τ_i as in (2.14),*

$$d_W(W, Z) \leq 2 \sum_{i=1}^n (E|\xi_i \eta_i \tau_i| + |E(\xi_i \eta_i)| E|\tau_i|) + \sum_{i=1}^n E|\xi_i \eta_i^2|. \quad (2.15)$$

In situations of dependent random variables, couplings turn out to be useful. Here the exchangeable pair coupling is briefly presented. For a given random variable of interest W , the exchangeable pair approach relies on a construction of a random variable W' such that (W, W') is an exchangeable pair; equivalently $(W, W') \stackrel{d}{=} (W', W)$, where $\stackrel{d}{=}$ denotes equality in distribution. An exchangeable pair is called a λ -Stein pair if it also satisfies

$$E(W'|W) = (1 - \lambda)W,$$

where $\lambda \in (0, 1)$. An example of such an approach for independent random variables is given in (Chen et al., 2011, p.23), where $W = \sum_{i=1}^n \xi_i$, with $\xi_i, i \in \{1, 2, \dots, n\}$, are independent random variables with zero mean and such that $\text{Var}(W) = 1$. Allowing $\{\xi'_i, i = 1, 2, \dots, n\}$ to be an independent copy of $\{\xi_i, i = 1, 2, \dots, n\}$ and I to have the uniform distribution on $\{1, 2, \dots, n\}$ we define $W' = W - \xi_I + \xi'_I$. By construction (W, W') and (W', W) have the same distribution and simple calculations show that

$$E(W'|W) = \left(1 - \frac{1}{n}\right)W.$$

Hence, (W, W') is a λ -Stein pair with $\lambda = \frac{1}{n}$. This approach can be extended to the multivariate setting, where W and W' are vectors. Such an approach will be used in Chapter 4 through a theorem by Reinert and Röllin (2009).

So far, this section considered univariate random variables. In this thesis Stein's method for multivariate normal approximation in \mathbb{R}^d is also used. We start with the multivariate Stein equation. For $g : \mathbb{R}^d \rightarrow \mathbb{R}$ being twice differentiable and $\nabla g, D^2g$ being

the gradient vector and Hessian matrix of g , respectively, we use the Stein equation

$$\text{Tr}(D^2 g(\mathbf{w})) - \mathbf{w} \nabla g(\mathbf{w}) = h(\mathbf{w}) - \mathbb{E}[h(\mathbf{Z})], \quad (2.16)$$

where $\mathbf{Z} \sim \mathbf{N}_d(\mathbf{0}, I_{d \times d})$ and $\text{Tr}(\cdot)$ denotes the trace of a matrix. The function $h: \mathbb{R}^d \rightarrow \mathbb{R}$, still plays the role of the test function. Lemma 2.6 in (Chen et al., 2011, p.17) shows that if h has three bounded derivatives then the unique bounded solution of (2.16) is

$$g_h(\mathbf{w}) = - \int_0^\infty \left[\mathbb{E} \left[h \left(\mathbf{w} e^{-u} + \sqrt{1 - e^{-2u}} \mathbf{Z} \right) \right] - \mathbb{E}[h(\mathbf{Z})] \right] du$$

and satisfies

$$\|g_h^{(\mathbf{k})}\| \leq \frac{1}{k} \|h^{(\mathbf{k})}\|.$$

Here for a vector $\mathbf{k} = (k_1, k_2, \dots, k_d)$ with $k_j \in \mathbb{Z}_0^+$, we have $h^{(\mathbf{k})}(\mathbf{w}) = \frac{\partial^{|\mathbf{k}|}}{\prod_{j=1}^d \partial w_{k_j}} h(\mathbf{w})$ and $|\mathbf{k}| = \sum_{j=1}^d k_j$. The idea for the construction of exchangeable pair couplings is similar as in the univariate case.

Significant progress on Stein's method beyond normal approximation has been achieved. General distributions can be treated via the generator approach of Barbour (1990), the density approach of Stein et al. (2004), further developed in Ley and Swan (2013b) and Ley and Swan (2013a), and the parsimonious Stein operator approach of Ley et al. (2014). The most worth-mentioning extension is Stein's method for the Poisson distribution. This method is also known as Stein-Chen method and was first presented in Chen (1975). An illustration through various examples of the broad applicability and power of the Stein-Chen method for Poisson approximation can be found in Arratia et al. (1990). In addition, using Stein-Chen method, Barbour and Holst (1989) give upper bounds on the distributional distance between the distribution of a sum of Bernoulli random variables and the Poisson distribution, while Barbour et al. (1992) treats in detail the Stein-Chen method.

Further distributions in which Stein's method has been adapted include the uniform (Diaconis, 1989), point mass (Reinert, 1995), geometric (Peköz, 1996), Beta (Döbler, 2015), (Goldstein and Reinert, 2013), Laplace (Pike and Ren, 2014), binomial (Ehm, 1991), negative binomial (Barbour et al., 2015), exponential (Chatterjee et al., 2011), (Peköz and Röllin, 2011) and gamma distributions (Luk, 1994), (Nourdin and Peccati, 2009), (Pickett, 2004) as well as the class of variance-Gamma distributions (Gaunt, 2013). Related to the normal distribution, Barbour (1990) and Götze (1991) study similar results for more general Gaussian distributions. Nourdin and Peccati (2012) generalise the results to an infinite-dimensional setting. The survey article Ross (2011) gives the main concepts of Stein's method for distributional approximation by the normal, Poisson, exponential and geometric distributions.

In recent years there has been a growing interest in assessing the quality of asymptotic results in statistical theory using Stein's method. Recent examples include Ley et al. (2015) where the impact of the choice of the prior distribution in Bayesian statistics is assessed and Gaunt et al. (2015) for the chi-square approximation to the Pearson Statistic and of course the quantitative results of this thesis for assessing the normal approximation of the MLE. Our bounds have already influenced Pinelis and Molzon (2016).

Chapter 3

Bounds on the distance to normal for the MLE

In this chapter we assess the normal approximation of the MLE. Section 3.1 first briefly relates the Kolmogorov and the bounded Wasserstein distance. In addition, insights to conservative confidence intervals are given. Section 3.2 gives the main upper bound on the distributional distance between the distribution of the MLE and the normal distribution in terms of the bounded Wasserstein distance.

We specify the result for the distribution of the MLE for single-parameter exponential family distributions. Section 3.3 provides an upper bound for the distance between the distribution of the MLE for exponential family distributions and the normal distribution using the results of Section 3.2. We treat the example of the exponential distribution; simulation-based results illustrate the behaviour of the bound. In Section 3.4 we explain how, in discrete distributions, a perturbation method solves the issue of the MLE having positive probability of being on the boundary of the parameter space. The Poisson distribution with parameter space $[0, \infty)$ serves as an example to illustrate the results. Finally, in Section 3.5, bounds are provided for cases where a closed-form expression

for the MLE is not available. To achieve this task we first give bounds on the mean squared error of the MLE. The results of this chapter form a paper with Gesine Reinert, which is to appear in *Bernoulli*.

3.1 The bounded Wasserstein and the Kolmogorov distance

For $Z \sim N(0, 1)$, the aim is to bound the bounded Wasserstein distance

$$d_{bW} \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right), \quad (3.1)$$

with $d_{bW}(\cdot, \cdot)$ as defined in (2.7). Recall that using $H_K = \{ \mathbb{1}_{[\cdot, \leq x]}, x \in \mathbb{R} \}$ as the class of functions in (2.3) yields the Kolmogorov distance,

$$d_K \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right).$$

The next proposition links our results to the Kolmogorov distance.

Proposition 3.1. *Let W be any real-valued random variable and $Z \sim N(0, 1)$. Assume that there exist $\delta_1 > 0$ and $\delta_2 \geq 0$, such that for any function h with $\|h'\|$ and $\|h\|$ bounded,*

$$|\mathbb{E}[h(W)] - \mathbb{E}[h(Z)]| \leq \delta_1 \|h'\| + \delta_2 \|h\|. \quad (3.2)$$

Then,

$$d_K(W, Z) \leq 2\sqrt{\delta_1} + \delta_2.$$

Proof. The proof of this proposition follows the approach explained in the proof of Theorem 3.3 of Chen et al. (2011), p.48, where a similar result is shown that links the

Wasserstein and the Kolmogorov distance. Let $z \in \mathbb{R}$ and for $\alpha = \sqrt{\delta_1}(2\pi)^{\frac{1}{4}}, z \in \mathbb{R}$, let

$$h_\alpha(w) = \begin{cases} 1, & \text{if } w \leq z, \\ 1 + \frac{z-w}{\alpha}, & \text{if } z < w \leq z + \alpha, \\ 0, & \text{if } w > z + \alpha, \end{cases}$$

so that h_α is bounded Lipschitz with $\|h_\alpha\| \leq 1$ and $\|h'_\alpha\| \leq \frac{1}{\alpha}$. The definition of $h_\alpha(\cdot)$ gives that

$$\begin{aligned} \mathbb{P}(W \leq z) - \mathbb{P}(Z \leq z) &\leq \mathbb{E}[h_\alpha(W)] - \mathbb{E}[h_\alpha(Z)] + \mathbb{E}[h_\alpha(Z)] - \mathbb{P}(Z \leq z) \\ &\leq |\mathbb{E}[h_\alpha(W)] - \mathbb{E}[h_\alpha(Z)]| + \mathbb{E}[h_\alpha(Z)] - \mathbb{P}(Z \leq z) \\ &\leq \delta_1 \|h'_\alpha\| + \delta_2 \|h_\alpha\| + \mathbb{P}(z < Z \leq z + \alpha) \\ &\leq \frac{\delta_1}{\alpha} + \delta_2 + \frac{\alpha}{\sqrt{2\pi}} \\ &\leq 2 \frac{\sqrt{\delta_1}}{(2\pi)^{\frac{1}{4}}} + \delta_2 \leq 2\sqrt{\delta_1} + \delta_2. \end{aligned}$$

Similarly, using now

$$h_\alpha^*(w) = \begin{cases} 0, & \text{if } w \leq z - \alpha, \\ \frac{w-(z-\alpha)}{\alpha}, & \text{if } z - \alpha < w \leq z, \\ 1, & \text{if } w > z, \end{cases}$$

we can show that $\mathbb{P}(W \leq z) - \mathbb{P}(Z \leq z) \geq -(2\sqrt{\delta_1} + \delta_2)$, which completes the proof. \square

The Kolmogorov distance relates directly to exact conservative confidence intervals. Our bounds for the MLE (before specialising on d_{bW}) are of the form of the expression in (3.2), where in our case $W = \sqrt{ni(\theta_0)}(\hat{\theta}_n(\mathbf{X}) - \theta_0)$, δ_1 and δ_2 are explicitly given

with $\delta_1 = \mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$ and $\delta_2 = \mathcal{O}\left(\frac{1}{n}\right)$. Using Proposition 3.1,

$$d_K\left(\sqrt{ni(\theta_0)}(\hat{\theta}_n(\mathbf{X}) - \theta_0), Z\right) \leq 2\sqrt{\delta_1} + \delta_2 =: B_K.$$

Therefore, for $y \in \mathbb{R}$:

$$\begin{aligned} & \left| \mathbb{P}\left(\sqrt{ni(\theta_0)}(\hat{\theta}_n(\mathbf{X}) - \theta_0) \leq y\right) - \mathbb{P}(Z \leq y) \right| \leq B_K \\ & \Leftrightarrow -B_K \leq \mathbb{P}\left(\sqrt{ni(\theta_0)}(\hat{\theta}_n(\mathbf{X}) - \theta_0) \leq y\right) - \mathbb{P}(Z \leq y) \leq B_K. \end{aligned} \quad (3.3)$$

For $\Phi^{-1}(\cdot)$ the quantile function for the standard normal distribution, applying (3.3) to $y = \Phi^{-1}\left(\frac{\alpha}{2} - B_K\right)$ and to $y = \Phi^{-1}\left(1 - \frac{\alpha}{2} + B_K\right)$ yields

$$\mathbb{P}\left(\Phi^{-1}\left(\frac{\alpha}{2} - B_K\right) \leq \sqrt{ni(\theta_0)}(\hat{\theta}_n(\mathbf{X}) - \theta_0) \leq \Phi^{-1}\left(1 - \frac{\alpha}{2} + B_K\right)\right) \geq 1 - \alpha.$$

Hence, if the expected Fisher Information number for one random variable, $i(\theta_0)$, is known, then

$$\left(\hat{\theta}_n(\mathbf{X}) - \frac{\Phi^{-1}\left(1 - \frac{\alpha}{2} + B_K\right)}{\sqrt{ni(\theta_0)}}, \hat{\theta}_n(\mathbf{X}) - \frac{\Phi^{-1}\left(\frac{\alpha}{2} - B_K\right)}{\sqrt{ni(\theta_0)}}\right)$$

is a conservative $100(1 - \alpha)\%$ confidence interval for θ_0 . Note that B_K may depend on θ_0 .

3.2 Bounds in terms of the bounded Wasserstein distance

The bounded Wasserstein distance links in well with Stein's method because the Lipschitz test functions are differentiable almost everywhere and so Taylor expansions can

be used. From now on, $\frac{d}{d\theta} \log f(X_1|\theta_0) := \frac{d}{d\theta} \log f(X_1|\theta) \Big|_{\theta=\theta_0}$. The next two results provide a bound for (a) and (b) in (2.2), respectively.

Proposition 3.2. *Suppose X_1, X_2, \dots, X_n are i.i.d. random variables with probability density (or mass) function $f(x_i|\theta)$. Assume that the regularity conditions (R1)-(R4) are satisfied and $E_{\theta_0} \left| \frac{d}{d\theta} \log f(X_1|\theta_0) \right|^3$ exists. Let $Z \sim N(0, 1)$. Then for $h: \mathbb{R} \rightarrow \mathbb{R}$, such that h is absolutely continuous and bounded*

$$\left| E_{\theta_0} \left[h \left(\frac{l'(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}} \right) \right] - E[h(Z)] \right| \leq \frac{\|h'\|}{\sqrt{n}} \left(2 + \frac{1}{[i(\theta_0)]^{\frac{3}{2}}} \left[E_{\theta_0} \left| \frac{d}{d\theta} \log f(X_1|\theta_0) \right|^3 \right] \right). \quad (3.4)$$

In particular,

$$d_{bW} \left(\frac{l'(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}}, Z \right) \leq \frac{1}{\sqrt{n}} \left(2 + \frac{1}{[i(\theta_0)]^{\frac{3}{2}}} \left[E_{\theta_0} \left| \frac{d}{d\theta} \log f(X_1|\theta_0) \right|^3 \right] \right). \quad (3.5)$$

Proof. Let $Y_i = Y_i(X_i; \theta_0) = \frac{\frac{d}{d\theta} \log f(X_i|\theta_0)}{\sqrt{i(\theta_0)}}, i = 1, 2, \dots, n$, which are i.i.d. random variables as X_1, X_2, \dots, X_n are i.i.d. The regularity conditions (R1)-(R4) ensure that $E_{\theta_0}[Y_i] = 0$ and $\text{Var}_{\theta_0}[Y_i] = 1$. Let $W = W(\mathbf{X}; \theta_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i = \frac{l'(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}}$, so that $E_{\theta_0}[W] = 0$ and $\text{Var}_{\theta_0}[W] = 1$. Applying Lemma 2.2 to $K = Z \sim N(0, 1)$ gives the result in (3.4). Using that $\|h'\| \leq 1$ for $h \in H_{bW}$, we obtain the bound in (3.5). \square

The following theorem is the main result of this chapter; it is a bound on the distributional distance between the exact distribution of the MLE and its limiting normal distribution for the case of i.i.d. random variables from a single-parameter distribution. We specify the result for the bounded Wasserstein distance. For the sake of presentation we drop the subscript θ_0 from the expectation and variance.

Theorem 3.1. *Let X_1, X_2, \dots, X_n be i.i.d. random variables with probability density (or mass) function $f(x_i|\theta)$ such that the MLE, $\hat{\theta}_n(\mathbf{X})$, exists and is unique and the regularity conditions (R1)-(R4) are satisfied. Assume that $E \left| \frac{d}{d\theta} \log f(X_1|\theta_0) \right|^3 < \infty$ and*

that $\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^4 \right] < \infty$. Let $0 < \varepsilon = \varepsilon(\theta_0) = \varepsilon_0$ be such that $|\theta - \theta_0| < \varepsilon$ with ε_0 as in (R3). Also, assume that $\mathbb{E} \left(\left(\sum_{i=1}^n M(X_i) \right)^2 \middle| |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right) < \infty$, with $M(X)$ as in (R3). Let $Z \sim \mathcal{N}(0, 1)$. Then for $h : \mathbb{R} \rightarrow \mathbb{R}$, such that h is absolutely continuous and bounded,

$$\begin{aligned}
& \left| \mathbb{E} \left[h \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) \right] - \mathbb{E}[h(Z)] \right| \\
& \leq \frac{\|h'\|}{\sqrt{n}} \left(2 + \frac{1}{[i(\theta_0)]^{\frac{3}{2}}} \left[\mathbb{E} \left| \frac{d}{d\theta} \log f(X_1 | \theta_0) \right|^3 \right] \right) \\
& + \frac{\|h'\|}{\sqrt{i(\theta_0)}} \sqrt{\text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right)} \sqrt{\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} \\
& + \frac{2\|h\|}{\varepsilon^2} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \\
& + \frac{\|h'\|}{2\sqrt{n i(\theta_0)}} \left[\mathbb{E} \left(\left(\sum_{i=1}^n M(X_i) \right)^2 \middle| |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right) \right]^{\frac{1}{2}} \left[\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^4 \right] \right]^{\frac{1}{2}}.
\end{aligned} \tag{3.6}$$

In particular, for $h \in H_{bW}$ as in (2.6), (3.6) gives

$$\begin{aligned}
d_{bW} \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) & \leq \frac{1}{\sqrt{n}} \left(2 + \frac{1}{[i(\theta_0)]^{\frac{3}{2}}} \left[\mathbb{E} \left| \frac{d}{d\theta} \log f(X_1 | \theta_0) \right|^3 \right] \right) \\
& + \frac{1}{\sqrt{i(\theta_0)}} \sqrt{\text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right)} \sqrt{\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} \\
& + \frac{2}{\varepsilon^2} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \\
& + \frac{1}{2\sqrt{n i(\theta_0)}} \left[\mathbb{E} \left(\left(\sum_{i=1}^n M(X_i) \right)^2 \middle| |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right) \right]^{\frac{1}{2}} \left[\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^4 \right] \right]^{\frac{1}{2}}.
\end{aligned} \tag{3.7}$$

The following lemma, often known as Chebyshev's other inequality (see Fink and Jodeit (1984) for more information), is useful to obtain (3.7).

Lemma 3.1. *Let $M \geq 0$ be a random variable and $\varepsilon > 0$. For every continuous increasing function $f : [0, \infty) \rightarrow [0, \infty)$,*

$$\mathbb{E}[f(M)|M < \varepsilon] \leq \mathbb{E}[f(M)].$$

Proof of Lemma 3.1. Let $\varepsilon > 0$ and f a continuous increasing function with $f(m) \geq 0$ for $m > 0$. Then,

$$\begin{aligned} \mathbb{E}[f(M)] &= \mathbb{E}[f(M)|M < \varepsilon] \mathbb{P}(M < \varepsilon) + \mathbb{E}[f(M)|M \geq \varepsilon] \mathbb{P}(M \geq \varepsilon) \\ &= \mathbb{E}[f(M)|M < \varepsilon] (1 - \mathbb{P}(M \geq \varepsilon)) + \mathbb{E}[f(M)|M \geq \varepsilon] \mathbb{P}(M \geq \varepsilon) \\ &= \mathbb{E}[f(M)|M < \varepsilon] + \mathbb{P}(M \geq \varepsilon) (\mathbb{E}[f(M)|M \geq \varepsilon] - \mathbb{E}[f(M)|M < \varepsilon]) \\ &\geq \mathbb{E}[f(M)|M < \varepsilon] \quad \text{as } f(m) \text{ is increasing.} \end{aligned}$$

□

Proof of Theorem 3.1. The regularity conditions and the definition of the MLE ensure that $0 = l'(\hat{\theta}_n(x); x)$. A second order Taylor expansion of $l'(\hat{\theta}_n(x); x)$ about θ_0 gives

$$l''(\theta_0; x) (\hat{\theta}_n(x) - \theta_0) = -l'(\theta_0; x) - R_1(\theta_0; x), \quad (3.8)$$

where

$$R_1(\theta_0; x) = \frac{1}{2} (\hat{\theta}_n(x) - \theta_0)^2 l^{(3)}(\theta^*; x) \quad (3.9)$$

is the remainder term with θ^* lying between $\hat{\theta}_n(x)$ and θ_0 . Rearranging (3.8) and adding zero gives

$$-n i(\theta_0) (\hat{\theta}_n(x) - \theta_0) = -l'(\theta_0; x) - R_1(\theta_0; x) - (\hat{\theta}_n(x) - \theta_0) [l''(\theta_0; x) + n i(\theta_0)].$$

As $i(\theta_0) \neq 0$,

$$\sqrt{ni(\theta_0)} (\hat{\theta}_n(x) - \theta_0) = \frac{l'(\theta_0; x) + R_1(\theta_0; x) + R_2(\theta_0, x)}{\sqrt{ni(\theta_0)}},$$

where

$$R_2(\theta_0, x) = (\hat{\theta}_n(x) - \theta_0) (l''(\theta_0; x) + ni(\theta_0)). \quad (3.10)$$

For $Z \sim N(0, 1)$ and $h \in H$ given in (2.6),

$$\begin{aligned} & \left| \mathbb{E} \left[h \left(\sqrt{ni(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) \right] - \mathbb{E}[h(Z)] \right| \leq \\ & \left| \mathbb{E} \left[h \left(\frac{l'(\theta_0; \mathbf{X}) + R_1(\theta_0; \mathbf{X}) + R_2(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}} \right) - h \left(\frac{l'(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}} \right) \right] \right| \end{aligned} \quad (3.11)$$

$$+ \left| \mathbb{E} \left[h \left(\frac{l'(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}} \right) \right] - \mathbb{E}[h(Z)] \right|. \quad (3.12)$$

The upper bound for (3.12) is given in Proposition 3.2. With regards to (3.11), let for ease of presentation

$$\begin{aligned} C_1 &= C_1(h, \theta_0; \mathbf{X}) = h \left(\frac{l'(\theta_0; \mathbf{X}) + R_1(\theta_0; \mathbf{X}) + R_2(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}} \right) - h \left(\frac{l'(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}} \right) \\ C_2 &= C_2(h, \theta_0; \mathbf{X}) = h \left(\frac{l'(\theta_0; \mathbf{X}) + R_1(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}} \right) - h \left(\frac{l'(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}} \right). \end{aligned} \quad (3.13)$$

Using this notation, the triangle inequality yields

$$\begin{aligned} (3.11) &= |\mathbb{E}[C_1]| \\ &= \left| \mathbb{E} \left[h \left(\frac{l'(\theta_0; \mathbf{X}) + R_1(\theta_0; \mathbf{X}) + R_2(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}} \right) - h \left(\frac{l'(\theta_0; \mathbf{X}) + R_1(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}} \right) + C_2 \right] \right| \\ &\leq \frac{\|h'\|}{\sqrt{ni(\theta_0)}} \mathbb{E} |R_2(\theta_0; \mathbf{X})| + \mathbb{E} |C_2|, \end{aligned} \quad (3.14)$$

where the last inequality is due to a first order Taylor expansion of $h \left(\frac{l'(\theta_0; \mathbf{X}) + R_1(\theta_0; \mathbf{X}) + R_2(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}} \right)$ about $\frac{l'(\theta_0; \mathbf{X}) + R_1(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}}$. The definition of $R_2(\theta_0; x)$ in

(3.10) and the Cauchy - Schwarz inequality give that

$$\begin{aligned}
 \mathbb{E}|R_2(\theta_0; \mathbf{X})| &= \mathbb{E} \left| (ni(\theta_0) + l''(\theta_0; \mathbf{X})) (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right| \\
 &\leq \sqrt{\mathbb{E}[ni(\theta_0) + l''(\theta_0; \mathbf{X})]^2} \sqrt{\mathbb{E}[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2]} \\
 &= \sqrt{n \text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right)} \sqrt{\mathbb{E}[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2]}. \quad (3.15)
 \end{aligned}$$

To bound $\mathbb{E}|C_2|$ with C_2 as in (3.13), note that the term $R_1(\theta_0; \mathbf{X})$ is in general not uniformly bounded. For all x the rather crude bound $|C_2| \leq 2\|h\|$ is valid. If $|\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon$ then a better bound is available. Hence, we condition on whether $|\hat{\theta}_n(\mathbf{X}) - \theta_0| \geq \varepsilon$ or $|\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon$, with $\varepsilon = \varepsilon_0 > 0$ as in condition (R3). By Markov's inequality

$$\mathbb{P}(|\hat{\theta}_n(\mathbf{X}) - \theta_0| \geq \varepsilon) \leq \frac{\mathbb{E}[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2]}{\varepsilon^2}. \quad (3.16)$$

Using the law of total expectation,

$$\begin{aligned}
 |\mathbb{E}[C_2]| &\leq \mathbb{E}|C_2| = \mathbb{E}(|C_2| | |\hat{\theta}_n(\mathbf{X}) - \theta_0| \geq \varepsilon) \mathbb{P}(|\hat{\theta}_n(\mathbf{X}) - \theta_0| \geq \varepsilon) \\
 &\quad + \mathbb{E}(|C_2| | |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon) \mathbb{P}(|\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon).
 \end{aligned}$$

Using (3.16) for the first term and a first order Taylor expansion of $h \left(\frac{l'(\theta_0; \mathbf{X}) + R_1(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}} \right)$ about $\frac{l'(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}}$ for the second term gives

$$\begin{aligned}
 |\mathbb{E}[C_2]| &\leq \frac{2\|h\|}{\varepsilon^2} \mathbb{E}[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2] \\
 &\quad + \left| \mathbb{E} \left(\frac{R_1(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}} h'(t(\mathbf{X})) \mid |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right) \right| \mathbb{P}(|\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon) \\
 &\leq \frac{2\|h\|}{\varepsilon^2} \mathbb{E}[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2] \\
 &\quad + \frac{\|h'\|}{2\sqrt{ni(\theta_0)}} \mathbb{E} \left((\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \left| l^{(3)}(\theta^*; \mathbf{X}) \right| \mid |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right),
 \end{aligned}$$

where $t(\mathbf{X})$ lies between $\frac{l'(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}}$ and $\frac{l'(\theta_0; \mathbf{X}) + R_1(\theta_0; \mathbf{X})}{\sqrt{ni(\theta_0)}}$. Using (3.14), (3.15) and the fact that for $|\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon$,

$$\begin{aligned} |R_1(\theta_0; \mathbf{X})| &\leq \frac{1}{2}(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \sup_{\theta: |\theta - \theta_0| < \varepsilon} |l^{(3)}(\theta; \mathbf{X})| \\ &= \frac{1}{2}(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \sup_{\theta: |\theta - \theta_0| < \varepsilon} \left| \sum_{i=1}^n \frac{d^3}{d\theta^3} \log f(X_i | \theta) \right| \\ &\leq \frac{1}{2}(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \sum_{i=1}^n \left\{ \sup_{\theta: |\theta - \theta_0| < \varepsilon} \left| \frac{d^3}{d\theta^3} \log f(X_i | \theta) \right| \right\} \\ &\leq \frac{1}{2}(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \sum_{i=1}^n M(X_i), \end{aligned} \quad (3.17)$$

then

$$\begin{aligned} |E[C_1]| &\leq \frac{2\|h\|}{\varepsilon^2} E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \\ &\quad + \frac{\|h'\|}{\sqrt{i(\theta_0)}} \sqrt{\text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right)} \sqrt{E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} \\ &\quad + \frac{\|h'\|}{2\sqrt{ni(\theta_0)}} E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \sum_{i=1}^n M(X_i) \middle| |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right]. \end{aligned}$$

The next step is based on the Cauchy - Schwarz inequality and the fact that

$$E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^4 \middle| |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right] \leq E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^4 \right], \quad (3.18)$$

due to Lemma 3.1, giving

$$\begin{aligned} |E[C_1]| &\leq \frac{2\|h\|}{\varepsilon^2} E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \\ &\quad + \frac{\|h'\|}{\sqrt{i(\theta_0)}} \sqrt{\text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right)} \sqrt{E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} \\ &\quad + \frac{\|h'\|}{2\sqrt{ni(\theta_0)}} \left[E \left(\left(\sum_{i=1}^n M(X_i) \right)^2 \middle| |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right) \right]^{\frac{1}{2}} \left[E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^4 \right] \right]^{\frac{1}{2}}. \end{aligned} \quad (3.19)$$

The first result of the theorem in (3.6) is obtained using (3.4) and (3.19). Taking $h \in H_{bW}$ as in (2.6), then $\|h\| \leq 1$ and $\|h'\| \leq 1$ leading to the result in (3.7). \square

Remark 3.1. (1) Using (3.10), if $l''(\theta_0; x) \equiv -n i(\theta_0)$ then $R_2(\theta_0; x) \equiv 0$ and the second term of the bounds given in Theorem 3.1 vanishes.

(2) For $M_1 = M_1(\theta_0)$ a constant that may depend on the unknown parameter θ_0 and for $0 < \varepsilon = \varepsilon(\theta_0) = \varepsilon_0$ as in (R3), if $\left| \frac{d^3}{d\theta^3} \log f(x|\theta) \right| \leq M_1$ then $M(x) = M_1$ and the bound is simpler as the conditional expectation in the last term of the bound vanishes. This is because step (3.18) in the proof is not any more necessary since

$$\begin{aligned} & \mathbb{E} \left[\sum_{i=1}^n M(X_i) (\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \mid |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right] \\ &= n M_1 \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \mid |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right] \leq n M_1 \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]. \end{aligned}$$

Thus, in this case the bound on the bounded Wasserstein distance becomes

$$\begin{aligned} d_{bW} \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) &\leq \frac{1}{\sqrt{n}} \left(2 + \frac{1}{[i(\theta_0)]^{\frac{3}{2}}} \left[\mathbb{E} \left| \frac{d}{d\theta} \log f(X_1 | \theta_0) \right|^3 \right] \right) \\ &+ \frac{1}{\sqrt{i(\theta_0)}} \sqrt{\text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right)} \sqrt{\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} \\ &+ \frac{2}{\varepsilon^2} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \\ &+ \frac{\sqrt{n} M_1}{2 \sqrt{i(\theta_0)}} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]. \end{aligned}$$

Furthermore, if the third derivative of the log-likelihood function is uniformly bounded by a constant, $\left| \frac{d^3}{d\theta^3} \log f(x|\theta) \right| \leq \tilde{B}$, then we do not need to use ε and with C_1 as in (3.13)

$$\begin{aligned} |\mathbb{E}[C_1]| &\leq \frac{\|h'\|}{\sqrt{i(\theta_0)}} \sqrt{\text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right)} \sqrt{\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} \\ &+ \tilde{B} \frac{\|h'\| \sqrt{n}}{2 \sqrt{i(\theta_0)}} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \end{aligned}$$

leading to

$$\begin{aligned} d_{bW} \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) &\leq \frac{1}{\sqrt{n}} \left(2 + \frac{1}{[i(\theta_0)]^{\frac{3}{2}}} \left[E \left| \frac{d}{d\theta} \log f(X_1 | \theta_0) \right|^3 \right] \right) \\ &+ \frac{1}{\sqrt{i(\theta_0)}} \sqrt{\text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right)} \sqrt{E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} \\ &+ \frac{\tilde{B}\sqrt{n}}{2\sqrt{i(\theta_0)}} E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]. \end{aligned}$$

(3) The rate of convergence of the mean squared error, $E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]$, is $\mathcal{O} \left(\frac{1}{n} \right)$.

This result is obtained using that

$$E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] = \text{Var} \left[\hat{\theta}_n(\mathbf{X}) \right] + \text{bias}^2 \left[\hat{\theta}_n(\mathbf{X}) \right]. \quad (3.20)$$

Under the standard asymptotics (from the regularity conditions (R1)-(R4)) the MLE is asymptotically efficient,

$$n \text{Var}[\hat{\theta}_n(\mathbf{X})] \xrightarrow{n \rightarrow \infty} [i(\theta_0)]^{-1},$$

and hence the variance of the MLE is of order $\frac{1}{n}$. In addition, from Theorem 2.2 the bias of the MLE is equal to zero up to the order $\frac{1}{\sqrt{n}}$; see also Cox and Snell (1968), where it is stated that the bias of the MLE is $o \left(\frac{1}{n} \right)$, but no explicit conditions are given. Combining these two results and using (3.20) shows that the mean squared error of the MLE is of order $\frac{1}{n}$. In the examples that follow, the remaining terms in the bound are of order at most $\frac{1}{\sqrt{n}}$.

(4) When the calculation of $E \left(\left| \frac{d}{d\theta} \log f(X_1 | \theta_0) \right|^3 \right)$ is awkward, Hölder's inequality can be used, giving $E \left(\left| \frac{d}{d\theta} \log f(X_1 | \theta_0) \right|^3 \right) \leq \left[E \left(\frac{d}{d\theta} \log f(X_1 | \theta_0) \right)^4 \right]^{\frac{3}{4}}$.

3.3 Single-parameter exponential families

This section specifies Theorem 3.1 for the distribution of the MLE for single-parameter exponential family distributions. Many popular distributions which have the same underlying structure based on relatively simple properties are exponential families, such as the normal, Gamma and Laplace distributions. Generalisations of exponential families can be found in Lauritzen (1988) and Berk (1972).

This section focuses on the simplest case of a scalar parameter. The distribution of a random variable, X , is said to be a *single-parameter exponential family distribution* if the probability density (or mass) function is of the form

$$f(x|\theta) = \exp \{k(\theta)T(x) - A(\theta) + S(x)\} \mathbb{1}_{\{x \in B\}}, \quad (3.21)$$

where the set $B = \{x : f(x|\theta) > 0\}$ is the support of X and does not depend on θ ; $k(\theta)$ and $A(\theta)$ are functions of the parameter; $T(x)$ and $S(x)$ are functions only of the data. The choice of the functions $k(\theta)$ and $T(X)$ is not unique. If $k(\theta) = \theta$ we have the so-called *canonical case*. In this case θ and $T(X)$ are called the *natural parameter* and *natural observation* (Casella and Berger, 2002). For x_1, x_2, \dots, x_n the data points, we make the following assumptions, where (Ass.Ex.1)-(Ass.Ex.3) are used to ensure the existence and uniqueness of the MLE and (A1)-(A4) follow from the regularity conditions in Section 2.2.

(Ass.Ex.1) $\Theta \subset \mathbb{R}$ is open and connected;

(Ass.Ex.2) $\lim_{\theta \rightarrow \partial\Theta} [k(\theta) \sum_{i=1}^n T(x_i) - nA(\theta) + \sum_{i=1}^n S(x_i)] = -\infty$;

(Ass.Ex.3) We have $k''(\theta) \sum_{i=1}^n T(x_i) - nA''(\theta) < 0$ at every point θ , for which

$$k'(\theta) \sum_{i=1}^n T(x_i) - nA'(\theta) = 0;$$

(A1) $k'(\theta) \neq 0, \forall \theta \in \Theta$ and $D(\theta) = \frac{A'(\theta)}{k'(\theta)}$ is invertible;

(A2) $l(\theta; x)$ is three times continuously differentiable with respect to θ , which is equivalent to both $k^{(3)}(\theta)$ and $A^{(3)}(\theta)$ existing and being continuous. In addition, integration of the density function over x and differentiation with respect to θ are three times interchangeable;

(A3) for any $\theta_0 \in \Theta$, there exists a positive number ε_0 and a function $M(x)$ (both of which may depend on θ_0) such that

$$\left| k^{(3)}(\theta)T(x) - A^{(3)}(\theta) \right| \leq M(x), \quad \forall x \in B, \quad \forall \theta \in \Theta : |\theta - \theta_0| < \varepsilon_0,$$

with $E[M(X)] < \infty$;

(A4) $\text{Var}[T(X)] > 0$;

(A5) $E|T(X) - D(\theta_0)|^3$ exists. This assumption is required for meaningful bounds.

Proposition 3.3. *In the case of the single-parameter exponential family, the assumptions (A1)-(A4) imply the regularity conditions (R1)-(R4) as expressed in Chapter 3.*

Proof. First, according to (Davison, 2008, p.172), an exponential family of order p is written in its minimal representation if the set $\{1, k_1(\theta), k_2(\theta), \dots, k_p(\theta)\}$ is linearly independent. From (A1), because $k'(\theta) \neq 0$, then $k(\theta)$ is not constant in θ . Therefore, a single-parameter exponential family is always in minimal representation and thus identifiable (Geyer, 2013). Continuing, (A2) and (A3) are just the versions of (R2) and (R3) in the case that the random variables follow a single-parameter exponential family distribution. Finally, using (3.21), simple steps yield

$$i(\theta_0) = \text{Var} \left(\frac{d}{d\theta} \log f(X|\theta_0) \right) = \text{Var}(k'(\theta_0)T(X) - A'(\theta_0)) = [k'(\theta_0)]^2 \text{Var}(T(X)),$$

which indicates that since, from (A1), $k'(\theta) \neq 0$, $\forall \theta \in \Theta$, then

$$i(\theta_0) > 0 \Leftrightarrow \text{Var}(T(X)) > 0.$$

Thus, (A4) also implies (R4). □

Since exponential family distributions have the same structure for the probability density (or mass) function, then it would make sense for the MLE to have a specific general form for these distributions. This is shown in the proposition that follows.

Proposition 3.4. *Let X_1, X_2, \dots, X_n be i.i.d. random variables from a single-parameter exponential family distribution. The probability density (or mass) function is given in (3.21). Assume that (A1), (A2) and (Ass.Ex.1)-(Ass.Ex.3) hold. The MLE in this case is unique and equal to*

$$\hat{\theta}_n(\mathbf{X}) = D^{-1} \left(\frac{1}{n} \sum_{i=1}^n T(X_i) \right). \quad (3.22)$$

Proof. From (A2) and (3.21),

$$l^{(j)}(\theta; \mathbf{x}) = k^{(j)}(\theta) \sum_{i=1}^n T(x_i) - nA^{(j)}(\theta), \text{ for } j = 1, 2, 3$$

and straightforward calculations show that $l'(\hat{\theta}_n(\mathbf{X}); \mathbf{X}) = 0$ has the unique solution $\hat{\theta}_n(\mathbf{X}) = D^{-1} \left(\frac{1}{n} \sum_{i=1}^n T(X_i) \right)$, where $D(\theta) = \frac{A'(\theta)}{k'(\theta)}$ is invertible by (A1). Using now (Ass.Ex.1) and (Ass.Ex.3) gives that $\hat{\theta}_n(\mathbf{X}) = D^{-1} \left(\frac{1}{n} \sum_{i=1}^n T(X_i) \right)$ is the unique MLE in the case of single-parameter exponential family distributions. □

The following corollary of Theorem 3.1 gives an upper bound on the bounded Wasserstein distance between the distribution of the MLE as in (3.22) and the normal distribution.

Corollary 3.1. *Let X_1, X_2, \dots, X_n be i.i.d. random variables with the probability density (or mass) function of a single-parameter exponential family distribution. Assume that*

(A1)-(A5) are satisfied and that (Ass.Ex.1)-(Ass.Ex.3) hold. With $Z \sim N(0, 1)$, $h \in H_{bW}$ as in (2.6), $0 < \varepsilon = \varepsilon(\theta_0) = \varepsilon_0$ such that $\theta \in (\theta_0 - \varepsilon, \theta_0 + \varepsilon) \cap \Theta$ where ε_0 is as in (A3) and for $E \left(\left(\sum_{i=1}^n M(X_i) \right)^2 \middle| \left| \hat{\theta}_n(\mathbf{X}) - \theta_0 \right| < \varepsilon \right) < \infty$, with $M(X)$ as in (A3), it holds that

$$\begin{aligned} d_{bW} \left(\sqrt{ni(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) &\leq \frac{1}{\sqrt{n}} \left(2 + \frac{E|T(X_1) - D(\theta_0)|^3}{[\text{Var}[T(X_1)]]^{\frac{3}{2}}} \right) \\ &+ \frac{|k''(\theta_0)|}{\sqrt{i(\theta_0)}} \sqrt{\text{Var}(T(X_1))} \sqrt{E[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2]} \\ &+ \frac{2}{\varepsilon^2} E[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2] \\ &+ \frac{1}{2\sqrt{ni(\theta_0)}} \left[E \left(\left(\sum_{i=1}^n M(X_i) \right)^2 \middle| \left| \hat{\theta}_n(\mathbf{X}) - \theta_0 \right| < \varepsilon \right) \right]^{\frac{1}{2}} \left[E[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^4] \right]^{\frac{1}{2}}. \end{aligned} \quad (3.23)$$

Proof. For the first term of the bound let $Y_i = Y_i(X_i; \theta_0) = \frac{\frac{d}{d\theta} \log f(X_i | \theta_0)}{\sqrt{i(\theta_0)}}$, $i = 1, 2, \dots, n$. Proposition 3.2 requires the calculation of $E|Y_1|^3$ and

$$\frac{d}{d\theta} \log f(X_i | \theta) \Big|_{\theta=\theta_0} = k'(\theta_0)T(X_i) - A'(\theta_0)$$

yields

$$\begin{aligned} E \left| \frac{d}{d\theta} \log f(X_i | \theta_0) \right|^3 &= E |k'(\theta_0)T(X_i) - A'(\theta_0)|^3 = E |k'(\theta_0)(T(X_i) - D(\theta_0))|^3 \\ &= |k'(\theta_0)|^3 E|T(X_i) - D(\theta_0)|^3, \quad \forall i \in \{1, 2, \dots, n\}. \end{aligned}$$

In addition, $i(\theta_0) = \text{Var} \left[\frac{d}{d\theta} \log f(X_i | \theta_0) \right] = [k'(\theta_0)]^2 \text{Var}[T(X_i)] > 0$ from (A1) and (A4). These quantities can now be applied to get the first term of the bound in (3.23). For the remaining terms, using (3.7),

$$\text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right) = \text{Var} (k''(\theta_0)T(X_1) - A''(\theta_0)) = [k''(\theta_0)]^2 \text{Var}(T(X_1)),$$

which yields the assertion of the corollary. \square

Remark 3.2. In the canonical case, $k''(\theta_0) \equiv 0$ and the second term of the bound vanishes. Also, $\frac{d^2}{d\theta^2} \log f(x|\theta) = -A''(\theta)$ and $i(\theta_0) = A''(\theta_0)$. In addition, $\frac{d^3}{d\theta^3} \log f(x|\theta) = -A^{(3)}(\theta)$ is independent of the random variables. Thus, using (2) of Remark 3.1 we get

$$d_{bW} \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) \leq \frac{1}{\sqrt{n}} \left(2 + \frac{E|T(X_1) - D(\theta_0)|^3}{[\text{Var}[T(X_1)]]^{\frac{3}{2}}} \right) + E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \left(\frac{2}{\varepsilon^2} + \frac{\sqrt{n}}{2\sqrt{A''(\theta_0)}} \sup_{\theta: |\theta - \theta_0| < \varepsilon} |A^{(3)}(\theta)| \right). \quad (3.24)$$

3.3.1 The exponential distribution

Here, we consider two examples using the exponential distribution, firstly, in its canonical form, and then under a change of parametrisation. In each case, we check if all necessary assumptions are satisfied in order to obtain the bound in (3.23).

The canonical case

In the case of X_1, X_2, \dots, X_n exponentially distributed, $\text{Exp}(\theta)$, i.i.d. random variables where $\theta > 0$, the probability density function is

$$f(x|\theta) = \theta \exp\{-\theta x\} = \exp\{\log \theta - \theta x\} = \exp\{k(\theta)T(x) - A(\theta) + S(x)\} \mathbb{1}_{\{x \in B\}}, \quad (3.25)$$

where $B = (0, \infty)$, $\theta \in \Theta = (0, \infty)$, $T(x) = -x$, $k(\theta) = \theta$, $A(\theta) = -\log \theta$ and $S(x) = 0$. Hence $\text{Exp}(\theta)$ is a single-parameter canonical exponential family distribution.

Corollary 3.2. Let X_1, X_2, \dots, X_n be i.i.d. random variables that follow the $\text{Exp}(\theta_0)$ distribution. The MLE exists, it is unique, equal to $\hat{\theta}_n(\mathbf{X}) = \frac{1}{\bar{X}}$ and (A1)-(A5) are satisfied. For $h \in H_{bW}$ as in (2.6),

$$d_{bW} \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) < \frac{4.41456}{\sqrt{n}} + \frac{8(n+2)}{(n-1)(n-2)} + \frac{8\sqrt{n}(n+2)}{(n-1)(n-2)}. \quad (3.26)$$

Proof. We check the existence and uniqueness of the MLE through the assumptions (Ass.Ex.1)-(Ass.Ex.3). (Ass.Ex.1) holds since the parameter space Θ is open and connected. For (Ass.Ex.2),

$$k(\theta) \sum_{i=1}^n T(x_i) - nA(\theta) + \sum_{i=1}^n S(x_i) = -\theta \sum_{i=1}^n x_i + n \log \theta \xrightarrow[\theta \rightarrow \infty]{\theta \rightarrow 0} -\infty.$$

For (Ass.Ex.3), simple steps give

$$l'(\theta; \mathbf{x}) = \frac{n}{\theta} - \sum_{i=1}^n x_i, \quad l''(\theta; \mathbf{x}) = -\frac{n}{\theta^2}. \quad (3.27)$$

Thus, $l'(\theta; \mathbf{X}) = 0$ if and only if $\theta = \hat{\theta}_n(\mathbf{X}) = \frac{n}{\sum_{i=1}^n X_i}$. In addition, $l''(\hat{\theta}_n(\mathbf{X}); \mathbf{X}) < 0$ and thus (Ass.Ex.3) holds. The MLE is $\hat{\theta}_n(\mathbf{X}) = \frac{n}{\sum_{i=1}^n X_i}$ and it is unique. To verify the assumptions (A1)-(A5), for (A1), $k(\theta) = \theta$ and thus $k'(\theta) = 1 \neq 0$, $\forall \theta \in \Theta$ and $D(\theta) = A'(\theta) = -\frac{1}{\theta}$ with $D^{-1}(\theta) = D(\theta) = -\frac{1}{\theta}$. For (A2), the first two derivatives for the sample are shown in (3.27). Also, $l^{(3)}(\theta; \mathbf{x}) = \frac{2n}{\theta^3}$. Using (3.25) gives

$$\begin{aligned} \int_0^\infty \frac{d}{d\theta} f(x|\theta) dx &= \int_0^\infty \exp\{-\theta x\} - \theta x \exp\{-\theta x\} dx = \frac{1}{\theta} - E(X) = 0 \\ \int_0^\infty \frac{d^2}{d\theta^2} f(x|\theta) dx &= \int_0^\infty -2x \exp\{-\theta x\} + \theta x^2 \exp\{-\theta x\} dx \\ &= -\frac{2}{\theta} E(X) + E(X^2) = 0 \\ \int_0^\infty \frac{d^3}{d\theta^3} f(x|\theta) dx &= \int_0^\infty 3x^2 \exp\{-\theta x\} - \theta x^3 \exp\{-\theta x\} dx \\ &= \frac{3}{\theta} E(X^2) - E(X^3) = 0. \end{aligned}$$

Now,

$$\left| k^{(3)}(\theta)T(x) - A^{(3)}(\theta) \right| = \left| -A^{(3)}(\theta) \right| = \left| \frac{2}{\theta^3} \right|, \quad (3.28)$$

which needs to get bounded in a neighbourhood of θ_0 in order for (A3) to hold. Let $\varepsilon = \varepsilon(\theta_0) = \varepsilon_0$ be a positive constant such that $\theta \in (\theta_0 - \varepsilon, \theta_0 + \varepsilon) \cap \Theta$. From (3.28), if $\varepsilon > \theta_0$ then $(\theta_0 - \varepsilon, \theta_0 + \varepsilon) \cap \Theta = (0, \theta_0 + \varepsilon)$ and the result in (3.28) can not be upper bounded in $(0, \theta_0 + \varepsilon)$. Therefore, $0 < \varepsilon < \theta_0$ with $(\theta_0 - \varepsilon, \theta_0 + \varepsilon) \cap \Theta = (\theta_0 - \varepsilon, \theta_0 + \varepsilon)$ and

$$\left| k^{(3)}(\theta)T(x) - A^{(3)}(\theta) \right| = \left| -A^{(3)}(\theta) \right| = \left| \frac{2}{\theta^3} \right| < \frac{2}{(\theta_0 - \varepsilon)^3} = M(x),$$

with $E[M(X)] = \frac{2}{(\theta_0 - \varepsilon)^3} < \infty$, and hence (A3) holds. We choose ε to be in the middle of the interval $(0, \theta_0)$ as there is a trade-off on its choice for the second and third term of the bound in (3.24). The second term is a product of the mean squared error with $\frac{2}{\varepsilon^2}$ indicating that we should choose ε away from zero. However, the third term of the bound includes $\sup_{\theta: |\theta - \theta_0| < \varepsilon} \left| A^{(3)}(\theta) \right| = \frac{2}{(\theta_0 - \varepsilon)^3}$ and a wise choice for ε would be such that it is not close to θ_0 . An optimisation process with respect to ε becomes quite tedious and therefore we choose ε to be the midpoint of $(0, \theta_0)$, which is away from both zero and θ_0 and also behaves very well. Now, since $T(x) = x$, we have that $\text{Var}(T(X)) = \text{Var}(X) = \frac{1}{\theta_0^2} > 0$ and (A4) is also satisfied. For (A5), basic calculations of integrals show that $E|T(X) - D(\theta_0)|^3 = E\left|\frac{1}{\theta_0} - X\right|^3 < \frac{2.41456}{\theta_0^3}$ and (A5) is satisfied. We have already showed that $\sup_{\theta: |\theta - \theta_0| < \varepsilon} \left| A^{(3)}(\theta) \right| = \frac{2}{(\theta_0 - \varepsilon)^3}$. For $\varepsilon = \frac{\theta_0}{2}$, $\sup_{\theta: |\theta - \theta_0| < \varepsilon} \left| A^{(3)}(\theta) \right| = \frac{16}{\theta_0^3}$. In addition, since $X_i \sim \text{Exp}(\theta)$, $\forall i \in \{1, 2, \dots, n\}$ we have that $\bar{X} \sim G(n, n\theta)$, with $G(\alpha, \beta)$ being the Gamma distribution with shape parameter α and rate parameter β . Using now the results in pp.70-73 of Kendall and Stuart (1969) gives that

$$\begin{aligned} E[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2] &= E[(\hat{\theta}_n(\mathbf{X})^2] - 2\theta_0 E[\hat{\theta}_n(\mathbf{X})] + \theta_0^2 \\ &= \frac{(n\theta_0)^2}{(n-1)(n-2)} - \frac{2n\theta_0^2}{n-1} + \theta_0^2 = \frac{(n+2)\theta_0^2}{(n-1)(n-2)}. \end{aligned}$$

Using (3.24) we obtain the result in (3.26). □

Remark 3.3. (1) The rate of convergence of the bound is $\frac{1}{\sqrt{n}}$. Note that the bound does not depend on the value of θ_0 .

(2) The calculation of $E \left| \frac{1}{\theta_0} - X \right|^3$ requires a significant amount of steps. Instead, one could use Hölder's inequality with $E \left| \frac{1}{\theta_0} - X \right|^3 \leq \left[E \left(\frac{1}{\theta_0} - X \right)^4 \right]^{\frac{3}{4}} = \frac{9^{\frac{3}{4}}}{\theta_0^3}$ using the results in pp.70-73 of Kendall and Stuart (1969).

The non-canonical case

Let X_1, X_2, \dots, X_n be i.i.d. random variables from $\text{Exp}(\frac{1}{\theta})$, with p.d.f.

$$\begin{aligned} f(x|\theta) &= \frac{1}{\theta} \exp \left\{ -\frac{1}{\theta} x \right\} = \exp \left\{ -\log \theta - \frac{1}{\theta} x \right\} \\ &= \exp \{ k(\theta) T(x) - A(\theta) + S(x) \} \mathbb{1}_{\{x \in B\}}, \end{aligned} \quad (3.29)$$

where $B = (0, \infty)$, $\theta \in \Theta = (0, \infty)$, $T(x) = -x$, $k(\theta) = \frac{1}{\theta}$, $A(\theta) = \log \theta$ and $S(x) = 0$. Thus, $\text{Exp}(\frac{1}{\theta})$ is a non-canonical exponential family distribution.

Corollary 3.3. Let X_1, X_2, \dots, X_n be i.i.d. random variables that follow the $\text{Exp}(\frac{1}{\theta_0})$ distribution. The MLE exists, it is unique, equal to $\hat{\theta}_n(\mathbf{X}) = \bar{X}$ and (A1)-(A5) are satisfied. For $h \in H_{bW}$ as in (2.6),

$$\begin{aligned} d_{bW} \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) &< \frac{4.41456}{\sqrt{n}} + \frac{8}{n} + \frac{2}{\sqrt{n}} \\ &+ \frac{1}{\sqrt{n}} \left(80 \left[3 \left(\frac{2}{n} + 1 \right) \right]^{\frac{1}{2}} \right). \end{aligned} \quad (3.30)$$

Proof. Corollary 3.2 and the invariance property of the MLE give that $\hat{\theta}_n(\mathbf{X}) = \bar{X}$. To check (A1)-(A5), condition (A1) holds since $k'(\theta) = -\frac{1}{\theta^2} \neq 0$, $\forall \theta \in \Theta$ and $D^{-1}(\theta) = D(\theta) = -\theta$. Equation (3.29) and simple calculations show that the log-likelihood is three times differentiable with respect to the parameter. For one observation, x , we have

that

$$\begin{aligned}
\int_0^\infty \frac{d}{d\theta} f(x|\theta) dx &= \int_0^\infty -\frac{1}{\theta^2} \exp\left\{-\frac{1}{\theta}x\right\} + \frac{x}{\theta^3} \exp\left\{-\frac{1}{\theta}x\right\} dx \\
&= -\frac{1}{\theta} + \frac{1}{\theta^2} E(X) = 0 \\
\int_0^\infty \frac{d^2}{d\theta^2} f(x|\theta) dx &= \int_0^\infty \exp\left\{-\frac{1}{\theta}x\right\} \left(\frac{2}{\theta^3} - \frac{4x}{\theta^4} + \frac{x^2}{\theta^5}\right) dx \\
&= \frac{2}{\theta^2} - \frac{4}{\theta^3} E(X) + \frac{1}{\theta^4} E(X^2) = 0 \\
\int_0^\infty \frac{d^3}{d\theta^3} f(x|\theta) dx &= \int_0^\infty \exp\left\{-\frac{1}{\theta}x\right\} \left(-\frac{6}{\theta^4} + \frac{18x}{\theta^5} - \frac{9x^2}{\theta^6} + \frac{x^3}{\theta^7}\right) dx \\
&= -\frac{6}{\theta^3} + \frac{18}{\theta^4} E(X) - \frac{9}{\theta^5} E(X^2) + \frac{1}{\theta^6} E(X^3) = 0
\end{aligned}$$

and thus (A2) is satisfied. Now, using the triangle inequality,

$$\left|k^{(3)}(\theta)T(x) - A^{(3)}(\theta)\right| = \left|\frac{6x}{\theta^4} - \frac{2}{\theta^3}\right| \leq \frac{6x}{\theta^4} + \frac{2}{\theta^3}. \quad (3.31)$$

To upper bound (3.31) in a neighbourhood of θ so that (A3) holds, we take as in the canonical case $0 < \varepsilon < \theta_0$, since if $\varepsilon > \theta_0$ then $(\theta_0 - \varepsilon, \theta_0 + \varepsilon) \cap \Theta = (0, \theta_0 + \varepsilon)$ and (3.31) cannot get bounded. Thus, $\varepsilon < \theta_0$ and $(\theta_0 - \varepsilon, \theta_0 + \varepsilon) \cap \Theta = (\theta_0 - \varepsilon, \theta_0 + \varepsilon)$ which leads to

$$\left|k^{(3)}(\theta)T(x) - A^{(3)}(\theta)\right| = \left|\frac{6x}{\theta^4} - \frac{2}{\theta^3}\right| \leq \frac{6x}{\theta^4} + \frac{2}{\theta^3} \leq \frac{6x}{(\theta_0 - \varepsilon)^4} + \frac{2}{(\theta_0 - \varepsilon)^3} = M(x), \quad (3.32)$$

with $E[M(X)] = \frac{6\theta_0}{(\theta_0 - \varepsilon)^4} + \frac{2}{(\theta_0 - \varepsilon)^3} < \infty$ showing that (A3) is satisfied. Since $T(x) = x$, then $\text{Var}(T(X)) = \text{Var}(X) = \theta_0^2$ which is positive and (A4) holds. For (A5),

$$E|T(X) - D(\theta_0)|^3 = E[\theta_0 - X]^3 < 2.41456\theta_0^3.$$

The mean squared error is $E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] = E \left[(\bar{X} - \theta_0)^2 \right] = \frac{\theta_0^2}{n}$. Using (3.32) we have that

$$\sum_{i=1}^n M(X_i) = \frac{6}{(\theta_0 - \varepsilon)^4} \sum_{i=1}^n X_i + \frac{2n}{(\theta_0 - \varepsilon)^3} = \frac{2n}{(\theta_0 - \varepsilon)^4} (3\hat{\theta}_n(\mathbf{X}) + \theta_0 - \varepsilon).$$

Therefore,

$$\begin{aligned} & \left[E \left(\left(\sum_{i=1}^n M(X_i) \right)^2 \middle| |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right) \right]^{\frac{1}{2}} \\ &= \left[E \left(\left(\frac{2n}{(\theta_0 - \varepsilon)^4} (3\hat{\theta}_n(\mathbf{X}) + \theta_0 - \varepsilon) \right)^2 \middle| |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right) \right]^{\frac{1}{2}} \\ &\leq \frac{2n}{(\theta_0 - \varepsilon)^4} \left[E \left((3|\hat{\theta}_n(\mathbf{X}) - \theta_0| + 4\theta_0 - \varepsilon)^2 \middle| |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right) \right]^{\frac{1}{2}} \\ &< \frac{2n}{(\theta_0 - \varepsilon)^4} \left[(2\varepsilon + 4\theta_0)^2 \right]^{\frac{1}{2}}, \quad \text{as } 4\theta_0 - \varepsilon > 0 \\ &= \frac{4n(2\theta_0 + \varepsilon)}{(\theta_0 - \varepsilon)^4}. \end{aligned}$$

The quantity $\left[E (\hat{\theta}_n(\mathbf{X}) - \theta_0)^4 \right]^{\frac{1}{2}}$ is calculated using the results in p.73 and the equations (3.38), p.70 of Kendall and Stuart (1969) along with the fact that $\hat{\theta}_n(\mathbf{X}) = \bar{X} \sim G\left(n, \frac{n}{\theta_0}\right)$. We obtain that $E(\hat{\theta}_n(\mathbf{X}) - \theta_0)^4 = \frac{3\theta_0^4}{n^2} \left(\frac{2}{n} + 1\right)$. Therefore,

$$\begin{aligned} & \left[E \left(\left(\sum_{i=1}^n M(X_i) \right)^2 \middle| |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right) \right]^{\frac{1}{2}} \left[E (\hat{\theta}_n(\mathbf{X}) - \theta_0)^4 \right]^{\frac{1}{2}} \\ &< \frac{4n(2\theta_0 + \varepsilon)}{(\theta_0 - \varepsilon)^4} \left[\frac{3\theta_0^4}{n^2} \left(\frac{2}{n} + 1\right) \right]^{\frac{1}{2}} = \frac{4(2\theta_0 + \varepsilon)}{(\theta_0 - \varepsilon)^4} \left[3\theta_0^4 \left(\frac{2}{n} + 1\right) \right]^{\frac{1}{2}}. \end{aligned} \quad (3.33)$$

We are now in a position to apply the general result in (3.1) to the specific case of i.i.d. random variables, which follow the exponential distribution with rate parameter $\frac{1}{\theta_0}$. There is, as in the canonical case, a trade-off on the choice of ε . On the one hand, the bound related to the third term in (3.23) requires ε to be large. On the other hand,

for the upper bound of $\left[\mathbb{E} \left((\sum_{i=1}^n M(X_i))^2 \mid |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right) \right]^{\frac{1}{2}}$ as shown in (3.33) it would be better to have ε away from θ_0 . It has already been explained that $\varepsilon \in (0, \theta_0)$ and we choose again $\varepsilon = \frac{\theta_0}{2}$, which yields the result in (3.30). \square

Remark 3.4. (1) *The rate of convergence related to the sample size of the above upper bound is $\frac{1}{\sqrt{n}}$ and the bound does not depend on θ_0 .*

(2) *The first term of the upper bound in (3.30) is the same as the one in (3.26). However, the rest of the bound is larger in (3.30) than in (3.26) $\forall n \in \mathbb{N}$.*

(3) *In the specific setting of independent, exponentially distributed random variables with rate parameter $\frac{1}{\theta_0}$, the MLE exists, it is unique and equal to $\hat{\theta}_n(\mathbf{X}) = \bar{X}$. Define $W = \frac{\sqrt{n}(\bar{X} - \theta_0)}{\theta_0} = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i$, where $Y_i = \frac{X_i - \theta_0}{\theta_0}$ are independent, zero mean and unit variance random variables. In addition, $\mathbb{E}(W) = 0$ and $\text{Var}(W) = \frac{n}{\theta_0^2} \text{Var}(\bar{X}) = \frac{1}{n\theta_0^2} \sum_{i=1}^n \text{Var}(X_i) = 1$. Therefore, the result in (2.12) gives*

$$d_{bW} \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) \leq \frac{1}{\sqrt{n}} \left(2 + \frac{1}{\theta_0^3} \mathbb{E}|X_1 - \theta_0|^3 \right) < \frac{4.41456}{\sqrt{n}}. \quad (3.34)$$

The upper bound given in (3.34) as a result of the direct use of Stein's method is smaller than the upper bound given in (3.30) using the general method from Section 3.2. However, in order to apply Stein's method directly, we require the quantity $\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{x}) - \theta_0)$ to be a sum of independent random variables. The general method, on the other hand, gives an upper bound for (3.1), whatever the MLE is, as long as the assumptions (A1)-(A5) hold.

In the next subsection results based on simulations are presented for both the canonical and the non-canonical parametrisation of the exponential distribution and for various sample sizes.

3.3.2 Empirical results

Here, we study the accuracy of our bounds by simulations. We start by generating 10,000 trials of n random independent observations, x , from the exponential distribution. The means for the canonical and the non-canonical case are equal to 1 and 2, respectively. The MLE, $\hat{\theta}_n(\mathbf{X})$, is evaluated in each trial, which in turn gives a vector of 10,000 values. After these values are standardised, we apply to them the function $h(x) = \frac{1}{x^2+2}$ and we calculate their sample mean, denoted by $\hat{E} \left[h \left(\sqrt{ni(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) \right]$. The function h is a member of the class H_{bW} as in (2.6) and $\|h\| = 0.5$, $\|h'\| = \frac{3\sqrt{1.5}}{16}$. We use these values to calculate the bound in (3.6). We define

$$Q_h(\theta_0) := \left| \hat{E} \left[h \left(\sqrt{ni(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) \right] - \tilde{E}[h(Z)] \right|,$$

where $\tilde{E}[h(Z)] = 0.379$ is the approximation of $E[h(Z)]$ up to three decimal places. We compare $Q_h(\theta_0)$ with the bound in (3.6), using the difference between their values as a measure of the error. The results presented in the following tables are based on this particular function h while the bounded Wasserstein metric is a supremum over a broader class of test functions, given in (2.6). The results from the simulations are shown in Tables 3.1 and 3.2.

Table 3.1: Simulation results from the Exp(1) distribution

n	$Q_h(\theta_0)$	Upper bound	Error
10	0.008	1.955	1.947
100	0.002	0.336	0.334
1000	0.001	0.094	0.093
10,000	0.0002	0.029	0.0288
100,000	0.0001	0.009	0.0089

Table 3.2: Simulation results from the $\text{Exp}(0.5)$ distribution treated as a non-canonical exponential family

n	$Q_h(\theta_0)$	Upper bound	Error	Bound from Lemma 2.2
10	0.004	11.888	11.884	0.321
100	0.002	3.401	3.399	0.101
1000	0.002	1.058	1.056	0.032
10,000	0.001	0.333	0.332	0.010
100,000	0.0005	0.105	0.1045	0.003

The tables indicate that $Q_h(\theta_0)$, the bound and the error, decrease as the sample size gets larger. All the values in Table 3.1 are smaller than the respective ones in Table 3.2, as expected from Remark 3.4. The bounds are not very good for n up to 100. For the non-canonical case the bounds using directly Lemma 2.2 are, as expected, much better than those from the general approach. The bounds are conceptual and better constants may be possible.

As a graphical tool, both the density plots and the normal Quantile-Quantile (Q-Q) plot (for $n = 1000$) for the two examples are used, in order to observe the behaviour of the standardised values of the MLE as the sample size varies. In addition, the bar plots give the magnitude of the error against the sample size in order to capture the performance of Stein's method for each example. Figures 3.1 and 3.2 show the density plots, while the Q-Q plots with the bar plots for the canonical and the non-canonical case are presented in Figures 3.3 and 3.4, respectively.

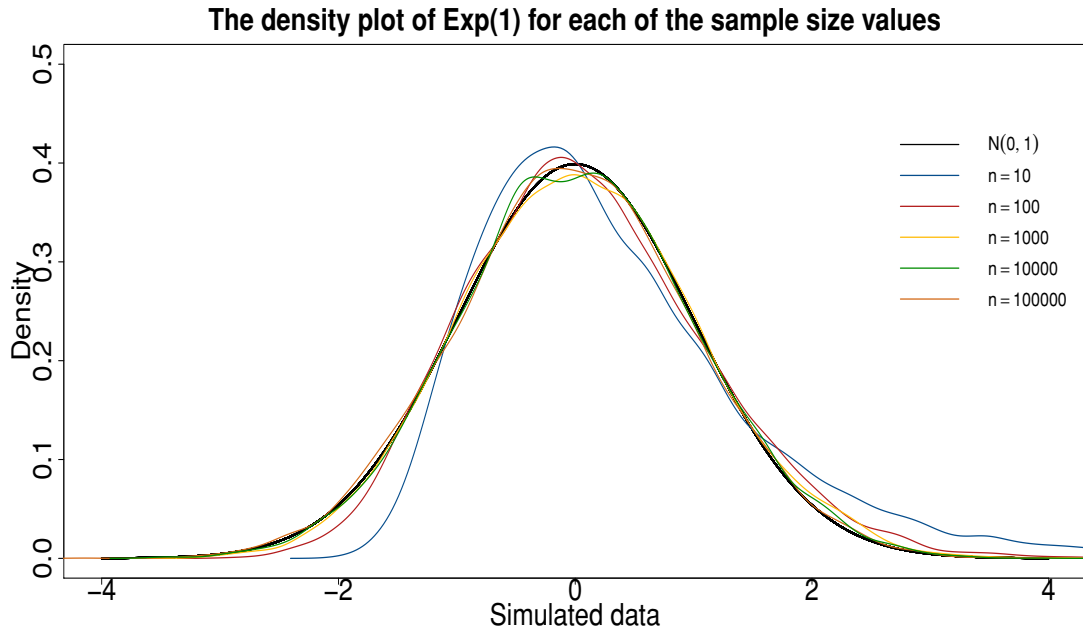


Fig. 3.1: The density function of $\sqrt{ni(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0)$ for various sample sizes. In this case X_1, X_2, \dots, X_n are i.i.d. random variables which follow the $\text{Exp}(1)$ distribution.

The density plots verify what has already been observed in Tables 3.1 and 3.2. The distribution of the standardised MLE gets closer to the standard normal as the sample size increases. Specifically, for both examples, the curves produced when $n = 10$ do not approximate the simulated standard normal distribution well. There is a significant but not adequate improvement when $n = 100$, since the density plots generated are closer to the density of the $N(0,1)$ distribution. Therefore, if the sample size gets greater than 100, then a significant drop in the difference between the distributional distance and the upper bound provided by Stein's method is expected. The 'Error' columns of the tables verify this conjecture.

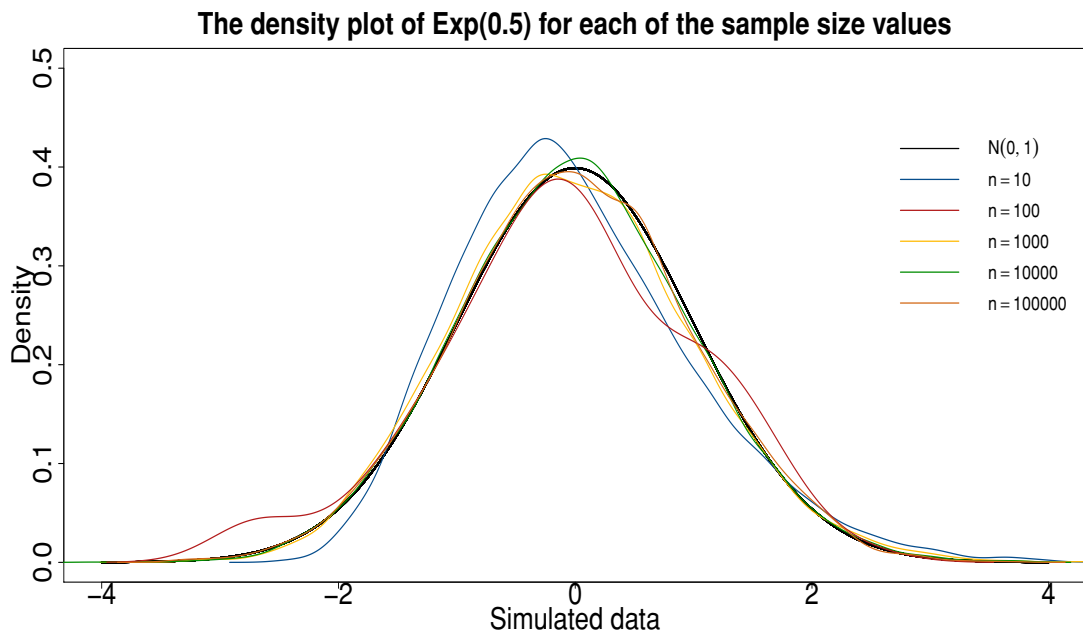


Fig. 3.2: The density function of $\sqrt{ni(\theta_0)}(\hat{\theta}_n(\mathbf{X}) - \theta_0)$ for various sample sizes. In this case X_1, X_2, \dots, X_n are i.i.d. random variables which follow the Exp(0.5) distribution.

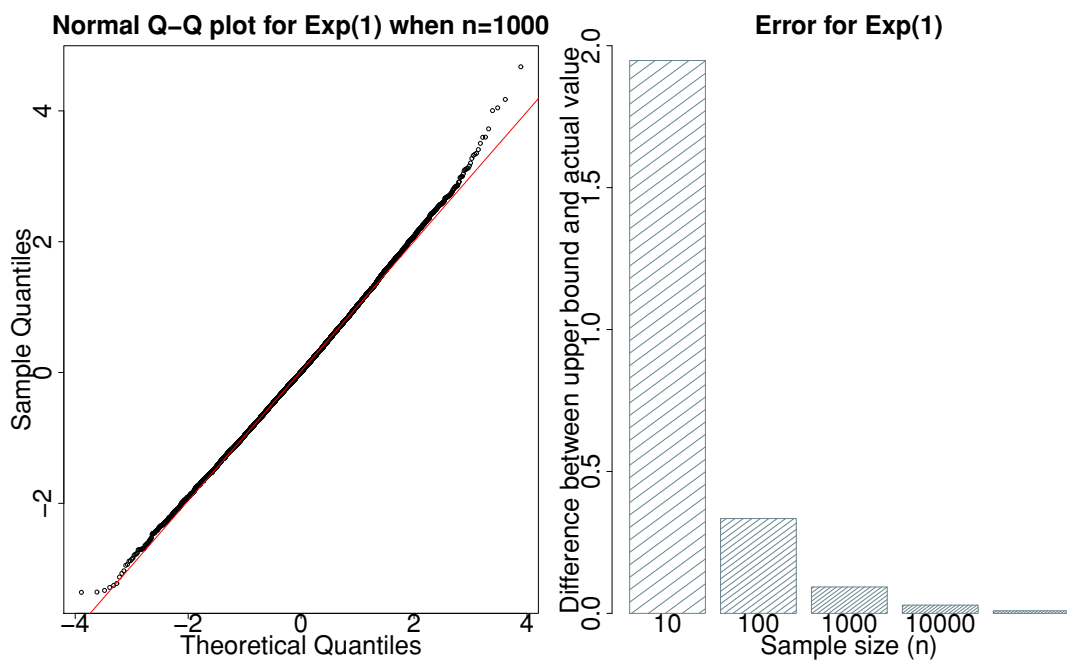


Fig. 3.3: The Q-Q plot, when $n = 1000$, and the bar plot of the error for different sample sizes for the exponential distribution example, Exp(1).

Figures 3.3 and 3.4 present two symmetric normal Q-Q plots, with the majority of the observations lying on the straight line that passes through the first and third quartiles.

There are a few extreme values in the tails of the plots, but the number of those cases is relatively small, compared to the total of 10,000 values for the standardised MLE. The bar plots presented in the figures show graphically what we concluded from the density plots. When the sample size, n , is greater than 100, the error drops very quickly.

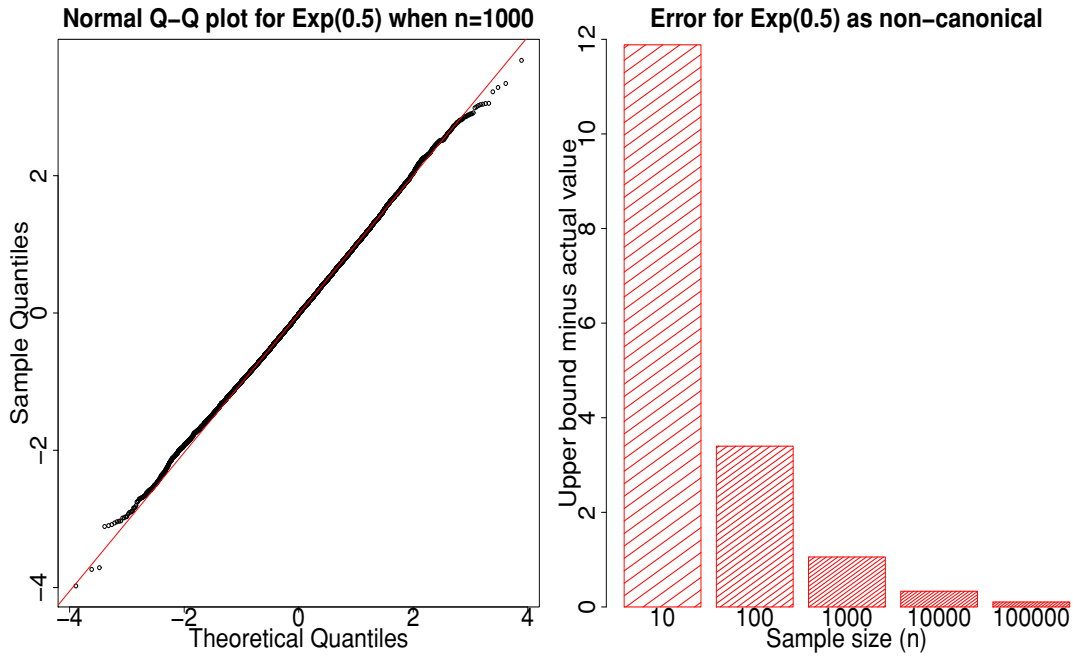


Fig. 3.4: The Q-Q plot, when $n = 1000$, and the bar plot of the error for different sample sizes for the exponential distribution example, $\text{Exp}(0.5)$.

3.4 Discrete distributions: The boundary issue

In this section, we consider discrete distributions that face the problem of the MLE having positive probability of being on the boundary of the parameter space. This issue causes complications for results which are based on asymptotic properties that require both the MLE and the parameter to be interior to the parameter space, such as Taylor expansions. The idea we follow is based on a perturbation method applied both on the parameter and the sample. The new, perturbed values of the parameter and the MLE will be interior to Θ , allowing us to apply part of the general results as explained

in Section 3.2. The above method is assessed through the Poisson distribution with parameter θ and parameter space $\Theta = [0, \infty)$.

3.4.1 The perturbation approach

A perturbation method based on a perturbation function, should be such that first of all, the function should perturb the quantity of interest in a way that ensures it will be interior to its domain. The second requirement is that the perturbed quantity should be as close as possible to the initial quantity. Let X be a random variable with support B , which is a connected, closed (or semi-closed) interval $[a, b]$ ($(a, b]$ or $[a, b)$), where $-\infty < a < b < \infty$. For $0 < \delta < \frac{b-a}{2}$, we are looking for a perturbation function, $q : B \rightarrow \overset{\circ}{B}$ (where in this case, $\overset{\circ}{B}$ denotes the interior of the set B) with $q(x) = kx + d$, such that:

$$(1) \quad q(a) = a + \delta \text{ and } q(b) = b - \delta.$$

$$(2) \quad \sup_{x \in B} |q(x) - x| \text{ is minimal.}$$

Solving this problem for k and d , gives $k = 1 - \frac{2\delta}{b-a}$ and $d = \delta + \frac{2a}{b-a}\delta$. There is only one solution, which minimises $\sup_{x \in B} |q(x) - x|$. Thus, the second requirement is also satisfied.

We choose $\delta = \delta(n) = \frac{c}{n}$ with $0 < c < \frac{n(b-a)}{2}$. Finally, the perturbation function is

$$q_c(x) = x + \frac{c}{n} - \frac{2c}{n} \left(\frac{x-a}{b-a} \right), \quad x \in B. \quad (3.35)$$

In the case where $B = (-\infty, b]$ or $B = [a, \infty)$, then $q_c(x) = x - \frac{c}{n}$ or $q_c(x) = x + \frac{c}{n}$, respectively.

Assuming existence and uniqueness of the MLE, $\hat{\theta}_n(\mathbf{X})$, for the parameter θ_0 , of a discrete distribution with parameter space as in the previous paragraph, the aim is to find an upper bound on

$$d_{bW}(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0), K),$$

where $K \sim N\left(0, \frac{1}{i(\theta_0)}\right)$. The quantity to bound is not exactly the one shown in (3.1), because the Expected Fisher Information number might not exist or not be finite when θ_0 lies on the boundary of the parameter space. For this purpose, the perturbation function in (3.35) is used for both the parameter and the data.

First we introduce some notation. For S being the discrete sample space, let $a := \inf \Theta$, $b := \sup \Theta$, $\underline{S} := \inf S$, $\bar{S} := \sup S$ and $0 < c_1 < \frac{n(b-a)}{2}$, $0 < c_2 < \frac{n(\bar{S}-\underline{S})}{2}$. In addition,

$$q_{c_1}(\theta_0) = \theta_0 + \frac{c_1}{n} - \frac{2c_1}{n} \left(\frac{\theta_0 - a}{b - a} \right) \quad (3.36)$$

is the perturbed parameter and

$$q_{c_2}(x_i) = x_i + \frac{c_2}{n} - \frac{2c_2}{n} \left(\frac{x_i - \underline{S}}{\bar{S} - \underline{S}} \right)$$

is the perturbed value for x_i . The perturbed data is

$$q_{c_2}(\mathbf{x}) = (q_{c_2}(x_1), q_{c_2}(x_2), \dots, q_{c_2}(x_n)). \quad (3.37)$$

We now assume that the domain of the MLE, $\hat{\theta}_n(\mathbf{x})$, can be extended so that $\hat{\theta}_n(q_{c_2}(\mathbf{x})) := \hat{\theta}_n(\mathbf{y})|_{\mathbf{y}=q_{c_2}(\mathbf{x})}$ is defined and gives values that belong to the parameter space Θ . The perturbed MLE is denoted by

$$\hat{\theta}_n^*(\mathbf{x}) := \hat{\theta}_n(q_{c_2}(\mathbf{x})) := \hat{\theta}_n(\mathbf{y})|_{\mathbf{y}=q_{c_2}(\mathbf{x})}. \quad (3.38)$$

Also,

$$\begin{aligned} l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{x})) &:= l'(\theta; \mathbf{y}) \Big|_{\substack{\theta=q_{c_1}(\theta_0), \\ \mathbf{y}=q_{c_2}(\mathbf{x})}}, & l'(\hat{\theta}_n^*(\mathbf{x}); q_{c_2}(\mathbf{x})) &:= l'(\theta; \mathbf{y}) \Big|_{\substack{\theta=\hat{\theta}_n^*(\mathbf{x}), \\ \mathbf{y}=q_{c_2}(\mathbf{x})}}, \\ l''(q_{c_1}(\theta_0); q_{c_2}(\mathbf{x})) &:= l''(\theta; \mathbf{y}) \Big|_{\substack{\theta=q_{c_1}(\theta_0), \\ \mathbf{y}=q_{c_2}(\mathbf{x})}}, & l^{(3)}(\theta; q_{c_2}(\mathbf{x})) &= l^{(3)}(\theta; \mathbf{y}) \Big|_{\mathbf{y}=q_{c_2}(\mathbf{x})}. \end{aligned}$$

For ease of presentation, abbreviate

$$Y_i = \frac{1}{\sqrt{ni}(q_{c_1}(\theta_0))} \left[\frac{d}{d\theta} \log f(U_i|\theta) \Big|_{\substack{\theta=q_{c_1}(\theta_0) \\ U_i=q_{c_2}(X_i)}} \right], \quad i \in \{1, 2, \dots, n\} \quad (3.39)$$

while

$$w_1 := w_1(n, c_1, c_2, \theta_0) \quad \text{and} \quad w_2 := w_2(n, c_1, c_2, \theta_0)$$

are its expectation and variance, respectively.

Theorem 3.2. *Let X_1, X_2, \dots, X_n be i.i.d. random variables which follow a single-parameter discrete distribution with parameter space the connected, closed or semi-closed interval $\Theta \subset \mathbb{R}$ and discrete sample space S . Assume that $i(\theta_0) > 0$ and let $\frac{1}{i(\theta_0)} = 0$ be the continuous extension of $\frac{1}{i(\theta)}$ to θ_0 when θ_0 is such that $i(\theta)$ approaches negative or positive infinity as θ approaches θ_0 . Let $h \in H_{bW}$ and $0 < \varepsilon = \varepsilon(q_{c_1}(\theta_0)) = \varepsilon_0$ such that as in (R3), $\left| \frac{d^3}{d\theta^3} \log f(q_{c_2}(x)|\theta) \right| \leq M(q_{c_2}(x)), \forall x \in S, \forall \theta \in \overset{\circ}{\Theta}$ such that $|\theta - q_{c_1}(\theta_0)| < \varepsilon_0$, where $\overset{\circ}{\Theta}$ denotes the interior of the parameter space Θ . Also, assume that $E \left(\left(\sum_{i=1}^n M(q_{c_2}(X_i)) \right)^2 \middle| \left| \hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0) \right| < \varepsilon \right) < \infty$. Let $K \sim N \left(0, \frac{1}{i(\theta_0)} \right)$. Then for $0 < c_1 < \frac{n(b-a)}{2}$, $0 < c_2 < \frac{n(\bar{S}-S)}{2}$, $q_{c_1}(\theta_0)$ as in (3.36) and $q_{c_2}(x)$ as in (3.37),*

$$\begin{aligned} d_{bW}(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0), K) &\leq \frac{c_1}{\sqrt{n}} \left| 1 - 2 \left(\frac{\theta_0 - a}{b - a} \right) \right| + \sqrt{n} E |\hat{\theta}_n(\mathbf{X}) - \hat{\theta}_n^*(\mathbf{X})| \\ &+ \left[\left| 1 - \frac{1}{\sqrt{w_2 n i(\theta_0)}} \right| \sqrt{n w_2 + (n w_1)^2} + \frac{\sqrt{n} |w_1|}{\sqrt{w_2 i(\theta_0)}} \right] \\ &+ \frac{1}{\sqrt{n}} \left(2 + \frac{1}{(w_2)^{\frac{3}{2}}} E |Y_1 - w_1|^3 \right) \mathbb{1} \left\{ \frac{1}{i(\theta_0)} > 0 \right\} + \frac{2}{\varepsilon^2} E \left[(\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0))^2 \right] \\ &+ \frac{1}{\sqrt{ni}(q_{c_1}(\theta_0))} \left\{ \sqrt{E[l''(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X})) + ni(q_{c_1}(\theta_0))]^2} \sqrt{E[(\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0))^2]} \right. \\ &\quad \left. + \frac{1}{2} \left[E \left(\left(\sum_{i=1}^n M(q_{c_2}(X_i)) \right)^2 \middle| \left| \hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0) \right| < \varepsilon \right) \right]^{\frac{1}{2}} \right. \\ &\quad \left. \times \left[E[(\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0))^4] \right]^{\frac{1}{2}} \right\}. \end{aligned} \quad (3.40)$$

Proof. For ease of presentation, the proof of the theorem is split into three steps.

Step 1: Perturbation of θ_0 . Using the triangle inequality and then a first order Taylor expansion of $h(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0))$ about $\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - q_{c_1}(\theta_0))$ gives

$$\begin{aligned} & |E[h(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0))] - E[h(K)]| \leq |E[h(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - q_{c_1}(\theta_0)))] - E[h(K)]| \\ & + |E[h(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0)) - h(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - q_{c_1}(\theta_0)))]| \\ & \leq |E[h(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - q_{c_1}(\theta_0)))] - E[h(K)]| + \sqrt{n}\|h'\|E|q_{c_1}(\theta_0) - \theta_0| \\ & = |E[h(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - q_{c_1}(\theta_0)))] - E[h(K)]| + \frac{\|h'\|c_1}{\sqrt{n}} \left| 1 - 2 \left(\frac{\theta_0 - a}{b - a} \right) \right|. \end{aligned} \quad (3.41)$$

Step 2: Perturbation of the MLE. To perturb the MLE, we perturb the data. This construction ensures that the MLE evaluated at $q_{c_2}(\mathbf{x})$ (as defined in (3.37)) is not on the boundary of the parameter space. Following the same process as in (3.41), using the triangle inequality and a first order Taylor expansion of $h(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - q_{c_1}(\theta_0)))$ about $\sqrt{n}(\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0))$, where $\hat{\theta}_n^*(\mathbf{X})$ is as in (3.38), gives

$$\begin{aligned} & |E[h(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - q_{c_1}(\theta_0)))] - E[h(K)]| \\ & \leq |E[h(\sqrt{n}(\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0)))] - E[h(K)]| \\ & + |E[h(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - q_{c_1}(\theta_0)) - h(\sqrt{n}(\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0)))]| \\ & \leq |E[h(\sqrt{n}(\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0)))] - E[h(K)]| + \sqrt{n}\|h'\|E|\hat{\theta}_n(\mathbf{X}) - \hat{\theta}_n^*(\mathbf{X})|. \end{aligned} \quad (3.42)$$

Step 3: The final bound. It remains to bound

$$|E[h(\sqrt{n}(\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0)))] - E[h(K)]|.$$

Since both $q_{c_1}(\theta_0)$ and $\hat{\theta}_n^*(x)$ are interior to Θ , a second-order Taylor expansion of $l'(\hat{\theta}_n^*(x); q_{c_2}(x)) = l'(\hat{\theta}_n(q_{c_2}(x)); q_{c_2}(x)) = 0$ about $q_{c_1}(\theta_0)$ yields

$$0 = l'(q_{c_1}(\theta_0); q_{c_2}(x)) + (\hat{\theta}_n^*(x) - q_{c_1}(\theta_0)) l''(q_{c_1}(\theta_0); q_{c_2}(x)) + R_1(q_{c_1}(\theta_0); q_{c_2}(x)), \quad (3.43)$$

where, similarly as in Section 3.2,

$$R_1(q_{c_1}(\theta_0); q_{c_2}(x)) = \frac{1}{2} (\hat{\theta}_n^*(x) - q_{c_1}(\theta_0))^2 l^{(3)}(\tilde{\theta}; q_{c_2}(x))$$

with

$$l^{(3)}(\tilde{\theta}; q_{c_2}(x)) = l^{(3)}(\theta; y) \Big|_{\substack{\theta=\tilde{\theta} \\ y=q_{c_2}(x)}}$$

for $\tilde{\theta}$ between $\hat{\theta}_n^*(x)$ and $q_{c_1}(\theta_0)$. A simple rearrangement of the terms in (3.43), leads to

$$\hat{\theta}_n^*(x) - q_{c_1}(\theta_0) = \frac{-l'(q_{c_1}(\theta_0); q_{c_2}(x)) - R_1(q_{c_1}(\theta_0); q_{c_2}(x))}{l''(q_{c_1}(\theta_0); q_{c_2}(x))}.$$

In general $l''(q_{c_1}(\theta_0); q_{c_2}(x)) \neq -ni(q_{c_1}(\theta_0))$ and similar steps to those used in the proof of Theorem 3.1 yield

$$\hat{\theta}_n^*(x) - q_{c_1}(\theta_0) = \frac{l'(q_{c_1}(\theta_0); q_{c_2}(x)) + R_1(q_{c_1}(\theta_0); q_{c_2}(x)) + R_2(q_{c_1}(\theta_0); q_{c_2}(x))}{ni(q_{c_1}(\theta_0))},$$

where

$$R_2(q_{c_1}(\theta_0); q_{c_2}(x)) = (\hat{\theta}_n^*(x) - q_{c_1}(\theta_0)) [l''(q_{c_1}(\theta_0); q_{c_2}(x)) + ni(q_{c_1}(\theta_0))].$$

Using that $q_{c_2}(X) = (q_{c_2}(X_1), q_{c_2}(X_2), \dots, q_{c_2}(X_n))$, the triangle inequality gives

$$\begin{aligned} & \left| \mathbb{E} \left[h \left(\sqrt{n} \left(\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0) \right) \right) \right] - \mathbb{E}[h(K)] \right| \\ & \leq \left| \mathbb{E} \left[h \left(\frac{l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X}))}{\sqrt{n} i(q_{c_1}(\theta_0))} \right) \right] - \mathbb{E}[h(K)] \right| \end{aligned} \quad (3.44)$$

$$\begin{aligned} & + \left| \mathbb{E} \left[h \left(\frac{l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X})) + R_1(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X})) + R_2(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X}))}{\sqrt{n} i(q_{c_1}(\theta_0))} \right) \right. \right. \\ & \quad \left. \left. - h \left(\frac{l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X}))}{\sqrt{n} i(q_{c_1}(\theta_0))} \right) \right] \right|. \end{aligned} \quad (3.45)$$

(A) To find an upper bound on (3.44) using Lemma 2.2, note that

$$\frac{l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X}))}{\sqrt{n} i(q_{c_1}(\theta_0))} = \sum_{i=1}^n Y_i,$$

where Y_i as in (3.39) and $w_1 := w_1(n, c_1, c_2, \theta_0)$, $w_2 := w_2(n, c_1, c_2, \theta_0)$ are its expectation and variance, respectively. These quantities depend on the sample size and on the perturbed values $(q_{c_1}(\theta_0)$ and $q_{c_2}(x_i)$). Define $\tilde{Y}_i = \frac{Y_i - w_1}{\sqrt{w_2 i(\theta_0)}}$, $\forall i \in \{1, 2, \dots, n\}$ with $\mathbb{E}(\tilde{Y}_i) = 0$ and $\text{Var}(\tilde{Y}_i) = \frac{1}{i(\theta_0)}$. As a consequence of X_1, X_2, \dots, X_n being i.i.d. random variables, $\tilde{Y}_1, \tilde{Y}_2, \dots, \tilde{Y}_n$ are i.i.d. random variables too. Using the triangle inequality and that

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{Y}_i = \frac{1}{\sqrt{w_2 n i(\theta_0)}} \left(\frac{l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X}))}{\sqrt{n} i(q_{c_1}(\theta_0))} - n w_1 \right) \quad (3.46)$$

gives

$$\begin{aligned} & \left| \mathbb{E} \left[h \left(\frac{l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X}))}{\sqrt{n} i(q_{c_1}(\theta_0))} \right) \right] - \mathbb{E}[h(K)] \right| \leq \\ & \left| \mathbb{E} \left[h \left(\frac{1}{\sqrt{w_2 n i(\theta_0)}} \left(\frac{l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X}))}{\sqrt{n} i(q_{c_1}(\theta_0))} - n w_1 \right) \right) \right] - \mathbb{E}[h(K)] \right| \\ & + \left| \mathbb{E} \left[h \left(\frac{l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X}))}{\sqrt{n} i(q_{c_1}(\theta_0))} \right) - h \left(\frac{1}{\sqrt{w_2 n i(\theta_0)}} \left(\frac{l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X}))}{\sqrt{n} i(q_{c_1}(\theta_0))} - n w_1 \right) \right) \right] \right|. \end{aligned} \quad (3.47)$$

The first term of the bound in (3.47) will be bounded using Lemma 2.2 with $W = \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{Y}_i$ as in (3.46). Thus,

$$\begin{aligned} & \left| \mathbb{E} \left[h \left(\frac{1}{\sqrt{w_2 n i(\theta_0)}} \left(\frac{l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X}))}{\sqrt{n} i(q_{c_1}(\theta_0))} - n w_1 \right) \right) \right] - \mathbb{E}[h(K)] \right| \\ & \leq \frac{\|h'\|}{\sqrt{n}} \left(2 + [i(\theta_0)]^{\frac{3}{2}} \mathbb{E}|\tilde{Y}_1|^3 \right) = \frac{\|h'\|}{\sqrt{n}} \left(2 + \frac{1}{(w_2)^{\frac{3}{2}}} \mathbb{E}|Y_1 - w_1|^3 \right). \end{aligned} \quad (3.48)$$

When $\frac{1}{i(\theta_0)} = 0$ then $\tilde{Y}_i = 0, \forall i \in \{1, 2, \dots, n\}$ and the first term of the bound in (3.47) is equal to zero. For the second term of the upper bound in (3.47) a first-order Taylor expansion and the Cauchy-Schwarz inequality yield

$$\begin{aligned} & \left| \mathbb{E} \left[h \left(\frac{l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X}))}{\sqrt{n} i(q_{c_1}(\theta_0))} \right) - h \left(\frac{1}{\sqrt{w_2 n i(\theta_0)}} \left(\frac{l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X}))}{\sqrt{n} i(q_{c_1}(\theta_0))} - n w_1 \right) \right) \right] \right| \\ & \leq \|h'\| \left| 1 - \frac{1}{\sqrt{w_2 n i(\theta_0)}} \right| \mathbb{E} \left| \frac{l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X}))}{\sqrt{n} i(q_{c_1}(\theta_0))} \right| + \frac{\|h'\| \sqrt{n} |w_1|}{\sqrt{w_2 i(\theta_0)}} \\ & \leq \|h'\| \left| 1 - \frac{1}{\sqrt{w_2 n i(\theta_0)}} \right| \left(\text{Var} \left(\frac{l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X}))}{\sqrt{n} i(q_{c_1}(\theta_0))} \right) + \frac{[\mathbb{E}(l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X})))]^2}{n [i(q_{c_1}(\theta_0))]^2} \right)^{\frac{1}{2}} \\ & \quad + \frac{\|h'\| \sqrt{n} |w_1|}{\sqrt{w_2 i(\theta_0)}} \\ & = \|h'\| \left[\left| 1 - \frac{1}{\sqrt{w_2 n i(\theta_0)}} \right| \sqrt{n w_2 + (n w_1)^2} + \frac{\sqrt{n} |w_1|}{\sqrt{w_2 i(\theta_0)}} \right]. \end{aligned} \quad (3.49)$$

(B) To complete the proof, it remains to find an upper bound for (3.45). The idea is the same as the one used for (3.19). We condition on whether $|\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0)| \geq \varepsilon$ or $|\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0)| < \varepsilon$, where $0 < \varepsilon = \varepsilon(q_{c_1}(\theta_0))$ such that $\left| \frac{d^3}{d\theta^3} \log f(u|\theta) \right|_{u=q_{c_2}(x)} \leq M(q_{c_2}(x)), \forall x \in S, \forall \theta \in \overset{\circ}{\Theta}$ with $|\theta - q_{c_1}(\theta_0)| < \varepsilon$. Following the same process as in Section 3.2 yields

$$\begin{aligned}
& \left| \mathbb{E} \left[h \left(\frac{l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X})) + R_1(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X})) + R_2(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X}))}{\sqrt{n}i(q_{c_1}(\theta_0))} \right) \right. \right. \\
& \quad \left. \left. - h \left(\frac{l'(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X}))}{\sqrt{n}i(q_{c_1}(\theta_0))} \right) \right] \right| \\
& \leq \frac{2\|h\|}{\varepsilon^2} \mathbb{E} \left[(\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0))^2 \right] \\
& \quad + \frac{\|h'\|}{\sqrt{n}i(q_{c_1}(\theta_0))} \left\{ \sqrt{\mathbb{E} [l''(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X})) + ni(q_{c_1}(\theta_0))]^2} \sqrt{\mathbb{E} \left[(\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0))^2 \right]} \right. \\
& \quad + \frac{1}{2} \left[\mathbb{E} \left(\left(\sum_{i=1}^n M(q_{c_2}(X_i)) \right)^2 \middle| |\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0)| < \varepsilon \right) \right]^{\frac{1}{2}} \\
& \quad \left. \times \left[\mathbb{E} \left[(\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0))^4 \right] \right]^{\frac{1}{2}} \right\}. \tag{3.50}
\end{aligned}$$

Combining (3.41), (3.42), (3.48), (3.49) and (3.50) and the fact that $\|h\| \leq 1$, $\|h'\| \leq 1$ gives the result in (3.40). \square

Remark 3.5. (1) In order for the above bound to approach zero as the sample size, n , increases we require that $\mathbb{E} |\hat{\theta}_n(\mathbf{X}) - \hat{\theta}_n^*(\mathbf{X})| = o\left(\frac{1}{\sqrt{n}}\right)$.

(2) When both endpoints of the parameter space are not finite, then parameter perturbation is not necessary. In the case where one of the two endpoints of the now semi-closed parameter space is infinite, then it suffices to change the form of the perturbed parameter, which now becomes

$$\begin{aligned}
q_{c_1}(\theta_0) &= \theta_0 - \frac{c_1}{n}, & \text{if the left endpoint is equal to } -\infty \\
q_{c_1}(\theta_0) &= \theta_0 + \frac{c_1}{n}, & \text{if the right endpoint is equal to } \infty.
\end{aligned}$$

The same remark holds regarding the sample space and the relevant perturbation of the data.

3.4.2 Example: The Poisson distribution

In this subsection the Poisson distribution with parameter $\theta \in \Theta = [0, \infty)$ is treated. The value $\theta = 0$ must be in the parameter space in order for the MLE, $\hat{\theta}_n(\mathbf{X}) = \bar{X}$, to exist and to be unique. The $\text{Poisson}(\theta)$ distribution with the aforementioned parameter space is not a single-parameter exponential family. When $\theta = 0$ is included in the parameter space the requirements of an exponential family are not satisfied as the set of values x for which the relevant probability mass function

$$f(x|\theta) = \frac{e^{-\theta} \theta^x}{x!}, \quad \theta \in [0, \infty), x \in \mathbb{Z}_0^+$$

is positive, is different for $\theta = 0$ than for any other value of the parameter θ ; the support of the distribution depends on the parameter. As the parameter space is neither an open and connected nor a compact subset of \mathbb{R} , we can not use Proposition 2.1 or Proposition 2.2, respectively. Therefore, the steps that prove existence and uniqueness of the MLE are presented below.

Let X_1, X_2, \dots, X_n be i.i.d. random variables from $\text{Poisson}(\theta)$, $\theta \in [0, \infty)$. The likelihood function is

$$L(\theta; \mathbf{x}) = \prod_{i=1}^n f(x_i|\theta) = \frac{\exp\{-n\theta\} \theta^{\sum_{i=1}^n x_i}}{\prod_{i=1}^n x_i!}.$$

To maximise the likelihood function we account for whether $\sum_{i=1}^n x_i = 0$ or not. Each case is examined separately.

Case 1: If $\sum_{i=1}^n x_i = 0$ (equivalently $x_1 = x_2 = \dots = x_n = 0$), then

$$L(\theta; \mathbf{0}) = \exp\{-n\theta\}$$

is a decreasing function of θ . It takes its maximum value on $[0, \infty)$ when θ attains its minimum; hence $\hat{\theta}_n(\mathbf{X}) = 0 = \bar{X}$.

Case 2: If $\sum_{i=1}^n x_i \neq 0$, we follow the usual concept on finding the MLE. The first derivative of the log-likelihood function with respect to θ is

$$l'(\theta; \mathbf{x}) = -n + \frac{\sum_{i=1}^n x_i}{\theta}. \quad (3.51)$$

We can see that $l'(\theta; \mathbf{x}) = 0 \Leftrightarrow \frac{\sum_{i=1}^n x_i}{\theta} = n \Leftrightarrow \theta = \hat{\theta}_n(\mathbf{X}) = \bar{X}$. In order to verify that $\hat{\theta}_n(\mathbf{X}) = \bar{X}$ is the MLE, $l''(\hat{\theta}_n(\mathbf{x}); \mathbf{x}) < 0$ is required. Simple calculations yield

$$l''(\theta; \mathbf{x}) = -\frac{\sum_{i=1}^n x_i}{\theta^2} \Rightarrow l''(\hat{\theta}_n(\mathbf{x}); \mathbf{x}) = -\frac{n\bar{x}}{\bar{x}^2} = -\frac{n}{\bar{x}} < 0.$$

In conclusion, using the results from Cases 1 and 2 above, we deduce that for X_1, X_2, \dots, X_n being i.i.d. random variables that follow the $\text{Poisson}(\theta)$, $\theta \in [0, \infty)$, the MLE exists, it is unique and equal to $\hat{\theta}_n(\mathbf{X}) = \bar{X}$.

Following the steps of the proof of Theorem 3.2, using also Hölder's inequality for the third absolute moment in the fourth term of the bound in (3.40) gives the following result.

Corollary 3.4. *Let X_1, X_2, \dots, X_n be i.i.d. random variables which follow the $\text{Poisson}(\theta_0)$ distribution, with $\theta_0 \in [0, \infty)$. For $K \sim N(0, \theta_0)$, $h \in H_{bW}$ as in (2.6) and $0 < c_1 \leq c_2$ positive constants and $C = c_2 - c_1 \geq 0$,*

1) *if $\theta_0 > 0$ then*

$$\begin{aligned} d_{bW}(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0), K) &< \frac{c_1}{\sqrt{n}} + \frac{c_2}{\sqrt{n}} + \frac{1}{\sqrt{n}} \left[2 + \frac{(3\theta_0 + 1)^{\frac{3}{4}}}{\theta_0^{\frac{3}{4}}} \right] \\ &+ \frac{C}{\sqrt{n}} + \frac{8}{n(\theta_0 + \frac{c_1}{n})^2} \left(\theta_0 + \frac{C^2}{n} \right) + \frac{1}{\sqrt{n}(\theta_0 + \frac{c_1}{n})} \left(\theta_0 + \frac{C^2}{n} \right) \\ &+ \frac{12}{\sqrt{n}(\theta_0 + \frac{c_1}{n})} \left[\frac{\theta_0}{n} (1 + 3n\theta_0) + \frac{4C\theta_0}{n} + \frac{6C^2\theta_0}{n} + \frac{C^4}{n^2} \right]^{\frac{1}{2}}. \end{aligned} \quad (3.52)$$

2) if $\theta_0 = 0$ then

$$d_{bW}(\sqrt{n}\hat{\theta}_n(\mathbf{X}), K) = 0.$$

Proof. We have already proved that the unique MLE in this case is $\hat{\theta}_n(\mathbf{X}) = \bar{X}$. Using the second result of Remark 3.5, the perturbed parameter is $q_{c_1}(\theta_0) = \theta_0 + \frac{c_1}{n}$, while the perturbed data are $q_{c_2}(x_i) = x_i + \frac{c_2}{n}$ so that the perturbed MLE is

$$\hat{\theta}_n^*(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n \left(X_i + \frac{c_2}{n} \right) = \hat{\theta}_n(\mathbf{X}) + \frac{c_2}{n},$$

where c_1 and c_2 are chosen such that $0 < c_1 \leq c_2$. It is easy to deduce that, in this specific example, the first term of the general bound in (3.40) becomes $\frac{c_1}{\sqrt{n}}$. For the second term,

$$\sqrt{n}E|\hat{\theta}_n(\mathbf{X}) - \hat{\theta}_n^*(\mathbf{X})| = \sqrt{n} \left| -\frac{c_2}{n} \right| = \frac{c_2}{\sqrt{n}}.$$

The expected Fisher information number is

$$i(\theta_0) = \text{Var} \left(\frac{d}{d\theta} \log f(X|\theta_0) \right) = \frac{1}{\theta_0^2} \text{Var}(X) = \frac{1}{\theta_0}.$$

To find $w_1 := w_1(n)$ and $w_2 := w_2(n)$ in the general bound, using (3.51)

$$Y_i = \frac{q_{c_1}(\theta_0)}{\sqrt{n}} \left(-1 + \frac{q_{c_2}(X_i)}{q_{c_1}(\theta_0)} \right) = \frac{q_{c_2}(X_i) - q_{c_1}(\theta_0)}{\sqrt{n}} = \frac{1}{\sqrt{n}} \left(X_i + \frac{c_2}{n} - \theta_0 - \frac{c_1}{n} \right),$$

results in $w_1 = \frac{c_2 - c_1}{n^{\frac{3}{2}}}$ and $w_2 = \frac{\theta_0}{n}$. Thus, for the third term of the general upper bound in (3.40),

$$\begin{aligned} & \left| 1 - \frac{1}{\sqrt{w_2 n i(\theta_0)}} \right| \sqrt{n w_2 + (n w_1)^2} + \frac{\sqrt{n} |w_1|}{\sqrt{w_2 i(\theta_0)}} \\ &= \left| 1 - \frac{\sqrt{\theta_0}}{\sqrt{\theta_0}} \right| \sqrt{\theta_0 + \frac{c_2 - c_1}{n^2}} + \frac{\sqrt{n} |c_2 - c_1|}{n} = \frac{|c_2 - c_1|}{\sqrt{n}}. \end{aligned}$$

For the fourth term in (3.40), the equations (3.38), p.70 of Kendall and Stuart (1969), and the results in p.73 of the book show that $E(X_i - \theta_0)^4 = \theta_0(1 + 3\theta_0)$. This result and Hölder's inequality yield

$$\begin{aligned} \frac{1}{\sqrt{n}} \left(2 + \frac{1}{(w_2)^{\frac{3}{2}}} E|Y_1 - w_1|^3 \right) &\leq \frac{1}{\sqrt{n}} \left(2 + \frac{1}{(w_2)^{\frac{3}{2}}} [E(Y_1 - w_1)^4]^{\frac{3}{4}} \right) \\ &= \frac{1}{\sqrt{n}} \left(2 + \left(\frac{n}{\theta_0} \right)^{\frac{3}{2}} \left[E \left(\frac{X_1 - \theta_0}{\sqrt{n}} \right)^4 \right]^{\frac{3}{4}} \right) = \frac{1}{\sqrt{n}} \left(2 + \left(\frac{3\theta_0 + 1}{\theta_0} \right)^{\frac{3}{4}} \right). \end{aligned}$$

Denoting by $C = c_2 - c_1 \geq 0$, then $E \left[(\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0))^2 \right] = \frac{\theta_0}{n} + \frac{C^2}{n^2}$. Also,

$$\begin{aligned} E \left[(\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0))^4 \right] &= E \left[\left(\hat{\theta}_n(\mathbf{X}) - \theta_0 + \frac{C}{n} \right)^4 \right] \\ &= E \left[\left(\frac{\sum_{i=1}^n X_i - n\theta_0 + C}{n} \right)^4 \right] \\ &= \frac{1}{n^4} E \left[\left(\sum_{i=1}^n X_i - n\theta_0 \right)^4 \right] + \frac{4C}{n^4} E \left[\left(\sum_{i=1}^n X_i - n\theta_0 \right)^3 \right] + \frac{6C^2}{n^4} E \left[\left(\sum_{i=1}^n X_i - n\theta_0 \right)^2 \right] \\ &\quad + \frac{4C^3}{n^4} E \left[\sum_{i=1}^n X_i - n\theta_0 \right] + \frac{C^4}{n^4}. \end{aligned}$$

Using again the results in Kendall and Stuart (1969) and the fact that $\sum_{i=1}^n X_i \sim \text{Poisson}(n\theta_0)$ yields

$$E \left[(\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0))^4 \right] = \frac{\theta_0}{n^3} (1 + 3n\theta_0) + \frac{4C\theta_0}{n^3} + \frac{6C^2\theta_0}{n^3} + \frac{C^4}{n^4}.$$

We have that $\left| \frac{d^3}{d\theta^3} \log f(U|\theta) \right|_{U=q_{c_2}(X)} = \left| \frac{2(X + \frac{c_2}{n})}{\theta^3} \right|$. To find the supremum in a neighbourhood of θ , we take $0 < \varepsilon = \varepsilon(q_{c_1}(\theta_0))$ such that $\varepsilon < q_{c_1}(\theta_0)$, since if $\varepsilon > q_{c_1}(\theta_0)$ then $(q_{c_1}(\theta_0) - \varepsilon, q_{c_1}(\theta_0) + \varepsilon) \cap \overset{\circ}{\Theta} = (0, q_{c_1}(\theta_0) + \varepsilon)$ and then $\sup_{\theta: |\theta - q_{c_1}(\theta_0)| < \varepsilon} \left| \frac{2(X_i + \frac{c_2}{n})}{\theta^3} \right|$ is not finite. Thus, for $0 < \varepsilon < q_{c_1}(\theta_0)$,

$$\sup_{\theta: |\theta - q_{c_1}(\theta_0)| < \varepsilon} \left| \frac{2(X_i + \frac{c_2}{n})}{\theta^3} \right| = \frac{2(X_i + \frac{c_2}{n})}{(q_{c_1}(\theta_0) - \varepsilon)^3} = M(q_{c_2}(X_i))$$

and

$$\begin{aligned}
& \left[\mathbb{E} \left(\left(\sum_{i=1}^n M(q_{c_2}(X_i)) \right)^2 \middle| |\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0)| < \varepsilon \right) \right]^{\frac{1}{2}} \\
&= \left[\mathbb{E} \left(\left(\frac{2n\hat{\theta}_n^*(\mathbf{X})}{(q_{c_1}(\theta_0) - \varepsilon)^3} \right)^2 \middle| |\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0)| < \varepsilon \right) \right]^{\frac{1}{2}} \\
&= \frac{2n}{(q_{c_1}(\theta_0) - \varepsilon)^3} \left[\mathbb{E} \left((\hat{\theta}_n^*(\mathbf{X}))^2 \middle| |\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0)| < \varepsilon \right) \right]^{\frac{1}{2}} \\
&= \frac{2n}{(q_{c_1}(\theta_0) - \varepsilon)^3} \left[\mathbb{E} \left((|\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0) + q_{c_1}(\theta_0)|)^2 \middle| |\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0)| < \varepsilon \right) \right]^{\frac{1}{2}} \\
&\leq \frac{2n}{(q_{c_1}(\theta_0) - \varepsilon)^3} \left[\mathbb{E} \left((|\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0)| + q_{c_1}(\theta_0))^2 \middle| |\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0)| < \varepsilon \right) \right]^{\frac{1}{2}} \\
&< \frac{2n}{(q_{c_1}(\theta_0) - \varepsilon)^3} \left[(\varepsilon + q_{c_1}(\theta_0))^2 \right]^{\frac{1}{2}} = \frac{2n}{(q_{c_1}(\theta_0) - \varepsilon)^3} (\varepsilon + q_{c_1}(\theta_0)).
\end{aligned}$$

There is one last term to bound in order to apply the results of Theorem 3.2 to this specific example of the Poisson distribution. Simple steps yield

$$\begin{aligned}
& \sqrt{\mathbb{E} \left[(l''(q_{c_1}(\theta_0); q_{c_2}(\mathbf{X})) + ni(q_{c_1}(\theta_0)))^2 \right]} \\
&= \sqrt{\mathbb{E} \left[\left(-\frac{\sum_{i=1}^n q_{c_2}(X_i)}{(q_{c_1}(\theta_0))^2} + \frac{n}{q_{c_1}(\theta_0)} \right)^2 \right]} \\
&= \frac{n}{(q_{c_1}(\theta_0))^2} \sqrt{\mathbb{E} \left[(\hat{\theta}_n^*(\mathbf{X}) - q_{c_1}(\theta_0))^2 \right]} = \frac{n}{(\theta_0 + \frac{c_1}{n})^2} \sqrt{\left(\frac{\theta_0}{n} + \frac{C^2}{n^2} \right)}.
\end{aligned}$$

Using all these results to (3.40), gives

$$\begin{aligned}
d_{bw}(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0), K) &< \frac{c_1}{\sqrt{n}} + \frac{c_2}{\sqrt{n}} + \frac{1}{\sqrt{n}} \left[2 + \frac{(3\theta_0 + 1)^{\frac{3}{4}}}{\theta_0^{\frac{3}{4}}} \right] \mathbb{1}\{\theta_0 > 0\} \\
&+ \frac{C}{\sqrt{n}} + \frac{2}{\varepsilon^2} \left(\frac{\theta_0}{n} + \frac{C^2}{n^2} \right) + \frac{\sqrt{n}}{(\theta_0 + \frac{c_1}{n})} \left(\frac{\theta_0}{n} + \frac{C^2}{n^2} \right) \\
&+ \frac{\theta_0 + \frac{c_1}{n}}{2\sqrt{n}} \left(\frac{2n(\varepsilon + q_{c_1}(\theta_0))}{(q_{c_1}(\theta_0) - \varepsilon)^3} \right) \left[\frac{\theta_0}{n^3} (1 + 3n\theta_0) + \frac{4C\theta_0}{n^3} + \frac{6C^2\theta_0}{n^3} + \frac{C^4}{n^4} \right]^{\frac{1}{2}}.
\end{aligned}$$

We have already explained that $\varepsilon \in (0, q_{c_1}(\theta_0))$. There is a trade-off on the choice of ε for the fifth and the seventh term of the general bound in (3.40). The fifth term indicates that as the value of ε gets nearer zero the quality of the bound is worse, while the seventh term requires ε to be as small as possible since $M(q_{c_2}(X_i)) = \frac{2(X_i + \frac{c_2}{n})}{(q_{c_1}(\theta_0) - \varepsilon)^3}$. An optimisation process is quite tedious and thus we choose $\varepsilon = \frac{q_{c_1}(\theta_0)}{2} = \frac{\theta_0 + \frac{c_1}{n}}{2}$. Giving this value to ε yields the result of the Corollary. \square

Remark 3.6. (1) *The upper bound in (3.52) for the distributional distance between the actual distribution of the MLE and the normal distribution in the case of i.i.d. random variables following the Poisson(θ) distribution, with $\theta \in [0, \infty)$ is of order $\frac{1}{\sqrt{n}}$.*

(2) *We notice that the bound in (3.52) achieves its minimum as a function of C when $C = 0$ (equivalently $c_1 = c_2 = c > 0$). This yields*

$$d_{bW}(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0), K) < \frac{2c}{\sqrt{n}} + \frac{1}{\sqrt{n}} \left[2 + \frac{(3\theta_0 + 1)^{\frac{3}{4}}}{\theta_0^{\frac{3}{4}}} \right] \\ + \frac{8\theta_0}{n(\theta_0 + \frac{c}{n})^2} + \frac{\theta_0}{\sqrt{n}(\theta_0 + \frac{c}{n})} + \frac{12}{\sqrt{n}(\theta_0 + \frac{c}{n})} \left[\frac{\theta_0}{n} + 3\theta_0^2 \right]^{\frac{1}{2}}.$$

(3) *Since the MLE is unique and equal to $\hat{\theta}_n(\mathbf{X}) = \bar{X}$, Lemma 2.2 could be used directly for \bar{X} . Define $W = \sqrt{n}(\bar{X} - \theta_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i$, where $Y_i = X_i - \theta_0$ are independent, zero mean random variables. Also, $E(W) = 0$ and $\text{Var}(W) = n\text{Var}(\bar{X}) = \frac{1}{n} \sum_{i=1}^n \text{Var}(X_i) = \theta_0$. Therefore, (2.12) for $K \sim N(0, \theta_0)$ and Hölder's inequality give for $\theta_0 > 0$*

$$d_{bW}(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0), K) \leq \frac{1}{\sqrt{n}} \left(2 + \frac{1}{\theta_0^{\frac{3}{2}}} [E(Y_1)^4]^{\frac{3}{4}} \right) = \frac{1}{\sqrt{n}} \left(2 + \frac{(3\theta_0 + 1)^{\frac{3}{4}}}{\theta_0^{\frac{3}{4}}} \right).$$

This bound, obtained by the direct application of Stein's method, is smaller than the bound given in Corollary 3.4. However the interest in the example treated in this section, where $\Theta = [0, \infty)$, is in adapting the approach to such cases where the MLE could be on the boundary of the parameter space and also it is not assumed that the MLE is a sum of random variables.

Empirical results

As in the exponential distribution example, the accuracy of the distributional upper bound given in (3.52) is studied using some simulated data. Similar steps as in Section 3.3 are followed. We produce 10,000 trials of n random independent observations x from the Poisson distribution with mean equal to 2. We then take the mean (which is the MLE) of the observations for each trial. However, these 10,000 values are not standardised now. We just subtract the value of the parameter and multiply by \sqrt{n} . Then, the function $h(x) = \frac{1}{x^2+2}, \forall x \in \mathbb{Z}_0^+$ is applied to these values and we take the sample mean, namely $\hat{E}[h(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0))]$. Since $K \sim N(0, 2)$, basic calculations give $\tilde{E}[h(K)] = 0.328$ which is the approximation of $E[h(K)]$ up to three decimal places. With

$$Q_2(h, \theta_0) := |\hat{E}[h(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0))] - \tilde{E}[h(K)]|,$$

we compare $Q_2(h, \theta_0)$ with the bound in (3.52). Table 3.3 summarises the results for each sample size, n . In the table we denote by B_c the bound for $c_1 = c_2 = c$, such that $C = 0$ in (3.52). We present the results for $c \in \{0.1, 1, 2\}$. The error column in the table is the difference between $Q_2(h, \theta_0)$ and the value of the bound for $c = 1$, denoted by B_1 .

Table 3.3: Simulation results from the Poisson(2) distribution

n	$Q_2(h, \theta_0)$	B_1	Error	$B_{0.1}$	B_2	Bound from Lemma 2.2
10	0.0008	1.2626	1.2618	1.1756	1.3642	0.4762
100	0.0006	0.3699	0.3693	0.3297	0.4147	0.1506
1000	0.0005	0.1130	0.1125	0.9999	0.1275	0.0476
10000	0.0004	0.0353	0.0349	0.0311	0.0399	0.0151
100000	0.0002	0.0111	0.0109	0.0098	0.0125	0.0048

Table 3.3 indicates that the quantity $Q_2(h, \theta_0)$, the upper bounds and the relevant error decrease as the sample size increases. In addition, the smaller the constant c is, the

smaller the bound. The error (for $c = 1$) is adequately small when the sample size is equal to 100. The upper bound obtained by applying directly Lemma 2.2 is, as expected, smaller than the one we get from the general perturbation method.

The accuracy of Stein's method in the case of i.i.d. random variables that follow the Poisson distribution is verified using some graphical tools. The behaviour of the MLE is shown in the density and the Q-Q plots presented in Figure 3.5 and 3.6, respectively. Both of these plots show that the distribution of the MLE is very close to the normal distribution. The bar plot in Figure 3.6 gives the magnitude of the error and its considerable decrease when the sample size gets greater than or equal to 100.

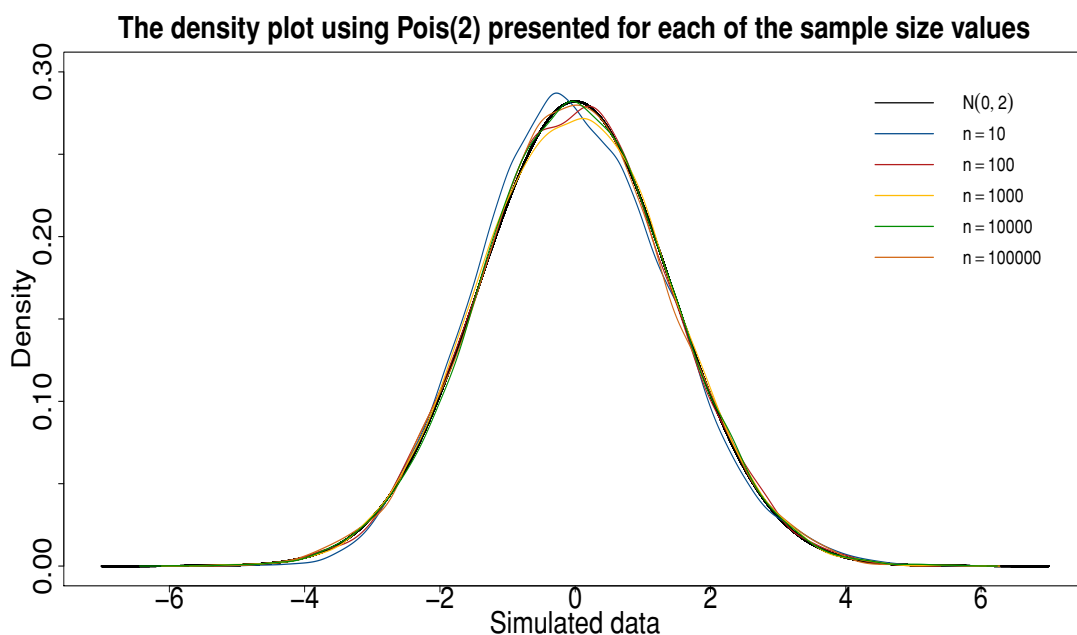


Fig. 3.5: The density function of $\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0)$ for various sample sizes. In this case X_1, X_2, \dots, X_n are i.i.d. random variables which follow the Poisson(2) distribution.

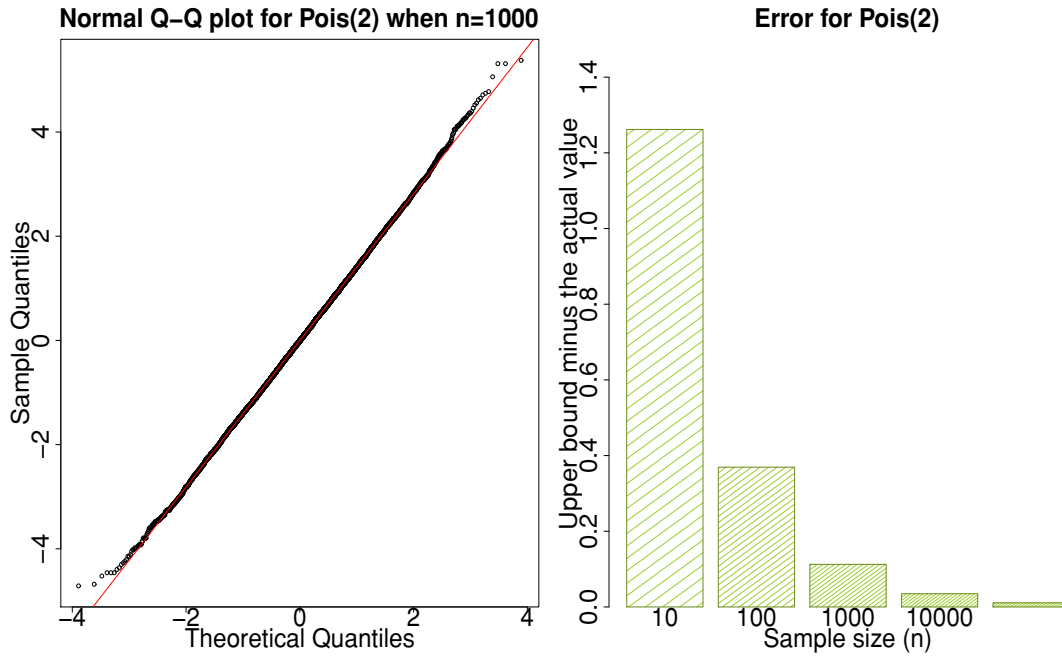


Fig. 3.6: The Q-Q plot, when $n = 1000$, and the bar plot of the error for different sample sizes for the Poisson distribution with mean equal to 2.

3.5 Bounds on the mean squared error

The previous sections considered models where the MLE exists, it is unique and can be found analytically as an explicit function of the random variables X_1, X_2, \dots, X_n . The general upper bound in Section 3.2 has been proven without assuming or using an explicit form of the MLE. The bound in (3.7) includes terms, such as the mean squared error, $E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]$, of which the calculation requires to know the expression of the MLE. This fact generates problems in models where no closed-form solution to the maximization problem is known or available; in these cases, a numerical method, such as Newton-Raphson algorithm, can often be used to approximate the MLE. Hence a normal approximation is of interest in these cases too.

This section focuses on finding upper bounds for situations where an analytic form for the MLE is not available. In the proof for the final upper bound in Theorem 3.1, an

explicit form of the MLE was not used. However, if the MLE is not known, then the mean squared error (MSE) of the MLE, appearing in the bound for (3.1) should be bounded by a quantity which is independent of $\hat{\theta}_n(\mathbf{X})$.

Let X_1, X_2, \dots, X_n be i.i.d. random variables. Apart from the regularity conditions, first defined in Section 2.2, we require the following further assumptions that make the steps and the calculations easier and ensure a meaningful upper bound:

(Fur.1) The support, S , is bounded;

(Fur.2) For all $\theta_0 \in \Theta$ there exists $\varepsilon_0 = \varepsilon_0(\theta_0) > 0$ such that $\left| \frac{d^3}{d\theta^3} \log f(x_1|\theta) \right| \leq M$ for all $\theta \in \Theta$ such that $|\theta - \theta_0| < \varepsilon_0$, where $M = M(\theta_0)$ is a constant that may depend on the unknown parameter θ_0 ;

(Fur.3) The sample size, n , satisfies

$$n > \frac{\|x\|^2 \left[M\varepsilon_0 + \sqrt{(M\varepsilon_0)^2 + 8[i(\theta_0)]^2} \right]^2}{4[i(\theta_0)]^3 \varepsilon_0^2}. \quad (3.53)$$

This ensures that $1 - 2 \frac{\|x^2\|}{ni(\theta_0)\varepsilon_0^2} - \frac{\|x\|M}{\sqrt{n}[i(\theta_0)]^{\frac{3}{2}}} > 0$ for ε_0 as in (Fur.2).

Towards the aim of bounding the MSE by a quantity which is independent of the MLE, we will use the smooth test function $h(x) = x^2$, which under (Fur.1) is bounded, in (3.6) of Theorem 3.1 and rearrange the bound. For ease of presentation, let

$$D_1 = D_1(\theta_0, x, n) = 1 - 2 \frac{\|x^2\|}{ni(\theta_0)\varepsilon_0^2} - \frac{\|x\|M}{\sqrt{n}[i(\theta_0)]^{\frac{3}{2}}}.$$

Theorem 3.3. *Let X_1, X_2, \dots, X_n be i.i.d. random variables with probability density (or mass) function $f(x_i|\theta)$ with bounded support. Assume that the regularity conditions (R1) - (R4), as well as the assumptions (Fur.1) - (Fur.3) are satisfied. Also assume that the MLE exists and that it is unique. Then $A_1 = A_1(\theta_0, n)$ is an upper bound for*

$$\sqrt{\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]}, \text{ where}$$

$$A_1 = \frac{1}{D_1 \sqrt{ni(\theta_0)}} \left\{ \frac{\|x\| \sqrt{\text{Var} \left[\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right]}}{\sqrt{ni(\theta_0)}} \right. \quad (3.54)$$

$$\left. + \left[\frac{\|x\|^2 \text{Var} \left[\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right]}{n[i(\theta_0)]^2} + D_1 \left[1 + \frac{2\|x\|}{\sqrt{n}} \left(2 + \frac{\mathbb{E} \left| \frac{d}{d\theta} \log f(X_1 | \theta_0) \right|^3}{[i(\theta_0)]^{\frac{3}{2}}} \right) \right] \right]^{\frac{1}{2}} \right\}.$$

Proof. Using the notation for the remainder terms, the triangle inequality, conditional expectations, Markov's inequality and Stein's method, as in Section 3.2, gives for $\varepsilon = \varepsilon(\theta_0) = \varepsilon_0$

$$\begin{aligned} \left| \mathbb{E} \left[h \left(\sqrt{ni(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) \right] - \mathbb{E}[h(Z)] \right| &\leq \frac{\|h'\|}{\sqrt{n}} \left(2 + \frac{\mathbb{E} \left| \frac{d}{d\theta} \log f(X_1 | \theta_0) \right|^3}{[i(\theta_0)]^{\frac{3}{2}}} \right) \\ &+ \frac{\|h'\|}{\sqrt{i(\theta_0)}} \sqrt{\text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right)} \sqrt{\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} \\ &+ \frac{2\|h\|}{\varepsilon^2} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \\ &+ \frac{\|h'\|}{2\sqrt{ni(\theta_0)}} \mathbb{E} \left((\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \sup_{\theta: |\theta - \theta_0| < \varepsilon} \left| l^{(3)}(\theta; \mathbf{X}) \right| \middle| |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right). \end{aligned}$$

Using now (Fur.2)

$$\begin{aligned} \left| \mathbb{E} \left[h \left(\sqrt{ni(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) \right] - \mathbb{E}[h(Z)] \right| &\leq \frac{\|h'\|}{\sqrt{n}} \left(2 + \frac{\mathbb{E} \left| \frac{d}{d\theta} \log f(X_1 | \theta_0) \right|^3}{[i(\theta_0)]^{\frac{3}{2}}} \right) \\ &+ \frac{\|h'\|}{\sqrt{i(\theta_0)}} \sqrt{\text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right)} \sqrt{\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} \\ &+ \frac{\|h'\| \sqrt{nM}}{2\sqrt{i(\theta_0)}} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] + \frac{2\|h\|}{\varepsilon^2} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]. \quad (3.55) \end{aligned}$$

Straightforward calculations and denoting by B_{x^2} the upper bound for (3.1) when $h(x) = x^2$, lead to

$$\begin{aligned} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] &= \frac{1}{ni(\theta_0)} \left| \mathbb{E} \left[\sqrt{ni(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right]^2 - \mathbb{E}(Z^2) + \mathbb{E}(Z^2) \right| \\ &\leq \frac{1}{ni(\theta_0)} (B_{x^2} + 1), \end{aligned} \quad (3.56)$$

where

$$\begin{aligned} B_{x^2} &\leq 2 \frac{\|x\|}{\sqrt{n}} \left(2 + \frac{\mathbb{E} \left| \frac{d}{d\theta} \log f(X_1 | \theta_0) \right|^3}{[i(\theta_0)]^{\frac{3}{2}}} \right) + \left[\frac{2\|x^2\|}{\varepsilon^2} + \frac{\|x\|\sqrt{nM}}{\sqrt{i(\theta_0)}} \right] \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \\ &\quad + 2 \frac{\|x\|}{\sqrt{i(\theta_0)}} \sqrt{\text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right)} \sqrt{\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]}. \end{aligned} \quad (3.57)$$

Combining the results in (3.56) and (3.57) yields

$$\begin{aligned} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] &\leq 2 \frac{\|x\|}{n^{\frac{3}{2}} i(\theta_0)} \left(2 + \frac{\mathbb{E} \left| \frac{d}{d\theta} \log f(X_1 | \theta_0) \right|^3}{[i(\theta_0)]^{\frac{3}{2}}} \right) \\ &\quad + \frac{1}{ni(\theta_0)} \left[\frac{2\|x^2\|}{\varepsilon^2} + \frac{\|x\|\sqrt{nM}}{\sqrt{i(\theta_0)}} \right] \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \\ &\quad + 2 \frac{\|x\|}{n[i(\theta_0)]^{\frac{3}{2}}} \sqrt{\text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right)} \sqrt{\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} + \frac{1}{ni(\theta_0)}. \end{aligned} \quad (3.58)$$

The next step is to solve the simple quadratic inequality in (3.58), with unknown $A = A(\theta_0) = \sqrt{\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]}$. We get

$$\begin{aligned} A^2 &\leq \frac{1}{ni(\theta_0)} \left(1 + \frac{2\|x\|}{\sqrt{n}} \left(2 + \frac{\mathbb{E} \left| \frac{d}{d\theta} \log f(X_1 | \theta_0) \right|^3}{[i(\theta_0)]^{\frac{3}{2}}} \right) + \left[\frac{2\|x^2\|}{\varepsilon^2} + \frac{\|x\|\sqrt{nM}}{\sqrt{i(\theta_0)}} \right] A^2 \right. \\ &\quad \left. + \frac{2\|x\| \sqrt{\text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right)}}{\sqrt{i(\theta_0)}} A \right), \end{aligned}$$

which yields

$$D_1 A^2 - 2 \frac{\|x\| \sqrt{\text{Var}\left(\frac{d^2}{d\theta^2} \log f(X_1|\theta_0)\right)}}{n[i(\theta_0)]^{\frac{3}{2}}} A - \frac{1}{ni(\theta_0)} \left[1 + 2 \frac{\|x\|}{\sqrt{n}} \left(2 + \frac{E\left|\frac{d}{d\theta} \log f(X_1|\theta_0)\right|^3}{[i(\theta_0)]^{\frac{3}{2}}} \right) \right] \leq 0. \quad (3.59)$$

Thus, the two solutions of the above quadratic inequality are given by

$$A_{1,2} = \frac{1}{D_1 \sqrt{ni(\theta_0)}} \left\{ \sqrt{\text{Var}\left(\frac{d^2}{d\theta^2} \log f(X_1|\theta_0)\right)} \frac{\|x\|}{\sqrt{ni(\theta_0)}} \pm \left[\frac{\|x\|^2 \text{Var}\left(\frac{d^2}{d\theta^2} \log f(X_1|\theta_0)\right)}{n[i(\theta_0)]^2} + D_1 \left(1 + \frac{2\|x\|}{\sqrt{n}} \left(2 + \frac{E\left|\frac{d}{d\theta} \log f(X_1|\theta_0)\right|^3}{[i(\theta_0)]^{\frac{3}{2}}} \right) \right) \right]^{\frac{1}{2}} \right\}. \quad (3.60)$$

Since $D_1 > 0$, the root with the plus sign in (3.60) is positive, while the other is negative.

Therefore the inequality in (3.59) is satisfied for $A = \sqrt{E\left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2\right]} \leq A_1$, with A_1 as in (3.54). \square

Remark 3.7. (1) Using Theorem 3.3, the final upper bound for (3.1) which is useful when no analytic expression of the MLE is available, becomes

$$d_{bw}\left(\sqrt{ni(\theta_0)}(\hat{\theta}_n(\mathbf{X}) - \theta_0), Z\right) \leq \frac{1}{\sqrt{n}} \left(2 + \frac{E\left|\frac{d}{d\theta} \log f(X_1|\theta_0)\right|^3}{[i(\theta_0)]^{\frac{3}{2}}} \right) + \left[\frac{2}{\varepsilon^2} + \frac{\sqrt{n}M}{2\sqrt{i(\theta_0)}} \right] (A_1)^2 + \frac{\sqrt{\text{Var}\left(\frac{d^2}{d\theta^2} \log f(X_1|\theta_0)\right)}}{\sqrt{i(\theta_0)}} A_1. \quad (3.61)$$

(2) The order of A_1 in terms of the sample size is $\frac{1}{\sqrt{n}}$. Hence, the order of the final upper bound in (3.61) is also $\frac{1}{\sqrt{n}}$.

(3) The result of the Theorem holds when (Fur.3) is satisfied, meaning that we have a lower bound on the sample size n , as can be seen from (3.53). This lower bound can be

quite large, and thus Theorem 3.3 is mainly of theoretical interest.

(4) The result holds only when the support of the distribution is bounded; see (Fur.1). An approach that holds when the support of the distribution is not necessarily bounded, but assumes boundedness of the parameter space, is explained in Theorem 6.5, where the random variables are not necessarily independent.

Example: The Beta distribution

Consider the example of i.i.d random variables from the Beta distribution with one of the two shape parameters being unknown. In this case, the MLE can only be expressed in terms of the inverse of the digamma function, $\Psi(\theta) = \frac{d}{d\theta} \log \Gamma(\theta)$. The general result in Theorem 3.3 is used, in order to obtain an upper bound for the MSE and use it to get an upper bound for (3.1). The following corollary gives the result.

Corollary 3.5. *Let X_1, X_2, \dots, X_n be i.i.d. random variables from the $\text{Beta}(\theta_0, \beta)$ distribution, where β is known and θ_0 is the unknown parameter. Let*

$$B_1 = B_1(\theta_0) = 8 \left(\Psi_3(\theta_0) + \Psi_3(\theta_0 + \beta) + 3[\Psi_1(\theta_0)]^2 + 3[\Psi_1(\theta_0 + \beta)]^2 \right),$$

where $\Psi_j(\theta)$, $j \in \mathbb{N}$ is the j^{th} derivative of the digamma function, $\Psi(\theta)$. Also, let

$$B_2 = B_2(\theta_0) = \frac{96\beta}{\theta_0^4} + \frac{\beta\pi^4}{15}, \quad D_{\Psi_1} = D_{\Psi_1}(\theta_0, \beta) = \Psi_1(\theta_0) - \Psi_1(\theta_0 + \beta)$$

and

$$B_3 = B_3(\theta_0, n) = \frac{\left[\left(1 + \frac{2}{\sqrt{n}} \left(2 + \frac{(B_1)^{\frac{3}{4}}}{(D_{\Psi_1})^{\frac{3}{2}}} \right) \right) \left(1 - \frac{8}{n\theta_0^2 D_{\Psi_1}} - \frac{B_2}{\sqrt{n}(D_{\Psi_1})^{\frac{3}{2}}} \right) \right]^{\frac{1}{2}}}{\sqrt{D_{\Psi_1}} - \frac{8}{n\theta_0^2 \sqrt{D_{\Psi_1}}} - \frac{B_2}{\sqrt{n} D_{\Psi_1}}}. \quad (3.62)$$

Let

$$n > \frac{\left[B_2 \frac{\theta_0}{2} + \sqrt{\frac{(B_2 \theta_0)^2}{4} + 8[\Psi_1(\theta_0) - \Psi_1(\theta_0 + \beta)]^2} \right]^2}{[\Psi_1(\theta_0) - \Psi_1(\theta_0 + \beta)]^3 \theta_0^2}.$$

Then for $Z \sim N(0, 1)$

$$\begin{aligned} d_{bW} \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) &\leq \frac{1}{\sqrt{n}} \left(2 + \frac{(B_1)^{\frac{3}{4}}}{[\Psi_1(\theta_0) - \Psi_1(\theta_0 + \beta)]^{\frac{3}{2}}} \right) \\ &\quad + \frac{8}{n\theta_0^2} (B_3)^2 + \frac{B_2(B_3)^2}{2\sqrt{n}[\Psi_1(\theta_0) - \Psi_1(\theta_0 + \beta)]^{\frac{1}{2}}}. \end{aligned} \quad (3.63)$$

Proof. The probability density function is

$$f(x|\theta) = \frac{\Gamma(\theta + \beta)}{\Gamma(\theta)\Gamma(\beta)} x^{\theta-1} (1-x)^{\beta-1}, \quad (3.64)$$

with $\theta > 0$ and $x \in [0, 1]$. Hence

$$\begin{aligned} l(\theta; \mathbf{x}) &= n[\log(\Gamma(\theta + \beta)) - \log(\Gamma(\theta)) - \log(\Gamma(\beta))] \\ &\quad + (\theta - 1) \sum_{i=1}^n \log x_i + (\beta - 1) \sum_{i=1}^n \log(1 - x_i) \end{aligned} \quad (3.65)$$

and

$$l'(\theta; \mathbf{x}) = n[\Psi(\theta + \beta) - \Psi(\theta)] + \sum_{i=1}^n \log x_i$$

$$l^{(j)}(\theta; \mathbf{x}) = n(\Psi_{j-1}(\theta + \beta) - \Psi_{j-1}(\theta)), \quad j \in \mathbb{N} \setminus \{1\}.$$

Now we show that the conditions (R1)-(R4) and the assumptions (Fur.1)-(Fur.3) are satisfied. For (R1) it is obvious. As for (R2), the three times differentiability of the density function can be verified from (3.65). In addition, using (3.64) and the expressions for the logarithmic expectations of a Beta distributed random variable, it is easy to verify $\int_0^1 \frac{d^j}{d\theta^j} f(x|\theta) dx = \frac{d^j}{d\theta^j} \int_0^1 f(x|\theta) dx = 0, j \in \{1, 2, 3\}$ for (R2). Now for (R3), a first order Taylor expansion yields

$$\left| \frac{d^3}{d\theta^3} \log f(x|\theta) \right| = |\Psi_2(\theta + \beta) - \Psi_2(\theta)| = \beta |\Psi_3(\theta^*)|, \quad (3.66)$$

for θ^* between θ and $\theta + \beta$. Using that

$$\Psi_m(z) = (-1)^{m+1} m! \sum_{k=0}^{\infty} \frac{1}{(z+k)^{m+1}}, \text{ for } z \in \mathbb{C} \setminus \{\mathbb{Z}_0^-\} \text{ and } m > 0, \quad (3.67)$$

gives

$$\Psi_3(z) = 6 \sum_{k=0}^{\infty} \frac{1}{(z+k)^4}, \quad (3.68)$$

with $\Psi_3(z)$ being a decreasing function of z . We need to upper bound (3.66) in a neighbourhood of θ in order for (R3) to hold. Let $\varepsilon = \varepsilon(\theta_0) = \varepsilon_0$ be a positive constant such that $\theta \in (\theta_0 - \varepsilon, \theta_0 + \varepsilon) \cap \Theta$ as in (R3). If $\varepsilon > \theta_0$ then $(\theta_0 - \varepsilon, \theta_0 + \varepsilon) \cap \Theta = (0, \theta_0 + \varepsilon)$ and the result in (3.66) can not be upper bounded in $(0, \theta_0 + \varepsilon)$ as $\Psi_3(\theta)$ is not bounded for $\theta \in (0, \theta_0 + \varepsilon)$. Therefore, $0 < \varepsilon < \theta_0$ and using that $\sum_{i=1}^{\infty} \frac{1}{i^4} = \frac{\pi^4}{90}$ gives that for θ^* between θ and $\theta + \beta$

$$\begin{aligned} \sup_{\theta: |\theta - \theta_0| < \varepsilon} \left| \frac{d^3}{d\theta^3} \log f(x|\theta) \right| &= \sup_{\theta: |\theta - \theta_0| < \varepsilon} |\Psi_2(\theta + \beta) - \Psi_2(\theta)| \\ &= \beta \sup_{\theta: |\theta - \theta_0| < \varepsilon} |\Psi_3(\theta^*)| \leq \beta \sup_{\theta: |\theta - \theta_0| < \varepsilon} |\Psi_3(\theta)|, \text{ as } \Psi_3(\cdot) \text{ is a decreasing function} \\ &\leq \beta |\Psi_3(\theta_0 - \varepsilon)| = 6\beta \sum_{k=0}^{\infty} \frac{1}{(\theta_0 - \varepsilon + k)^4} \\ &\leq 6\beta \left[\frac{1}{(\theta_0 - \varepsilon)^4} + \sum_{k=1}^{\infty} \frac{1}{k^4} \right] \\ &= \frac{6\beta}{(\theta_0 - \varepsilon)^4} + \frac{\beta\pi^4}{15} = M =: M(\varepsilon), \end{aligned} \quad (3.69)$$

with $M < \infty$. Hence, (R3) holds as well. Also, $i(\theta_0) = \Psi_1(\theta_0) - \Psi_1(\theta_0 + \beta)$ which is positive since it is obvious from (3.67) that $\Psi_1(z)$ is a decreasing function. The assumption (Fur.1) holds with $\|x\| \leq 1$. The assumption (Fur.2) is satisfied from (3.69). Now, since $i(\theta_0) = \Psi_1(\theta_0) - \Psi_1(\theta_0 + \beta)$ take

$$n > \frac{\left[M\varepsilon + \sqrt{(M\varepsilon)^2 + 8[\Psi_1(\theta_0) - \Psi_1(\theta_0 + \beta)]^2} \right]^2}{4\varepsilon^2[\Psi_1(\theta_0) - \Psi_1(\theta_0 + \beta)]^3}$$

in order for (Fur.3) to be satisfied. To find A_1 as in (3.54), firstly, as $E \left| \frac{d}{d\theta} \log f(X_1|\theta_0) \right|^3$ is not straightforward to evaluate due to the absolute value in the expectation, it is easily seen that using Hölder's inequality $E \left| \frac{d}{d\theta} \log f(X_1|\theta_0) \right|^3 \leq \left[E \left[\left(\frac{d}{d\theta} \log f(X_1|\theta_0) \right)^4 \right] \right]^{\frac{3}{4}}$ we find an upper bound for

$$\begin{aligned} E \left[\left(\frac{d}{d\theta} \log f(X_1|\theta_0) \right)^4 \right] &= E \left[(\log X_1 + \Psi(\theta_0 + \beta) - \Psi(\theta_0))^4 \right] \\ &= E \left[(\log X_1 - E(\log X_1))^4 \right]. \end{aligned}$$

If $G_1 \sim \Gamma(\theta_0, \lambda)$ and $G_2 \sim \Gamma(\beta, \lambda)$ are independent, then $\frac{G_1}{G_1+G_2} \sim \text{Beta}(\theta_0, \beta)$. Thus, with $X_1 = \frac{G_1}{G_1+G_2}$

$$\begin{aligned} &E \left[\left(\frac{d}{d\theta} \log f(X_1|\theta_0) \right)^4 \right] \\ &= E \left[((\log G_1 - E[\log G_1]) + (E(\log(G_1 + G_2)) - \log(G_1 + G_2)))^4 \right] \\ &\leq 8 \left\{ E \left[(\log G_1 - E(\log G_1))^4 \right] + E \left[(\log(G_1 + G_2) - E(\log(G_1 + G_2)))^4 \right] \right\}. \end{aligned} \tag{3.70}$$

Now the fourth central moment of the logarithm of a Gamma distributed random variable is calculated using that $\int_0^\infty \frac{z^{\alpha-1} e^{-z} (\log z)^k}{\Gamma(\alpha)} dz = \frac{\Gamma^{(k)}(\alpha)}{\Gamma(\alpha)}$, for any $\alpha > 0$ and $k \in \mathbb{N}$. For $Y \sim \Gamma(\alpha, \lambda)$ and letting $z = \lambda y$

$$E(\log Y) = \Psi(\alpha) - \log \lambda.$$

Using again $z = \lambda y$,

$$\begin{aligned} E \left[(\log Y - E(\log Y))^4 \right] &= \int_0^\infty \frac{z^{\alpha-1} e^{-z}}{\Gamma(\alpha)} \left(\log \left(\frac{z}{\lambda} \right) - E \left(\log \left(\frac{Z}{\lambda} \right) \right) \right)^4 dz \\ &= \int_0^\infty \frac{z^{\alpha-1} e^{-z}}{\Gamma(\alpha)} (\log z - E(\log Z))^4 dz \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{\Gamma(\alpha)} \sum_{k=0}^4 \binom{4}{k} (-1)^k [\Psi(\alpha)]^{4-k} \int_0^\infty z^{\alpha-1} e^{-z} (\log z)^k dz \\
&= -3[\Psi(\alpha)]^4 + 6[\Psi(\alpha)]^2 \frac{\Gamma''(\alpha)}{\Gamma(\alpha)} - 4\Psi(\alpha) \frac{\Gamma^{(3)}(\alpha)}{\Gamma(\alpha)} + \frac{\Gamma^{(4)}(\alpha)}{\Gamma(\alpha)}.
\end{aligned}$$

At this point, the digamma function can be used in order to simplify the expression above. Following simple steps it can be easily verified that

$$\begin{aligned}
\frac{\Gamma''(\alpha)}{\Gamma(\alpha)} &= \Psi_1(\alpha) + [\Psi(\alpha)]^2, & \frac{\Gamma^{(3)}(\alpha)}{\Gamma(\alpha)} &= \Psi_2(\alpha) + 3\Psi(\alpha)\Psi_1(\alpha) + [\Psi(\alpha)]^3, \\
\frac{\Gamma^{(4)}(\alpha)}{\Gamma(\alpha)} &= \Psi_3(\alpha) + 4\Psi_2(\alpha)\Psi(\alpha) + 6\Psi_1(\alpha)[\Psi(\alpha)]^2 + 3[\Psi_1(\alpha)]^2 + [\Psi(\alpha)]^4.
\end{aligned}$$

Hence for $Y \sim \Gamma(\alpha, \lambda)$

$$\mathbb{E} \left[(\log Y - \mathbb{E}(\log Y))^4 \right] = \Psi_3(\alpha) + 3[\Psi_1(\alpha)]^2 \quad (3.71)$$

and therefore, from (3.70),

$$\begin{aligned}
&\mathbb{E} \left[\left(\frac{d}{d\theta} \log f(X_1 | \theta_0) \right)^4 \right] \\
&\leq 8 \left(\Psi_3(\theta_0) + \Psi_3(\theta_0 + \beta) + 3[\Psi_1(\theta_0)]^2 + 3[\Psi_1(\theta_0 + \beta)]^2 \right) = B_1.
\end{aligned}$$

With $M = M(\varepsilon)$ as in (3.69), taking $\varepsilon = \frac{\theta_0}{2}$, we conclude that

$$\sup_{\theta: |\theta - \theta_0| < \varepsilon} \left| \frac{d^3}{d\theta^3} \log f(x_1 | \theta) \right| \leq \frac{96\beta}{\theta_0^4} + \frac{\beta\pi^4}{15} = B_2.$$

Using (3.65), gives

$$\text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right) = \text{Var}(\Psi_1(\theta_0 + \beta) - \Psi_1(\theta_0)) = 0.$$

Having found all the necessary quantities, the upper bound in (3.54) is calculated and multiplied by \sqrt{n} . This is equal to B_3 shown in (3.62), which is an upper bound for

$\sqrt{n\mathbb{E}\left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2\right]}$ in the specific case of i.i.d. random variables from the Beta distribution. Using this bound in (3.55) gives the result in (3.63). \square

Now, we study the accuracy of our bound for the MSE of the MLE by simulations. For the simulations, $\theta_0 = 1.5$, $\beta = 1$ and in this case of β being equal to 1, the MLE is analytically known with $\hat{\theta}_n(\mathbf{X}) = -\frac{n}{\sum_{i=1}^n \log X_i}$. We find that $n \geq 7398$, in order for (Fur.3) to be satisfied. The process to simulate is quite simple. Let $n \in \{7398, 7399, \dots, 8397\}$ and for each n , start by generating 10,000 trials of n random independent observations, x , from the Beta distribution with parameter values as above. The MLE, $\hat{\theta}_n(\mathbf{X})$, is evaluated in each trial, which in turn gives a vector of 10,000 values. Thus, for each n from 7398 to 8397, we evaluate the sample MSE, $\hat{\mathbb{E}}\left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2\right] = \frac{1}{10,000} \sum_{i=1}^{10,000} (\hat{\theta}_n(x)_i - \theta_0)^2$ and compare it with its upper bound, $\left(\frac{B_3}{\sqrt{n}}\right)^2$, where B_3 is given in (3.62). The difference between their values measures the error of our bound on the MSE. Some simulation results are shown in Table 3.4.

Table 3.4: Simulation results from the Beta(1.5, 1) distribution

n	$\hat{\mathbb{E}}\left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2\right]$	Upper bound	Error
7500	0.0002	0.0999	0.0997
7700	0.0002	0.0333	0.0331
7900	0.0002	0.0197	0.0195
8100	0.0002	0.0139	0.0137
8300	0.0002	0.0107	0.0105

The table indicates that the bound and the error decrease as the sample size increases, as expected, since the order of the upper bound for the MSE is $\frac{1}{n}$. In addition, the smaller the sample size, the larger the bound is. The bounds are considerably larger than the estimated MSE and they are not numerically sharp. In addition, because of the relatively strong requirement that $n \geq 7398$, these bounds on the MSE are mainly of theoretical interest.

Chapter 4

Multi-parameter distributions

This chapter extends the work presented in Chapter 3 by focusing on the multi-parameter case in the presence of independent but not necessarily identically distributed (i.n.i.d.) and possibly high-dimensional random vectors. The interest is on assessing the asymptotic normality of the vector MLE. The approach we follow is based on Stein's method under a multivariate setting. Let

$$H = \left\{ h : \mathbb{R}^d \rightarrow \mathbb{R} : h \text{ is three times differentiable with bounded derivatives} \right\} \quad (4.1)$$

be the class of test functions we use in this chapter. We abbreviate $\|h\|_1 := \sup_i \left\| \frac{\partial}{\partial x_i} h \right\|$, $\|h\|_2 := \sup_{i,j} \left\| \frac{\partial^2}{\partial x_i \partial x_j} h \right\|$ and $\|h\|_3 := \sup_{i,j,k} \left\| \frac{\partial^3}{\partial x_i \partial x_j \partial x_k} h \right\|$. For upper bounds for the distributional distance of interest we employ Lemma 4.1 below, which is based on the notion of exchangeability. For two random vectors \mathbf{W} and \mathbf{W}' in \mathbb{R}^d , $d \in \mathbb{Z}^+$, we say that $(\mathbf{W}, \mathbf{W}')$ is an exchangeable pair if and only if $(\mathbf{W}, \mathbf{W}') \stackrel{\text{d}}{=} (\mathbf{W}', \mathbf{W})$, where $\stackrel{\text{d}}{=}$ denotes equality in distribution.

Lemma 4.1. (*Reinert and Röllin, 2009*). Assume that $(\mathbf{W}, \mathbf{W}')$ is an exchangeable pair of \mathbb{R}^d -valued random vectors such that $E(\mathbf{W}) = \mathbf{0}$, $E(\mathbf{W}\mathbf{W}^\top) = \Sigma$, with $\Sigma \in \mathbb{R}^{d \times d}$

symmetric and positive definite. Suppose further that

$$\mathbb{E}(\mathbf{W}' - \mathbf{W} | \mathbf{W}) = -\Lambda \mathbf{W} + \mathbf{R} \quad (4.2)$$

for an invertible $d \times d$ matrix Λ and a $\sigma(\mathbf{W})$ -measurable random vector \mathbf{R} . Then if \mathbf{Z} has the d -dimensional standard normal distribution, we have for every $h \in H$ as in (4.1),

$$\left| \mathbb{E}[h(\mathbf{W})] - \mathbb{E}\left[h\left(\Sigma^{\frac{1}{2}} \mathbf{Z}\right)\right] \right| \leq \frac{\|h\|_2}{4} A + \frac{\|h\|_3}{12} B + \left(\|h\|_1 + \frac{1}{2} d \|\Sigma\|^{\frac{1}{2}} \|h\|_2 \right) C,$$

where, with $\lambda^{(i)} := \sum_{m=1}^d \left| (\Lambda^{-1})_{m,i} \right|$,

$$\begin{aligned} A &= \sum_{i=1}^d \sum_{j=1}^d \lambda^{(i)} \sqrt{\text{Var}\left(\mathbb{E}\left((W'_i - W_i)(W'_j - W_j) \mid \mathbf{W}\right)\right)}, \\ B &= \sum_{i=1}^d \sum_{j=1}^d \sum_{k=1}^d \lambda^{(i)} \mathbb{E}\left|(W'_i - W_i)(W'_j - W_j)(W'_k - W_k)\right|, \\ C &= \sum_{i=1}^d \lambda^{(i)} \sqrt{\text{Var}(R_i)}. \end{aligned}$$

Reinert and Röllin (2009) were not the first to use Stein's original approach of exchangeable pairs for multivariate normal approximation. This goes back to Chatterjee and Meckes (2008) where three abstract normal approximation theorems using exchangeable pairs are presented that frequently produce total variation distance bounds of the optimal order in the univariate case. Chatterjee and Meckes (2008) also explain that because the differential equation (2.8) in the univariate case is first order rather than second, it allows for reducing the degree of smoothness of the test function needed by one, over what is required in the multivariate case. This lead to obtaining bounds on the Kolmogorov distance, as explained in Section 3.1. The multivariate Stein equation in (2.16) is a second order differential equation and an improvement on bounds for the total variation distance in the multivariate case is not available. Although Chatterjee

and Meckes (2008) give in general sharper bounds, the results of Reinert and Röllin (2009) are more widely applicable; see Meckes (2009) for a comparison of these two different approaches. In addition, bounds for the so-called Fisher information distance for sums of i.i.d. random vectors are given in Nourdin et al. (2014).

The chapter is organised as follows. Section 4.1 treats the case of i.n.i.d. random vectors. After presenting the main result for the asymptotic normality of the MLE under sufficient regularity conditions, we give an upper bound on the distributional distance between the distribution of the vector MLE and the multivariate normal distribution. Special attention is given to linear regression models with an application to the simplest case of the straight-line model. In Section 4.2, under weaker regularity conditions, we present the upper bound for the case of i.i.d. random vectors. Specific results for independent random variables that follow the normal distribution with unknown mean and variance are also given. Subsection 4.2.2 contains an upper bound on the distributional distance of interest, which holds even in cases where no analytic expression of the vector MLE is available. We illustrate the results through the Beta distribution with both shape parameters unknown. The results of this chapter have been submitted to the *Journal of Multivariate Analysis* as a single-author paper.

4.1 Non-identically distributed random vectors

In this section we examine the case of i.n.i.d. t -dimensional random vectors, for $t \in \mathbb{N}$. Apart from the assumptions set in Propositions 2.1 and 2.2 for the existence and uniqueness of the MLE, we use some regularity conditions, first stated in Hoadley (1971), in order to establish the asymptotic normality of the MLE.

4.1.1 A general bound

For $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ being i.n.i.d. random vectors, we denote by $f_i(\mathbf{x}, \boldsymbol{\theta})$ the probability density (or mass) function for \mathbf{X}_i . Thus, the likelihood function is $L(\boldsymbol{\theta}, \mathbf{x}) = \prod_{i=1}^n f_i(\mathbf{x}_i | \boldsymbol{\theta})$, with its logarithm still being denoted by $l(\boldsymbol{\theta}; \mathbf{x})$. Assuming that the parameter space is an open subset of \mathbb{R}^d , Proposition 2.1 is sufficient for the existence and uniqueness of $\hat{\boldsymbol{\theta}}_n(\mathbf{X})$. In addition, we work under the following regularity conditions for the asymptotic normality of the MLE to hold, (Hoadley, 1971):

- (N1) $\hat{\boldsymbol{\theta}}_n(\mathbf{X}) \xrightarrow{\mathbb{P}} \boldsymbol{\theta}_0$, as $n \rightarrow \infty$, where $\boldsymbol{\theta}_0$ is the true parameter value;
- (N2) the gradient vector $\nabla(\log(f_k(\mathbf{X}_k | \boldsymbol{\theta}))) \in \mathbb{R}^{d \times 1}$ and the Hessian matrix $J_k(\mathbf{X}_k, \boldsymbol{\theta}) = \left\{ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log(f_k(\mathbf{X}_k | \boldsymbol{\theta})) \right\}_{i,j=1,2,\dots,d} \in \mathbb{R}^{d \times d}$ exist almost surely $\forall k \in \{1, 2, \dots, n\}$ with respect to the probability measure \mathbb{P} , where $\mathbb{P}[\cdot]$ denotes probability with respect to $\boldsymbol{\theta}_0$;
- (N3) $J_k(\mathbf{X}_k, \boldsymbol{\theta})$ is a continuous function of $\boldsymbol{\theta}$, uniformly in k , almost surely with respect to \mathbb{P} and is a measurable function of \mathbf{X}_k ;
- (N4) $\mathbb{E}[\nabla(\log(f_k(\mathbf{X}_k | \boldsymbol{\theta})))] = \mathbf{0}$, $k = 1, 2, \dots, n$;
- (N5) with $[\nabla(\log(f_k(\mathbf{X}_k | \boldsymbol{\theta})))]^\top \in \mathbb{R}^{1 \times d}$ denoting the transpose of $\nabla(\log(f_k(\mathbf{X}_k | \boldsymbol{\theta})))$,

$$I_k(\boldsymbol{\theta}) = \mathbb{E}[[\nabla(\log(f_k(\mathbf{X}_k | \boldsymbol{\theta})))][\nabla(\log(f_k(\mathbf{X}_k | \boldsymbol{\theta})))^\top] = -\mathbb{E}[J_k(\mathbf{X}_k, \boldsymbol{\theta})];$$

(N6) for

$$\bar{I}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{j=1}^n I_j(\boldsymbol{\theta}), \quad (4.3)$$

there exists a matrix $\bar{I}(\boldsymbol{\theta}) \in \mathbb{R}^{d \times d}$ such that $\bar{I}_n(\boldsymbol{\theta}) \xrightarrow[n \rightarrow \infty]{} \bar{I}(\boldsymbol{\theta})$. In addition, $\bar{I}_n(\boldsymbol{\theta}), \bar{I}(\boldsymbol{\theta})$ are symmetric matrices for all $\boldsymbol{\theta}$ and $\bar{I}(\boldsymbol{\theta})$ is positive definite;

- (N7) for some $\delta > 0$, $\frac{\sum_k \mathbb{E}[\lambda^\top \nabla(\log(f_k(\mathbf{X}_k | \boldsymbol{\theta}_0)))]^2 + \delta}{n^{\frac{2+\delta}{2}}} \xrightarrow[n \rightarrow \infty]{} 0$ for all $\boldsymbol{\lambda} \in \mathbb{R}^d$;

(N8) with $\|\cdot\|$ the ordinary Euclidean norm on \mathbb{R}^d , then for $k, i, j \in \{1, 2, \dots, d\}$ there exist $\varepsilon > 0$, $K > 0$, $\delta > 0$ and random variables $B_{k,ij}(\mathbf{X}_k)$ such that

- (i) $\sup \left\{ \left| \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log(f_k(\mathbf{X}_k | \mathbf{t})) \right| : \|\mathbf{t} - \boldsymbol{\theta}_0\| \leq \varepsilon \right\} \leq B_{k,ij}(\mathbf{X}_k);$
- (ii) $E |B_{k,ij}(\mathbf{X}_k)|^{1+\delta} \leq K.$

Assuming that $\hat{\boldsymbol{\theta}}_n(\mathbf{X})$ exists and is unique, the following theorem gives the result for the asymptotic normality of the MLE in the i.n.i.d. case, in a slightly different way than the result in Hoadley (1971).

Theorem 4.1. *Let $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ be independent random vectors with probability density (or mass) functions $f_i(\mathbf{x}_i | \boldsymbol{\theta})$, where $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^d$. Assume that the MLE exists and it is unique and that the regularity conditions (N1)-(N8) hold. Also let $\mathbf{Z} \sim N_d(\mathbf{0}, I_{d \times d})$, where $\mathbf{0}$ is the $d \times 1$ zero vector and $I_{d \times d}$ is the $d \times d$ identity matrix. Then,*

$$\left[\sum_{k=1}^n I_k(\boldsymbol{\theta}_0) \right]^{\frac{1}{2}} (\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0) \xrightarrow[n \rightarrow \infty]{d} \mathbf{Z}. \quad (4.4)$$

Proof. Hoadley (1971) proves in Theorem 2, p.1983 that under the regularity conditions (N1)-(N8)

$$\sqrt{n} (\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0) \xrightarrow[n \rightarrow \infty]{d} [\bar{I}(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \mathbf{Z}.$$

Using this result and (N6) we obtain that

$$\left[\frac{1}{n} \sum_{k=1}^n I_k(\boldsymbol{\theta}_0) \right]^{\frac{1}{2}} \sqrt{n} (\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0) \xrightarrow[n \rightarrow \infty]{d} [\bar{I}(\boldsymbol{\theta}_0)]^{\frac{1}{2}} [\bar{I}(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \mathbf{Z} = \mathbf{Z},$$

which is the result of the theorem. □

Our motivation remains to assess the quality of the normal approximation in (4.4) through explicit upper bounds on the distributional distance of interest. From now on, unless otherwise stated, $\bar{I}_n(\boldsymbol{\theta})$ is as in (4.3). Let the subscript (m) denote an index for

which $\left| \hat{\theta}_n(\mathbf{x})_{(m)} - \theta_{0(m)} \right|$ is the largest among the d components;

$(m) \in \{1, 2, \dots, d\}$ is such that

$$\left| \hat{\theta}_n(\mathbf{x})_{(m)} - \theta_{0(m)} \right| \geq \left| \hat{\theta}_n(\mathbf{x})_j - \theta_{0j} \right|, \quad \forall j \in \{1, 2, \dots, d\}$$

and also, for ease of presentation, let

$$C_{(m)} = C_{(m)}(\mathbf{X}, \boldsymbol{\theta}_0) := \hat{\theta}_n(\mathbf{X})_{(m)} - \theta_{0(m)}. \quad (4.5)$$

Our main result is as follows.

Theorem 4.2. *Let $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ be i.n.i.d. \mathbb{R}^t -valued, $t \in \mathbb{Z}^+$, random vectors with probability density (or mass) function $f_i(\mathbf{x}_i | \boldsymbol{\theta})$, for which the parameter space Θ is an open subset of \mathbb{R}^d . Assume that (N1)-(N8) are satisfied and that the MLE exists and it is unique. In addition, assume that for any $\boldsymbol{\theta}_0 \in \Theta$ there exists $0 < \varepsilon = \varepsilon(\boldsymbol{\theta}_0)$ and functions $M_{kjl}(\mathbf{x})$, $\forall k, j, l \in \{1, 2, \dots, d\}$ such that $\left| \frac{\partial^3}{\partial \theta_k \partial \theta_j \partial \theta_l} l(\boldsymbol{\theta}, \mathbf{x}) \right| \leq M_{kjl}(\mathbf{x})$ for all $\boldsymbol{\theta} \in \Theta$ such that $|\theta_j - \theta_{0j}| < \varepsilon \forall j \in \{1, 2, \dots, d\}$. Also, for $C_{(m)}$ as in (4.5), assume that $\mathbb{E} \left((M_{k_{jv}}(\mathbf{X}))^2 \middle| |C_{(m)}| < \varepsilon \right) < \infty$. Let $\{\mathbf{X}'_i, i = 1, 2, \dots, n\}$ be an independent copy of $\{\mathbf{X}_i, i = 1, 2, \dots, n\}$. For $h \in H$, with H as in (4.1) and $\mathbf{Z} \sim \mathbf{N}_d(\mathbf{0}, I_{d \times d})$ it holds that*

$$\begin{aligned} & \left| \mathbb{E} \left[h \left(\sqrt{n} [\bar{I}_n(\boldsymbol{\theta}_0)]^{\frac{1}{2}} \left(\hat{\theta}_n(\mathbf{X}) - \boldsymbol{\theta}_0 \right) \right) \right] - \mathbb{E}[h(\mathbf{Z})] \right| \\ & \leq \frac{\|h\|_1}{\sqrt{n}} K_1(\boldsymbol{\theta}_0) + \frac{\|h\|_2}{\sqrt{n}} K_2(\boldsymbol{\theta}_0) + \frac{\|h\|_3}{\sqrt{n}} K_3(\boldsymbol{\theta}_0) + \frac{2\|h\|}{\varepsilon^2} \mathbb{E} \left(\sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right), \end{aligned} \quad (4.6)$$

where

$$\begin{aligned} K_1(\boldsymbol{\theta}_0) &= \sum_{k=1}^d \sum_{l=1}^d \left| \left[[\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{lk} \right| \\ & \quad \times \sum_{j=1}^d \left[\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right] \mathbb{E} \left[\left(\frac{\partial^2}{\partial \theta_j \partial \theta_k} l(\boldsymbol{\theta}_0; \mathbf{X}) + n [\bar{I}_n(\boldsymbol{\theta}_0)]_{kj} \right)^2 \right] \right]^{\frac{1}{2}} \end{aligned}$$

$$\begin{aligned}
& + \frac{1}{2} \left\{ \sum_{k=1}^d \sum_{l=1}^d \left| [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right|_{lk} \left| \sum_{j=1}^d \sum_{v=1}^d \left[\mathbb{E} \left((\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \boldsymbol{\theta}_{0_j})^2 (\hat{\boldsymbol{\theta}}_n(\mathbf{X})_v - \boldsymbol{\theta}_{0_v})^2 \right) \right]^{\frac{1}{2}} \right. \right. \\
& \quad \left. \left. \times \left[\mathbb{E} \left((M_{k j v}(\mathbf{X}))^2 \right) |C_{(m)}| < \varepsilon \right) \right]^{\frac{1}{2}} \right\}, \tag{4.7}
\end{aligned}$$

$$\begin{aligned}
K_2(\boldsymbol{\theta}_0) &= \frac{1}{4\sqrt{n}} \sum_{j=1}^d \left[\sum_{i=1}^n \text{Var} \left(\left(\sum_{k=1}^d [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{jk} \frac{\partial}{\partial \boldsymbol{\theta}_k} \log(f_i(\mathbf{X}_i | \boldsymbol{\theta}_0)) \right)^2 \right]^{\frac{1}{2}} \\
&+ \frac{1}{2\sqrt{n}} \sum_{k=1}^{d-1} \sum_{j>k}^d \left[\sum_{i=1}^n \text{Var} \left(\sum_{q=1}^d \sum_{v=1}^d [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{jq} [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{kv} \\
&\quad \times \frac{\partial}{\partial \boldsymbol{\theta}_q} \log(f_i(\mathbf{X}_i | \boldsymbol{\theta}_0)) \frac{\partial}{\partial \boldsymbol{\theta}_v} \log(f_i(\mathbf{X}_i | \boldsymbol{\theta}_0)) \right]^{\frac{1}{2}} \tag{4.8}
\end{aligned}$$

and

$$\begin{aligned}
K_3(\boldsymbol{\theta}_0) &= \frac{1}{12n} \sum_{i=1}^n \mathbb{E} \left(\sum_{m=1}^d \left| \sum_{l=1}^d [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right|_{ml} \right. \\
&\quad \left. \times \left(\frac{\partial}{\partial \boldsymbol{\theta}_l} \{ \log(f_i(\mathbf{X}'_i | \boldsymbol{\theta}_0)) - \log(f_i(\mathbf{X}_i | \boldsymbol{\theta}_0)) \} \right) \right|^3 \Bigg) . \tag{4.9}
\end{aligned}$$

The following lemma can be seen as a multivariate extension of Chebyshev's other inequality in Lemma 3.1 and it is useful for bounding conditional expectations, which sometimes can be difficult to derive.

Lemma 4.2. *Let $\mathbf{M} \in \mathbb{R}^d$ be a random vector with $M_i > 0 \forall i = 1, 2, \dots, d$ and $\varepsilon > 0$. For every continuous function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ such that $f(\mathbf{m})$ is increasing and $f(\mathbf{m}) \geq 0$, for $m_i > 0 \forall i \in \{1, 2, \dots, d\}$, where $\mathbf{m} = (m_1, m_2, \dots, m_d)$,*

$$\mathbb{E}[f(\mathbf{M}) | M_i < \varepsilon \forall i = 1, 2, \dots, d] \leq \mathbb{E}[f(\mathbf{M})].$$

Proof of Lemma 4.2. Let $k \in \{1, 2, \dots, d\}$. We set $M_{d+1} = 0$. It will be shown that for $k = 1, 2, \dots, d$ we have that

$$\mathbb{E}[f(\mathbf{M}) | M_i < \varepsilon, i = k, \dots, d] \leq \mathbb{E}[f(\mathbf{M}) | M_i < \varepsilon, i = k+1, \dots, d].$$

From the law of total expectation,

$$\begin{aligned} & \mathbb{E}[f(\mathbf{M}) | M_i < \varepsilon, i = k+1, \dots, d] \\ &= \mathbb{E}[f(\mathbf{M}) | M_i < \varepsilon, i = k, \dots, d] \mathbb{P}[M_k < \varepsilon | M_i < \varepsilon, i = k+1, \dots, d] \\ &+ \mathbb{E}[f(\mathbf{M}) | M_i < \varepsilon, i = k+1, \dots, d, M_k \geq \varepsilon] \mathbb{P}[M_k \geq \varepsilon | M_i < \varepsilon, i = k+1, \dots, d]. \end{aligned}$$

Using that

$$\mathbb{P}[M_k < \varepsilon | M_i < \varepsilon, i = k+1, \dots, d] = 1 - \mathbb{P}[M_k \geq \varepsilon | M_i < \varepsilon, i = k+1, \dots, d]$$

yields

$$\begin{aligned} & \mathbb{E}[f(\mathbf{M}) | M_i < \varepsilon, i = k+1, \dots, d] \\ &= \mathbb{E}[f(\mathbf{M}) | M_i < \varepsilon, i = k, \dots, d] \\ &+ \mathbb{P}[M_k \geq \varepsilon | M_i < \varepsilon, i = k+1, \dots, d] \left\{ \mathbb{E}[f(\mathbf{M}) | M_i < \varepsilon, i = k+1, \dots, d, M_k \geq \varepsilon] \right. \\ &\quad \left. - \mathbb{E}[f(\mathbf{M}) | M_i < \varepsilon, i = k, \dots, d] \right\}. \end{aligned} \tag{4.10}$$

Since $f(\mathbf{m})$ is an increasing function,

$$\mathbb{E}[f(\mathbf{M}) | M_i < \varepsilon, i = k+1, \dots, d, M_k \geq \varepsilon] - \mathbb{E}[f(\mathbf{M}) | M_i < \varepsilon, i = k, \dots, d] \geq 0. \tag{4.11}$$

Applying this to (4.10) gives that

$$\mathbb{E}[f(\mathbf{M})|M_i < \varepsilon, i = k, \dots, d] \leq \mathbb{E}[f(\mathbf{M})|M_i < \varepsilon, i = k+1, \dots, d].$$

A simple iteration over k gives that

$$\mathbb{E}[f(\mathbf{M})|M_i < \varepsilon \forall i = 1, 2, \dots, d] \leq \mathbb{E}[f(\mathbf{M})], \quad (4.12)$$

which is the result of the lemma. \square

Proof of Theorem 4.2. The regularity conditions and the definition of the MLE give that

$$\frac{\partial}{\partial \theta_k} l(\hat{\theta}_n(\mathbf{x}); \mathbf{x}) = 0 \quad \forall k \in \{1, 2, \dots, d\}.$$

Using a second-order Taylor expansion of $\frac{\partial}{\partial \theta_k} l(\hat{\theta}_n(\mathbf{x}); \mathbf{x})$ about θ_0 ,

$$\begin{aligned} 0 &= \frac{\partial}{\partial \theta_k} l(\theta_0; \mathbf{x}) + \sum_{j=1}^d (\hat{\theta}_n(\mathbf{x})_j - \theta_{0j}) \left(\frac{\partial^2}{\partial \theta_k \partial \theta_j} l(\theta_0; \mathbf{x}) \right) \\ &\quad + \frac{1}{2} \sum_{j=1}^d \sum_{q=1}^d (\hat{\theta}_n(\mathbf{x})_j - \theta_{0j}) (\hat{\theta}_n(\mathbf{x})_q - \theta_{0q}) \left(\frac{\partial^3}{\partial \theta_k \partial \theta_j \partial \theta_q} l(\theta; \mathbf{x}) \Big|_{\theta=\theta_0^*} \right), \end{aligned}$$

where θ_0^* is between $\hat{\theta}_n(\mathbf{x})$ and θ_0 . Rearranging,

$$\begin{aligned} &\sum_{j=1}^d (\hat{\theta}_n(\mathbf{x})_j - \theta_{0j}) \left(\frac{\partial^2}{\partial \theta_k \partial \theta_j} l(\theta_0; \mathbf{x}) \right) \\ &= -\frac{\partial}{\partial \theta_k} l(\theta_0; \mathbf{x}) - \frac{1}{2} \sum_{j=1}^d \sum_{q=1}^d (\hat{\theta}_n(\mathbf{x})_j - \theta_{0j}) (\hat{\theta}_n(\mathbf{x})_q - \theta_{0q}) \left(\frac{\partial^3}{\partial \theta_k \partial \theta_j \partial \theta_q} l(\theta; \mathbf{x}) \Big|_{\theta=\theta_0^*} \right), \end{aligned}$$

which gives that for all $k = 1, 2, \dots, d$,

$$\begin{aligned} &\sum_{j=1}^d (\hat{\theta}_n(\mathbf{x})_j - \theta_{0j}) \left(-n[\bar{I}_n(\theta_0)]_{kj} + \frac{\partial^2}{\partial \theta_k \partial \theta_j} l(\theta_0; \mathbf{x}) + n[\bar{I}_n(\theta_0)]_{kj} \right) \\ &= -\frac{\partial}{\partial \theta_k} l(\theta_0; \mathbf{x}) - \frac{1}{2} \sum_{j=1}^d \sum_{q=1}^d (\hat{\theta}_n(\mathbf{x})_j - \theta_{0j}) (\hat{\theta}_n(\mathbf{x})_q - \theta_{0q}) \left(\frac{\partial^3}{\partial \theta_k \partial \theta_j \partial \theta_q} l(\theta; \mathbf{x}) \Big|_{\theta=\theta_0^*} \right) \end{aligned}$$

so that

$$\begin{aligned}
& \sum_{j=1}^d n[\bar{I}_n(\boldsymbol{\theta}_0)]_{kj}(\hat{\boldsymbol{\theta}}_n(\mathbf{x})_j - \theta_{0j}) = \frac{\partial}{\partial \theta_k} l(\boldsymbol{\theta}_0; \mathbf{x}) \\
& + \sum_{j=1}^d (\hat{\boldsymbol{\theta}}_n(\mathbf{x})_j - \theta_{0j}) \left(\frac{\partial^2}{\partial \theta_k \partial \theta_j} l(\boldsymbol{\theta}_0; \mathbf{x}) + n[\bar{I}_n(\boldsymbol{\theta}_0)]_{kj} \right) \\
& + \frac{1}{2} \sum_{j=1}^d \sum_{q=1}^d (\hat{\boldsymbol{\theta}}_n(\mathbf{x})_j - \theta_{0j})(\hat{\boldsymbol{\theta}}_n(\mathbf{x})_q - \theta_{0q}) \left(\frac{\partial^3}{\partial \theta_k \partial \theta_j \partial \theta_q} l(\boldsymbol{\theta}; \mathbf{x}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0^*} \right). \quad (4.13)
\end{aligned}$$

Using (4.13), which holds $\forall k \in \{1, 2, \dots, d\}$,

$$\begin{aligned}
& \sqrt{n}[\bar{I}_n(\boldsymbol{\theta}_0)]^{\frac{1}{2}}(\hat{\boldsymbol{\theta}}_n(\mathbf{x}) - \boldsymbol{\theta}_0) \\
& = \frac{1}{\sqrt{n}}[\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \left\{ \nabla(l(\boldsymbol{\theta}_0; \mathbf{x})) \right. \\
& \quad \left. + \sum_{j=1}^d (\hat{\boldsymbol{\theta}}_n(\mathbf{x})_j - \theta_{0j}) \left(\nabla \left(\frac{\partial}{\partial \theta_j} l(\boldsymbol{\theta}_0; \mathbf{x}) \right) + n[\bar{I}_n(\boldsymbol{\theta}_0)]_{[j]} \right) \right\} \\
& + \frac{1}{2\sqrt{n}}[\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \left\{ \sum_{j=1}^d \sum_{q=1}^d (\hat{\boldsymbol{\theta}}_n(\mathbf{x})_j - \theta_{0j})(\hat{\boldsymbol{\theta}}_n(\mathbf{x})_q - \theta_{0q}) \right. \\
& \quad \left. \times \left(\nabla \left(\frac{\partial^2}{\partial \theta_j \partial \theta_q} l(\boldsymbol{\theta}; \mathbf{x}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0^*} \right) \right) \right\},
\end{aligned}$$

where $[\bar{I}_n(\boldsymbol{\theta}_0)]_{[j]}$ is the j^{th} column of the matrix $\bar{I}_n(\boldsymbol{\theta}_0)$. The triangle inequality gives

$$\begin{aligned}
& \left| \mathbb{E} \left[h \left(\sqrt{n}[\bar{I}_n(\boldsymbol{\theta}_0)]^{\frac{1}{2}}(\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0) \right) \right] - \mathbb{E}[h(\mathbf{Z})] \right| \\
& \leq \left| \mathbb{E} \left[h \left(\frac{1}{\sqrt{n}}[\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \nabla(l(\boldsymbol{\theta}_0; \mathbf{X})) \right) \right] - \mathbb{E}[h(\mathbf{Z})] \right| \quad (4.14)
\end{aligned}$$

$$\begin{aligned}
& + \left| \mathbb{E} \left[h \left(\sqrt{n}[\bar{I}_n(\boldsymbol{\theta}_0)]^{\frac{1}{2}}(\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0) \right) - h \left(\frac{1}{\sqrt{n}}[\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \nabla(l(\boldsymbol{\theta}_0; \mathbf{X})) \right) \right] \right|. \quad (4.15)
\end{aligned}$$

Step 1: Upper bound for (4.14). First of all, $\nabla(l(\boldsymbol{\theta}_0; \mathbf{x})) = \sum_{i=1}^n \nabla(\log(f_i(\mathbf{x}_i|\boldsymbol{\theta}_0)))$ due to independence. The results of Lemma 4.1 will be used for

$$\mathbf{W} = \frac{1}{\sqrt{n}} [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \sum_{i=1}^n \nabla(\log(f_i(\mathbf{X}_i|\boldsymbol{\theta}_0))) = (W_1, W_2, \dots, W_d)^\top \in \mathbb{R}^{d \times 1}. \quad (4.16)$$

From (4.16) we have that for all $k \in \{1, 2, \dots, d\}$, $W_k = \sum_{i=1}^n \xi_{ik}$, with

$$\xi_{ik} = \frac{1}{\sqrt{n}} \sum_{j=1}^d \left[[\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{kj} \frac{\partial}{\partial \theta_j} \log(f_i(\mathbf{X}_i|\boldsymbol{\theta}_0)). \quad (4.17)$$

From the regularity conditions, $E(\nabla(l(\boldsymbol{\theta}_0; \mathbf{X}))) = \mathbf{0}$ and thus $E(\mathbf{W}) = \mathbf{0}$. Also, $\bar{I}_n(\boldsymbol{\theta}_0)$ is symmetric. Therefore, $[\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}}$ is also symmetric. From the regularity conditions $\sum_{i=1}^n \text{Cov}[\nabla(\log(f_i(\mathbf{X}_i|\boldsymbol{\theta}_0)))] = n\bar{I}_n(\boldsymbol{\theta}_0)$ and basic calculations show that $\text{Cov}(\mathbf{W}) = I_{d \times d}$. Since $E(\mathbf{W}) = \mathbf{0}$ and $E(\mathbf{W}\mathbf{W}^\top) = I_{d \times d}$, the first assumption of Lemma 4.1 is satisfied.

The next step is to define \mathbf{W}' such that $(\mathbf{W}, \mathbf{W}')$ is an exchangeable pair satisfying (4.2). Let $\{\mathbf{X}'_i, i = 1, 2, \dots, n\}$ be an independent copy of $\{\mathbf{X}_i, i = 1, 2, \dots, n\}$ and let the index $I \in \{1, 2, \dots, n\}$ follow the uniform distribution on $\{1, 2, \dots, n\}$, independently of $\{\mathbf{X}_i, \mathbf{X}'_i, i = 1, 2, \dots, n\}$. Let

$$\xi'_{ik} = \frac{1}{\sqrt{n}} \sum_{j=1}^d \left[[\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{kj} \frac{\partial}{\partial \theta_j} \log(f_i(\mathbf{X}'_i|\boldsymbol{\theta}_0))$$

and

$$W'_k = W_k - \xi_{Ik} + \xi'_{Ik}, \quad \forall k \in \{1, 2, \dots, d\},$$

with $E(W'_k - W_k | \mathbf{W}) = E(\xi'_{Ik} - \xi_{Ik} | \mathbf{W}) = -E(\xi_{Ik} | \mathbf{W}) = -\frac{1}{n} \sum_{i=1}^n E(\xi_{ik} | \mathbf{W}) = -\frac{W_k}{n}$.

Hence (4.2) is satisfied with $\Lambda = \frac{1}{n} I_{d \times d}$ and $\mathbf{R} = \mathbf{0}$. We conclude that the assumptions

of Lemma 4.1 are satisfied and Lemma 4.1 gives that

$$|\mathbb{E}[h(\mathbf{W})] - \mathbb{E}[h(\mathbf{Z})]| \leq n \left(\frac{\|h\|_2}{4} \sum_{i=1}^d \sum_{j=1}^d [\text{Var}(\mathbb{E}[(W'_i - W_i)(W'_j - W_j) | \mathbf{W}])]^{\frac{1}{2}} \right) \quad (4.18)$$

$$+ n \left(\frac{\|h\|_3}{12} \sum_{i=1}^d \sum_{j=1}^d \sum_{k=1}^d \mathbb{E}|(W'_i - W_i)(W'_j - W_j)(W'_k - W_k)| \right). \quad (4.19)$$

To bound the variance of the conditional expectations in (4.18), let $\mathcal{A} = \sigma(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n)$. Since $\sigma(\mathbf{W}) \subset \mathcal{A}$, for any random variable Y , $\text{Var}(\mathbb{E}[Y | \mathbf{W}]) \leq \text{Var}(\mathbb{E}[Y | \mathcal{A}])$. Simple steps yield

$$(4.18) \leq n \frac{\|h\|_2}{4} \left\{ \sum_{j=1}^d [\text{Var}(\mathbb{E}[(\xi'_{Ij} - \xi_{Ij})^2 | \mathcal{A}])]^{\frac{1}{2}} + 2 \sum_{j>k}^d \sum_{k=1}^{d-1} [\text{Var}(\mathbb{E}[(\xi'_{Ik} - \xi_{Ik})(\xi'_{Ij} - \xi_{Ij}) | \mathcal{A}])]^{\frac{1}{2}} \right\}. \quad (4.20)$$

After straightforward calculations and using that $\{\mathbf{X}'_i, i = 1, 2, \dots, n\}$ is an independent copy of $\{\mathbf{X}_i, i = 1, 2, \dots, n\}$ and ξ'_{ik} is independent of \mathcal{A} we get that

$$(4.20) = n \frac{\|h\|_2}{4} \left\{ \sum_{j=1}^d [\text{Var}(\mathbb{E}[(\xi'_{Ij})^2] - 2\mathbb{E}[\xi'_{Ij}]\mathbb{E}[\xi_{Ij} | \mathcal{A}] + \mathbb{E}[\xi_{Ij}^2 | \mathcal{A}])]^{\frac{1}{2}} + 2 \sum_{k=1}^{d-1} \sum_{j>k}^d [\text{Var}(\mathbb{E}[\xi'_{Ik}\xi'_{Ij}] - \mathbb{E}[\xi'_{Ij}]\mathbb{E}[\xi_{Ik} | \mathcal{A}] - \mathbb{E}[\xi'_{Ik}]\mathbb{E}[\xi_{Ij} | \mathcal{A}] + \mathbb{E}[\xi_{Ik}\xi_{Ij} | \mathcal{A}])]^{\frac{1}{2}} \right\}. \quad (4.21)$$

Using that $\mathbb{E}[\xi'_{ik}] = 0$,

$$\begin{aligned}
(4.21) &= n \frac{\|h\|_2}{4} \left\{ \sum_{j=1}^d [\text{Var}(\mathbb{E}[(\xi'_{Ij})^2] + \mathbb{E}[\xi_{Ij}^2|\mathcal{A}])]^{\frac{1}{2}} \right. \\
&\quad \left. + 2 \sum_{k=1}^{d-1} \sum_{j>k}^d [\text{Var}(\mathbb{E}[\xi'_{Ik}\xi'_{Ij}] + \mathbb{E}[\xi_{Ik}\xi_{Ij}|\mathcal{A}])]^{\frac{1}{2}} \right\} \\
&= n \frac{\|h\|_2}{4} \left\{ \sum_{j=1}^d \left[\frac{1}{n^2} \text{Var} \left(\sum_{i=1}^n \mathbb{E}[\xi_{ij}^2|\mathcal{A}] \right) \right]^{\frac{1}{2}} + 2 \sum_{j>k}^d \sum_{k=1}^{d-1} \left[\frac{1}{n^2} \text{Var} \left(\sum_{i=1}^n \mathbb{E}[\xi_{ik}\xi_{ij}|\mathcal{A}] \right) \right]^{\frac{1}{2}} \right\} \\
&= \frac{\|h\|_2}{4} \left\{ \sum_{j=1}^d \left[\text{Var} \left(\sum_{i=1}^n \xi_{ij}^2 \right) \right]^{\frac{1}{2}} + 2 \sum_{j>k}^d \sum_{k=1}^{d-1} \left[\text{Var} \left(\sum_{i=1}^n \xi_{ik}\xi_{ij} \right) \right]^{\frac{1}{2}} \right\} = \frac{\|h\|_2}{\sqrt{n}} K_2(\boldsymbol{\theta}_0),
\end{aligned}$$

with $K_2(\boldsymbol{\theta}_0)$ in (4.8). For (4.19), using (4.17), after basic calculations we obtain that

$$(4.19) \leq \frac{\|h\|_3}{\sqrt{n}} K_3(\boldsymbol{\theta}_0),$$

with $K_3(\boldsymbol{\theta}_0)$ as in (4.9). Thus,

$$(4.14) \leq \frac{\|h\|_2}{\sqrt{n}} K_2(\boldsymbol{\theta}_0) + \frac{\|h\|_3}{\sqrt{n}} K_3(\boldsymbol{\theta}_0). \quad (4.22)$$

Step 2: Upper bound for (4.15). We will use multivariate Taylor expansions, conditional expectations and Markov's inequality. The approach is similar to the one in the single-parameter case of Chapter 3. For ease of presentation, let

$$\begin{aligned}
\mathbf{R}_1(\boldsymbol{\theta}_0; \mathbf{x}) &= \frac{1}{2\sqrt{n}} [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \\
&\quad \times \sum_{j=1}^d \sum_{q=1}^d (\hat{\boldsymbol{\theta}}_n(\mathbf{x})_j - \boldsymbol{\theta}_{0j}) (\hat{\boldsymbol{\theta}}_n(\mathbf{x})_q - \boldsymbol{\theta}_{0q}) \left(\nabla \left(\frac{\partial^2}{\partial \boldsymbol{\theta}_j \partial \boldsymbol{\theta}_q} l(\boldsymbol{\theta}; \mathbf{x}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0^*} \right) \right) \\
T_1 &= T_1(\boldsymbol{\theta}_0; \mathbf{X}, h) := h \left(\sqrt{n} [\bar{I}_n(\boldsymbol{\theta}_0)]^{\frac{1}{2}} (\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0) \right) \\
&\quad - h \left(\frac{1}{\sqrt{n}} [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} (\nabla(l(\boldsymbol{\theta}_0; \mathbf{X}))) \right) \\
T_2 &= T_2(\boldsymbol{\theta}_0; \mathbf{X}, h) := h \left(\frac{1}{\sqrt{n}} [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} (\nabla(l(\boldsymbol{\theta}_0; \mathbf{x}))) + \mathbf{R}_1(\boldsymbol{\theta}_0; \mathbf{x}) \right) \\
&\quad - h \left(\frac{1}{\sqrt{n}} [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} (\nabla(l(\boldsymbol{\theta}_0; \mathbf{X}))) \right). \quad (4.23)
\end{aligned}$$

Using the above notation,

$$\begin{aligned}
 (4.15) &= |\mathbb{E}[T_1]| \\
 &= \left| \mathbb{E} \left[h \left(\sqrt{n} [\bar{I}_n(\boldsymbol{\theta}_0)]^{\frac{1}{2}} \left(\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0 \right) \right) \right. \right. \\
 &\quad \left. \left. - h \left(\frac{1}{\sqrt{n}} [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} (\nabla(l(\boldsymbol{\theta}_0; \mathbf{x}))) + \mathbf{R}_1(\boldsymbol{\theta}_0; \mathbf{X}) \right) + T_2 \right] \right|.
 \end{aligned}$$

Now, the triangle inequality gives that

$$\begin{aligned}
 (4.15) &\leq \mathbb{E} \left| h \left(\sqrt{n} [\bar{I}_n(\boldsymbol{\theta}_0)]^{\frac{1}{2}} \left(\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0 \right) \right) \right. \\
 &\quad \left. - h \left(\frac{1}{\sqrt{n}} [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} (\nabla(l(\boldsymbol{\theta}_0; \mathbf{x}))) + \mathbf{R}_1(\boldsymbol{\theta}_0; \mathbf{X}) \right) \right| \quad (4.24)
 \end{aligned}$$

$$+ \mathbb{E}|T_2|. \quad (4.25)$$

The next step is based on a first order multivariate Taylor expansion. With $A_{[j]}$ denoting the j^{th} row of a matrix A , we get that

$$\begin{aligned}
 h \left(\sqrt{n} [\bar{I}_n(\boldsymbol{\theta}_0)]^{\frac{1}{2}} \left(\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0 \right) \right) &= h \left(\frac{1}{\sqrt{n}} [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} (\nabla(l(\boldsymbol{\theta}_0; \mathbf{x}))) + \mathbf{R}_1(\boldsymbol{\theta}_0; \mathbf{X}) \right) \\
 &+ \sum_{j=1}^d \left(\sqrt{n} [\bar{I}_n(\boldsymbol{\theta}_0)]^{\frac{1}{2}} \right)_{[j]} \left(\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0 \right) - \frac{1}{\sqrt{n}} [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \Big|_{[j]} (\nabla(l(\boldsymbol{\theta}_0; \mathbf{X}))) \\
 &- \frac{1}{2\sqrt{n}} [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \Big|_{[j]} \left\{ \sum_{k=1}^d \sum_{q=1}^d (\hat{\boldsymbol{\theta}}_n(\mathbf{x})_k - \boldsymbol{\theta}_{0_k}) (\hat{\boldsymbol{\theta}}_n(\mathbf{x})_q - \boldsymbol{\theta}_{0_q}) \right. \\
 &\quad \left. \times \left(\nabla \left(\frac{\partial^2}{\partial \boldsymbol{\theta}_k \partial \boldsymbol{\theta}_q} l(\boldsymbol{\theta}; \mathbf{x}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0^*} \right) \right) \right\} \frac{\partial}{\partial x_j} h(\mathbf{t}(\mathbf{X})),
 \end{aligned}$$

where $\mathbf{t}(\mathbf{X})$ is between $\frac{1}{\sqrt{n}} [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} (\nabla(l(\boldsymbol{\theta}_0; \mathbf{X}))) + \mathbf{R}_1(\boldsymbol{\theta}_0; \mathbf{X})$ and $\sqrt{n} [\bar{I}_n(\boldsymbol{\theta}_0)]^{\frac{1}{2}} \left(\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0 \right)$. Thus,

$$\begin{aligned}
& \left| h \left(\sqrt{n} [\bar{I}_n(\boldsymbol{\theta}_0)]^{\frac{1}{2}} \left(\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0 \right) \right) - h \left(\frac{1}{\sqrt{n}} [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} (\nabla(l(\boldsymbol{\theta}_0; \mathbf{x}))) + \mathbf{R}_1(\boldsymbol{\theta}_0; \mathbf{X}) \right) \right| \\
& \leq \|h\|_1 \left| \sum_{j=1}^d \left(\sqrt{n} [\bar{I}_n(\boldsymbol{\theta}_0)]^{\frac{1}{2}} \right)_{[j]} \left(\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0 \right) - \frac{1}{\sqrt{n}} [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right|_{[j]} \nabla(l(\boldsymbol{\theta}_0; \mathbf{X})) \\
& \quad - \frac{1}{2\sqrt{n}} [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right|_{[j]} \left\{ \sum_{k=1}^d \sum_{q=1}^d (\hat{\boldsymbol{\theta}}_n(\mathbf{x})_k - \theta_{0_k}) (\hat{\boldsymbol{\theta}}_n(\mathbf{x})_q - \theta_{0_q}) \right. \\
& \quad \left. \times \left(\nabla \left(\frac{\partial^2}{\partial \theta_k \partial \theta_q} l(\boldsymbol{\theta}; \mathbf{x}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0^*} \right) \right) \right\} \right|.
\end{aligned}$$

Using (4.13) component-wise, we conclude that

$$\begin{aligned}
(4.24) & \leq \frac{\|h\|_1}{\sqrt{n}} \sum_{k=1}^d \sum_{l=1}^d \left| [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right|_{lk} \left| \sum_{j=1}^d \mathbb{E} \left((\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \theta_{0_j}) \right. \right. \\
& \quad \left. \left. \times \left(\frac{\partial^2}{\partial \theta_j \partial \theta_k} l(\boldsymbol{\theta}_0; \mathbf{X}) + n [\bar{I}_n(\boldsymbol{\theta}_0)]_{kj} \right) \right) \right| \\
& \leq \frac{\|h\|_1}{\sqrt{n}} \sum_{k=1}^d \sum_{l=1}^d \left| [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right|_{lk} \left| \sum_{j=1}^d \left[\mathbb{E} \left[(\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \theta_{0_j})^2 \right] \right. \right. \\
& \quad \left. \left. \times \mathbb{E} \left[\left(\frac{\partial^2}{\partial \theta_j \partial \theta_k} l(\boldsymbol{\theta}_0; \mathbf{X}) + n [\bar{I}_n(\boldsymbol{\theta}_0)]_{kj} \right)^2 \right] \right] \right|^{\frac{1}{2}},
\end{aligned} \tag{4.26}$$

where for the second inequality we used Cauchy-Schwarz inequality. To bound now $\mathbb{E}|T_2|$, with T_2 as in (4.23), we need to take into account that $\frac{\partial^3}{\partial \theta_k \partial \theta_q \partial \theta_j} l(\boldsymbol{\theta}; \mathbf{x}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0^*}$ is not uniformly bounded. Therefore, with $\varepsilon > 0$, the law of total expectation yields

$$\mathbb{E}|T_2| \leq 2\|h\| \mathbb{P} \left(\left| \hat{\boldsymbol{\theta}}_n(\mathbf{X})_{(m)} - \boldsymbol{\theta}_{0(m)} \right| \geq \varepsilon \right) + \mathbb{E} \left(|T_2| \mid \left| \hat{\boldsymbol{\theta}}_n(\mathbf{X})_{(m)} - \boldsymbol{\theta}_{0(m)} \right| < \varepsilon \right).$$

Using Chebyshev's inequality,

$$\mathbb{P} \left(\left| \hat{\boldsymbol{\theta}}_n(\mathbf{X})_{(m)} - \boldsymbol{\theta}_{0(m)} \right| \geq \varepsilon \right) \leq \frac{1}{\varepsilon^2} \mathbb{E} \left(\sum_{j=1}^d (\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \theta_{0_j})^2 \right). \tag{4.27}$$

To bound $E \left(|T_2| \left| \left| \hat{\theta}_n(\mathbf{X})_{(m)} - \theta_{0(m)} \right| < \varepsilon \right)$, we use a first-order Taylor expansion and

$$\begin{aligned} h \left(\frac{1}{\sqrt{n}} [\bar{I}_n(\theta_0)]^{-\frac{1}{2}} (\nabla(l(\theta_0; \mathbf{x}))) + \mathbf{R}_1(\theta_0; \mathbf{x}) \right) &= h \left(\frac{1}{\sqrt{n}} [\bar{I}_n(\theta_0)]^{-\frac{1}{2}} (\nabla(l(\theta_0; \mathbf{X}))) \right) \\ &+ \sum_{j=1}^d \left(\frac{1}{2\sqrt{n}} [\bar{I}_n(\theta_0)]^{-\frac{1}{2}} \right)_{[j]} \left\{ \sum_{k=1}^d \sum_{q=1}^d (\hat{\theta}_n(\mathbf{x})_k - \theta_{0_k}) (\hat{\theta}_n(\mathbf{x})_q - \theta_{0_q}) \right. \\ &\quad \left. \times \left(\nabla \left(\frac{\partial^2}{\partial \theta_j \partial \theta_q} l(\theta; \mathbf{x}) \Big|_{\theta=\theta_0^*} \right) \right) \right\} \frac{\partial}{\partial x_j} h(\tilde{\mathbf{t}}(\mathbf{X})), \end{aligned}$$

where now $\tilde{\mathbf{t}}(\mathbf{X})$ is between $\frac{1}{\sqrt{n}} [\bar{I}_n(\theta_0)]^{-\frac{1}{2}} (\nabla(l(\theta_0; \mathbf{x}))) + \mathbf{R}_1(\theta_0; \mathbf{x})$ and $\frac{1}{\sqrt{n}} [\bar{I}_n(\theta_0)]^{-\frac{1}{2}} (\nabla(l(\theta_0; \mathbf{x})))$. Using again (4.13),

$$\begin{aligned} |T_2| &\leq \frac{\|h\|_1}{2\sqrt{n}} \sum_{k=1}^d \sum_{l=1}^d \left| [\bar{I}_n(\theta_0)]^{-\frac{1}{2}} \right|_{lk} \\ &\quad \times \left\{ \sum_{j=1}^d \sum_{v=1}^d \left| (\hat{\theta}_n(\mathbf{X})_j - \theta_{0_j}) (\hat{\theta}_n(\mathbf{X})_v - \theta_{0_v}) \frac{\partial^3}{\partial \theta_k \partial \theta_j \partial \theta_v} l(\theta; \mathbf{X}) \Big|_{\theta=\theta_0^*} \right| \right\}. \end{aligned} \quad (4.28)$$

Therefore, from (4.27) and (4.28) we have for $C_{(m)}$ as in (4.5) that

$$\begin{aligned} E|T_2| &\leq \frac{2\|h\|}{\varepsilon^2} E \left(\sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0_j})^2 \right) \\ &+ \frac{\|h\|_1}{2\sqrt{n}} \sum_{k=1}^d \sum_{l=1}^d \left| [\bar{I}_n(\theta_0)]^{-\frac{1}{2}} \right|_{lk} E \left(\sum_{j=1}^d \sum_{v=1}^d \left| (\hat{\theta}_n(\mathbf{X})_j - \theta_{0_j}) (\hat{\theta}_n(\mathbf{X})_v - \theta_{0_v}) \right. \right. \\ &\quad \left. \left. \times \frac{\partial^3}{\partial \theta_k \partial \theta_j \partial \theta_v} l(\theta; \mathbf{X}) \Big|_{\theta=\theta_0^*} \right| \Big| C_{(m)} < \varepsilon \right). \end{aligned}$$

The Cauchy-Schwarz inequality and Lemma 4.2 yield

$$\begin{aligned}
\mathbb{E}|T_2| &\leq \frac{2\|h\|}{\varepsilon^2} \mathbb{E} \left(\sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right) \\
&+ \frac{\|h\|_1}{2\sqrt{n}} \left\{ \sum_{k=1}^d \sum_{l=1}^d \left| [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right|_{lk} \left| \sum_{j=1}^d \sum_{v=1}^d \left[\mathbb{E} \left((\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 (\hat{\theta}_n(\mathbf{X})_v - \theta_{0v})^2 \right) \right]^{\frac{1}{2}} \right. \right. \\
&\quad \times \left. \left[\mathbb{E} \left(\left(\sup_{\substack{\boldsymbol{\theta}: |\theta_m - \theta_{0m}| < \varepsilon \\ \forall m \in \{1, 2, \dots, d\}}} \left| \frac{\partial^3}{\partial \theta_k \partial \theta_j \partial \theta_v} l(\boldsymbol{\theta}; \mathbf{X}) \right| \right)^2 \right) \right]^{\frac{1}{2}} \right\} \\
&\leq \frac{2\|h\|}{\varepsilon^2} \mathbb{E} \left(\sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right) \\
&+ \frac{\|h\|_1}{2\sqrt{n}} \left\{ \sum_{k=1}^d \sum_{l=1}^d \left| [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right|_{lk} \left| \sum_{j=1}^d \sum_{v=1}^d \left[\mathbb{E} \left((\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 (\hat{\theta}_n(\mathbf{X})_v - \theta_{0v})^2 \right) \right]^{\frac{1}{2}} \right. \right. \\
&\quad \times \left. \left[\mathbb{E} \left((M_{k j v}(\mathbf{X}))^2 \right) \right]^{\frac{1}{2}} \right\}. \tag{4.29}
\end{aligned}$$

Therefore, from (4.26) and (4.29) we obtain that

$$(4.15) \leq \frac{2\|h\|}{\varepsilon^2} \mathbb{E} \left(\sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right) + \frac{\|h\|_1}{\sqrt{n}} K_1(\boldsymbol{\theta}_0), \tag{4.30}$$

where $K_1(\boldsymbol{\theta}_0)$ is as in (4.7). Using now (4.22) and (4.30) gives the assertion. \square

Remark 4.1. (I) Assuming that $\bar{I}_n(\boldsymbol{\theta}_0) = \mathcal{O}(1)$ in (4.3) and using Theorem 4.1 yields, for fixed d , $\mathbb{E} \left(\sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right) = \mathcal{O} \left(\frac{1}{n} \right)$. To see this, use that from the asymptotic normality of the MLE as expressed in Theorem 4.1, $\sqrt{n} \mathbb{E} \left(\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0 \right) \xrightarrow[n \rightarrow \infty]{} \mathbf{0}$ and thus

$$\mathbb{E} \left(\hat{\theta}_n(\mathbf{X})_j - \theta_{0j} \right) = o \left(\frac{1}{\sqrt{n}} \right), \quad \forall j \in \{1, 2, \dots, d\}. \tag{4.31}$$

Theorem 4.1 shows that $\text{Cov} \left(\sqrt{n} [\bar{I}_n(\boldsymbol{\theta}_0)]^{\frac{1}{2}} \left(\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0 \right) \right) \xrightarrow{n \rightarrow \infty} I_{d \times d}$. Therefore,

$$n [\bar{I}_n(\boldsymbol{\theta}_0)]^{\frac{1}{2}} \text{Cov} \left(\hat{\boldsymbol{\theta}}_n(\mathbf{X}) \right) [\bar{I}_n(\boldsymbol{\theta}_0)]^{\frac{1}{2}} \xrightarrow{n \rightarrow \infty} I_{d \times d}. \quad (4.32)$$

Assuming that the matrix $\bar{I}_n(\boldsymbol{\theta}_0)$ as defined in (4.3) is $\mathcal{O}(1)$, it follows from (4.32) that

$$\text{Var} \left(\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j \right) = \mathcal{O} \left(\frac{1}{n} \right), \quad \forall j \in \{1, 2, \dots, d\}.$$

Combining these results,

$$\mathbb{E} \left[\left(\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \boldsymbol{\theta}_{0j} \right)^2 \right] = \text{Var} \left(\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j \right) + \left[\mathbb{E} \left(\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \boldsymbol{\theta}_{0j} \right) \right]^2 = \mathcal{O} \left(\frac{1}{n} \right). \quad (4.33)$$

Furthermore, using (4.31), (4.33) then if $\bar{I}_n(\boldsymbol{\theta}_0) = \mathcal{O}(1)$ it can be deduced that

$$\begin{aligned} K_1(\boldsymbol{\theta}_0) &= \mathcal{O}(1) \\ K_2(\boldsymbol{\theta}_0) &= \mathcal{O}(1) \\ K_3(\boldsymbol{\theta}_0) &= \mathcal{O}(1), \end{aligned} \quad (4.34)$$

where $K_1(\boldsymbol{\theta}_0), K_2(\boldsymbol{\theta}_0), K_3(\boldsymbol{\theta}_0)$ are as in (4.7), (4.8), (4.9), respectively. Hence, using (4.33) and (4.34), if $\bar{I}_n(\boldsymbol{\theta}_0) = \mathcal{O}(1)$ then the upper bound in Theorem 4.2 is $\mathcal{O} \left(\frac{1}{\sqrt{n}} \right)$.

(2) Often Hölder's inequality will be used to bound the third term as the calculation of absolute third moments can be quite complicated, even for simple multi-parameter distributions.

(3) In terms of the dimensionality d of the parameter, $K_1(\boldsymbol{\theta}_0) = \mathcal{O}(d^4)$, $K_2(\boldsymbol{\theta}_0) = \mathcal{O}(d^4)$ and $K_3(\boldsymbol{\theta}_0) = \mathcal{O}(d^8)$ as can be deduced from (4.7), (4.8) and (4.9), respectively. The last term of the bound in (4.6) is of order d in terms of the dimensionality of the parameter. Thus, for $d \gg n$ the bound does not behave well, but d could grow moderately with n . For example $d = o(n^\alpha)$, $0 < \alpha < \frac{1}{16}$ would still yield a bound which goes to zero as n goes to infinity.

4.1.2 Linear regression

This subsection calculates the bound in (4.6) for linear regression models. The asymptotic normality of the MLE in linear regression models has been proven in Fahrmeir and Kaufmann (1985). We give the example of a straight-line regression and the bound turns out to be of order $\mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$, where n is the sample size. The following notation is used throughout this subsection. The vector $\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)^\top \in \mathbb{R}^{n \times 1}$ denotes the response variable for the linear regression, while $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_d)^\top \in \mathbb{R}^{d \times 1}$ is the vector of the d parameters and $\boldsymbol{\epsilon} = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)^\top \in \mathbb{R}^{n \times 1}$ is the vector of the error terms, which are i.i.d. random variables with $\epsilon_i \sim N(0, \sigma^2) \forall i \in \{1, 2, \dots, n\}$. The true value of the unknown parameter $\boldsymbol{\beta}$ is denoted by $\boldsymbol{\beta}_0 = (\beta_{0_1}, \beta_{0_2}, \dots, \beta_{0_d})^\top \in \mathbb{R}^{d \times 1}$. The design matrix is

$$X = \begin{pmatrix} 1 & x_{1,2} & \dots & x_{1,d} \\ 1 & x_{2,2} & \dots & x_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,2} & \dots & x_{n,d} \end{pmatrix}.$$

For the model

$$\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

the aim is to find bounds on the distributional distance between the distribution of the MLE, $\hat{\boldsymbol{\beta}}$, and the normal distribution. In the case where $\beta_j, j = 1, 2, \dots, d$ are independent random variables, then we have linear random-effects models and bounds related to the normal approximation of variance component maximum likelihood estimators are given in Dicker and Erdogdu (2016+).

The probability density function for Y_i is

$$f(y_i|\boldsymbol{\beta}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (y_i - X_{[i]}\boldsymbol{\beta})^2 \right\}, \quad (4.35)$$

where $X_{[i]}$ denotes the i^{th} row of the design matrix. The parameter space $\Theta = \mathbb{R}^d$ is open and if $X^\top X$ is of full rank, the matrix $X^\top X$ is invertible and the vector MLE is

$$\hat{\beta} = (X^\top X)^{-1} X^\top Y. \quad (4.36)$$

We now bound the corresponding distributional distance.

Corollary 4.1. *Let $Y_i, i \in \{1, 2, \dots, n\}$ be independent normal random variables with*

$$Y_i \sim N(X_{[i]}\beta_0, \sigma^2),$$

where $X_{[i]}$ denotes the i^{th} row of the design matrix X and σ^2 is known. Assume that the $d \times d$ matrix $X^\top X$ is of full rank. Let $\{Y'_i, i = 1, 2, \dots, n\}$ be an independent copy of $\{Y_i, i = 1, 2, \dots, n\}$ and $\mathbf{Z} \sim N_d(\mathbf{0}, I_{d \times d})$ and $\bar{I}_n(\beta) = \frac{1}{n} \sum_{j=1}^n I_j(\beta)$. Then for $h : \mathbb{R}^d \rightarrow \mathbb{R}$, a three times differentiable function with bounded derivatives,

$$\begin{aligned} & \left| E \left[h \left(\sqrt{n} [\bar{I}_n(\beta)]^{\frac{1}{2}} (\hat{\beta} - \beta_0) \right) \right] - E[h(\mathbf{Z})] \right| \\ & \leq \frac{\|h\|_2}{4} \sum_{j=1}^d \left[\sum_{i=1}^n \text{Var} \left(\left(\sum_{k=1}^d \frac{X_{ik}}{\sigma} [X^\top X]^{-\frac{1}{2}} \right)_{jk} \left(Y_i - \sum_{m=1}^d X_{im} \beta_{0_m} \right) \right)^2 \right]^{\frac{1}{2}} \\ & + \frac{\|h\|_2}{2} \sum_{j>k}^d \sum_{k=1}^{d-1} \left[\sum_{i=1}^n \text{Var} \left(\sum_{q=1}^d \sum_{v=1}^d \frac{X_{iq} X_{iv}}{\sigma^2} [X^\top X]^{-\frac{1}{2}} \right)_{jq} [X^\top X]^{-\frac{1}{2}} \right]_{kv} \\ & \quad \times \left(Y_i - \sum_{m=1}^d X_{im} \beta_{0_m} \right)^2 \right]^{\frac{1}{2}} \\ & + \frac{\|h\|_3}{12} \sum_{i=1}^n E \left(\sum_{m=1}^d \left| \sum_{l=1}^d \frac{X_{il}}{\sigma} [X^\top X]^{-\frac{1}{2}} \right|_{ml} (Y_i - Y'_i) \right)^3. \end{aligned} \quad (4.37)$$

Proof. Using (4.35) and a vector representation, we get the following expressions for the log-likelihood and its first derivative with respect to the parameter vector:

$$\begin{aligned}
l(\beta; \mathbf{y}) &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - X\beta)^\top (\mathbf{y} - X\beta) \\
\frac{d}{d\beta} l(\beta; \mathbf{y}) &= -\frac{1}{2\sigma^2} \frac{d}{d\beta} \{(\mathbf{y}^\top - \beta^\top X^\top)(\mathbf{y} - X\beta)\} \\
&= -\frac{1}{2\sigma^2} \frac{d}{d\beta} \{\mathbf{y}^\top \mathbf{y} - \beta^\top X^\top \mathbf{y} - \mathbf{y}^\top X \beta + \beta^\top X^\top X \beta\} \\
&= \frac{1}{\sigma^2} (X^\top \mathbf{y} - X^\top X \beta).
\end{aligned}$$

The Hessian matrix for the log-likelihood function,

$$H(\beta; \mathbf{y}) = \left\{ \frac{\partial^2}{\partial \beta_i \partial \beta_j} l(\beta; \mathbf{y}) \right\}_{i,j=1,\dots,d} = -\frac{1}{\sigma^2} X^\top X,$$

does not depend on \mathbf{y} . Thus, $\bar{I}_n(\beta_0) = \frac{1}{n\sigma^2} X^\top X$ and so $[\bar{I}_n(\beta)]^{-\frac{1}{2}} = \sigma\sqrt{n} [X^\top X]^{-\frac{1}{2}}$. The result in (4.36) gives that

$$\begin{aligned}
\sqrt{n} [\bar{I}_n(\beta_0)]^{\frac{1}{2}} (\hat{\beta} - \beta_0) &= \frac{1}{\sigma} \left\{ [X^\top X]^{-\frac{1}{2}} X^\top \mathbf{Y} - [X^\top X]^{\frac{1}{2}} \beta_0 \right\} \\
&= \frac{1}{\sqrt{n}} \left[\sigma\sqrt{n} [X^\top X]^{-\frac{1}{2}} \right] \frac{1}{\sigma^2} (X^\top \mathbf{Y} - X^\top X \beta_0) \\
&= \frac{1}{\sqrt{n}} [\bar{I}_n(\beta_0)]^{-\frac{1}{2}} \frac{d}{d\beta} l(\beta; \mathbf{y}) \Big|_{\beta=\beta_0}. \tag{4.38}
\end{aligned}$$

The expression in (4.38) is the same as W in (4.16) and therefore the quantity of interest $\left| E \left[h \left(\sqrt{n} [\bar{I}_n(\beta_0)]^{\frac{1}{2}} (\hat{\beta} - \beta_0) \right) \right] - E[h(\mathbf{Z})] \right|$ is equal to (4.14) with (4.15) being equal to zero for this specific case of the linear regression model. Therefore, using (4.22) and

$$\frac{\partial}{\partial \beta_k} \log(f_i(Y_i | \beta_0)) = \frac{X_{ik}}{\sigma^2} \left(Y_i - \sum_{m=1}^d X_{im} \beta_{0m} \right)$$

in Theorem 4.2 yields the result of the Corollary. \square

Example: The simple linear model ($d=2$)

Here, we apply the results of (4.37) to the case of a straight-line regression with two unknown parameters. The model is

$$Y_i = \beta_1 + \beta_2(x_i - \bar{x}) + \varepsilon_i, \quad \forall i \in \{1, 2, \dots, n\}.$$

The unknown parameters β_1 and β_2 are the intercept and slope of the regression, respectively. The i.i.d. random variables $\varepsilon_i \sim N(0, \sigma^2)$, $\forall i \in \{1, 2, \dots, n\}$ are as before. It is well known that the MLE exists, it is unique and equal to $\hat{\beta} = \left(\bar{Y}, \frac{\sum_{i=1}^n (x_i - \bar{x}) Y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)^\top$. The following corollary of Theorem 4.2 gives, for this specific example, an upper bound on the distributional distance between the exact distribution of $\hat{\beta}$ and the normal distribution.

Corollary 4.2. *Let Y_1, Y_2, \dots, Y_n be independent random variables with $Y_i \sim N(\beta_1 + \beta_2(x_i - \bar{x}), \sigma^2)$, with the case of $x_i = x_j$ for $i \neq j$, $i, j \in \{1, 2, \dots, n\}$ being excluded. For $\mathbf{Z} \sim N_2(\mathbf{0}, I_{2 \times 2})$ and $h \in H$ as in (4.1),*

$$\begin{aligned} \left| \mathbb{E} \left[h \left(\sqrt{n} [\bar{I}_n(\beta)]^{\frac{1}{2}} (\hat{\beta} - \beta_0) \right) \right] - \mathbb{E}[h(\mathbf{Z})] \right| &\leq \frac{\|h\|_2}{4} \left(\sqrt{\frac{2}{n}} + \frac{\sqrt{2 \sum_{i=1}^n (x_i - \bar{x})^4}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) + \frac{\|h\|_2}{\sqrt{2n}} \\ &\quad + \frac{8\|h\|_3}{3\sqrt{\pi}} \left(\frac{1}{\sqrt{n}} + \frac{\sum_{i=1}^n |x_i - \bar{x}|^3}{[\sum_{i=1}^n (x_i - \bar{x})^2]^{\frac{3}{2}}} \right). \end{aligned}$$

Proof. We have that

$$X = \begin{pmatrix} 1 & x_1 - \bar{x} \\ 1 & x_2 - \bar{x} \\ \vdots & \vdots \\ 1 & x_n - \bar{x} \end{pmatrix}, \quad X^\top X = \begin{pmatrix} n & 0 \\ 0 & \sum_{i=1}^n (x_i - \bar{x})^2 \end{pmatrix}, \quad [X^\top X]^{-\frac{1}{2}} = \begin{pmatrix} \frac{1}{\sqrt{n}} & 0 \\ 0 & \frac{1}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \end{pmatrix}, \quad (4.39)$$

which shows that $X^\top X$ is invertible if and only if $\sum_{i=1}^n (x_i - \bar{x})^2 \neq 0$, which holds if x_i 's are not all identical. The quantities of the bound in (4.37) are calculated for this specific

case. We use that $Y_i - \beta_1 - (x_i - \bar{x})\beta_2 \stackrel{d}{=} \sigma Z_i$, where $Z_i \sim N(0, 1)$. For the first term in (4.37) we obtain that

$$\begin{aligned}
& \sum_{j=1}^2 \left[\sum_{i=1}^n \text{Var} \left(\left(\sum_{k=1}^2 \frac{X_{ik}}{\sigma} [X^\top X]^{-\frac{1}{2}} \right)_{jk} \left(Y_i - \sum_{m=1}^2 X_{im} \beta_m \right) \right)^2 \right]^{\frac{1}{2}} \\
&= \sum_{j=1}^2 \left[\sum_{i=1}^n \text{Var} \left(\left(\left(\frac{X_{i1}}{\sigma} [X^\top X]^{-\frac{1}{2}} \right)_{j1} + \frac{X_{i2}}{\sigma} [X^\top X]^{-\frac{1}{2}} \right)_{j2} (\sigma Z_i) \right)^2 \right]^{\frac{1}{2}} \\
&= \left[\sum_{i=1}^n \text{Var} \left(\left(\frac{1}{\sqrt{n}} Z_i \right)^2 \right) \right]^{\frac{1}{2}} + \left[\sum_{i=1}^n \text{Var} \left(\left(\frac{x_i - \bar{x}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} Z_i \right)^2 \right) \right]^{\frac{1}{2}} \\
&= \frac{1}{n} \left[\sum_{i=1}^n \text{Var} (Z_i^2) \right]^{\frac{1}{2}} + \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \left[\sum_{i=1}^n (x_i - \bar{x})^4 \text{Var} (Z_i^2) \right]^{\frac{1}{2}} \\
&= \sqrt{\frac{2}{n}} + \frac{\sqrt{2 \sum_{i=1}^n (x_i - \bar{x})^4}}{\sum_{i=1}^n (x_i - \bar{x})^2}. \tag{4.40}
\end{aligned}$$

For an upper bound for the second term of (4.37), since $d = 2$ we have that $k = 1, j = 2$ leading to

$$\begin{aligned}
& \left[\sum_{i=1}^n \text{Var} \left(\sum_{q=1}^2 \sum_{v=1}^2 \frac{X_{iq} X_{iv}}{\sigma^2} [X^\top X]^{-\frac{1}{2}} \right)_{2q} [X^\top X]^{-\frac{1}{2}} \right)_{1v} \left(Y_i - \sum_{m=1}^2 X_{im} \beta_m \right)^2 \right]^{\frac{1}{2}} \\
&= \left[\sum_{i=1}^n \text{Var} \left(\frac{X_{i2} X_{i1}}{\sigma^2} [X^\top X]^{-\frac{1}{2}} \right)_{22} [X^\top X]^{-\frac{1}{2}} \right)_{11} (\sigma Z_i)^2 \right]^{\frac{1}{2}} \\
&= \left[\sum_{i=1}^n \text{Var} \left(\frac{x_i - \bar{x}}{\sqrt{n \sum_{i=1}^n (x_i - \bar{x})^2}} Z_i^2 \right) \right]^{\frac{1}{2}} = \frac{1}{\sqrt{n \sum_{i=1}^n (x_i - \bar{x})^2}} \left[\sum_{i=1}^n \text{Var} (Z_i^2) (x_i - \bar{x})^2 \right]^{\frac{1}{2}} \\
&= \sqrt{\frac{2}{n}}. \tag{4.41}
\end{aligned}$$

For the final term of (4.37), because Y'_i is an independent copy of Y_i , then $Y'_i - Y_i \sim N(0, 2\sigma^2)$, with $E|Y'_i - Y_i|^3 = 8 \frac{\sigma^3}{\sqrt{\pi}}$. Using that

$$(|a| + |b|)^3 \leq 4(|a|^3 + |b|^3), \quad a, b \in \mathbb{R} \tag{4.42}$$

yields

$$\begin{aligned}
& \sum_{i=1}^n \mathbb{E} \left(\sum_{m=1}^2 \left| \sum_{l=1}^2 \frac{X_{il}}{\sigma} [X^\top X]^{-\frac{1}{2}} \right]_{ml} (Y_i - Y'_i) \right|^3 \\
&= \sum_{i=1}^n \mathbb{E} \left(\left| \left(\frac{X_{i1}}{\sigma} [X^\top X]^{-\frac{1}{2}} \right]_{11} + \frac{X_{i2}}{\sigma} [X^\top X]^{-\frac{1}{2}} \right]_{22} \right| (Y_i - Y'_i) \right|^3 \\
&\leq \sum_{i=1}^n \mathbb{E} \left(\left(\frac{1}{\sigma \sqrt{n}} + \frac{|x_i - \bar{x}|}{\sigma \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \right) |Y'_i - Y_i| \right)^3 \\
&\leq 4 \sum_{i=1}^n \left(\frac{8}{n^{\frac{3}{2}} \sqrt{\pi}} + \frac{8|x_i - \bar{x}|^3}{[\sum_{i=1}^n (x_i - \bar{x})^2]^{\frac{3}{2}} \sqrt{\pi}} \right) = \frac{32}{\sqrt{\pi}} \left(\frac{1}{\sqrt{n}} + \frac{\sum_{i=1}^n |x_i - \bar{x}|^3}{[\sum_{i=1}^n (x_i - \bar{x})^2]^{\frac{3}{2}}} \right). \quad (4.43)
\end{aligned}$$

Summarizing, in the case of Y_1, Y_2, \dots, Y_n being independent random variables with $Y_i \sim N(\beta_1 + \beta_2(x_i - \bar{x}), \sigma^2)$, we apply to (4.37) the results of (4.40), (4.41) and (4.43) to obtain the assertion of the corollary. \square

Remark 4.2. The bound is $\mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$.

4.2 Identically distributed random vectors

In this section we use weaker regularity conditions than (N1)-(N8) in order to find an upper bound in the case of independent and identically distributed random vectors. Following Davison (2008), we make the following assumptions:

- (R.C.1) The densities defined by any two different values of θ are distinct;
- (R.C.2) the log-likelihood function is three times differentiable with respect to the unknown vector parameter, θ , and the third partial derivatives are continuous in θ ;
- (R.C.3) for any $\theta_0 \in \Theta$ and for \mathbb{X} denoting the support of the data $\mathbf{x} = (x_1, x_2, \dots, x_n)$, there exists $\varepsilon_0 > 0$ and functions $M_{rst}(\mathbf{x})$ (they can depend on θ_0), such that

for $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_d)$ and $r, s, t, j = 1, \dots, d$,

$$\frac{1}{n} \left| \frac{\partial^3}{\partial \theta_r \partial \theta_s \partial \theta_t} l(\boldsymbol{\theta}; \mathbf{x}) \right| \leq M_{rst}(\mathbf{x}), \forall \mathbf{x} \in \mathbb{X}, |\theta_j - \theta_{0j}| < \varepsilon_0,$$

with $E(M_{rst}(\mathbf{X})) < \infty$;

(R.C.4) the expected Fisher information matrix $I(\boldsymbol{\theta})$ is finite, symmetric and positive definite. For $r, s = 1, 2, \dots, d$, its elements satisfy

$$[I(\boldsymbol{\theta})]_{rs} = E \left\{ \frac{\partial}{\partial \theta_r} l(\boldsymbol{\theta}; \mathbf{X}) \frac{\partial}{\partial \theta_s} l(\boldsymbol{\theta}; \mathbf{X}) \right\} = E \left\{ - \frac{\partial^2}{\partial \theta_r \partial \theta_s} l(\boldsymbol{\theta}; \mathbf{X}) \right\}.$$

This condition implies that $I(\boldsymbol{\theta})$ is the covariance matrix of the score vector.

These regularity conditions in the multi-parameter case resemble those in Section 2.2 where it is assumed that the parameter is scalar. From now on, unless otherwise stated, the notation $I(\boldsymbol{\theta})$ stands for the expected Fisher information matrix for one random vector. Under (R.C.1)-(R.C.4), (Davison, 2008, p.118) shows that

$$\sqrt{n} [I(\boldsymbol{\theta}_0)]^{\frac{1}{2}} (\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0) \xrightarrow[n \rightarrow \infty]{d} N_d(\mathbf{0}, I_{d \times d}).$$

The upper bound on the distributional distance between the distribution of a vector MLE and the multivariate normal in the case of i.i.d. random vectors is the same as the bound in Theorem 4.2 and thus it is not given again. The bound can be simplified due to the fact that in the i.i.d. case $\bar{I}_n(\boldsymbol{\theta}_0) = I(\boldsymbol{\theta}_0)$ and $f_i(\mathbf{x}_i) = f(\mathbf{x}_i)$, $\forall i \in \{1, 2, \dots, n\}$. In the next example of independent random variables from the normal distribution with both mean and variance unknown the bound can be easily calculated and it is, as expected, of the order $\frac{1}{\sqrt{n}}$.

4.2.1 The normal distribution

Here, we apply Theorem 4.2 in the case of X_1, X_2, \dots, X_n independent and identically distributed random variables from $N(\mu, \sigma^2)$ with $\theta_0 = (\mu, \sigma^2)$. It is well-known that the MLE exists, it is unique and equal to $\hat{\theta}_n(\mathbf{X}) = (\hat{\mu}, \hat{\sigma}^2)^\top = (\bar{X}, \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2)^\top$; see for example (Davison, 2008, p.116). In addition, the regularity conditions (R.C.1)-(R.C.4) are satisfied.

Corollary 4.3. *Let X_1, X_2, \dots, X_n be i.i.d. random variables that follow the $N(\mu, \sigma^2)$ distribution. For $\mathbf{Z} \sim N_d(\mathbf{0}, I_{d \times d})$ and $h \in H$ as defined in (4.1),*

$$\begin{aligned} \left| \mathbb{E} \left[h \left(\sqrt{n} [I(\theta_0)]^{\frac{1}{2}} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) \right] - \mathbb{E}[h(\mathbf{Z})] \right| &\leq \frac{5}{2} \frac{\|h\|_2}{\sqrt{n}} + 19 \frac{\|h\|_3}{\sqrt{n}} \\ &+ \frac{8\|h\|}{n\sigma^2} (1 + 2\sigma^2) + 4\sqrt{2} \frac{\|h\|_1}{\sqrt{n}} + \frac{4\|h\|_1}{\sqrt{n}} \left[3 + 16\sqrt{2} \sqrt{\frac{1}{n} + \frac{\sigma^2}{4}} \right] \\ &+ \frac{32\|h\|_1}{\sqrt{n}} \left[1 + 648 \left[\left(\frac{3}{2} + \frac{\sigma^2}{4} \right)^2 + \frac{3}{n^2} \right] \right]^{\frac{1}{2}}. \end{aligned} \quad (4.44)$$

Proof. The first and second-order partial derivatives are

$$\begin{aligned} \frac{\partial}{\partial \mu} l(\theta_0; \mathbf{x}) &= \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu), & \frac{\partial}{\partial \sigma^2} l(\theta_0; \mathbf{x}) &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2, \\ \frac{\partial^2}{\partial \mu^2} l(\theta_0; \mathbf{x}) &= -\frac{n}{\sigma^2}, & \frac{\partial^2}{\partial \sigma^4} l(\theta_0; \mathbf{x}) &= \frac{n}{2\sigma^4} - \frac{1}{\sigma^6} \sum_{i=1}^n (x_i - \mu)^2, \\ \frac{\partial^2}{\partial \mu \partial \sigma^2} l(\theta_0; \mathbf{x}) &= \frac{\partial^2}{\partial \sigma^2 \partial \mu} l(\theta_0; \mathbf{x}) = -\frac{n}{\sigma^4} (\bar{X} - \mu). \end{aligned} \quad (4.45)$$

Hence the expected Fisher Information matrix for one random variable is

$$I(\theta_0) = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{2\sigma^4} \end{pmatrix}, \text{ so that } [I(\theta_0)]^{-\frac{1}{2}} = \begin{pmatrix} \sigma & 0 \\ 0 & \sqrt{2}\sigma^2 \end{pmatrix}. \quad (4.46)$$

Using (4.16), let

$$\mathbf{W} = \frac{1}{\sqrt{n}} \left(\frac{1}{\sigma} \sum_{i=1}^n (X_i - \mu), -\frac{n}{\sqrt{2}} + \frac{1}{\sqrt{2}\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 \right)^\top = (W_1, W_2)^\top,$$

where $W_i = \sum_{j=1}^n \xi_{ji}$ with $\xi_{i1} = \frac{X_i - \mu}{\sqrt{n}\sigma}$ and $\xi_{i2} = \frac{(X_i - \mu)^2 - \sigma^2}{\sigma^2\sqrt{2n}}$, $i = 1, 2, \dots, n$. Now we bound the terms in Theorem 4.2 in order of appearance. We start with the term $\frac{\|h\|_1}{\sqrt{n}} K_1(\boldsymbol{\theta}_0)$, where $K_1(\boldsymbol{\theta}_0)$ is given in (4.7). For the first quantity, using that $E[(\bar{X} - \mu)^2] = \frac{\sigma^2}{n}$ and $E[(\hat{\sigma}^2 - \sigma^2)^2] = \frac{\sigma^4}{n} (2 - \frac{1}{n})$ yields

$$\begin{aligned} & \frac{\|h\|_1}{\sqrt{n}} \sum_{k=1}^2 \sum_{l=1}^2 \left| [\bar{I}_n(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right|_{lk} \\ & \quad \times \sum_{j=1}^2 \left[E[(\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \boldsymbol{\theta}_{0j})^2] E \left[\left(\frac{\partial^2}{\partial \boldsymbol{\theta}_j \partial \boldsymbol{\theta}_k} l(\boldsymbol{\theta}_0; \mathbf{X}) + n[\bar{I}_n(\boldsymbol{\theta}_0)]_{kj} \right)^2 \right] \right]^{\frac{1}{2}} \\ & = \frac{\sigma \|h\|_1}{\sqrt{n}} \left[E[(\hat{\sigma}^2 - \sigma^2)^2] E \left[\left(\frac{-n(\bar{X} - \mu)}{\sigma^4} \right)^2 \right] \right]^{\frac{1}{2}} \\ & \quad + \frac{\sqrt{2}\sigma^2 \|h\|_1}{\sqrt{n}} \left\{ E[(\bar{X} - \mu)^2] E \left[\left(\frac{-n(\bar{X} - \mu)}{\sigma^4} \right)^2 \right] \right]^{\frac{1}{2}} \\ & \quad + \left[E[(\hat{\sigma}^2 - \sigma^2)^2] E \left[\left(\frac{n}{\sigma^4} - \frac{1}{\sigma^6} \sum_{i=1}^n (X_i - \mu)^2 \right)^2 \right] \right]^{\frac{1}{2}} \right\} \\ & \leq \frac{\sqrt{n} \|h\|_1}{\sigma^3} \sqrt{E[(\hat{\sigma}^2 - \sigma^2)^2]} \sqrt{E[(\bar{X} - \mu)^2]} + \frac{\sqrt{2n} \|h\|_1}{\sigma^2} E[(\bar{X} - \mu)^2] \\ & \quad + \frac{\sqrt{2}\sigma^2 \|h\|_1}{\sqrt{n}} \sqrt{E[(\hat{\sigma}^2 - \sigma^2)^2] E \left[\left(\frac{n}{\sigma^4} - \frac{1}{\sigma^6} \sum_{i=1}^n (X_i - \mu)^2 \right)^2 \right]} \\ & \leq \frac{\|h\|_1}{\sqrt{n}} \sqrt{2 - \frac{1}{n}} + \frac{\sqrt{2} \|h\|_1}{\sqrt{n}} \\ & \quad + \frac{\sqrt{2}\sigma^4 \|h\|_1}{n} \sqrt{\left(2 - \frac{1}{n} \right) \left(\frac{n^2}{\sigma^8} - \frac{2n}{\sigma^{10}} E \left[\sum_{i=1}^n (X_i - \mu)^2 \right] + \frac{1}{\sigma^{12}} \left[E \left[\sum_{i=1}^n (X_i - \mu)^2 \right] \right]^2 \right)} \\ & = \frac{\|h\|_1}{\sqrt{n}} \sqrt{2 - \frac{1}{n}} + \frac{\sqrt{2} \|h\|_1}{\sqrt{n}} + \frac{\sqrt{2}\sigma^4 \|h\|_1}{n} \sqrt{\left(2 - \frac{1}{n} \right) \left(\frac{n^2}{\sigma^8} - \frac{2n^2}{\sigma^8} + \frac{n}{\sigma^8} (2 + n) \right)} \end{aligned}$$

$$= \frac{\|h\|_1}{\sqrt{n}} \left(\sqrt{2 - \frac{1}{n}} + \sqrt{2} \right) + \frac{\sqrt{2}\|h\|_1}{\sqrt{n}} \sqrt{2 \left(2 - \frac{1}{n} \right)} \leq 4\sqrt{2} \frac{\|h\|_1}{\sqrt{n}}. \quad (4.47)$$

For the second quantity in (4.7), we find an upper bound for

$$\sup_{\theta: |\theta_m - \theta_{0m}| < \varepsilon} \left| \frac{\partial^3}{\partial \theta_k \partial \theta_j \partial \theta_i} l(\theta; \mathbf{X}) \right| \leq nM_{kji}(\mathbf{X}),$$

where $m, k, j, i \in \{1, 2\}$ and $M_{kji}(\mathbf{x})$ is as in (R.C.3). Below, $\theta = (\theta_1, \theta_2)$ is the vector parameter and $\theta_0 = (\theta_{01}, \theta_{02}) = (\mu, \sigma^2)$ is the true, unknown value. We have

$$\sup_{\theta: |\theta_m - \theta_{0m}| < \varepsilon} \left| \frac{\partial^3}{\partial \theta_1^3} l(\theta; \mathbf{X}) \right| = 0 =: M_{111}(\mathbf{X})$$

as well as

$$\begin{aligned} \sup_{\theta: |\theta_m - \theta_{0m}| < \varepsilon} \left| \frac{\partial^3}{\partial \theta_2^3} l(\theta; \mathbf{X}) \right| &= \sup_{\theta: |\theta_m - \theta_{0m}| < \varepsilon} \left| -\frac{n}{\theta_2^3} + \frac{3}{\theta_2^4} \sum_{i=1}^n (X_i - \theta_1)^2 \right| \\ &< \frac{n}{(\sigma^2 - \varepsilon)^3} + \frac{3}{(\sigma^2 - \varepsilon)^4} \sup_{\theta: |\theta_m - \theta_{0m}| < \varepsilon} \left| \sum_{i=1}^n (X_i - \bar{X} + \bar{X} - \mu + \mu - \theta_1)^2 \right| \\ &\leq \frac{n}{(\sigma^2 - \varepsilon)^3} + \frac{3}{(\sigma^2 - \varepsilon)^4} \sup_{\theta: |\theta_m - \theta_{0m}| < \varepsilon} \left| 3 \sum_{i=1}^n \left[(X_i - \bar{X})^2 + (\bar{X} - \mu)^2 + (\mu - \theta_1)^2 \right] \right| \\ &< \frac{n}{(\sigma^2 - \varepsilon)^3} + \frac{9n}{(\sigma^2 - \varepsilon)^4} \left(\hat{\sigma}^2 + (\bar{X} - \mu)^2 + \varepsilon^2 \right) =: nM_{222}(\mathbf{X}). \end{aligned} \quad (4.48)$$

Moreover,

$$\begin{aligned} \sup_{\theta: |\theta_m - \theta_{0m}| < \varepsilon} \left| \frac{\partial^3}{\partial \theta_1 \partial \theta_2^2} l(\theta; \mathbf{X}) \right| &= \sup_{\theta: |\theta_m - \theta_{0m}| < \varepsilon} \left| \frac{\partial^3}{\partial \theta_2^2 \partial \theta_1} l(\theta; \mathbf{X}) \right| \\ &= \sup_{\theta: |\theta_m - \theta_{0m}| < \varepsilon} \left| \frac{2n}{\theta_2^3} (\bar{X} - \theta_1) \right| \\ &= \sup_{\theta: |\theta_m - \theta_{0m}| < \varepsilon} \left| \frac{2n}{\theta_2^3} (\bar{X} - \mu + \mu - \theta_1) \right| \\ &< \frac{2n}{(\sigma^2 - \varepsilon)^3} (|\bar{X} - \mu| + \varepsilon) =: nM_{122}(\mathbf{X}) \end{aligned} \quad (4.49)$$

and

$$\begin{aligned} \sup_{\theta: |\theta_m - \theta_{0m}| < \varepsilon} \left| \frac{\partial^3}{\partial \theta_1^2 \partial \theta_2} l(\theta; \mathbf{X}) \right| &= \sup_{\theta: |\theta_m - \theta_{0m}| < \varepsilon} \left| \frac{\partial^3}{\partial \theta_2 \partial \theta_1^2} l(\theta; \mathbf{X}) \right| \\ &= \sup_{\theta: |\theta_m - \theta_{0m}| < \varepsilon} \left| \frac{n}{\theta_2^2} \right| < \frac{n}{(\sigma^2 - \varepsilon)^2} =: nM_{112}(\mathbf{X}). \end{aligned} \quad (4.50)$$

In addition, we calculate $E\left((\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 (\hat{\theta}_n(\mathbf{X})_i - \theta_{0i})^2\right), \forall i, j \in \{1, 2\}$. For this purpose we use that in the case of i.i.d. random variables from the Normal distribution,

$$\text{Cov}\left(\bar{X} - \mu, \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 - \sigma^2\right) = 0 \quad (\text{Casella and Berger, 2002, p.218}). \quad (4.51)$$

We have that

$$E(\bar{X} - \mu)^4 = \frac{3\sigma^4}{n^2}$$

and

$$E\left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 - \sigma^2\right)^4 = \frac{\sigma^8}{n^2} \left(12 + \frac{4}{n} - \frac{15}{n^2}\right) < 16 \frac{\sigma^8}{n^2}, \text{ for all } n \in \mathbb{N}$$

and with $G \sim \chi_{n-1}^2$, so that $E(G) = n - 1$ and $\text{Var}(G) = 2(n - 1)$,

$$\begin{aligned} &E\left[(\bar{X} - \mu)^2 \left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 - \sigma^2\right)^2\right] \\ &= E[(\bar{X} - \mu)^2] E\left[\left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 - \sigma^2\right)^2\right] \quad (\text{using (4.51)}) \\ &= E[(\bar{X} - \mu)^2] \frac{\sigma^4}{n^2} E\left[\left(\sum_{i=1}^n \left(\frac{X_i - \bar{X}}{\sigma}\right)^2 - n\right)^2\right] \\ &= \frac{\sigma^6}{n^3} E[(G - (n - 1) - 1)^2] \\ &= \frac{\sigma^6}{n^3} (E[(G - (n - 1))^2] - 2E[G - (n - 1)] + 1) = \frac{\sigma^6}{n^3} (\text{Var}[G] + 1) \\ &= \frac{\sigma^6}{n^3} (2(n - 1) + 1) = \frac{\sigma^6}{n^2} \left(2 - \frac{1}{n}\right) < 2 \frac{\sigma^6}{n^2}. \end{aligned}$$

For $C_{(m)}$ as in (4.5) and $[I(\theta_0)]^{-\frac{1}{2}}$ in (4.46), the second quantity in $K_1(\theta_0)$ becomes

$$\begin{aligned}
& \frac{\|h\|_1}{2\sqrt{n}} \left\{ \sum_{k=1}^2 \sum_{l=1}^2 \left| [I(\theta_0)]^{-\frac{1}{2}} \right|_{lk} \left| \sum_{j=1}^2 \sum_{i=1}^2 \left[E \left((\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 (\hat{\theta}_n(\mathbf{X})_i - \theta_{0i})^2 \right) \right]^{\frac{1}{2}} \right. \right. \\
& \quad \times \left. \left[E \left((nM_{kji}(\mathbf{X}))^2 \middle| |C_{(m)}| < \varepsilon \right) \right]^{\frac{1}{2}} \right\} \\
&= \frac{\|h\|_1}{2\sqrt{n}} \left\{ \left| [I(\theta_0)]^{-\frac{1}{2}} \right|_{11} \left| \sum_{j=1}^2 \sum_{i=1}^2 \left[E \left((\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 (\hat{\theta}_n(\mathbf{X})_i - \theta_{0i})^2 \right) \right]^{\frac{1}{2}} \right. \right. \\
& \quad \times \left. \left[E \left((nM_{1ji}(\mathbf{X}))^2 \middle| |C_{(m)}| < \varepsilon \right) \right]^{\frac{1}{2}} \right. \\
& \quad + \left| [I(\theta_0)]^{-\frac{1}{2}} \right|_{22} \left| \sum_{j=1}^2 \sum_{i=1}^2 \left[E \left((\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 (\hat{\theta}_n(\mathbf{X})_i - \theta_{0i})^2 \right) \right]^{\frac{1}{2}} \right. \\
& \quad \times \left. \left[E \left((nM_{2ji}(\mathbf{X}))^2 \middle| |C_{(m)}| < \varepsilon \right) \right]^{\frac{1}{2}} \right\} \\
&= \frac{\|h\|_1 \sigma}{2\sqrt{n}} \left\{ \frac{2n}{(\sigma^2 - \varepsilon)^2} \left[E \left((\bar{X} - \mu)^2 (\hat{\sigma}^2 - \sigma^2)^2 \right) \right]^{\frac{1}{2}} \right. \\
& \quad + \left. \left[E \left((\hat{\sigma}^2 - \sigma^2)^4 \right) \right]^{\frac{1}{2}} \frac{2n}{(\sigma^2 - \varepsilon)^3} \left[E \left((|\bar{X} - \mu| + \varepsilon)^2 \middle| |C_{(m)}| < \varepsilon \right) \right]^{\frac{1}{2}} \right\} \\
&+ \frac{\|h\|_1 \sqrt{2} \sigma^2}{2\sqrt{n}} \left\{ \frac{n}{(\sigma^2 - \varepsilon)^2} \left[E (\bar{X} - \mu)^4 \right]^{\frac{1}{2}} \right. \\
& \quad + \frac{4n}{(\sigma^2 - \varepsilon)^3} \left[E \left((\bar{X} - \mu)^2 (\hat{\sigma}^2 - \sigma^2)^2 \right) \right]^{\frac{1}{2}} \left[E \left((|\bar{X} - \mu| + \varepsilon)^2 \middle| |C_{(m)}| < \varepsilon \right) \right]^{\frac{1}{2}} \\
& \quad + n \left[E \left((\hat{\sigma}^2 - \sigma^2)^4 \right) \right]^{\frac{1}{2}} \left[E \left((M_{222}(\mathbf{X}))^2 \middle| |C_{(m)}| < \varepsilon \right) \right]^{\frac{1}{2}} \right\} \\
&\leq \frac{\|h\|_1 \sigma}{2\sqrt{n}} \left\{ \frac{2\sqrt{2} \sigma^3}{(\sigma^2 - \varepsilon)^2} + \frac{8\sigma^4}{(\sigma^2 - \varepsilon)^3} \left[E \left((|\bar{X} - \mu| + \varepsilon)^2 \middle| |C_{(m)}| < \varepsilon \right) \right]^{\frac{1}{2}} \right\} \\
&+ \frac{\|h\|_1 \sigma^2}{\sqrt{2n}} \left\{ \frac{\sqrt{3} \sigma^2}{(\sigma^2 - \varepsilon)^2} + 4\sqrt{2} \frac{\sigma^3}{(\sigma^2 - \varepsilon)^3} \left[E \left((|\bar{X} - \mu| + \varepsilon)^2 \middle| |C_{(m)}| < \varepsilon \right) \right]^{\frac{1}{2}} \right. \\
& \quad \left. + 4\sigma^4 \left[E \left((M_{222}(\mathbf{X}))^2 \middle| |C_{(m)}| < \varepsilon \right) \right]^{\frac{1}{2}} \right\}. \tag{4.52}
\end{aligned}$$

We next bound $\left[\mathbb{E} \left((M_{222}(\mathbf{X}))^2 \middle| |C_{(m)}| < \varepsilon \right) \right]^{\frac{1}{2}}$ and $\left[\mathbb{E} \left((|\bar{X} - \mu| + \varepsilon)^2 \middle| |C_{(m)}| < \varepsilon \right) \right]^{\frac{1}{2}}$. For $\left[\mathbb{E} \left((|\bar{X} - \mu| + \varepsilon)^2 \middle| |C_{(m)}| < \varepsilon \right) \right]^{\frac{1}{2}}$, we have

$$\begin{aligned} \left[\mathbb{E} \left((|\bar{X} - \mu| + \varepsilon)^2 \middle| |C_{(m)}| < \varepsilon \right) \right]^{\frac{1}{2}} &\leq \left[2\mathbb{E} \left((\bar{X} - \mu)^2 \middle| |C_{(m)}| < \varepsilon \right) + 2\varepsilon^2 \right]^{\frac{1}{2}} \\ &\leq \sqrt{2} \left[\mathbb{E} \left[(\bar{X} - \mu)^2 \right] + \varepsilon^2 \right]^{\frac{1}{2}} \quad \text{using Lemma 4.2} \\ &= \sqrt{2} \sqrt{\frac{\sigma^2}{n} + \varepsilon^2}. \end{aligned} \quad (4.53)$$

In addition, simple steps yield

$$\begin{aligned} M_{222}(\mathbf{X}) &= \frac{1}{(\sigma^2 - \varepsilon)^3} + \frac{9}{(\sigma^2 - \varepsilon)^4} \left(\hat{\sigma}^2 + (\bar{X} - \mu)^2 + \varepsilon^2 \right) \\ &= \frac{1}{(\sigma^2 - \varepsilon)^3} + \frac{9}{(\sigma^2 - \varepsilon)^4} \left(\hat{\sigma}^2 - \sigma^2 + (\bar{X} - \mu)^2 + \varepsilon^2 + \sigma^2 \right) \\ &\leq \frac{1}{(\sigma^2 - \varepsilon)^3} + \frac{9}{(\sigma^2 - \varepsilon)^4} \left(\left| \hat{\sigma}^2 - \sigma^2 \right| + (\bar{X} - \mu)^2 + \varepsilon^2 + \sigma^2 \right), \end{aligned}$$

which leads to

$$\begin{aligned} [M_{222}(\mathbf{X})]^2 &\leq 2 \left[\frac{1}{(\sigma^2 - \varepsilon)^6} + \frac{81}{(\sigma^2 - \varepsilon)^8} \left(\left| \hat{\sigma}^2 - \sigma^2 \right| + (\bar{X} - \mu)^2 + \varepsilon^2 + \sigma^2 \right)^2 \right] \\ &\leq 2 \left[\frac{1}{(\sigma^2 - \varepsilon)^6} + \frac{162}{(\sigma^2 - \varepsilon)^8} \left[\left(\left| \hat{\sigma}^2 - \sigma^2 \right| + \varepsilon^2 + \sigma^2 \right)^2 + (\bar{X} - \mu)^4 \right] \right]. \end{aligned} \quad (4.54)$$

Using the result in (4.54) and Lemma 4.2 yields

$$\begin{aligned} &\left[\mathbb{E} \left((M_{222}(\mathbf{X}))^2 \middle| |C_{(m)}| < \varepsilon \right) \right]^{\frac{1}{2}} \\ &\leq \sqrt{2} \left[\frac{1}{(\sigma^2 - \varepsilon)^6} + \frac{162}{(\sigma^2 - \varepsilon)^8} \left[(\varepsilon + \varepsilon^2 + \sigma^2)^2 + 3\frac{\sigma^4}{n^2} \right] \right]^{\frac{1}{2}}. \end{aligned} \quad (4.55)$$

Using (4.52), (4.53) and (4.55), the second term in (4.7) multiplied by $\frac{\|h\|_1}{\sqrt{n}}$ becomes

$$\begin{aligned}
& \frac{\|h\|_1}{2\sqrt{n}} \left\{ \sum_{k=1}^d \sum_{l=1}^d \left| [I(\theta_0)]^{-\frac{1}{2}} \right|_{lk} \left| \sum_{j=1}^d \sum_{i=1}^d \left[E \left((\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 (\hat{\theta}_n(\mathbf{X})_i - \theta_{0i})^2 \right) \right]^{\frac{1}{2}} \right. \right. \\
& \quad \left. \left. \times \left[E \left((nM_{kji}(\mathbf{X}))^2 \mid |C_{(m)}| < \varepsilon \right) \right]^{\frac{1}{2}} \right\} \\
& \leq \frac{\|h\|_1 \sigma^4}{\sqrt{n}(\sigma^2 - \varepsilon)^2} \left[\sqrt{2} + \frac{4\sqrt{2}\sigma}{\sigma^2 - \varepsilon} \sqrt{\frac{\sigma^2}{n} + \varepsilon^2} \right] \\
& \quad + \frac{\|h\|_1 \sigma^4}{\sqrt{2n}(\sigma^2 - \varepsilon)^2} \left[\sqrt{3} + \frac{8\sigma}{\sigma^2 - \varepsilon} \sqrt{\frac{\sigma^2}{n} + \varepsilon^2} \right] \\
& \quad + \frac{4\|h\|_1 \sigma^6}{\sqrt{n}(\sigma^2 - \varepsilon)^3} \left[1 + \frac{162}{(\sigma^2 - \varepsilon)^2} \left[(\varepsilon + \varepsilon^2 + \sigma^2)^2 + 3\frac{\sigma^4}{n^2} \right] \right]^{\frac{1}{2}} \\
& = \frac{\|h\|_1 \sigma^4}{\sqrt{n}(\sigma^2 - \varepsilon)^2} \left[\sqrt{2} + \sqrt{\frac{3}{2}} + \frac{4\sigma}{\sigma^2 - \varepsilon} \sqrt{\frac{\sigma^2}{n} + \varepsilon^2} \left(\sqrt{2} + \frac{2}{\sqrt{2}} \right) \right] \\
& \quad + \frac{4\|h\|_1 \sigma^6}{\sqrt{n}(\sigma^2 - \varepsilon)^3} \left[1 + \frac{162}{(\sigma^2 - \varepsilon)^2} \left[(\varepsilon + \varepsilon^2 + \sigma^2)^2 + 3\frac{\sigma^4}{n^2} \right] \right]^{\frac{1}{2}} \\
& \leq \frac{\|h\|_1 \sigma^4}{\sqrt{n}(\sigma^2 - \varepsilon)^2} \left[3 + \frac{8\sqrt{2}\sigma}{\sigma^2 - \varepsilon} \sqrt{\frac{\sigma^2}{n} + \varepsilon^2} \right] \\
& \quad + \frac{4\|h\|_1 \sigma^6}{\sqrt{n}(\sigma^2 - \varepsilon)^3} \left[1 + \frac{162}{(\sigma^2 - \varepsilon)^2} \left[(\varepsilon + \varepsilon^2 + \sigma^2)^2 + 3\frac{\sigma^4}{n^2} \right] \right]^{\frac{1}{2}}. \tag{4.56}
\end{aligned}$$

Taking $\varepsilon = \frac{\sigma^2}{2}$ yields

$$(4.56) = \frac{4\|h\|_1}{\sqrt{n}} \left[3 + 16\sqrt{2} \sqrt{\frac{1}{n} + \frac{\sigma^2}{4}} \right] + \frac{32\|h\|_1}{\sqrt{n}} \left[1 + 648 \left[\left(\frac{3}{2} + \frac{\sigma^2}{4} \right)^2 + \frac{3}{n^2} \right] \right]^{\frac{1}{2}}. \tag{4.57}$$

Combining the results of (4.47) and (4.57) yields to a bound on the first term, $\frac{\|h\|_1}{\sqrt{n}} K_1(\theta_0)$, of the general upper bound in (4.6). The second term, $\frac{\|h\|_2}{\sqrt{n}} K_2(\theta_0)$, is a sum of two quantities as (4.8) shows. For the first quantity,

$$\begin{aligned}
& \frac{\|h\|_2}{4\sqrt{n}} \sum_{j=1}^d \left[\text{Var} \left(\left(\sum_{k=1}^d [I(\theta_0)]^{-\frac{1}{2}} \right]_{jk} \frac{\partial}{\partial \theta_k} \log f(X_1 | \theta_0) \right)^2 \right]^{\frac{1}{2}} \\
&= \frac{\|h\|_2}{4\sqrt{n}} \left\{ \left[\text{Var} \left(\frac{(X_1 - \mu)^2}{\sigma^2} \right) \right]^{\frac{1}{2}} + \left[\text{Var} \left(\frac{(X_1 - \mu)^4}{2\sigma^4} - \frac{(X_1 - \mu)^2}{\sigma^2} \right) \right]^{\frac{1}{2}} \right\} \\
&= \frac{\|h\|_2}{4\sqrt{n}} \left\{ \sqrt{2} + \left[\frac{1}{4} \text{Var} \left[\left(\frac{X_1 - \mu}{\sigma} \right)^4 \right] + \text{Var} \left[\left(\frac{X_1 - \mu}{\sigma} \right)^2 \right] \right. \right. \\
&\quad \left. \left. - \text{Cov} \left(\left(\frac{X_1 - \mu}{\sigma} \right)^4, \left(\frac{X_1 - \mu}{\sigma} \right)^2 \right) \right]^{\frac{1}{2}} \right\}. \tag{4.58}
\end{aligned}$$

Since $\frac{X_1 - \mu}{\sigma} \sim N(0, 1)$, simple steps yield

$$\begin{aligned}
\text{Var} \left[\left(\frac{X_1 - \mu}{\sigma} \right)^4 \right] &= \text{E} \left[\left(\frac{X_1 - \mu}{\sigma} \right)^8 \right] - \left[\text{E} \left[\left(\frac{X_1 - \mu}{\sigma} \right)^4 \right] \right]^2 = 96 \\
\text{Var} \left[\left(\frac{X_1 - \mu}{\sigma} \right)^2 \right] &= \text{E} \left[\left(\frac{X_1 - \mu}{\sigma} \right)^4 \right] - \left[\text{E} \left[\left(\frac{X_1 - \mu}{\sigma} \right)^2 \right] \right]^2 = 2 \\
\text{Cov} \left(\left(\frac{X_1 - \mu}{\sigma} \right)^4, \left(\frac{X_1 - \mu}{\sigma} \right)^2 \right) &= \text{E} \left[\left(\frac{X_1 - \mu}{\sigma} \right)^6 \right] \\
&\quad - \text{E} \left[\left(\frac{X_1 - \mu}{\sigma} \right)^4 \right] \text{E} \left[\left(\frac{X_1 - \mu}{\sigma} \right)^2 \right] = 12.
\end{aligned}$$

Applying the above results to (4.58),

$$\frac{\|h\|_2}{4\sqrt{n}} \sum_{j=1}^d \left[\text{Var} \left(\left(\sum_{k=1}^d [I(\theta_0)]^{-\frac{1}{2}} \right]_{jk} \frac{\partial}{\partial \theta_k} \log f(X_1 | \theta_0) \right)^2 \right]^{\frac{1}{2}} = \frac{\|h\|_2}{4\sqrt{n}} (\sqrt{2} + \sqrt{14}). \tag{4.59}$$

The second quantity in $\frac{\|h\|_2}{\sqrt{n}} K_2(\theta_0)$ is also easily calculated. Since $[I(\theta_0)]^{-\frac{1}{2}}$ is diagonal, we obtain that

$$\begin{aligned}
& \frac{\|h\|_2}{2\sqrt{n}} \sum_{j>k}^d \sum_{k=1}^{d-1} \left[\text{Var} \left(\sum_{q=1}^d \sum_{v=1}^d \left[[I(\theta_0)]^{-\frac{1}{2}} \right]_{jq} \frac{\partial}{\partial \theta_q} \log f(X_1 | \theta_0) \right. \right. \\
& \quad \left. \left. \times \left[[I(\theta_0)]^{-\frac{1}{2}} \right]_{kv} \frac{\partial}{\partial \theta_v} \log f(X_1 | \theta_0) \right) \right]^{\frac{1}{2}} \\
&= \frac{\|h\|_2}{2\sqrt{n}} \left[\text{Var} \left(\left[[I(\theta_0)]^{-\frac{1}{2}} \right]_{11} \frac{\partial}{\partial \mu} \log f(X_1 | \theta_0) \left[[I(\theta_0)]^{-\frac{1}{2}} \right]_{22} \frac{\partial}{\partial \sigma^2} \log f(X_1 | \theta_0) \right) \right]^{\frac{1}{2}}.
\end{aligned} \tag{4.60}$$

Using (4.45) for one random variable and (4.46) yields

$$\begin{aligned}
(4.60) &= \frac{\|h\|_2}{2\sqrt{n}} \left[\text{Var} \left(\left(\frac{X_1 - \mu}{\sigma} \right) \left(\frac{(X_1 - \mu)^2}{\sqrt{2}\sigma^2} - \frac{1}{\sqrt{2}} \right) \right) \right]^{\frac{1}{2}} \\
&= \frac{\|h\|_2}{2\sqrt{2n}} \left[\text{Var} \left(\left(\frac{X_1 - \mu}{\sigma} \right)^3 - \left(\frac{X_1 - \mu}{\sigma} \right) \right) \right]^{\frac{1}{2}} \\
&= \frac{\|h\|_2}{2\sqrt{2n}} \left[\text{Var} \left(\left(\frac{X_1 - \mu}{\sigma} \right)^3 \right) + \text{Var} \left(\frac{X_1 - \mu}{\sigma} \right) - 2\text{Cov} \left(\left(\frac{X_1 - \mu}{\sigma} \right)^3, \frac{X_1 - \mu}{\sigma} \right) \right]^{\frac{1}{2}} \\
&= \frac{\|h\|_2}{2\sqrt{2n}} \left[\mathbb{E} \left[\left(\frac{X_1 - \mu}{\sigma} \right)^6 \right] + 1 - 2\mathbb{E} \left[\left(\frac{X_1 - \mu}{\sigma} \right)^4 \right] \right]^{\frac{1}{2}} \\
&= \frac{\|h\|_2}{2\sqrt{2n}} \sqrt{10} = \frac{\sqrt{5}\|h\|_2}{2\sqrt{n}}.
\end{aligned} \tag{4.61}$$

For an upper bound for the third term $\frac{\|h\|_3}{\sqrt{n}} K_3(\theta_0)$ in (4.6), with $K_3(\theta_0)$ as in (4.9), we use that X'_1 is an independent copy of X_1 . The triangle inequality and (4.42) give that

$$\begin{aligned}
& \frac{\|h\|_3}{12\sqrt{n}} \mathbb{E} \left(\sum_{i=1}^d \left| \sum_{l=1}^d \left[[I(\theta_0)]^{-\frac{1}{2}} \right]_{il} \left(\frac{\partial}{\partial \theta_l} \log f(X'_1 | \theta_0) - \frac{\partial}{\partial \theta_l} \log f(X_1 | \theta_0) \right) \right| \right)^3 \\
&= \frac{\|h\|_3}{12\sqrt{n}} \mathbb{E} \left(\left| \frac{X'_1 - \mu}{\sigma} - \frac{X_1 - \mu}{\sigma} \right| + \frac{1}{\sqrt{2}\sigma^2} |(X'_1 - \mu)^2 - (X_1 - \mu)^2| \right)^3 \\
&\leq \frac{\|h\|_3}{3\sqrt{n}} \left[\mathbb{E} \left| \frac{X'_1 - X_1}{\sigma} \right|^3 + \mathbb{E} \left| \frac{(X'_1 - \mu)^2 - (X_1 - \mu)^2}{\sqrt{2}\sigma^2} \right|^3 \right] \\
&\leq \frac{8\|h\|_3}{3\sqrt{n}} \left[\frac{\mathbb{E}|X_1|^3}{\sigma^3} + \frac{\mathbb{E}(X_1 - \mu)^6}{2\sqrt{2}\sigma^6} \right] = \frac{8\|h\|_3}{3\sqrt{n}} \left[\frac{2\sqrt{2}}{\sqrt{\pi}} + \frac{15}{2\sqrt{2}} \right] \leq \frac{19}{\sqrt{n}} \|h\|_3.
\end{aligned} \tag{4.62}$$

For the last term of (4.6), $E(\bar{X} - \mu)^2 = \frac{\sigma^2}{n}$ and $E\left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 - \sigma^2\right)^2 = \frac{\sigma^4}{n} \left(2 - \frac{1}{n}\right)$. For the choice of $\varepsilon = \varepsilon_0$ as in (R.C.3), (4.48), (4.49) and (4.50) require that $0 < \varepsilon < \sigma^2$. Again, due to a trade-off between the first and the last term of the bound in (4.6), we choose $\varepsilon = \frac{\sigma^2}{2}$ which gives that

$$\begin{aligned} \frac{2\|h\|}{\varepsilon^2} \left(\sum_{j=1}^d E \left[\left(\hat{\theta}_n(X)_j - \theta_{0j} \right)^2 \right] \right) &= \frac{8\|h\|}{\sigma^4} \left(\frac{\sigma^2}{n} + \frac{\sigma^4}{n} \left(2 - \frac{1}{n} \right) \right) \\ &= \frac{8\|h\|}{n\sigma^2} \left(1 + \sigma^2 \left(2 - \frac{1}{n} \right) \right) \\ &\leq \frac{8\|h\|}{n\sigma^2} (1 + 2\sigma^2). \end{aligned} \quad (4.63)$$

Using the results in (4.47), (4.57), (4.59), (4.61), (4.62) and (4.63) we get the result of the Corollary. \square

Remark 4.3. (1) The rate of convergence of the upper bound in (4.44) is $\frac{1}{\sqrt{n}}$.

(2) There might be cases where the parameters depend on the sample size, so that $\mu = \mu(n)$ and $\sigma^2 = \sigma^2(n)$. The bound in (4.44) does not depend on $\mu(n)$ and goes to zero as long as

$$\begin{aligned} (i) \quad &\frac{1}{n\sigma^2(n)} \xrightarrow{n \rightarrow \infty} 0, \\ (ii) \quad &\frac{\sigma^2(n)}{\sqrt{n}} \xrightarrow{n \rightarrow \infty} 0, \end{aligned}$$

are both satisfied. From (i), the order of $\sigma^2(n)$ should not be less than or equal to $\frac{1}{n}$, while from (ii) we see that $\sigma^2(n)$ should be of order smaller than \sqrt{n} . For instance, $\sigma^2(n) = cn^{\frac{1}{4}}$, where $c \in \mathbb{R}$ is a constant, satisfies the above limits. The bound in (4.44) is then of order $\frac{1}{n^{\frac{1}{4}}}$.

4.2.2 Bounds when the MLE is not known explicitly

In the single-parameter case, it was possible to find an upper bound for the mean squared error and then use it to get upper bounds on the distributional distance of interest which can be applied when the MLE is not expressed in a closed-form. In this subsection, we give similar bounds for the multi-parameter case with multivariate (dimensionality $t \in \mathbb{N}$) i.i.d. random vectors $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$, with $\mathbf{X}_i = (X_{i1}, X_{i2}, \dots, X_{it})^\top$. We make some extra assumptions,

(Con.1) $\forall j \in \{1, 2, \dots, t\}$, the support S_j of X_{ij} is a bounded interval in \mathbb{R} ; let $s_j :=$

$$\sup_{x_{ij} \in S_j} \{|x_{ij}|\} \text{ and } s := \max \{s_1, s_2, \dots, s_t\};$$

(Con.2) for all $\boldsymbol{\theta}_0 \in \Theta$, where Θ is the open parameter space, there exists an $\varepsilon_0 =$

$$\varepsilon_0(\boldsymbol{\theta}_0) > 0 \text{ such that for all } \boldsymbol{\theta} \in \Theta \text{ with } |\theta_j - \theta_{0j}| < \varepsilon_0, \forall j = 1, 2, \dots, d$$

$$\sup_{\substack{\boldsymbol{\theta}: |\theta_q - \theta_{0q}| < \varepsilon_0 \\ \forall q \in \{1, 2, \dots, d\}}} \left| \frac{\partial^3}{\partial \theta_k \partial \theta_j \partial \theta_i} \log f(\mathbf{x}_1 | \boldsymbol{\theta}) \right| \leq M_{kji},$$

where $M_{kji} = M_{kji}(\boldsymbol{\theta}_0)$ is a constant that may depend only on the unknown parameter $\boldsymbol{\theta}_0$;

(Con.3) for $M = \sup_{i,j} \left\{ \left| [I(\boldsymbol{\theta}_0)]^{-1} \right|_{ij} \right\}$, the sample size satisfies

$$n > \frac{s^2 d^2}{4 \varepsilon_0^2} \left(M \varepsilon_0 \sum_{l=1}^d \sum_{k=1}^d \left| [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right|_{lk} \sum_{m=1}^d \sum_{i=1}^d M_{kim} \right. \\ \left. + \left[M^2 \varepsilon_0^2 \left(\sum_{l=1}^d \sum_{k=1}^d \left| [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right|_{lk} \sum_{m=1}^d \sum_{i=1}^d M_{kim} \right)^2 + 8M \right]^{\frac{1}{2}} \right)^2.$$

If (Con.1)-(Con.3) hold, then

$$2d^2 s^2 M + dsM \sqrt{n} \varepsilon_0^2 \sum_{l=1}^d \sum_{k=1}^d \left| [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right|_{lk} \sum_{m=1}^d \sum_{i=1}^d M_{kim} - n \varepsilon_0^2 < 0$$

holds with ε_0 as in (Con.2). The upper bound in (4.6) can be split into terms coming from Lemma 4.1, and terms due to Taylor expansions and conditional expectations. For ease of presentation, since the terms coming from Lemma 4.1 are not related to $\hat{\theta}_n(\mathbf{X})$, from now on, unless otherwise stated, we abbreviate

$$\begin{aligned}
 D := D(\theta_0, h, \mathbf{X}) &:= \frac{\|h\|_2}{4\sqrt{n}} \sum_{j=1}^d \left[\text{Var} \left(\left(\sum_{k=1}^d [I(\theta_0)]^{-\frac{1}{2}} \right]_{jk} \frac{\partial}{\partial \theta_k} \log f(\mathbf{X}_1 | \theta_0) \right)^2 \right]^{\frac{1}{2}} \\
 &+ \frac{\|h\|_2}{2\sqrt{n}} \sum_{k=1}^{d-1} \sum_{j>k}^d \left[\text{Var} \left(\sum_{q=1}^d \sum_{v=1}^d [I(\theta_0)]^{-\frac{1}{2}} \right]_{jq} \frac{\partial}{\partial \theta_q} \log f(\mathbf{X}_1 | \theta_0) \right. \\
 &\quad \left. \times [I(\theta_0)]^{-\frac{1}{2}} \right]_{kv} \frac{\partial}{\partial \theta_v} \log f(\mathbf{X}_1 | \theta_0) \right]^{\frac{1}{2}} \\
 &+ \frac{\|h\|_3}{12\sqrt{n}} \mathbb{E} \left(\sum_{i=1}^d \left| \sum_{l=1}^d [I(\theta_0)]^{-\frac{1}{2}} \right]_{il} \left(\frac{\partial}{\partial \theta_l} \log f(\mathbf{X}'_1 | \theta_0) - \frac{\partial}{\partial \theta_l} \log f(\mathbf{X}_1 | \theta_0) \right) \right|^3.
 \end{aligned} \tag{4.64}$$

In order to give an upper bound when $\hat{\theta}_n(\mathbf{X})$ is not known explicitly, we bound $\sqrt{\mathbb{E} \left[\sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right]}$ by a quantity which does not require knowledge of the MLE. The result is given in the following theorem.

Theorem 4.3. *Let $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ be i.i.d. \mathbb{R}^t -valued random elements, for $t \in \mathbb{N}$, with probability density (or mass) function $f(\mathbf{x}_i | \theta)$, where θ is the d -valued vector parameter. Assume that (R.C.1)-(R.C.4) and also (Con.1)-(Con.3) are satisfied. Also $\mathbf{Z} \sim \mathcal{N}_d(\mathbf{0}, I_{d \times d})$ and we assume existence and uniqueness of the MLE, $\hat{\theta}_n(\mathbf{X})$. For $\varepsilon = \varepsilon_0$ as in (Con.2) and using the notation*

$$\begin{aligned}
\gamma &= \sum_{j=1}^d \left| \left[[I(\boldsymbol{\theta}_0)]^{-1} \right]_{jj} \right| + \frac{M}{2\sqrt{n}} \sum_{j=1}^d \sqrt{\text{Var} \left(\left(\sum_{k=1}^d \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{jk} \frac{\partial}{\partial \boldsymbol{\theta}_k} \log f(\mathbf{X}_1 | \boldsymbol{\theta}_0) \right)^2 \right)} \\
&+ \frac{M}{\sqrt{n}} \sum_{k=1}^{d-1} \sum_{j>k}^d \left[\text{Var} \left(\sum_{q=1}^d \sum_{v=1}^d \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{jq} \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{kv} \right. \right. \\
&\quad \left. \left. \times \frac{\partial}{\partial \boldsymbol{\theta}_q} \log f(\mathbf{X}_1 | \boldsymbol{\theta}_0) \frac{\partial}{\partial \boldsymbol{\theta}_v} \log f(\mathbf{X}_1 | \boldsymbol{\theta}_0) \right) \right]^{\frac{1}{2}} \\
\omega &= 1 - 2 \frac{d^2 s^2 M}{n \varepsilon^2} - \frac{dsM}{\sqrt{n}} \sum_{l=1}^d \sum_{k=1}^d \left| \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{lk} \right| \sum_{m=1}^d \sum_{i=1}^d M_{kmi} \\
v &= 2d^{\frac{3}{2}} sM \sum_{l=1}^d \sum_{k=1}^d \left| \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{lk} \right| \sqrt{\sum_{i=1}^d \text{Var} \left(\frac{\partial^2}{\partial \boldsymbol{\theta}_k \partial \boldsymbol{\theta}_i} \log f(\mathbf{X}_1 | \boldsymbol{\theta}_0) \right)} \quad (4.65)
\end{aligned}$$

it holds that

$$\sqrt{\mathbb{E} \left[\sum_{j=1}^d (\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \boldsymbol{\theta}_{0j})^2 \right]} \leq \frac{1}{\sqrt{n}} \left(\frac{\frac{v}{\sqrt{n}} + \sqrt{\frac{v^2}{n} + 4\omega\gamma}}{2\omega} \right). \quad (4.66)$$

Remark 4.4. The quantities γ , ω and v are $\mathcal{O}(1)$.

Proof. Using (4.6) and (4.64) gives that for $C_{(m)}$ as in (4.5),

$$\begin{aligned}
&\left| \mathbb{E} \left[h \left(\sqrt{n} [I(\boldsymbol{\theta}_0)]^{\frac{1}{2}} (\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0) \right) \right] - \mathbb{E}[h(\mathbf{Z})] \right| \leq D + \frac{2\|h\|}{\varepsilon^2} \mathbb{E} \left[\sum_{j=1}^d (\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \boldsymbol{\theta}_{0j})^2 \right] \\
&+ \frac{\|h\|_1}{\sqrt{n}} \sum_{k=1}^d \sum_{l=1}^d \left| \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{lk} \right| \\
&\quad \times \sum_{j=1}^d \left[\mathbb{E} \left[(\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \boldsymbol{\theta}_{0j})^2 \right] \mathbb{E} \left[\left(\frac{\partial^2}{\partial \boldsymbol{\theta}_j \partial \boldsymbol{\theta}_k} l(\boldsymbol{\theta}_0; \mathbf{X}) + n [I(\boldsymbol{\theta}_0)]_{kj} \right)^2 \right] \right]^{\frac{1}{2}} \quad (4.67) \\
&+ \frac{\|h\|_1}{2\sqrt{n}} \left\{ \sum_{k=1}^d \sum_{l=1}^d \left| \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{lk} \right| \right. \\
&\quad \times \mathbb{E} \left(\left| \sum_{j=1}^d \sum_{i=1}^d (\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \boldsymbol{\theta}_{0j}) (\hat{\boldsymbol{\theta}}_n(\mathbf{X})_i - \boldsymbol{\theta}_{0i}) \frac{\partial^3}{\partial \boldsymbol{\theta}_k \partial \boldsymbol{\theta}_j \partial \boldsymbol{\theta}_i} l(\boldsymbol{\theta}_0^*; \mathbf{X}) \right| \middle| |C_{(m)}| < \varepsilon \right) \left. \right\}. \quad (4.68)
\end{aligned}$$

Step 1: Upper bound for (4.67). Using that $E\left(\frac{\partial^2}{\partial\theta_j\partial\theta_k}l(\theta_0; \mathbf{X}) + n[I(\theta_0)]_{kj}\right) = 0$, $\forall j, k \in \{1, 2, \dots, d\}$,

$$\begin{aligned}
 (4.67) &= \|h\|_1 \sum_{k=1}^d \sum_{l=1}^d \left| \left[[I(\theta_0)]^{-\frac{1}{2}} \right]_{lk} \right| \\
 &\quad \times \sum_{j=1}^d \sqrt{E\left[(\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2\right]} \sqrt{\text{Var}\left(\frac{\partial^2}{\partial\theta_k\partial\theta_j} \log f(\mathbf{X}_1|\theta_0)\right)} \\
 &\leq \|h\|_1 \sum_{k=1}^d \sum_{l=1}^d \left| \left[[I(\theta_0)]^{-\frac{1}{2}} \right]_{lk} \right| \\
 &\quad \times \sum_{j=1}^d \sqrt{E\left[(\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2\right]} \sqrt{\sum_{i=1}^d \text{Var}\left(\frac{\partial^2}{\partial\theta_k\partial\theta_i} \log f(\mathbf{X}_1|\theta_0)\right)}, \quad (4.69)
 \end{aligned}$$

where for the last inequality the trivial bound

$$\text{Var}\left(\frac{\partial^2}{\partial\theta_k\partial\theta_j} \log f(\mathbf{X}_1|\theta_0)\right) \leq \sum_{i=1}^d \text{Var}\left(\frac{\partial^2}{\partial\theta_k\partial\theta_i} \log f(\mathbf{X}_1|\theta_0)\right)$$

was used, since the variance of a random variable is always non-negative. Now, using that $\left(\sum_{j=1}^d \alpha_j\right)^2 \leq d \left(\sum_{j=1}^d \alpha_j^2\right)$ for $\alpha_j \in \mathbb{R}$, yields

$$\left(\sum_{j=1}^d \sqrt{E\left[(\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2\right]}\right)^2 \leq d \sum_{j=1}^d E\left[(\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2\right].$$

Taking square roots in both sides of the above inequality and applying this result to (4.69) gives

$$\begin{aligned}
 (4.67) &\leq \|h\|_1 \sqrt{d} \sum_{k=1}^d \sum_{l=1}^d \left| \left[[I(\theta_0)]^{-\frac{1}{2}} \right]_{lk} \right| \\
 &\quad \times \sqrt{\sum_{i=1}^d \text{Var}\left(\frac{\partial^2}{\partial\theta_k\partial\theta_i} \log f(\mathbf{X}_1|\theta_0)\right)} \sqrt{E\left[\sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2\right]}. \quad (4.70)
 \end{aligned}$$

Step 2: Upper bound for (4.68). Notice that from (Con.2), if $|\hat{\theta}_n(\mathbf{X})_j - \theta_{0j}| < \varepsilon$, $\forall j \in \{1, 2, \dots, d\}$, then $\left|\frac{\partial^3}{\partial\theta_k\partial\theta_j\partial\theta_i}l(\theta_0^*; \mathbf{x})\right| = \left|\sum_{l=1}^n \frac{\partial^3}{\partial\theta_k\partial\theta_j\partial\theta_i} \log f(\mathbf{x}_l|\theta_0^*)\right| \leq nM_{kji}$. In

addition,

$$\begin{aligned} & \sum_{j=1}^d \sum_{i=1}^d |(\hat{\theta}_n(\mathbf{X})_j - \theta_{0j}) (\hat{\theta}_n(\mathbf{X})_i - \theta_{0i})| M_{kji} \\ &= \sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 M_{kjj} + 2 \sum_{j>i}^d \sum_{i=1}^{d-1} |\hat{\theta}_n(\mathbf{X})_j - \theta_{0j}| |\hat{\theta}_n(\mathbf{X})_i - \theta_{0i}| M_{kij}. \end{aligned}$$

Using now that $2\alpha\beta \leq \alpha^2 + \beta^2, \forall \alpha, \beta \in \mathbb{R}$,

$$\begin{aligned} & \sum_{j=1}^d \sum_{i=1}^d |(\hat{\theta}_n(\mathbf{X})_j - \theta_{0j}) (\hat{\theta}_n(\mathbf{X})_i - \theta_{0i})| M_{kji} \\ & \leq \sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 M_{kjj} + \sum_{j>i}^d \sum_{i=1}^{d-1} [(\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 + (\hat{\theta}_n(\mathbf{X})_i - \theta_{0i})^2] M_{kji} \\ &= \sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \sum_{i=1}^d M_{kji} \\ & \leq \sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \sum_{m=1}^d \sum_{i=1}^d M_{kmi}. \end{aligned} \tag{4.71}$$

Using (4.71) yields

$$\begin{aligned} (4.68) & \leq \frac{\|h\|_1}{2\sqrt{n}} \sum_{k=1}^d \sum_{l=1}^d \left| [I(\theta_0)]^{-\frac{1}{2}} \right|_{lk} \\ & \quad \times \mathbb{E} \left(n \sum_{j=1}^d \sum_{i=1}^d M_{kij} |\hat{\theta}_n(\mathbf{X})_j - \theta_{0j}| |\hat{\theta}_n(\mathbf{X})_i - \theta_{0i}| \middle| |C_{(m)}| < \varepsilon \right) \\ & \leq \frac{\|h\|_1 \sqrt{n}}{2} \sum_{k=1}^d \sum_{l=1}^d \left| [I(\theta_0)]^{-\frac{1}{2}} \right|_{lk} \sum_{m=1}^d \sum_{i=1}^d M_{kmi} \mathbb{E} \left[\sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right]. \end{aligned} \tag{4.72}$$

Hence, from (4.70) and (4.72),

$$\begin{aligned}
& \left| \mathbb{E} \left[h \left(\sqrt{n} [I(\theta_0)]^{\frac{1}{2}} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) \right] - \mathbb{E}[h(\mathbf{Z})] \right| \leq D + \frac{2\|h\|}{\varepsilon^2} \mathbb{E} \left[\sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right] \\
& + \|h\|_1 \sqrt{d} \sum_{k=1}^d \sum_{l=1}^d \left| \left[[I(\theta_0)]^{-\frac{1}{2}} \right]_{lk} \right| \\
& \quad \times \sqrt{\sum_{i=1}^d \text{Var} \left(\frac{\partial^2}{\partial \theta_k \partial \theta_i} \log f(\mathbf{X}_1 | \theta_0) \right)} \sqrt{\mathbb{E} \left[\sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right]} \\
& + \frac{\|h\|_1 \sqrt{n}}{2} \sum_{k=1}^d \sum_{l=1}^d \left| \left[[I(\theta_0)]^{-\frac{1}{2}} \right]_{lk} \right| \sum_{m=1}^d \sum_{i=1}^d M_{kmi} \mathbb{E} \left[\sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right]. \tag{4.73}
\end{aligned}$$

Since D as defined in (4.64), is not related to the MLE, the upper bound in (4.73) depends on $\hat{\theta}_n(\mathbf{X})$ only through $\mathbb{E} \left[\sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right]$. Our purpose is to find a bound for $\mathbb{E} \left[\sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right]$ that does not depend on $\hat{\theta}_n(\mathbf{X})$.

Step 3: The MSE test function. To this purpose define the test function

$$h = h_{\theta_0} : \mathbb{R}^d \rightarrow \mathbb{R} : h(\mathbf{x}) = \mathbf{x}^\top [I(\theta_0)]^{-1} \mathbf{x}. \tag{4.74}$$

Then,

$$\begin{aligned}
& h \left(\sqrt{n} [I(\theta_0)]^{\frac{1}{2}} (\hat{\theta}_n(\mathbf{x}) - \theta_0) \right) \\
& = \left[\sqrt{n} [I(\theta_0)]^{\frac{1}{2}} (\hat{\theta}_n(\mathbf{x}) - \theta_0) \right]^\top [I(\theta_0)]^{-1} \sqrt{n} [I(\theta_0)]^{\frac{1}{2}} (\hat{\theta}_n(\mathbf{x}) - \theta_0) \\
& = n (\hat{\theta}_n(\mathbf{X}) - \theta_0)^\top (\hat{\theta}_n(\mathbf{X}) - \theta_0),
\end{aligned}$$

since $\left([I(\theta_0)]^{\frac{1}{2}} \right)^\top = [I(\theta_0)]^{\frac{1}{2}}$ as $[I(\theta_0)]^{\frac{1}{2}}$ is symmetric. Thus,

$$h \left(\sqrt{n} [I(\theta_0)]^{\frac{1}{2}} (\hat{\theta}_n(\mathbf{x}) - \theta_0) \right) = n \sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2.$$

Hence,

$$\begin{aligned}
\mathbb{E}[h(\mathbf{Z})] &= \sum_{j=1}^d \left[[I(\boldsymbol{\theta}_0)]^{-1} \right]_{jj} \mathbb{E}[Z_j^2] + 2 \sum_{\substack{k=1 \\ k < i}}^d \sum_{i=1}^d \left[[I(\boldsymbol{\theta}_0)]^{-1} \right]_{ki} \mathbb{E}[Z_k] \mathbb{E}[Z_i] \\
&= \sum_{j=1}^d \left[[I(\boldsymbol{\theta}_0)]^{-1} \right]_{jj}.
\end{aligned}$$

Denoting with B_h the bound in (4.73) for the test function h as in (4.74),

$$\begin{aligned}
&\mathbb{E} \left[\sum_{j=1}^d (\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \boldsymbol{\theta}_{0j})^2 \right] \\
&= \frac{1}{n} \left| \mathbb{E} \left[h \left(\sqrt{n} [I(\boldsymbol{\theta}_0)]^{\frac{1}{2}} (\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0) \right) \right] - \mathbb{E}[h(\mathbf{Z})] + \mathbb{E}[h(\mathbf{Z})] \right| \\
&\leq \frac{1}{n} \left(B_h + \sum_{j=1}^d \left| \left[[I(\boldsymbol{\theta}_0)]^{-1} \right]_{jj} \right| \right). \tag{4.75}
\end{aligned}$$

For the calculation of B_h , note that with s as in (Con.1) and M as in (Con.3),

$$\begin{aligned}
\|h\| &= d^2 s^2 M, & \|h\|_1 &= \sup_j \left| \frac{\partial}{\partial x_j} h(\mathbf{x}) \right| = 2dMs, \\
\|h\|_2 &= \sup_{i,j} \left| \frac{\partial^2}{\partial x_i \partial x_j} h(\mathbf{x}) \right| = 2M, & \|h\|_3 &= \sup_{i,j,k} \left| \frac{\partial^3}{\partial x_i \partial x_j \partial x_k} h(\mathbf{x}) \right| = 0.
\end{aligned}$$

Using (4.64) and (4.73) yields

$$\begin{aligned}
B_h &= \frac{M}{2\sqrt{n}} \sum_{j=1}^d \left[\text{Var} \left(\left(\sum_{k=1}^d \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{jk} \frac{\partial}{\partial \boldsymbol{\theta}_k} \log f(\mathbf{X}_1 | \boldsymbol{\theta}_0) \right)^2 \right) \right]^{\frac{1}{2}} \\
&+ \frac{M}{\sqrt{n}} \sum_{k=1}^{d-1} \sum_{j>k}^d \left[\text{Var} \left(\sum_{q=1}^d \sum_{v=1}^d \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{jq} \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{kv} \right. \right. \\
&\quad \left. \left. \times \frac{\partial}{\partial \boldsymbol{\theta}_q} \log f(\mathbf{X}_1 | \boldsymbol{\theta}_0) \frac{\partial}{\partial \boldsymbol{\theta}_v} \log f(\mathbf{X}_1 | \boldsymbol{\theta}_0) \right) \right]^{\frac{1}{2}} \\
&+ 2 \frac{d^2 s^2 M}{\varepsilon^2} \mathbb{E} \left(\sum_{j=1}^d (\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \boldsymbol{\theta}_{0j})^2 \right) \\
&+ \sqrt{nd} s M \sum_{l=1}^d \sum_{k=1}^d \left| \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{lk} \right| \sum_{m=1}^d \sum_{i=1}^d M_{kmi} \mathbb{E} \left(\sum_{j=1}^d (\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \boldsymbol{\theta}_{0j})^2 \right)
\end{aligned}$$

$$\begin{aligned}
& + 2d^{\frac{3}{2}}sM \sum_{l=1}^d \sum_{k=1}^d \left| \left[I(\boldsymbol{\theta}_0)^{-\frac{1}{2}} \right]_{lk} \right| \\
& \quad \times \sqrt{\sum_{i=1}^d \text{Var} \left(\frac{\partial^2}{\partial \theta_k \partial \theta_i} \log f(\mathbf{X}_1 | \boldsymbol{\theta}_0) \right)} \sqrt{\mathbb{E} \left(\sum_{j=1}^d (\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \theta_{0j})^2 \right)}.
\end{aligned} \tag{4.76}$$

Let $U := \sqrt{\mathbb{E} \left[\sum_{j=1}^d (\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \theta_{0j})^2 \right]}$. The results in (4.75) and (4.76) give

$$\begin{aligned}
0 \leq U^2 & \left(2 \frac{d^2 s^2 M}{n \varepsilon^2} + \frac{dsM}{\sqrt{n}} \sum_{l=1}^d \sum_{k=1}^d \left| \left[I(\boldsymbol{\theta}_0)^{-\frac{1}{2}} \right]_{lk} \right| \sum_{m=1}^d \sum_{i=1}^d M_{kmi} - 1 \right) \\
& + U \left(\frac{2}{n} d^{\frac{3}{2}} sM \sum_{l=1}^d \sum_{k=1}^d \left| \left[I(\boldsymbol{\theta}_0)^{-\frac{1}{2}} \right]_{lk} \right| \sqrt{\sum_{i=1}^d \text{Var} \left(\frac{\partial^2}{\partial \theta_k \partial \theta_i} \log f(\mathbf{X}_1 | \boldsymbol{\theta}_0) \right)} \right) + \frac{\gamma}{n},
\end{aligned} \tag{4.77}$$

with γ as in (4.65). Solving the quadratic inequality in (4.77) (with unknown U) and using (Con.3) related to the sample size, n , we obtain for v and ω as in (4.65) that

$$U \leq U_1 = \frac{1}{\sqrt{n}} \left(\frac{\frac{v}{\sqrt{n}} + \sqrt{\frac{v^2}{n} + 4\omega\gamma}}{2\omega} \right), \tag{4.78}$$

proving the result of the theorem. \square

Remark 4.5. As the bound (4.66) does not include $\hat{\boldsymbol{\theta}}_n(\mathbf{X})$, in cases where a closed-form expression for the vector MLE is not available, we can still get an upper bound on the distributional distance between the distribution of the MLE and the d -variate standard normal,

$$\begin{aligned}
& \left| \mathbb{E} \left[h \left(\sqrt{n} [I(\boldsymbol{\theta}_0)]^{\frac{1}{2}} (\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0) \right) \right] - \mathbb{E}[h(\mathbf{Z})] \right| \\
& \leq \frac{\|h\|_2}{4\sqrt{n}} \sum_{j=1}^d \left[\text{Var} \left(\left(\sum_{k=1}^d [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{jk} \frac{\partial}{\partial \theta_k} \log f(\mathbf{X}_1 | \boldsymbol{\theta}_0) \right)^2 \right]^{\frac{1}{2}} \\
& \quad + \frac{\|h\|_2}{2\sqrt{n}} \sum_{k < j} \sum_{j=1}^d \left[\text{Var} \left(\sum_{q=1}^d \sum_{v=1}^d [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{jq} [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{kv} \\
& \quad \quad \quad \times \frac{\partial}{\partial \theta_q} \log f(\mathbf{X}_1 | \boldsymbol{\theta}_0) \frac{\partial}{\partial \theta_v} \log f(\mathbf{X}_1 | \boldsymbol{\theta}_0) \right]^{\frac{1}{2}} \\
& \quad + \frac{\|h\|_3}{12\sqrt{n}} \mathbb{E} \left(\sum_{i=1}^d \left| \sum_{l=1}^d [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{il} \left(\frac{\partial}{\partial \theta_l} \log f(\mathbf{X}'_1 | \boldsymbol{\theta}_0) - \frac{\partial}{\partial \theta_l} \log f(\mathbf{X}_1 | \boldsymbol{\theta}_0) \right) \right|^3 \\
& \quad + 2 \frac{\|h\|}{\varepsilon^2} U_1^2 \\
& \quad + \|h\|_1 \sqrt{d} U_1 \sum_{k=1}^d \sum_{l=1}^d \left| [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{lk} \sqrt{\sum_{i=1}^d \text{Var} \left(\frac{\partial^2}{\partial \theta_k \partial \theta_i} \log f(\mathbf{X}_1 | \boldsymbol{\theta}_0) \right)} \\
& \quad + \frac{\|h\|_1 \sqrt{n}}{2} U_1^2 \sum_{k=1}^d \sum_{l=1}^d \left| [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{lk} \sum_{m=1}^d \sum_{i=1}^d M_{kmi}. \tag{4.79}
\end{aligned}$$

Example: The Beta distribution

Here, we find an upper bound for the specific example of i.i.d. random variables from the Beta distribution with both shape parameters being unknown. An analytic expression for the MLE is not available. Applying the result of (4.78) to bound $\mathbb{E} \left(\sum_{j=1}^2 (\hat{\boldsymbol{\theta}}_n(\mathbf{X})_j - \theta_{0j})^2 \right)$, gives an upper bound for the distributional distance of interest. Some useful notations are now presented. Firstly, as in Section 3.5, $\Psi_j(\cdot)$ is the j^{th} derivative of the digamma function Ψ , with $\Psi(z) = \frac{\Gamma'(z)}{\Gamma(z)}$, $z > 0$. From (3.67), $\Psi_j(z)$ can be defined through a sum. For $\alpha, \beta, x, y > 0$ and $0 < \varepsilon < \min\{x, y\}$, let

$$\begin{aligned}
\delta_I &:= \Psi_1(\alpha)\Psi_1(\beta) - \Psi_1(\alpha + \beta)(\Psi_1(\alpha) + \Psi_1(\beta)), \\
C_1(x, y) &:= \Psi_3(x) + \Psi_3(x + y) + 3[\Psi_1(x)]^2 + 3[\Psi_1(x + y)]^2, \\
C_2(x, y) &:= \Psi_1(x) - \Psi_1(x + y) + \sqrt{\delta_I}, \\
C_3(x, y) &:= C_1(x, y)[C_2(y, x)]^2, \\
C_4(x, y) &:= \frac{6x}{(y - \varepsilon)^4} + \frac{x\pi^4}{15} + \frac{6}{(x + y - \varepsilon)^3} + 7.26, \\
M_B &:= \frac{1}{\delta_I} \sup \{ \Psi_1(\alpha + \beta), \Psi_1(\min \{ \alpha, \beta \}) - \Psi_1(\alpha + \beta) \}, \\
\gamma_B &:= \frac{4M_B}{\sqrt{n}\delta_I [C_2(\alpha, \beta) + C_2(\beta, \alpha)]} \left\{ \left[(C_2(\beta, \alpha))^4 C_1(\alpha, \beta) + [\Psi_1(\alpha + \beta)]^4 C_1(\beta, \alpha) \right]^{\frac{1}{2}} \right. \\
&\quad \left. + \left[(C_2(\alpha, \beta))^4 C_1(\beta, \alpha) + [\Psi_1(\alpha + \beta)]^4 C_1(\alpha, \beta) \right]^{\frac{1}{2}} \right\} \\
&\quad + \frac{M_B \sqrt{24}}{\sqrt{n}\delta_I [C_2(\alpha, \beta) + C_2(\beta, \alpha)]} \left\{ [\Psi_1(\alpha + \beta)]^2 (C_3(\alpha, \beta) + C_3(\beta, \alpha)) \right. \\
&\quad \left. + 2\sqrt{C_1(\alpha, \beta)C_1(\beta, \alpha)} \left[[\Psi_1(\alpha + \beta)]^4 + [C_2(\alpha, \beta)]^2 [C_2(\beta, \alpha)]^2 \right] \right\}^{\frac{1}{2}} \\
&\quad + \frac{\Psi_1(\beta) + \Psi_1(\alpha) - 2\Psi_1(\alpha + \beta)}{\delta_I}, \\
\omega_B &:= 1 - 8 \frac{M_B}{n\varepsilon^2} - \frac{2M_B \{ (\Psi_1(\beta) + \sqrt{\delta_I}) C_4(\beta, \alpha) + (\Psi_1(\alpha) + \sqrt{\delta_I}) C_4(\alpha, \beta) \}}{\sqrt{n}\delta_I (C_2(\alpha, \beta) + C_2(\beta, \alpha))}.
\end{aligned} \tag{4.80}$$

Corollary 4.4. Let X_1, X_2, \dots, X_n be i.i.d. random variables from the $\text{Beta}(\alpha, \beta)$ distribution with $\theta_0 = (\alpha, \beta)$. Let $m = \min \{ \alpha, \beta \}$, $\varepsilon = \frac{m}{2} > 0$ and

$$\begin{aligned}
n &\geq \frac{4}{m^2} \left(\frac{mM_B (\Psi_1(\beta) + \sqrt{\delta_I}) C_4(\beta, \alpha) + (\Psi_1(\alpha) + \sqrt{\delta_I}) C_4(\alpha, \beta)}{2\sqrt{\delta_I} (C_2(\alpha, \beta) + C_2(\beta, \alpha))} \right. \\
&\quad \left. + \left[\left(\frac{mM_B (\Psi_1(\beta) + \sqrt{\delta_I}) C_4(\beta, \alpha) + (\Psi_1(\alpha) + \sqrt{\delta_I}) C_4(\alpha, \beta)}{2\sqrt{\delta_I} (C_2(\alpha, \beta) + C_2(\beta, \alpha))} \right)^2 + 8M_B \right]^{\frac{1}{2}} \right)^2.
\end{aligned} \tag{4.81}$$

Then,

a) When n satisfies (4.81), $E \left[\sum_{j=1}^2 (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right] \leq \sqrt{\frac{\gamma_B}{n\omega_B}}.$

b) For $\mathbf{Z} \sim N_2(\mathbf{0}, I_{2 \times 2})$, we obtain

$$\begin{aligned}
 & \left| E \left[h \left(\sqrt{n} [I(\theta_0)]^{\frac{1}{2}} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) \right] - E[h(\mathbf{Z})] \right| \leq \\
 & \frac{2\|h\|_2}{\sqrt{n}\delta_I [C_2(\alpha, \beta) + C_2(\beta, \alpha)]} \left\{ \left[(C_2(\beta, \alpha))^4 C_1(\alpha, \beta) + [\Psi_1(\alpha + \beta)]^4 C_1(\beta, \alpha) \right]^{\frac{1}{2}} \right. \\
 & \quad \left. + \left[(C_2(\alpha, \beta))^4 C_1(\beta, \alpha) + [\Psi_1(\alpha + \beta)]^4 C_1(\alpha, \beta) \right]^{\frac{1}{2}} \right\} \\
 & + \frac{\|h\|_2 \sqrt{6}}{\sqrt{n}\delta_I [C_2(\alpha, \beta) + C_2(\beta, \alpha)]} \left\{ [\Psi_1(\alpha + \beta)]^2 (C_3(\alpha, \beta) + C_3(\beta, \alpha)) \right. \\
 & \quad \left. + 2\sqrt{C_1(\alpha, \beta)C_1(\beta, \alpha)} \left[[\Psi_1(\alpha + \beta)]^4 + [C_2(\alpha, \beta)]^2 [C_2(\beta, \alpha)]^2 \right] \right\}^{\frac{1}{2}} \\
 & + \frac{32\|h\|_3 8^{\frac{3}{4}}}{3\sqrt{n} [\delta_I [C_2(\alpha, \beta) + C_2(\beta, \alpha)]]^{\frac{3}{2}}} \left\{ \left[[C_2(\beta, \alpha)]^3 + [\Psi_1(\alpha + \beta)]^3 \right] [C_1(\alpha, \beta)]^{\frac{3}{4}} \right. \\
 & \quad \left. + \left[[C_2(\alpha, \beta)]^3 + [\Psi_1(\alpha + \beta)]^3 \right] [C_1(\beta, \alpha)]^{\frac{3}{4}} \right\} \\
 & + \frac{8\|h\|\gamma_B}{nm^2\omega_B} + \frac{\|h\|_1 \gamma_B \{ (\Psi_1(\beta) + \sqrt{\delta_I}) C_4(\beta, \alpha) + (\Psi_1(\alpha) + \sqrt{\delta_I}) C_4(\alpha, \beta) \}}{2\sqrt{n}\omega_B \sqrt{\delta_I (C_2(\alpha, \beta) + C_2(\beta, \alpha))}}.
 \end{aligned} \tag{4.82}$$

The following lemma gives a useful approach to finding a square root of a matrix. The proof of the lemma is given in Somayya (1997).

Lemma 4.3. Let $M = \begin{pmatrix} A & B \\ C & D \end{pmatrix}$, where A, B, C, D may be real or complex numbers. In addition, let τ and δ be the trace and determinant of M , respectively. Let s and t , such

that $s^2 = \delta$ and $t^2 = \tau + 2s$. Then, if $t \neq 0$, a square root of M is

$$R = \frac{1}{t} \begin{pmatrix} A+s & B \\ C & D+s \end{pmatrix}.$$

Remark 4.6. The lemma holds for $s = \pm\sqrt{\delta}$ and $t = \pm\sqrt{\tau+2s}$. Using different signs for those quantities gives different roots of the matrix M . From now on, unless otherwise stated, we take the positive roots for s and t .

Proof of Corollary 4.4.

Part a). The probability density function is

$$f(x|\boldsymbol{\theta}) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1},$$

with $\alpha, \beta > 0$ and $x \in [0, 1]$. Hence, for $j, k \in \mathbb{Z}^+$

$$\begin{aligned} \frac{\partial}{\partial \alpha} \log f(x|\boldsymbol{\theta}) &= \Psi(\alpha + \beta) - \Psi(\alpha) + \log(x) \\ \frac{\partial}{\partial \beta} \log f(x|\boldsymbol{\theta}) &= \Psi(\alpha + \beta) - \Psi(\beta) + \log(1-x) \\ \frac{\partial^{j+1}}{\partial \alpha^{j+1}} \log f(x|\boldsymbol{\theta}) &= \Psi_j(\alpha + \beta) - \Psi_j(\alpha) \\ \frac{\partial^{j+1}}{\partial \beta^{j+1}} \log f(x|\boldsymbol{\theta}) &= \Psi_j(\alpha + \beta) - \Psi_j(\beta) \\ \frac{\partial^{k+j}}{\partial \alpha^k \partial \beta^j} \log f(x|\boldsymbol{\theta}) &= \Psi_{k+j-1}(\alpha + \beta). \end{aligned} \tag{4.83}$$

The assumptions (Con.1)-(Con.3) are satisfied. For (Con.1) notice that $x \in [0, 1]$. For (Con.2), let $\varepsilon_0 = \varepsilon = \varepsilon(\boldsymbol{\theta}_0) : 0 < \varepsilon < \min\{\alpha, \beta\}$. In addition, $\sum_{i=1}^{\infty} \frac{1}{i^4} = \frac{\pi^4}{90}$ and $\sum_{i=1}^{\infty} \frac{1}{i^3} = \zeta(3)$, where $\zeta(\cdot)$ is the Riemann zeta function and $\zeta(3)$ is known as Apéry's constant, which is an irrational number (≈ 1.202). For the calculations that follow,

$\zeta(3) < 1.21$ is used. With $\theta = (\theta_1, \theta_2)$, using (3.67), (3.68) and (3.69) yields

$$\begin{aligned}
 \sup_{\theta: |\theta_1 - \alpha| < \varepsilon, |\theta_2 - \beta| < \varepsilon} \left| \frac{\partial^3}{\partial \theta_1^3} \log f(x|\theta) \right| &\leq \frac{6\beta}{(\alpha - \varepsilon)^4} + \frac{\beta\pi^4}{15} := M_{111} \\
 \sup_{\theta: |\theta_1 - \alpha| < \varepsilon, |\theta_2 - \beta| < \varepsilon} \left| \frac{\partial^3}{\partial \theta_2^3} \log f(x|\theta) \right| &\leq \frac{6\alpha}{(\beta - \varepsilon)^4} + \frac{\alpha\pi^4}{15} := M_{222} \\
 \sup_{\theta: |\theta_1 - \alpha| < \varepsilon, |\theta_2 - \beta| < \varepsilon} \left| \frac{\partial^3}{\partial \theta_1^2 \partial \theta_2} \log f(x|\theta) \right| &= \sup_{\theta: |\theta_1 - \alpha| < \varepsilon, |\theta_2 - \beta| < \varepsilon} \left| \frac{\partial^3}{\partial \theta_1 \partial \theta_2^2} \log f(x|\theta) \right| \\
 &\leq \frac{2}{(\alpha + \beta - 2\varepsilon)^3} + 2\zeta(3) < \frac{2}{(\alpha + \beta - 2\varepsilon)^3} + 2.42 \\
 &:= M_{112} = M_{121} = M_{211} = M_{122} = M_{212} = M_{221}. \quad (4.84)
 \end{aligned}$$

Therefore, (Con.2) is also satisfied. Using (4.83), it is straightforward that the expected Fisher Information matrix in this case becomes

$$I(\theta_0) = \begin{pmatrix} \Psi_1(\alpha) - \Psi_1(\alpha + \beta) & -\Psi_1(\alpha + \beta) \\ -\Psi_1(\alpha + \beta) & \Psi_1(\beta) - \Psi_1(\alpha + \beta) \end{pmatrix}.$$

The inverse of this 2×2 matrix is

$$[I(\theta_0)]^{-1} = \frac{1}{\delta_I} \begin{pmatrix} \Psi_1(\beta) - \Psi_1(\alpha + \beta) & \Psi_1(\alpha + \beta) \\ \Psi_1(\alpha + \beta) & \Psi_1(\alpha) - \Psi_1(\alpha + \beta) \end{pmatrix}$$

and using Lemma 4.3

$$[I(\theta_0)]^{-\frac{1}{2}} = \frac{1}{[\delta_I (C_2(\alpha, \beta) + C_2(\beta, \alpha))]^{\frac{1}{2}}} \begin{pmatrix} C_2(\beta, \alpha) & \Psi_1(\alpha + \beta) \\ \Psi_1(\alpha + \beta) & C_2(\alpha, \beta) \end{pmatrix}. \quad (4.85)$$

For $M = M_B$, as defined in (Con.3), we have that

$$M_B = \frac{1}{\delta_I} \sup \{ \Psi_1(\alpha + \beta), \Psi_1(\alpha) - \Psi_1(\alpha + \beta), \Psi_1(\beta) - \Psi_1(\alpha + \beta) \}$$

$$= \frac{1}{\delta_I} \sup \{ \Psi_1(\alpha + \beta), \Psi_1(\min \{ \alpha, \beta \}) - \Psi_1(\alpha + \beta) \},$$

as $\Psi_1(\cdot)$ is a decreasing function. Having found M_B and $[I(\theta_0)]^{-\frac{1}{2}}$ and using (4.84) with the notation (4.80), (Con.3) is satisfied for

$$n \geq \frac{4}{m^2} \left(\frac{mM_B (\Psi_1(\beta) + \sqrt{\delta_I}) C_4(\beta, \alpha) + (\Psi_1(\alpha) + \sqrt{\delta_I}) C_4(\alpha, \beta)}{2\sqrt{\delta_I} (C_2(\alpha, \beta) + C_2(\beta, \alpha))} \right. \\ \left. + \left[\left(\frac{mM_B (\Psi_1(\beta) + \sqrt{\delta_I}) C_4(\beta, \alpha) + (\Psi_1(\alpha) + \sqrt{\delta_I}) C_4(\alpha, \beta)}{2\sqrt{\delta_I} (C_2(\alpha, \beta) + C_2(\beta, \alpha))} \right)^2 + 8M_B \right]^{\frac{1}{2}} \right)^2,$$

with $C_4(x, y)$ as in (4.80). We proceed with the bound for $E \left[\sum_{j=1}^2 (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right]$. Firstly, as the second-order partial derivatives of the log-likelihood function for the Beta distribution are not random,

$$\text{Var} \left(\frac{\partial^2}{\partial \alpha^2} \log f(X|\theta_0) \right) = \text{Var} \left(\frac{\partial^2}{\partial \beta^2} \log f(X|\theta_0) \right) = \text{Var} \left(\frac{\partial^2}{\partial \alpha \partial \beta} \log f(X|\theta_0) \right) = 0$$

and hence $v = 0$. Therefore, for the Beta distribution

$$E \left(\sum_{j=1}^2 (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right) \leq U_1 = \frac{\sqrt{4\omega\gamma}}{2\sqrt{n\omega}} = \sqrt{\frac{\gamma}{n\omega}}.$$

Our focus now turns to find an upper bound for γ . The first quantity of the result for γ in (4.65) is

$$\sum_{j=1}^d \left| [I(\theta_0)]^{-1} \right|_{jj} = \frac{1}{\delta_I} [\Psi_1(\beta) + \Psi_1(\alpha) - 2\Psi_1(\alpha + \beta)]. \quad (4.86)$$

For the second term, using equations (3.70) and (3.71),

$$E \left[\left(\frac{\partial}{\partial \alpha} \log f(X_1|\theta_0) \right)^4 \right] \leq 8C_1(\alpha, \beta) \\ E \left[\left(\frac{\partial}{\partial \beta} \log f(X_1|\theta_0) \right)^4 \right] \leq 8C_1(\beta, \alpha), \quad (4.87)$$

with $C_1(x, y)$ as in (4.80). Using (4.87) and (4.85) gives

$$\begin{aligned}
& \left[\text{Var} \left(\left(\sum_{k=1}^2 \left[I(\boldsymbol{\theta}_0) \right]^{-\frac{1}{2}} \right)_{1k} \frac{\partial}{\partial \boldsymbol{\theta}_k} \log f(X_1 | \boldsymbol{\theta}_0) \right)^2 \right]^{\frac{1}{2}} \\
& + \left[\text{Var} \left(\left(\sum_{k=1}^2 \left[I(\boldsymbol{\theta}_0) \right]^{-\frac{1}{2}} \right)_{2k} \frac{\partial}{\partial \boldsymbol{\theta}_k} \log f(X_1 | \boldsymbol{\theta}_0) \right)^2 \right]^{\frac{1}{2}} \\
& \leq \left[\text{E} \left(\left[I(\boldsymbol{\theta}_0) \right]^{-\frac{1}{2}} \right)_{11} \frac{\partial}{\partial \alpha} \log f(X_1 | \boldsymbol{\theta}_0) + \left[I(\boldsymbol{\theta}_0) \right]^{-\frac{1}{2}}_{12} \frac{\partial}{\partial \beta} \log f(X_1 | \boldsymbol{\theta}_0) \right)^4 \right]^{\frac{1}{2}} \\
& + \left[\text{E} \left(\left[I(\boldsymbol{\theta}_0) \right]^{-\frac{1}{2}} \right)_{12} \frac{\partial}{\partial \alpha} \log f(X_1 | \boldsymbol{\theta}_0) + \left[I(\boldsymbol{\theta}_0) \right]^{-\frac{1}{2}}_{22} \frac{\partial}{\partial \beta} \log f(X_1 | \boldsymbol{\theta}_0) \right)^4 \right]^{\frac{1}{2}} \\
& \leq \sqrt{8} \left[\left(\left[I(\boldsymbol{\theta}_0) \right]^{-\frac{1}{2}} \right)_{11} \right)^4 \text{E} \left(\frac{\partial}{\partial \alpha} \log f(X_1 | \boldsymbol{\theta}_0) \right)^4 \right. \\
& \quad \left. + \left(\left[I(\boldsymbol{\theta}_0) \right]^{-\frac{1}{2}} \right)_{12} \right)^4 \text{E} \left(\frac{\partial}{\partial \beta} \log f(X_1 | \boldsymbol{\theta}_0) \right)^4 \right]^{\frac{1}{2}} \\
& + \sqrt{8} \left[\left(\left[I(\boldsymbol{\theta}_0) \right]^{-\frac{1}{2}} \right)_{12} \right)^4 \text{E} \left(\frac{\partial}{\partial \alpha} \log f(X_1 | \boldsymbol{\theta}_0) \right)^4 \right. \\
& \quad \left. + \left(\left[I(\boldsymbol{\theta}_0) \right]^{-\frac{1}{2}} \right)_{22} \right)^4 \text{E} \left(\frac{\partial}{\partial \beta} \log f(X_1 | \boldsymbol{\theta}_0) \right)^4 \right]^{\frac{1}{2}} \\
& \leq \frac{8}{\delta_I [C_2(\alpha, \beta) + C_2(\beta, \alpha)]} \left\{ \left[(C_2(\beta, \alpha))^4 C_1(\alpha, \beta) + (\Psi_1(\alpha + \beta))^4 C_1(\beta, \alpha) \right]^{\frac{1}{2}} \right. \\
& \quad \left. + \left[(\Psi_1(\alpha + \beta))^4 C_1(\alpha, \beta) + (C_2(\alpha, \beta))^4 C_1(\beta, \alpha) \right]^{\frac{1}{2}} \right\}.
\end{aligned} \tag{4.88}$$

For the third term of the expression for γ as in (4.65) since $d = 2$,

$$\begin{aligned}
& \left[\text{Var} \left(\sum_{q=1}^2 \sum_{v=1}^2 \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{2q} \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{1v} \frac{\partial}{\partial \theta_q} \log f(X_1 | \boldsymbol{\theta}_0) \frac{\partial}{\partial \theta_v} \log f(X_1 | \boldsymbol{\theta}_0) \right) \right]^{\frac{1}{2}} \\
& \leq \left[\text{E} \left(\sum_{q=1}^2 \sum_{v=1}^2 \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{2q} \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{1v} \frac{\partial}{\partial \theta_q} \log f(X_1 | \boldsymbol{\theta}_0) \frac{\partial}{\partial \theta_v} \log f(X_1 | \boldsymbol{\theta}_0) \right)^2 \right]^{\frac{1}{2}} \\
& = \left[\text{E} \left(\left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{12} \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{11} \left(\frac{\partial}{\partial \alpha} \log f(X_1 | \boldsymbol{\theta}_0) \right)^2 \right. \right. \\
& \quad + \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{12} \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{22} \left(\frac{\partial}{\partial \beta} \log f(X_1 | \boldsymbol{\theta}_0) \right)^2 \\
& \quad + \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{22} \left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{11} \frac{\partial}{\partial \alpha} \log f(X_1 | \boldsymbol{\theta}_0) \frac{\partial}{\partial \beta} \log f(X_1 | \boldsymbol{\theta}_0) \\
& \quad \left. \left. + \left(\left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{12} \right)^2 \frac{\partial}{\partial \alpha} \log f(X_1 | \boldsymbol{\theta}_0) \frac{\partial}{\partial \beta} \log f(X_1 | \boldsymbol{\theta}_0) \right)^2 \right]^{\frac{1}{2}}. \quad (4.89)
\end{aligned}$$

Using now the known inequality, $(\sum_{i=1}^n a_i)^2 \leq n \sum_{i=1}^n a_i^2$ for $a_i \in \mathbb{R}$, and also the Cauchy-Schwarz inequality yields

$$\begin{aligned}
(4.89) & \leq \left[3 \left(\left(\left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{12} \right)^2 \left(\left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{11} \right)^2 \text{E} \left(\frac{\partial}{\partial \alpha} \log f(X_1 | \boldsymbol{\theta}_0) \right)^4 \right. \right. \\
& \quad + \left(\left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{22} \right)^2 \left(\left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{12} \right)^2 \text{E} \left(\frac{\partial}{\partial \beta} \log f(X_1 | \boldsymbol{\theta}_0) \right)^4 \\
& \quad + 2 \left(\left(\left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{22} \right)^2 \left(\left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{11} \right)^2 + \left(\left[[I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{12} \right)^4 \right) \\
& \quad \left. \times \sqrt{\text{E} \left(\frac{\partial}{\partial \alpha} \log f(X_1 | \boldsymbol{\theta}_0) \right)^4 \text{E} \left(\frac{\partial}{\partial \beta} \log f(X_1 | \boldsymbol{\theta}_0) \right)^4} \right]^{\frac{1}{2}} \\
& \leq \frac{\sqrt{24}}{\delta_I [C_2(\alpha, \beta) + C_2(\beta, \alpha)]} \left\{ (\Psi_1(\alpha + \beta))^2 [C_3(\alpha, \beta) + C_3(\beta, \alpha)] \right. \\
& \quad \left. + 2 [(C_2(\alpha, \beta))^2 (C_2(\beta, \alpha))^2 + (\Psi_1(\alpha + \beta))^4] \sqrt{C_1(\alpha, \beta) C_1(\beta, \alpha)} \right\}^{\frac{1}{2}}. \quad (4.90)
\end{aligned}$$

The inequalities in (4.86), (4.88) and (4.90) show that $\gamma \leq \gamma_B$ as in (4.80). To calculate ω ,

$$\begin{aligned}
& \sum_{l=1}^2 \sum_{k=1}^2 \left| \left[I(\theta_0) \right]^{-\frac{1}{2}} \right|_{lk} \sum_{m=1}^2 \sum_{i=1}^2 M_{kim} \\
&= \sum_{k=1}^2 \left[\left| \left[I(\theta_0) \right]^{-\frac{1}{2}} \right|_{1k} + \left| \left[I(\theta_0) \right]^{-\frac{1}{2}} \right|_{2k} \right] \sum_{m=1}^2 \sum_{i=1}^2 M_{kim} \\
&= \frac{\Psi_1(\beta) + \sqrt{\delta_I}}{\sqrt{\delta_I (C_2(\alpha, \beta) + C_2(\beta, \alpha))}} \sum_{m=1}^2 \sum_{i=1}^2 M_{1im} + \frac{\Psi_1(\alpha) + \sqrt{\delta_I}}{\sqrt{\delta_I (C_2(\alpha, \beta) + C_2(\beta, \alpha))}} \sum_{m=1}^2 \sum_{i=1}^2 M_{2im} \\
&= \frac{\Psi_1(\beta) + \sqrt{\delta_I}}{\sqrt{\delta_I (C_2(\alpha, \beta) + C_2(\beta, \alpha))}} \left(\frac{6\beta}{(\alpha - \varepsilon)^4} + \frac{\beta\pi^4}{15} + \frac{6}{(\alpha + \beta - \varepsilon)^3} + 7.26 \right) \\
&\quad + \frac{\Psi_1(\alpha) + \sqrt{\delta_I}}{\sqrt{\delta_I (C_2(\alpha, \beta) + C_2(\beta, \alpha))}} \left(\frac{6\alpha}{(\beta - \varepsilon)^4} + \frac{\alpha\pi^4}{15} + \frac{6}{(\alpha + \beta - \varepsilon)^3} + 7.26 \right) \\
&= \frac{(\Psi_1(\beta) + \sqrt{\delta_I}) C_4(\beta, \alpha) + (\Psi_1(\alpha) + \sqrt{\delta_I}) C_4(\alpha, \beta)}{\sqrt{\delta_I (C_2(\alpha, \beta) + C_2(\beta, \alpha))}} \tag{4.91}
\end{aligned}$$

shows that $\omega = \omega_B$ as defined in (4.80). Since $\omega_B > 0$ for n satisfying (Con.3), we conclude that

$$\mathbb{E} \left(\sum_{j=1}^2 \left(\hat{\theta}_n(\mathbf{X})_j - \theta_{0j} \right)^2 \right) \leq \frac{\sqrt{4\omega_B \gamma_B}}{2\sqrt{n}\omega_B} = \frac{\sqrt{\omega_B \gamma_B}}{\sqrt{n}\omega_B} = \sqrt{\frac{\gamma_B}{n\omega_B}}. \tag{4.92}$$

Part **b**). For (4.82) we use the general expression of the bound in (4.79). Upper bounds for the first two terms of (4.79) have already been obtained in (4.88) and (4.90). Now, for ease of presentation let

$$\begin{aligned}
K_1(X_1, X'_1) &:= \frac{\partial}{\partial \alpha} \log f(X'_1 | \theta_0) - \frac{\partial}{\partial \alpha} \log f(X_1 | \theta_0) \\
K_2(X_1, X'_1) &:= \frac{\partial}{\partial \beta} \log f(X'_1 | \theta_0) - \frac{\partial}{\partial \beta} \log f(X_1 | \theta_0). \tag{4.93}
\end{aligned}$$

Using (4.93) and the inequality in (4.42), we have for the third term in (4.79),

$$\begin{aligned}
& \mathbb{E} \left(\sum_{i=1}^2 \left| \sum_{l=1}^2 [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{il} \left(\frac{\partial}{\partial \boldsymbol{\theta}_l} \log f(X'_1 | \boldsymbol{\theta}_0) - \frac{\partial}{\partial \boldsymbol{\theta}_l} \log f(X_1 | \boldsymbol{\theta}_0) \right) \right|^3 \\
&= \mathbb{E} \left(\left| \sum_{l=1}^2 [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{1l} \left(\frac{\partial}{\partial \boldsymbol{\theta}_l} \log f(X'_1 | \boldsymbol{\theta}_0) - \frac{\partial}{\partial \boldsymbol{\theta}_l} \log f(X_1 | \boldsymbol{\theta}_0) \right) \right|^3 \\
&\quad + \left| \sum_{l=1}^2 [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{2l} \left(\frac{\partial}{\partial \boldsymbol{\theta}_l} \log f(X'_1 | \boldsymbol{\theta}_0) - \frac{\partial}{\partial \boldsymbol{\theta}_l} \log f(X_1 | \boldsymbol{\theta}_0) \right) \right|^3 \\
&\leq 4\mathbb{E} \left(\left| \sum_{l=1}^2 [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{1l} \left(\frac{\partial}{\partial \boldsymbol{\theta}_l} \log f(X'_1 | \boldsymbol{\theta}_0) - \frac{\partial}{\partial \boldsymbol{\theta}_l} \log f(X_1 | \boldsymbol{\theta}_0) \right) \right|^3 \\
&\quad + \left| \sum_{l=1}^2 [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{2l} \left(\frac{\partial}{\partial \boldsymbol{\theta}_l} \log f(X'_1 | \boldsymbol{\theta}_0) - \frac{\partial}{\partial \boldsymbol{\theta}_l} \log f(X_1 | \boldsymbol{\theta}_0) \right) \right|^3 \\
&= 4\mathbb{E} \left(\left| [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{11} (K_1(X_1, X'_1)) + [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{12} (K_2(X_1, X'_1)) \right|^3 \\
&\quad + \left| [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{21} (K_1(X_1, X'_1)) + [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{22} (K_2(X_1, X'_1)) \right|^3 \\
&\leq 16 \left(\left| [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{11} \right|^3 + \left| [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{12} \right|^3 \mathbb{E} (|K_1(X_1, X'_1)|^3) \\
&\quad + \left| [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{21} \right|^3 + \left| [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{22} \right|^3 \mathbb{E} (|K_2(X_1, X'_1)|^3) \\
&\leq 128 \left(\left| [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{11} \right|^3 + \left| [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{12} \right|^3 8^{\frac{3}{4}} [C_1(\boldsymbol{\alpha}, \boldsymbol{\beta})]^{\frac{3}{4}} \\
&\quad + \left| [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{21} \right|^3 + \left| [I(\boldsymbol{\theta}_0)]^{-\frac{1}{2}} \right]_{22} \right|^3 8^{\frac{3}{4}} [C_1(\boldsymbol{\beta}, \boldsymbol{\alpha})]^{\frac{3}{4}} \\
&= \frac{128}{[\delta_I(C_2(\boldsymbol{\alpha}, \boldsymbol{\beta}) + C_2(\boldsymbol{\beta}, \boldsymbol{\alpha}))]^{\frac{3}{2}}} \left\{ 8^{\frac{3}{4}} [C_1(\boldsymbol{\alpha}, \boldsymbol{\beta})]^{\frac{3}{4}} ([C_2(\boldsymbol{\beta}, \boldsymbol{\alpha})]^3 + [\Psi_1(\boldsymbol{\alpha} + \boldsymbol{\beta})]^3) \right. \\
&\quad \left. + 8^{\frac{3}{4}} [C_1(\boldsymbol{\beta}, \boldsymbol{\alpha})]^{\frac{3}{4}} ([C_2(\boldsymbol{\alpha}, \boldsymbol{\beta})]^3 + [\Psi_1(\boldsymbol{\alpha} + \boldsymbol{\beta})]^3) \right\}.
\end{aligned}$$

(4.94)

The above inequalities are a result of (4.42). To be more specific, for the last inequality in (4.94) we use that

$$\begin{aligned} & \left| \frac{\partial}{\partial \theta_j} \log f(X'_1 | \theta_0) - \frac{\partial}{\partial \theta_j} \log f(X_1 | \theta_0) \right|^3 \\ & \leq \left(\left| \frac{\partial}{\partial \theta_j} \log f(X'_1 | \theta_0) \right| + \left| \frac{\partial}{\partial \theta_j} \log f(X_1 | \theta_0) \right| \right)^3 \\ & \leq 4 \left(\left| \frac{\partial}{\partial \theta_j} \log f(X'_1 | \theta_0) \right|^3 + \left| \frac{\partial}{\partial \theta_j} \log f(X_1 | \theta_0) \right|^3 \right) \end{aligned}$$

and therefore,

$$\begin{aligned} & \mathbb{E} \left(\left| \frac{\partial}{\partial \theta_j} \log f(X'_1 | \theta_0) - \frac{\partial}{\partial \theta_j} \log f(X_1 | \theta_0) \right|^3 \right) \\ & \leq 4 \left(\mathbb{E} \left(\left| \frac{\partial}{\partial \theta_j} \log f(X'_1 | \theta_0) \right|^3 \right) + \mathbb{E} \left(\left| \frac{\partial}{\partial \theta_j} \log f(X_1 | \theta_0) \right|^3 \right) \right) \\ & = 8 \mathbb{E} \left(\left| \frac{\partial}{\partial \theta_j} \log f(X_1 | \theta_0) \right|^3 \right) \leq 8 \left[\mathbb{E} \left(\left| \frac{\partial}{\partial \theta_j} \log f(X_1 | \theta_0) \right|^4 \right) \right]^{\frac{3}{4}}. \end{aligned}$$

Thus, applying the results of (4.88), (4.90), (4.91), (4.92) and (4.94) to (4.79) Corollary 4.4 follows. \square

Remark 4.7. Using the notation (4.80), it is straightforward that the first three terms of the bound are $\mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$. In addition, since γ_B and ω_B are $\mathcal{O}(1)$, the fourth and the fifth term of the bound are of order $\frac{1}{n}$ and $\frac{1}{\sqrt{n}}$, respectively. Combining these results for the order of each of the terms, the order of the bound (4.82) is $\frac{1}{\sqrt{n}}$.

Chapter 5

An approach using the Delta method

The work presented in this chapter is based on a paper with Christophe Ley that has been submitted to the *Latin American Journal of Probability and Mathematical Statistics (ALEA)*. We return to the single-parameter setting and our aim is still to provide sharp explicit upper bounds on Zolotarev-type distances between the distribution of the MLE and the normal distribution. Our approach is based on a combination of the Delta method, Stein's method, Taylor expansions and conditional expectations. The bounds presented in this chapter are slightly restrictive in the sense that they are applicable only in situations where a function of the MLE can be expressed as a sum of independent and identically distributed terms which are functions of the observations. In these cases, such bounds are at least as sharp as those presented in Chapter 3.

The chapter is organised as follows. In Section 5.1 the new bound in the case where a function of the MLE is a sum of i.i.d. terms is described, proved and compared to the bound given in Theorem 3.1. In Section 5.2, the results are applied to the class of single-parameter exponential family distributions and some examples are treated in detail.

5.1 New bounds on the distance to normal for the MLE

The aim of the section is to find an upper bound for $d_{bW} \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right)$. The tools used to reach this result are now the Delta method, Stein's method (through Lemma 2.2) and, as in Chapter 3, Taylor expansions and conditional expectations.

Consider a sample of i.i.d. random variables X_1, X_2, \dots, X_n from a single-parameter distribution and, as in Chapter 3, $\hat{\theta}_n(\mathbf{X})$ denotes the MLE for the parameter $\theta \in \Theta$, which is as usual to exist. We are interested in settings where there exists a one-to-one twice differentiable function $q : \Theta \rightarrow \mathbb{R}$ such that

$$q(\hat{\theta}_n(\mathbf{X})) = \frac{1}{n} \sum_{i=1}^n g(X_i) \quad (5.1)$$

for some $g : \mathbb{R} \rightarrow \mathbb{R}$. Some classical examples include

- 1) The normal distribution with density $f(x|\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\}$, $x \in \mathbb{R}$, for which $\mu \in \mathbb{R}$ is an unknown parameter, whereas $\sigma > 0$ is considered to be known. The MLE for μ is

$$\hat{\theta}_n(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n X_i;$$

in this case $q(\theta) = \theta$ and $g(x) = x$.

- 2) The normal distribution, where now the mean μ is known and σ^2 represents the unknown parameter, with

$$\hat{\theta}_n(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2;$$

in this case $q(\theta) = \theta$ and $g(x) = (x - \mu)^2$.

- 3) The Weibull distribution with density $f(x|\alpha, \sigma) = \frac{\alpha}{\sigma} \left(\frac{x}{\sigma}\right)^{\alpha-1} \exp\left\{-\left(\frac{x}{\sigma}\right)^\alpha\right\}$, $x \geq 0$, where σ is the unknown scale parameter and $\alpha > 0$ is fixed. The MLE for σ is

$$\hat{\theta}_n(\mathbf{X}) = \left(\frac{1}{n} \sum_{i=1}^n X_i^\alpha \right)^{\frac{1}{\alpha}};$$

in this case, it is easily deduced that $q(\theta) = \theta^\alpha$ and $g(x) = x^\alpha$.

- 4) The Laplace scale model with density $f(x|\sigma) = \frac{1}{2\sigma} \exp\left\{-\frac{|x|}{\sigma}\right\}$, $\sigma > 0$, over \mathbb{R} , for which

$$\hat{\theta}_n(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n |X_i|;$$

in this case $q(\theta) = \theta$ and $g(x) = |x|$.

Moreover, the broad single-parameter exponential families, discussed in Section 3.3 satisfy condition (5.1); see Proposition 5.1 for details. Therefore, the results presented in this section apply to many of the well-known distributions.

Our strategy consists in benefiting from the special form of $q(\hat{\theta}_n(\mathbf{X}))$, which is a sum of random variables, and thus allows us to use the sharp bound of Lemma 2.2. It is precisely at this point that the Delta method comes into play: instead of comparing $\hat{\theta}_n(\mathbf{X})$ to $Z \sim N(0, 1)$ we rather compare $q(\hat{\theta}_n(\mathbf{X}))$ to Z , and then bound the distance between $\hat{\theta}_n(\mathbf{X})$ and $q(\hat{\theta}_n(\mathbf{X}))$. The univariate Delta method is given in the following theorem.

Theorem 5.1. *Let θ be the parameter with parameter space $\Theta \subset \mathbb{R}$. Also, let Y_n be a sequence of random variables which satisfies*

$$\sqrt{n}(Y_n - \theta) \xrightarrow[n \rightarrow \infty]{d} \sigma Z.$$

Then, for any function $g : \mathbb{R} \rightarrow \mathbb{R}$ such that $g'(\theta)$ exists and is not equal to zero we obtain

$$\sqrt{n}(g(Y_n) - g(\theta)) \xrightarrow[n \rightarrow \infty]{d} g'(\theta)\sigma Z. \quad (5.2)$$

For a proof of the theorem see (Casella and Berger, 2002, p.243). The following lemma is useful for the proof of the main result of this section.

Lemma 5.1. *Let X_1, X_2, \dots, X_n be i.i.d. random variables following a distribution with unknown parameter θ . Its true value is denoted by θ_0 . Assume that (R1)-(R4) of Section 2.2 are satisfied and that the MLE $\hat{\theta}_n(\mathbf{X})$ exists and is unique. Let $q : \Theta \rightarrow \mathbb{R}$ be a differentiable function such that $q(\hat{\theta}_n(\mathbf{X})) = \frac{1}{n} \sum_{i=1}^n g(X_i)$, for some function $g : \mathbb{R} \rightarrow \mathbb{R}$. With $E[\cdot]$ and $\text{Var}[\cdot]$ the expectation and variance under θ_0 , then if $E[g(X_1)]$ and $\text{Var}[g(X_1)]$ exist, we have that*

$$E[g(X_1)] = q(\theta_0) \quad \text{and} \quad \text{Var}[g(X_1)] = \frac{[q'(\theta_0)]^2}{i(\theta_0)}.$$

Proof. Since $q(\hat{\theta}_n(\mathbf{X})) = \frac{1}{n} \sum_{i=1}^n g(X_i)$, we have

$$\begin{aligned} E[q(\hat{\theta}_n(\mathbf{X}))] &= \frac{1}{n} \sum_{i=1}^n E[g(X_i)] = E[g(X_1)] \\ \text{Var}[\sqrt{n}q(\hat{\theta}_n(\mathbf{X}))] &= \frac{1}{n} \sum_{i=1}^n \text{Var}[g(X_i)] = \text{Var}[g(X_1)]. \end{aligned} \quad (5.3)$$

In addition, since the regularity conditions (R1)-(R4) are satisfied then the result (b) in (2.2) holds and the Delta method result in (5.2) yields

$$\sqrt{n}(q(\hat{\theta}_n(\mathbf{X})) - q(\theta_0)) \xrightarrow[n \rightarrow \infty]{d} \frac{q'(\theta_0)}{\sqrt{i(\theta_0)}} Z,$$

which indicates that

$$\sqrt{n}(E(q(\hat{\theta}_n(\mathbf{X}))) - q(\theta_0)) \xrightarrow[n \rightarrow \infty]{} 0 \quad \text{and} \quad \text{Var}(\sqrt{n}q(\hat{\theta}_n(\mathbf{X}))) \xrightarrow[n \rightarrow \infty]{} \frac{[q'(\theta_0)]^2}{i(\theta_0)}. \quad (5.4)$$

Combining the results in (5.3) and (5.4) we summarise that $E[q(\hat{\theta}_n(\mathbf{X}))] = E[g(X_1)]$ and as n goes to infinity the limit of $E(q(\hat{\theta}_n(\mathbf{X})))$ is $q(\theta_0)$, meaning that $E[g(X_1)] = q(\theta_0)$. Applying the same idea for $\text{Var}[\sqrt{n}q(\hat{\theta}_n(\mathbf{X}))]$ which is equal to $\text{Var}[g(X_1)]$ and converges to $\frac{[q'(\theta_0)]^2}{i(\theta_0)}$, yields the result of the Lemma. \square

We are now ready to give an upper bound on the distributional distance of interest using the Delta method from Theorem 5.1.

Theorem 5.2. *Let X_1, \dots, X_n be i.i.d. random variables with probability density (or mass) function $f(x_i|\theta)$ and let $Z \sim N(0, 1)$. Assume that the regularity conditions (R1)-(R4) are satisfied and that the MLE, $\hat{\theta}_n(\mathbf{X})$, exists and is unique. Furthermore let $q : \Theta \rightarrow \mathbb{R}$ be a one-to-one twice differentiable function with $q'(\theta) \neq 0 \forall \theta \in \Theta$ such that $q(\hat{\theta}_n(\mathbf{X})) = \frac{1}{n} \sum_{i=1}^n g(X_i)$, where the mapping $g : \mathbb{R} \rightarrow \mathbb{R}$ is such that $E[|g(X_1) - q(\theta_0)|^3] < \infty$. Then,*

(1) *if $q''(\cdot)$ is not uniformly bounded but for any $\theta_0 \in \Theta$ there exists $0 < \varepsilon = \varepsilon(\theta_0)$*

such that $\sup_{\theta: |\theta - \theta_0| < \varepsilon} |q''(\theta)| < \infty$, then

$$d_{bW} \left(\sqrt{ni(\theta_0)}(\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) \leq \frac{1}{\sqrt{n}} \left(2 + \frac{[i(\theta_0)]^{\frac{3}{2}}}{|q'(\theta_0)|^3} E[|g(X_1) - q(\theta_0)|^3] \right) + E[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2] \mathbb{1}_{\{q(\theta) \neq \theta\}} \left(\frac{2}{\varepsilon^2} + \frac{\sqrt{ni(\theta_0)}}{2|q'(\theta_0)|} \sup_{\theta: |\theta - \theta_0| < \varepsilon} |q''(\theta)| \right). \quad (5.5)$$

(2) *If $q''(\cdot)$ is uniformly bounded, with $|q''(\theta)| \leq B^*, \forall \theta \in \Theta$, then*

$$d_{bW} \left(\sqrt{ni(\theta_0)}(\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) \leq \frac{1}{\sqrt{n}} \left(2 + \frac{[i(\theta_0)]^{\frac{3}{2}}}{|q'(\theta_0)|^3} E[|g(X_1) - q(\theta_0)|^3] \right) + B^* \frac{\sqrt{ni(\theta_0)}}{2|q'(\theta_0)|} E[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2] \mathbb{1}_{\{q(\theta) \neq \theta\}}. \quad (5.6)$$

Proof. The asymptotic normality of the MLE is explicitly stated in Theorem 2.2.

Applying (5.2) to this result in combination with the requirement $q'(\theta_0) \neq 0$ yields

$$\frac{\sqrt{ni(\theta_0)}}{q'(\theta_0)} (q(\hat{\theta}_n(\mathbf{X})) - q(\theta_0)) \xrightarrow[n \rightarrow \infty]{d} Z,$$

with $q(\hat{\theta}_n(\mathbf{X})) = \frac{1}{n} \sum_{i=1}^n g(X_i)$. Using the triangle inequality,

$$\begin{aligned} & \left| \mathbb{E} \left[h \left(\sqrt{ni(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) \right] - \mathbb{E}[h(Z)] \right| \\ & \leq \left| \mathbb{E} \left[h \left(\frac{\sqrt{ni(\theta_0)}}{q'(\theta_0)} (q(\hat{\theta}_n(\mathbf{X})) - q(\theta_0)) \right) \right] - \mathbb{E}[h(Z)] \right| \end{aligned} \quad (5.7)$$

$$+ \left| \mathbb{E} \left[h \left(\sqrt{ni(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) - h \left(\frac{\sqrt{ni(\theta_0)}}{q'(\theta_0)} (q(\hat{\theta}_n(\mathbf{X})) - q(\theta_0)) \right) \right] \right|. \quad (5.8)$$

To find an upper bound for (5.7), some simple rewriting yields

$$\begin{aligned} \frac{\sqrt{ni(\theta_0)}}{q'(\theta_0)} (q(\hat{\theta}_n(\mathbf{X})) - q(\theta_0)) &= \frac{\sqrt{ni(\theta_0)}}{q'(\theta_0)} \left(\frac{1}{n} \sum_{i=1}^n g(X_i) - q(\theta_0) \right) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \left\{ \frac{\sqrt{i(\theta_0)}}{q'(\theta_0)} (g(X_i) - q(\theta_0)) \right\} = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i, \end{aligned}$$

where $Y_i = \frac{\sqrt{i(\theta_0)}}{q'(\theta_0)} (g(X_i) - q(\theta_0))$, $i = 1, \dots, n$ which are i.i.d. random variables. Using the result of Lemma 5.1 we deduce that $\mathbb{E}[Y_i] = 0$ and $\text{Var}[Y_i] = 1$, as Lemma 2.2 requires, which for $\sigma^2 = 1$ gives

$$d_{bW} \left(\frac{\sqrt{ni(\theta_0)}}{q'(\theta_0)} (q(\hat{\theta}_n(\mathbf{X})) - q(\theta_0)), Z \right) \leq \frac{1}{\sqrt{n}} \left(2 + \frac{[i(\theta_0)]^{\frac{3}{2}}}{|q'(\theta_0)|^3} \mathbb{E} \left[|g(X_1) - q(\theta_0)|^3 \right] \right). \quad (5.9)$$

For an upper bound on (5.8), since the case $q(\theta) = \theta$ is obvious, assume from now on that $q(\theta) \neq \theta$. We denote by

$$A := A(h, q, \theta_0, \mathbf{X}) := h \left(\sqrt{ni(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) - h \left(\frac{\sqrt{ni(\theta_0)}}{q'(\theta_0)} (q(\hat{\theta}_n(\mathbf{X})) - q(\theta_0)) \right).$$

To find an upper bound for $|E(A)|$, the approach is similar to the one developed in the proof of Theorem 3.1.

Case 1: $\forall \theta \in \Theta$ there exists $0 < \varepsilon = \varepsilon(\theta_0)$ such that $\sup_{\theta: |\theta - \theta_0| < \varepsilon} |q''(\theta)| < \infty$. Then, using the law of total expectation related to conditioning on $|\hat{\theta}_n(\mathbf{X}) - \theta_0| \geq \varepsilon$ or $|\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon$ yields

$$\begin{aligned} |E[A]| &\leq E[|A|] = E[|A| \mid |\hat{\theta}_n(\mathbf{X}) - \theta_0| \geq \varepsilon] \mathbb{P}(|\hat{\theta}_n(\mathbf{X}) - \theta_0| \geq \varepsilon) \\ &\quad + E[|A| \mid |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon] \mathbb{P}(|\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon). \end{aligned}$$

Markov's inequality and the elementary results of $\mathbb{P}(|\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon) \leq 1$ and $|A| \leq 2\|h\|$ further yield

$$|E(A)| \leq 2\|h\| \frac{E[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2]}{\varepsilon^2} + E[|A| \mid |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon]. \quad (5.10)$$

We now focus on the conditional expectation on the right-hand side of (5.10). A second-order Taylor expansion of $q(\hat{\theta}_n(x))$ about θ_0 gives

$$q(\hat{\theta}_n(x)) = q(\theta_0) + (\hat{\theta}_n(x) - \theta_0) q'(\theta_0) + \frac{1}{2} (\hat{\theta}_n(x) - \theta_0)^2 q''(\theta^*), \quad (5.11)$$

for θ^* between $\hat{\theta}_n(x)$ and θ_0 . Since $q'(\theta) \neq 0 \forall \theta \in \Theta$ is assumed, we can multiply both sides in (5.11) with $\frac{\sqrt{ni(\theta_0)}}{q'(\theta_0)}$. Rearranging the terms,

$$\begin{aligned} \frac{\sqrt{ni(\theta_0)} (q(\hat{\theta}_n(x)) - q(\theta_0))}{q'(\theta_0)} &= \sqrt{ni(\theta_0)} (\hat{\theta}_n(x) - \theta_0) \\ &\quad + \frac{\sqrt{ni(\theta_0)}}{2q'(\theta_0)} q''(\theta^*) (\hat{\theta}_n(x) - \theta_0)^2. \end{aligned}$$

The above result along with another first-order Taylor expansion gives

$$\begin{aligned} & h\left(\sqrt{ni(\theta_0)}(\hat{\theta}_n(x) - \theta_0)\right) - h\left(\frac{\sqrt{ni(\theta_0)}}{q'(\theta_0)}(q(\hat{\theta}_n(x)) - q(\theta_0))\right) \\ &= -\frac{\sqrt{ni(\theta_0)}}{2q'(\theta_0)}q''(\theta^*)(\hat{\theta}_n(x) - \theta_0)^2 h'(t(x)), \end{aligned} \quad (5.12)$$

where $t(x)$ is between $\sqrt{ni(\theta_0)}(\hat{\theta}_n(x) - \theta_0)$ and $\frac{\sqrt{ni(\theta_0)}}{q'(\theta_0)}(q(\hat{\theta}_n(x)) - q(\theta_0))$. Equality (5.12) combined with Lemma 3.1 related to conditional expectations yield

$$\begin{aligned} & \mathbb{E}[|A| \mid |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon] \\ &= \mathbb{E}\left[\left| -\frac{\sqrt{ni(\theta_0)}}{2q'(\theta_0)}q''(\theta^*)(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 h'(t(\mathbf{X})) \right| \mid |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon\right] \\ &\leq \frac{\|h'\|\sqrt{ni(\theta_0)}}{2|q'(\theta_0)|} \mathbb{E}\left[|q''(\theta^*)|(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \mid |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon\right] \\ &\leq \frac{\|h'\|\sqrt{ni(\theta_0)}}{2|q'(\theta_0)|} \sup_{\theta: |\theta - \theta_0| < \varepsilon} |q''(\theta)| \mathbb{E}\left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \mid |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon\right] \\ &\leq \frac{\|h'\|\sqrt{ni(\theta_0)}}{2|q'(\theta_0)|} \sup_{\theta: |\theta - \theta_0| < \varepsilon} |q''(\theta)| \mathbb{E}\left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2\right]. \end{aligned} \quad (5.13)$$

Combining the bounds in (5.9), (5.10) and (5.13) we get, for $h \in H_{bW}$ as in (2.6), the result in (5.5).

Case 2: $q''(\cdot)$ is uniformly bounded with $|q''(\theta)| \leq B^*$, $\forall \theta \in \Theta$. In this case, we do not need to take conditional expectations. The result in (5.12) gives

$$\begin{aligned} \mathbb{E}[|A|] &\leq \frac{\|h'\|\sqrt{ni(\theta_0)}}{2|q'(\theta_0)|} \mathbb{E}\left[|q''(\theta^*)|(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2\right] \\ &\leq B^* \frac{\|h'\|\sqrt{ni(\theta_0)}}{2|q'(\theta_0)|} \mathbb{E}\left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2\right]. \end{aligned}$$

The above result and the bound in (5.9) give for $h \in H_{bW}$ as in (2.6) the bound in (5.6) for the case where $q''(\cdot)$ is uniformly bounded in θ . \square

Remark 5.1. (1) *The convergence of the second and third terms in (5.5) is governed by the asymptotic behaviour of $E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]$, whose rate of convergence is $\mathcal{O} \left(\frac{1}{n} \right)$; see Remark 3.1.*

(2) *Note that if $q''(\cdot)$ is uniformly bounded, then $0 < \varepsilon = \varepsilon(\theta_0)$ is not needed leading to a simpler bound; see (5.6).*

(3) *In the simplest possible situation where $\hat{\theta}_n(\mathbf{X})$ is already a sum of i.i.d. terms, $q(x) = x$ and hence the upper bounds in (5.5) and (5.6) simplify to*

$$d_{bW} \left(\sqrt{ni(\theta_0)}(\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) \leq \frac{1}{\sqrt{n}} \left(2 + [i(\theta_0)]^{\frac{3}{2}} E \left[|g(X_1) - \theta_0|^3 \right] \right),$$

where Lemma 2.2 has been directly used.

We now compare this bound to (3.7). The bound using Delta method is simpler. A drawback is that in order to apply the Delta method approach an analytic expression for the MLE is required which in addition needs to be written as a function of a sum of independent terms. On the other hand, the approach through the score function, as explained in Proposition 3.2 and in the proof of Theorem 3.1 holds whatever the form of the MLE is, or even when an explicit analytic expression of the MLE is not known; we showed that in Section 3.5. Let us now comment on the two bounds in (3.7) and (5.5) term by term.

- For the first term of the bounds, the different positioning of the expected Fisher information number is explained by the fact that we apply Lemma 2.2 to the standardised version of $g(X_1), g(X_2), \dots, g(X_n)$, which have variance $\frac{[q'(\theta_0)]^2}{i(\theta_0)}$ (note that $i(\theta_0)$ is in the denominator), while in Theorem 3.1 the result was obtained by applying the lemma after standardising the first derivative of the log-likelihood function for each random variable. The variance in that case is equal to $i(\theta_0)$;
- the second and the third terms of the bound in (5.5) vanish when $q(\theta) = \theta$. For $q(\theta) \neq \theta$ the second term is the same in both bounds, whereas the third term in

(5.5) reads

$$\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \frac{\sqrt{ni(\theta_0)}}{2|q'(\theta_0)|} \sup_{\theta: |\theta - \theta_0| < \varepsilon} |q''(\theta)| \quad (5.14)$$

and is to be compared to

$$\begin{aligned} & \frac{1}{\sqrt{i(\theta_0)}} \sqrt{\text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right)} \sqrt{\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} \\ & + \frac{1}{2\sqrt{ni(\theta_0)}} \left[\mathbb{E} \left(\left(\sum_{i=1}^n M(X_i) \right)^2 \middle| |\hat{\theta}_n(\mathbf{X}) - \theta_0| < \varepsilon \right) \right]^{\frac{1}{2}} \\ & \times \left[\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^4 \right] \right]^{\frac{1}{2}}. \end{aligned} \quad (5.15)$$

The second derivative, $q''(\theta)$ in (5.14) plays the role of $l^{(3)}(\theta; \mathbf{X})$, up to a difference: in a neighbourhood of θ_0 , $l^{(3)}(\theta; \mathbf{X})$ is bounded by $\sum_{i=1}^n M(X_i)$. Consequently, the second term of (5.15) has \sqrt{n} in its numerator, exactly as in (5.14). The distinct positioning of the information number $i(\theta_0)$ has the same reason as for the first term. The bound in (5.14) is clearly simpler (has one less term) and sharper than (5.15) at the level of moments of $\hat{\theta}_n(\mathbf{X}) - \theta_0$ since

$$\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \leq \left[\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^4 \right] \right]^{\frac{1}{2}}$$

by the Cauchy-Schwarz inequality.

Therefore, in the simple, but still widely applicable setting where $\hat{\theta}_n(\mathbf{X})$ is a function of a sum of i.i.d. terms, the bound using the Delta method is simpler than the one in Theorem 3.1 and also of the same order. As already mentioned, an obvious advantage of the approach followed in Chapter 3 is its wider applicability as it works for all MLE settings and when no closed-form expression of the MLE is available.

5.2 Comparison in exponential families

In this section, we give the bound using the Delta method in the case of single-parameter exponential family distributions. The Generalised-Gamma distribution is also discussed and an upper bound is given. To allow for a further comparison between the two bounds, the exponential distribution both in the canonical and the non-canonical form as in Section 3.3 is used.

5.2.1 Bounds for single-parameter exponential families

For the single-parameter exponential families framework, where the probability density (or mass) function of a random variable is as in (3.21), it is straightforward to see that the function $q(\cdot)$ is defined in general.

Proposition 5.1. *Suppose X_1, \dots, X_n are i.i.d. random variables with probability density (or mass) function that can be expressed in the form of (3.21). Assume that (Ass.Ex.1)-(Ass.Ex.3) and (A1)-(A5) (Section 3.3) are satisfied. Then, with $D(\theta) = \frac{A'(\theta)}{k'(\theta)}$ as in (A1), we obtain that $q(\theta_0) = D(\theta_0)$ for $q : \Theta \rightarrow \mathbb{R}$ as in Theorem 5.2.*

Proof. Using (3.21), we have that

$$l'(\theta; \mathbf{x}) = k'(\theta) \sum_{i=1}^n T(x_i) - nA'(\theta) = 0 \iff D(\theta) = \frac{1}{n} \sum_{i=1}^n T(x_i).$$

Since $D(\cdot)$ is invertible (see condition (A1)), $\hat{\theta}_n(\mathbf{X}) = D^{-1} \left(\frac{1}{n} \sum_{i=1}^n T(X_i) \right)$. The claim readily follows. \square

The broad single-parameter exponential families satisfy (5.1). Consequently, Theorem 5.2 can be applied to (3.21), resulting in

Corollary 5.1. *Let X_1, \dots, X_n be i.i.d. random variables with the probability density (or mass) function of a single-parameter exponential family. Assume that (Ass.Ex.1)-(Ass.Ex.3) and (A1)-(A5) are satisfied. With $Z \sim N(0, 1)$, we obtain*

(1) *if $D''(\cdot)$ is not uniformly bounded but for any $\theta_0 \in \Theta$ there exists $0 < \varepsilon = \varepsilon(\theta_0)$*

such that $\sup_{\theta: |\theta - \theta_0| < \varepsilon} |D''(\theta)| < \infty$, then

$$\begin{aligned} d_{bW} \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) &\leq \frac{1}{\sqrt{n}} \left(2 + \frac{|k'(\theta_0)|^3 \mathbb{E} \left[|T(X_1) - D(\theta_0)|^3 \right]}{|A''(\theta_0) - k''(\theta_0)D(\theta_0)|^{\frac{3}{2}}} \right) \\ &+ \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \mathbb{1}_{\{D(\theta) \neq \theta\}} \left(\frac{2}{\varepsilon^2} \right. \\ &\quad \left. + \frac{\sqrt{n} |k'(\theta_0)|}{2 \sqrt{|A''(\theta_0) - k''(\theta_0)D(\theta_0)|}} \sup_{\theta: |\theta - \theta_0| < \varepsilon} |D''(\theta)| \right). \end{aligned} \quad (5.16)$$

(2) *If $D''(\cdot)$ is uniformly bounded, with $|D''(\theta)| \leq B^*, \forall \theta \in \Theta$, then*

$$\begin{aligned} d_{bW} \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) &\leq \frac{1}{\sqrt{n}} \left(2 + \frac{|k'(\theta_0)|^3 \mathbb{E} \left[|T(X_1) - D(\theta_0)|^3 \right]}{|A''(\theta_0) - k''(\theta_0)D(\theta_0)|^{\frac{3}{2}}} \right) \\ &+ B^* \frac{\sqrt{n} |k'(\theta_0)|}{2 \sqrt{|A''(\theta_0) - k''(\theta_0)D(\theta_0)|}} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \mathbb{1}_{\{D(\theta) \neq \theta\}}. \end{aligned} \quad (5.17)$$

Proof. We readily have

$$\begin{aligned} i(\theta_0) &= \mathbb{E} \left[-\frac{d^2}{d\theta^2} \log f(X_1 | \theta_0) \right] = A''(\theta_0) - k''(\theta_0) \mathbb{E}[T(X_1)] \\ &= \frac{A''(\theta_0)k'(\theta_0) - k''(\theta_0)A'(\theta_0)}{k'(\theta_0)} \\ q'(\theta_0) &= \frac{A''(\theta_0)k'(\theta_0) - k''(\theta_0)A'(\theta_0)}{[k'(\theta_0)]^2}. \end{aligned}$$

Therefore,

$$\frac{\sqrt{i(\theta_0)}}{|q'(\theta_0)|} = \frac{|k'(\theta_0)|^{\frac{3}{2}}}{\sqrt{|A''(\theta_0)k'(\theta_0) - k''(\theta_0)A'(\theta_0)|}} = \frac{|k'(\theta_0)|}{\sqrt{|A''(\theta_0) - k''(\theta_0)D(\theta_0)|}}.$$

This result, along with the fact that $g(x) = T(x)$ and $q(\theta) = D(\theta)$ by definition, allows to deduce the upper bounds from Theorem 5.2. \square

Remark 5.2. *It is particularly interesting to spell out this bound in the canonical case ($k(\theta) = \theta$). Since $k'(\theta) = 1$, then $D(\theta) = A'(\theta)$ and thus*

$$\begin{aligned} d_{bW} \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) &\leq \frac{1}{\sqrt{n}} \left(2 + \frac{\mathbb{E} \left[|T(X_1) - A'(\theta_0)|^3 \right]}{|A''(\theta_0)|^{\frac{3}{2}}} \right) \\ &+ \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \mathbb{1}_{\{A'(\theta) \neq \theta\}} \left(\frac{2}{\varepsilon^2} + \frac{\sqrt{n}}{2\sqrt{|A''(\theta_0)|}} \sup_{\theta: |\theta - \theta_0| < \varepsilon} |A^{(3)}(\theta)| \right), \end{aligned}$$

which is exactly the same as the bound in (3.24). If $|A^{(3)}(\theta)|$ is uniformly bounded by a constant B^* then

$$\begin{aligned} d_{bW} \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) &\leq \frac{1}{\sqrt{n}} \left(2 + \frac{\mathbb{E} \left[|T(X_1) - A'(\theta_0)|^3 \right]}{|A''(\theta_0)|^{\frac{3}{2}}} \right) \\ &+ B^* \frac{\sqrt{n}}{2\sqrt{|A''(\theta_0)|}} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \mathbb{1}_{\{A'(\theta) \neq \theta\}}. \end{aligned}$$

We now proceed to give the upper bound for the generalised Gamma distribution, which is a generalisation of the Gamma distribution with three parameters $\theta, p, d > 0$ and from now is denoted by $GG(\theta, d, p)$. With $\Gamma(\cdot)$ denoting as usual the gamma function, the probability density function for $x > 0$ is

$$f(x|\theta) = \frac{px^{d-1} \exp \left\{ -\left(\frac{x}{\theta}\right)^p \right\}}{\theta^d \Gamma\left(\frac{d}{p}\right)}. \quad (5.18)$$

The generalised gamma distribution includes many other known distributions as special cases. If $d = p$ then the generalized gamma distribution becomes the Weibull distribu-

tion. Furthermore, if $p = 1$, we get the gamma distribution (if in addition $d = 1$, then the exponential distribution with mean θ is obtained). When d and p are known and the scale θ_0 is the unknown parameter of interest, the MLE for θ_0 exists, it is unique and $\hat{\theta}_n(\mathbf{X}) = \left(\frac{p}{nd} \sum_{i=1}^n x_i^p\right)^{\frac{1}{p}}$.

Corollary 5.2. *Let X_1, \dots, X_n be i.i.d. random variables from the generalized gamma $GG(\theta, d, p)$ distribution, where the shape parameters $d, p > 0$ are known and the scale parameter $\theta_0 > 0$ is the unknown parameter of interest. Then,*

$$\begin{aligned} d_{bW} \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) &\leq \frac{1}{\sqrt{n}} \left(2 + \left(3 + 6 \frac{p}{d} \right)^{\frac{3}{4}} \right) \\ &+ \left(1 - 2 \left(\frac{p}{nd} \right)^{\frac{1}{p}} \frac{\Gamma\left(\frac{nd+1}{p}\right)}{\Gamma\left(\frac{nd}{p}\right)} + \left(\frac{p}{nd} \right)^{\frac{2}{p}} \frac{\Gamma\left(\frac{nd+2}{p}\right)}{\Gamma\left(\frac{nd}{p}\right)} \right) \mathbb{1}_{\{\{d \neq 1\} \cup \{p \neq 1\}\}} \\ &\times \left[8 + \frac{\sqrt{ndp} |p-1|}{2} \left(\frac{1}{2^{p-2}} \mathbb{1}_{\{p < 2\}} + \left(\frac{3}{2} \right)^{p-2} \mathbb{1}_{\{p \geq 2\}} \right) \right]. \end{aligned} \quad (5.19)$$

Proof. Rearranging the terms in (5.18) yields

$$f(x|\theta) = \exp \left\{ - \left(\frac{x}{\theta} \right)^p + \log p - d \log \theta + (d-1) \log x - \log \left(\Gamma \left(\frac{d}{p} \right) \right) \right\}$$

and in the terminology of single-parameter exponential families, $B = (0, \infty)$, $\Theta = (0, \infty)$, $T(x) = x^p$, $k(\theta) = -\frac{1}{\theta^p}$, $A(\theta) = d \log \theta$ and $S(x) = \log p + (d-1) \log x - \log \left(\Gamma \left(\frac{d}{p} \right) \right)$.

Simple steps yield

$$l(\theta_0; \mathbf{x}) = -\frac{1}{\theta_0^p} \sum_{i=1}^n x_i^p + n \log p - nd \log \theta_0 + (d-1) \log \left(\prod_{i=1}^n x_i \right) - n \log \left(\Gamma \left(\frac{d}{p} \right) \right)$$

and

$$l'(\theta_0; \mathbf{x}) = \frac{p}{\theta_0^{p+1}} \sum_{i=1}^n x_i^p - n \frac{d}{\theta_0}.$$

Thus, $l'(\theta; \mathbf{x}) = 0$ if and only if $\theta = \hat{\theta}_n(\mathbf{x}) = \left(\frac{p}{nd} \sum_{i=1}^n x_i^p\right)^{\frac{1}{p}}$. Also, $l''(\hat{\theta}_n(\mathbf{x}); \mathbf{x}) = -n \frac{pd}{[\hat{\theta}_n(\mathbf{x})]^2} < 0$, which shows that the MLE exists and is unique. To obtain the bound,

the results of Corollary 5.1 are employed. For the generalised gamma distribution, $D(\theta_0) = q(\theta_0) = \frac{d}{p}\theta_0^p$ and thus

$$\mathbb{E} \left[|T(X) - D(\theta_0)|^3 \right] = \mathbb{E} \left[\left| X^p - \frac{d}{p}\theta_0^p \right|^3 \right].$$

This third absolute moment is complicated to calculate. Using the change-of-variable technique, it can be easily shown that if $X \sim \text{GG}(\theta_0, d, p)$, then $X^p \sim \text{Gamma}\left(\frac{d}{p}, \frac{1}{\theta_0^p}\right)$. Therefore, Hölder's inequality gives

$$\begin{aligned} \mathbb{E} \left[\left| X^p - \frac{d}{p}\theta_0^p \right|^3 \right] &\leq \left[\mathbb{E} \left[\left(X^p - \frac{d}{p}\theta_0^p \right)^4 \right] \right]^{\frac{3}{4}} \\ &= \left[\mathbb{E}(X^{4p}) + \left(\frac{d\theta_0^p}{p} \right)^4 + 6 \left(\frac{d\theta_0^p}{p} \right)^2 \mathbb{E}(X^{2p}) - 4 \left(\frac{d\theta_0^p}{p} \right)^3 \mathbb{E}(X^p) - 4 \frac{d}{p}\theta_0^p \mathbb{E}(X^{3p}) \right]^{\frac{3}{4}} \\ &= \left[\theta_0^{4p} \frac{d}{p} \left(6 + 3 \frac{d}{p} \right) \right]^{\frac{3}{4}} = \theta_0^{3p} \left(\frac{d}{p} \right)^{\frac{3}{4}} \left(6 + 3 \frac{d}{p} \right)^{\frac{3}{4}}. \end{aligned} \quad (5.20)$$

In addition, since $k(\theta) = -\frac{1}{\theta^p}$, $A(\theta) = d \log \theta$ and $D(\theta) = \frac{d}{p}\theta^p$

$$\begin{aligned} \frac{|k'(\theta_0)|}{\sqrt{|A''(\theta_0) - k''(\theta_0)D(\theta_0)|}} &= \frac{p\theta_0^{-(p+1)}}{\sqrt{\left| -\frac{d}{\theta_0^2} + p(p+1)\theta_0^{-(p+2)}\frac{d}{p}\theta_0^p \right|}} \\ &= \frac{p\theta_0^{-(p+1)}}{\sqrt{\left| -\frac{d}{\theta_0^2} + \frac{d(p+1)}{\theta_0^2} \right|}} = \frac{p\theta_0^{-(p+1)}}{\sqrt{\frac{dp}{\theta_0^2}}} = \frac{\sqrt{p}}{\sqrt{d}\theta_0^p}. \end{aligned} \quad (5.21)$$

Using that $X_i^p \sim \text{Gamma}\left(\frac{d}{p}, \frac{1}{\theta_0^p}\right) \Rightarrow \sum_{i=1}^n X_i^p \sim \text{Gamma}\left(n\frac{d}{p}, \frac{1}{\theta_0^p}\right)$, we get

$$\begin{aligned} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] &= \left(\frac{p}{nd} \right)^{\frac{2}{p}} \mathbb{E} \left[\left(\sum_{i=1}^n X_i^p \right)^{\frac{2}{p}} \right] + \theta_0^2 - 2 \left(\frac{p}{nd} \right)^{\frac{1}{p}} \theta_0 \mathbb{E} \left[\left(\sum_{i=1}^n X_i^p \right)^{\frac{1}{p}} \right] \\ &= \theta_0^2 \left(\frac{p}{nd} \right)^{\frac{2}{p}} \frac{\Gamma\left(\frac{nd+2}{p}\right)}{\Gamma\left(\frac{nd}{p}\right)} + \theta_0^2 - 2\theta_0^2 \left(\frac{p}{nd} \right)^{\frac{1}{p}} \frac{\Gamma\left(\frac{nd+1}{p}\right)}{\Gamma\left(\frac{nd}{p}\right)} \end{aligned}$$

$$= \theta_0^2 \left(1 - 2 \left(\frac{p}{nd} \right)^{\frac{1}{p}} \frac{\Gamma\left(\frac{nd+1}{p}\right)}{\Gamma\left(\frac{nd}{p}\right)} + \left(\frac{p}{nd} \right)^{\frac{2}{p}} \frac{\Gamma\left(\frac{nd+2}{p}\right)}{\Gamma\left(\frac{nd}{p}\right)} \right). \quad (5.22)$$

We now check whether $D''(\theta)$ is uniformly bounded. Simple steps yield

$$|D''(\theta)| = |d(p-1)\theta^{p-2}|,$$

which is not bounded for $\theta \in (0, \infty)$. However for any $\theta_0 \in \Theta$ there exists $0 < \varepsilon = \varepsilon(\theta_0)$

such that $\sup_{\theta: |\theta - \theta_0| < \varepsilon} |D''(\theta)| < \infty$. Therefore, this example falls under the first case of

Corollary 5.1. Below, we see that the supremum depends on the value of p :

$$\begin{aligned} \sup_{\theta: |\theta - \theta_0| < \varepsilon} |d(p-1)\theta^{p-2}| &= d|p-1| \sup_{\theta: |\theta - \theta_0| < \varepsilon} |\theta^{p-2}| \\ &= d|p-1| \begin{cases} (\theta_0 - \varepsilon)^{p-2}, & \text{if } 0 < p < 2 \\ (\theta_0 + \varepsilon)^{p-2}, & \text{if } p \geq 2. \end{cases} \end{aligned} \quad (5.23)$$

We now explain how $0 < \varepsilon$ has been chosen such that $|\theta^{p-2}|$ is bounded for all $\theta \in (\theta_0 - \varepsilon, \theta_0 + \varepsilon) \cap \Theta$. We have that $|\theta^{p-2}|$ in (5.23) is bounded for $p < 2$ when $0 < \varepsilon < \theta_0$. In both cases of (5.23), $\sup_{\theta: |\theta - \theta_0| < \varepsilon} |d(p-1)\theta^{p-2}|$ gets smaller as ε approaches zero. On the other hand, the second term of the bound in (5.16) becomes smaller as ε gets larger. An optimisation process with respect to ε is quite complicated and we choose $\varepsilon = \frac{\theta_0}{2}$, which is far away from both zero and θ_0 . Using this value of ε and applying the results of (5.20), (5.21), (5.22) and (5.23) to (5.16), we obtain the result in (5.19). \square

Remark 5.3. (1) The bound in (5.19) is $\mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$. This is not obvious due to the term $\left(1 - 2 \left(\frac{p}{nd}\right)^{\frac{1}{p}} \frac{\Gamma\left(\frac{nd+1}{p}\right)}{\Gamma\left(\frac{nd}{p}\right)} + \left(\frac{p}{nd}\right)^{\frac{2}{p}} \frac{\Gamma\left(\frac{nd+2}{p}\right)}{\Gamma\left(\frac{nd}{p}\right)}\right)$. Using the following Taylor expansion for a

ratio of Gamma functions (see Tricomi and Erdélyi (1951))

$$\frac{\Gamma(z+a)}{\Gamma(z+b)} = z^{a-b} \left(1 + \frac{(a-b)(a+b-1)}{2z} + \mathcal{O}(|z|^{-2}) \right)$$

for large z (here, $\frac{nd}{p}$) and bounded a and b , we can see that this term is of order $\frac{1}{n}$, leading to the overall order of $\frac{1}{\sqrt{n}}$.

(2) The indicator function in (5.19) comes from the fact that $q(\theta) = \theta \Leftrightarrow d, p = 1$.

5.2.2 Bounds for the exponential distribution

In this subsection, we consider as an example one special case of the generalised gamma distribution; the exponential distribution. This choice of distribution allows us to compare the results of this section with those in Section 3.3. First, the canonical form of the distribution is treated and then we will change the parameterisation.

The canonical case: $\text{Exp}(\theta)$

We start with X_1, \dots, X_n exponentially distributed i.i.d. random variables with scale parameter $\theta > 0$. The probability density function is given in (3.25). With this in hand, it is easily noticed that $D(\theta_0) = q(\theta_0) = \frac{A'(\theta_0)}{k'(\theta_0)} = -\frac{1}{\theta_0}$. Equation (5.21) shows that for $p = d = 1$

$$\frac{|k'(\theta_0)|}{\sqrt{|A''(\theta_0) - k''(\theta_0)D(\theta_0)|}} = \theta_0.$$

The rest of the quantities in the bound of Corollary 5.1 have already been calculated in Section 3.3 and for $\varepsilon = \frac{\theta_0}{2}$, we obtain

$$\begin{aligned} d_{bw} \left(\sqrt{n i(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) &< \frac{4.41456}{\sqrt{n}} + 8 \frac{(n+2)}{(n-1)(n-2)} \\ &+ 8 \frac{\sqrt{n}(n+2)}{(n-1)(n-2)} = \mathcal{O} \left(\frac{1}{\sqrt{n}} \right). \end{aligned} \quad (5.24)$$

Remark 5.4. (1) The bound is $\mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$.

(2) The result in (5.24) is the same as in (3.26), which is not a surprise in view of Remark 5.2.

The non-canonical case: $\text{Exp}\left(\frac{1}{\theta}\right)$

We now examine again the case where X_1, \dots, X_n are i.i.d. random variables from $\text{Exp}\left(\frac{1}{\theta}\right)$. The probability density function is found in (3.29). It is straightforward that $D(\theta) = q(\theta) = \theta$ and $|D''(\theta)| = 0$, thus uniformly bounded, meaning that we are under the second case of Corollary 5.1 and we will use the bound in (5.17) which does not require to define $\varepsilon = \varepsilon(\theta_0) > 0$. The fact that $D(\theta) = \theta$ makes the last term of the bound in (5.17) vanish. Having also that $E|T(X) - D(\theta_0)|^3 = E|X - \theta_0|^3 < 2.41456\theta_0^3$ and

$$\frac{|k'(\theta_0)|}{\sqrt{|A''(\theta_0) - k''(\theta_0)D(\theta_0)|}} = \frac{1}{\theta_0}$$

yields

$$d_{bW}\left(\sqrt{n}i(\theta_0)(\hat{\theta}_n(\mathbf{X}) - \theta_0), Z\right) < \frac{4.41456}{\sqrt{n}}. \quad (5.25)$$

Remark 5.5. (1) The order of the bound in terms of the sample size is, as expected, $\frac{1}{\sqrt{n}}$.

(2) For the same distribution, Corollary 3.3 gives

$$\frac{4.41456}{\sqrt{n}} + \frac{8}{n} + \frac{2}{\sqrt{n}} + \frac{1}{\sqrt{n}} \left(80 \left[3 \left(\frac{2}{n} + 1 \right) \right]^{\frac{1}{2}} \right),$$

showing that the bound using the Delta method is an improvement of this bound. In Remark 3.4, it is also made clear that because in this specific example the MLE is a sum of random variables then someone can use directly Lemma 2.2 and bound the distributional distance of interest, which also gives the bound in (5.25).

Chapter 6

Locally dependent random variables

In this chapter, we assess the normal approximation of the MLE under the presence of a local dependence structure between the random variables. Stein's method for normal approximation is employed, but this time we use the results of Lemma 2.3 since they are appropriate due to the structure of the random variables. We work under assumptions that allow us to express the bound in terms of the Wasserstein distance which is defined in (2.5). A general approach is first developed to get upper bounds on the Wasserstein distance between the distribution of the suitably scaled MLE and the standard normal distribution. An example of normally distributed locally dependent random variables serves as an illustration of these results. After that, the asymptotic normality result is given under rather restrictive but sufficient regularity conditions and a new bound is obtained. Bounds are also derived when the MLE can not be expressed in a closed form.

The notion of an m -dependent sequence of random variables $\{X_i, i \in \mathbb{N}\}$ has already been introduced in Section 2.3. The sets $\{X_j, j \leq i\}$ and $\{X_j, j > i + m\}$ are independent. The joint density function of X_1, X_2, \dots, X_n is

$$\begin{aligned}
f(\mathbf{x}|\boldsymbol{\theta}) &= L(\boldsymbol{\theta}; \mathbf{x}) = f(x_1; \boldsymbol{\theta})f(x_2|x_1; \boldsymbol{\theta}) \dots f(x_n|x_{n-1}, \dots, x_{n-m}; \boldsymbol{\theta}) \\
&= f(x_1; \boldsymbol{\theta}) \prod_{i=2}^n f(x_i|x_{i-1}, \dots, x_{m_i^*}; \boldsymbol{\theta}),
\end{aligned} \tag{6.1}$$

where $m_i^* = \max\{i - m, 1\}$.

In Section 6.1 we explain the process of finding an upper bound on the Wasserstein distance between the distribution of the suitably scaled MLE and the standard normal distribution, that holds in a general framework. The quantity we are interested in is split into two terms with the one being bounded using Stein's method and the other using alternative techniques based mainly on Taylor expansions. These results are applied in Section 6.2 to the example of locally dependent normally distributed random variables. In Section 6.3, we discuss on a set of sufficient regularity conditions for the asymptotic normality of the MLE to hold. Under these, rather strict, conditions a different bound is obtained for the Wasserstein distance between the distribution of the MLE and its limiting normal distribution. Under further assumptions, Section 6.4 gives upper bounds for situations where the MLE can not be defined analytically.

6.1 A general approach

Under the assumption of the parameter space being an open and connected interval, Proposition 2.1 is sufficient for the existence and uniqueness of $\hat{\boldsymbol{\theta}}_n(\mathbf{X})$. Our purpose is to obtain an upper bound on the Wasserstein distance between the distribution of an appropriately scaled MLE and the standard normal distribution. The results of Lemma 2.3 will be applied to a sequence $\{\xi_i, i = 1, 2, \dots, n\}$ of $2m$ -dependent random variables. We denote by

$$\begin{aligned}
M_{1j} &:= \max\{1, j - 2m\} & M_{2j} &:= \min\{n, j + 2m\} \\
K_{1j} &:= \max\{1, j - 4m\} & K_{2j} &:= \min\{n, j + 4m\}.
\end{aligned}$$

In addition, the dependency neighbourhoods, A_j and B_j , as defined in (LD) on p.20 are

$$\begin{aligned} A_j &= \{M_{1j}, M_{1j} + 1, \dots, M_{2j} - 1, M_{2j}\} \\ B_j &= \{K_{1j}, K_{1j} + 1, \dots, K_{2j} - 1, K_{2j}\}. \end{aligned} \quad (6.2)$$

Having that $\forall i \in \{1, 2, \dots, n\}$, $|A_i|$ and $|B_i|$ denote the number of elements in the sets A_i and B_i , respectively, then

$$|A_i| \leq 4m + 1, \quad |B_i| \leq 8m + 1.$$

Figure 6.1 below is a graphical representation of our notation related to the $2m$ -dependence structure.

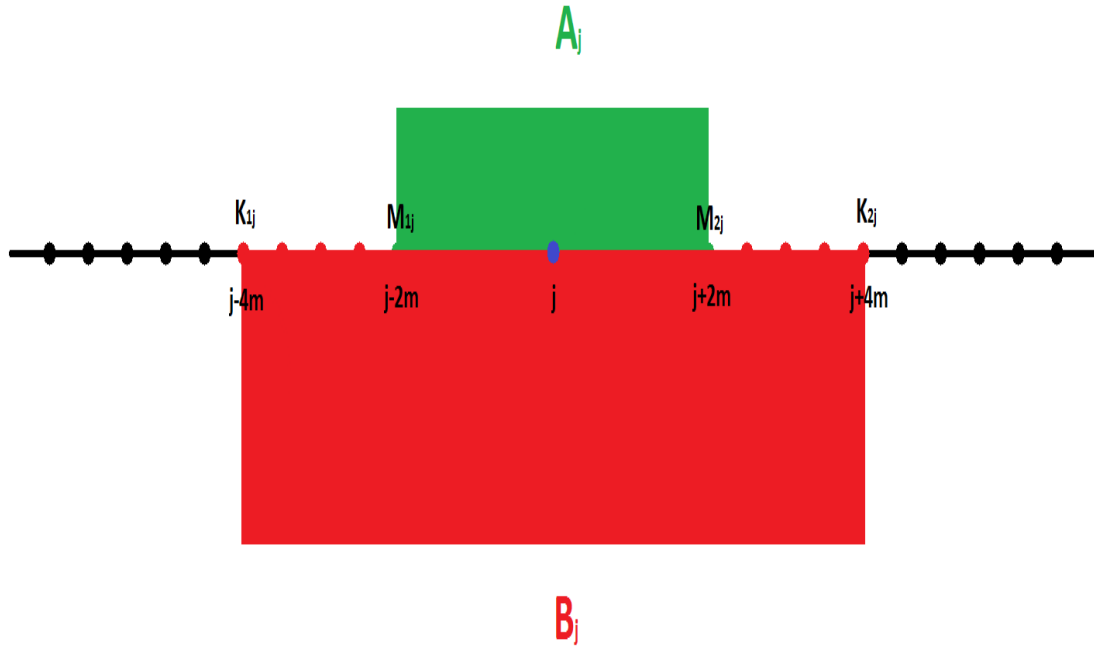


Fig. 6.1: Structure of a $2m$ -dependent sequence.

In this section, for $\{X_i, i = 1, 2, \dots, n\}$ being an m -dependent sequence, we work under the following assumptions:

(A.D.1) The log-likelihood function is three times differentiable with uniformly bounded third derivative in $\theta \in \Theta$, $(x_1, x_2, \dots, x_n) \in S$. We denote the supremum by

$$S_d(n) := \sup_{\substack{\theta \in \Theta \\ x \in S}} |l^{(3)}(\theta; x)| < \infty. \quad (6.3)$$

(A.D.2) $E \left[\frac{d}{d\theta} \log f(X_1 | \theta) \right] = E \left[\frac{d}{d\theta} \log f(X_i | X_{i-1}, \dots, X_{i-m}; \theta) \right] = 0$, for $i = 2, 3, \dots, n$.

(A.D.3) With θ_0 , as usual, denoting the true value of the unknown parameter,

$$\sqrt{n} E [\hat{\theta}_n(\mathbf{X}) - \theta_0] \xrightarrow{n \rightarrow \infty} 0. \quad (6.4)$$

(A.D.4) The limit of the reciprocal of $n \text{Var}(\hat{\theta}_n(\mathbf{X}))$ exists and from now on, unless otherwise stated,

$$0 < i_2(\theta_0) = \lim_{n \rightarrow \infty} \frac{1}{n \text{Var}(\hat{\theta}_n(\mathbf{X}))}. \quad (6.5)$$

The assumption (A.D.1) is stronger than (R3) of Section 2.2 which was used in the case of independent random variables. Using uniform boundedness of the third derivative of the log-likelihood function in (A.D.1) allows us to get bounds on the Wasserstein distance related to the MLE. The following theorem gives the bound.

Theorem 6.1. *Let $\{X_i, i = 1, 2, \dots, n\}$ be an m -dependent sequence of identically distributed random variables with probability density (or mass) function $f(x_i | x_{i-1}, \dots, x_{i-m}; \theta)$, where $\theta \in \Theta$ and $(x_1, x_2, \dots, x_n) \in S$, the support of the joint probability density (or mass) function. Assume that $\hat{\theta}_n(\mathbf{X})$ exists and is unique. In addition, assume that (A.D.1)-(A.D.4) hold and that $\text{Var}[l'(\theta_0; \mathbf{X})] > 0$. Let*

$$\alpha := \alpha(\theta_0, n) := \sqrt{\frac{\text{Var}(l'(\theta_0; \mathbf{X}))}{\text{Var}(\hat{\theta}_n(\mathbf{X}))}}, \quad (6.6)$$

which is assumed to be finite and not equal to zero. In addition, we denote by

$$\xi_1 = \frac{d}{d\theta} \log f(X_1 | \theta) \Big|_{\theta=\theta_0} \sqrt{\frac{n}{\text{Var}(l'(\theta_0; \mathbf{X}))}}$$

and for $i = 2, 3, \dots, n$,

$$\xi_i = \frac{d}{d\theta} \log f(X_i | X_{i-1}, \dots, X_{i-m}; \theta) \Big|_{\theta=\theta_0} \sqrt{\frac{n}{\text{Var}(l'(\theta_0; \mathbf{X}))}}.$$

For $Z \sim N(0, 1)$,

$$\begin{aligned} & d_W \left(\sqrt{n i_2(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) \\ & \leq \frac{2}{n^{\frac{3}{2}}} \sum_{i=1}^n \sum_{j \in A_i} \sum_{k \in B_i} \left[E \left((\xi_i)^4 \right) E \left((\xi_j)^4 \right) E \left((\xi_k)^4 \right) \right]^{\frac{1}{4}} \\ & \quad + \frac{2}{n^{\frac{3}{2}}} \sum_{i=1}^n \sum_{j \in A_i} \sum_{k \in B_i} \left[E \left((\xi_i)^2 \right) E \left((\xi_j)^2 \right) E \left((\xi_k)^2 \right) \right]^{\frac{1}{2}} \\ & \quad + \frac{1}{n^{\frac{3}{2}}} \sum_{i=1}^n |A_i| \sum_{j \in A_i} \left[E \left((\xi_i)^2 \right) E \left((\xi_j)^4 \right) \right]^{\frac{1}{2}} \\ & \quad + \left| \frac{\sqrt{n i_2(\theta_0)} \text{Var}[l'(\theta_0; \mathbf{X})]}{\alpha} - 1 \right| \\ & \quad + \frac{S_d(n) \sqrt{n i_2(\theta_0)}}{2\alpha} E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \\ & \quad + \frac{\sqrt{n i_2(\theta_0)}}{\alpha} \sqrt{E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} \sqrt{E \left[(l''(\theta_0; \mathbf{X}) + \alpha)^2 \right]}. \end{aligned} \quad (6.7)$$

Proof. By the definition of the MLE and (A.D.1), $l'(\hat{\theta}_n(\mathbf{x}); \mathbf{x}) = 0$. With $R_1(\theta_0; \mathbf{X})$ as in (3.9), a second order Taylor expansion gives that

$$\begin{aligned} & (\hat{\theta}_n(\mathbf{X}) - \theta_0) l''(\theta_0; \mathbf{X}) = -l'(\theta_0; \mathbf{X}) - R_1(\theta_0; \mathbf{X}) \\ & \Rightarrow (\hat{\theta}_n(\mathbf{X}) - \theta_0) (\alpha + l''(\theta_0; \mathbf{X}) - \alpha) = -l'(\theta_0; \mathbf{X}) - R_1(\theta_0; \mathbf{X}) \\ & \Rightarrow -\alpha (\hat{\theta}_n(\mathbf{X}) - \theta_0) = -l'(\theta_0; \mathbf{X}) - R_1(\theta_0; \mathbf{X}) - (\hat{\theta}_n(\mathbf{X}) - \theta_0) (l''(\theta_0; \mathbf{X}) + \alpha) \end{aligned}$$

Multiplying both sides by $-\frac{\sqrt{ni_2(\theta_0)}}{\alpha}$,

$$\begin{aligned} \sqrt{ni_2(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) &= \frac{\sqrt{ni_2(\theta_0)}}{\alpha} \left[l'(\theta_0; \mathbf{X}) + R_1(\theta_0; \mathbf{X}) \right. \\ &\quad \left. + (\hat{\theta}_n(\mathbf{X}) - \theta_0) (l''(\theta_0; \mathbf{X}) + \alpha) \right]. \end{aligned} \quad (6.8)$$

Applying the triangle inequality,

$$\begin{aligned} &\left| \mathbb{E} \left[h \left(\sqrt{ni_2(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) \right] - \mathbb{E}[h(Z)] \right| \\ &\leq \left| \mathbb{E} \left[h \left(\frac{\sqrt{ni_2(\theta_0)} l'(\theta_0; \mathbf{X})}{\alpha} \right) \right] - \mathbb{E}[h(Z)] \right| \end{aligned} \quad (6.9)$$

$$+ \left| \mathbb{E} \left[h \left(\sqrt{ni_2(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) - h \left(\frac{\sqrt{ni_2(\theta_0)} l'(\theta_0; \mathbf{X})}{\alpha} \right) \right] \right|. \quad (6.10)$$

Step 1: Bound for (6.9). Let, for ease of presentation $l'(\theta_0; \mathbf{X}) = \sum_{i=1}^n \tilde{\xi}_i$, where

$$\begin{aligned} \tilde{\xi}_1 &= \frac{d}{d\theta} \log f(X_1 | \theta) \Big|_{\theta=\theta_0}, \\ \tilde{\xi}_i &= \frac{d}{d\theta} \log f(X_i | X_{i-1}, \dots, X_{i-m}; \theta) \Big|_{\theta=\theta_0} \text{ for } i = 2, 3, \dots, n. \end{aligned}$$

Assumption (A.D.2) ensures that $\tilde{\xi}_i, i = 1, 2, \dots, n$ have mean zero. Furthermore, for some function $g : \mathbb{R}^{m+1} \rightarrow \mathbb{R}$, it holds that $\tilde{\xi}_i = g(X_i, X_{i-1}, \dots, X_{i-m})$ and taking into account that $\{X_i, i = 1, 2, \dots, n\}$ is an m -dependent sequence, we conclude that $\{\tilde{\xi}_i, i = 1, 2, \dots, n\}$ forms a $2m$ -dependent sequence. Define now

$$W := \frac{l'(\theta_0; \mathbf{X})}{\sqrt{\text{Var}[l'(\theta_0; \mathbf{X})]}} = \sum_{i=1}^n \left(\frac{\tilde{\xi}_i}{\sqrt{n}} \right), \quad (6.11)$$

with

$$\xi_i = \tilde{\xi}_i \sqrt{\frac{n}{\text{Var}[l'(\theta_0; \mathbf{X})]}}, \quad \forall i \in \{1, 2, \dots, n\}.$$

It follows that $\left\{\frac{\xi_i}{\sqrt{n}}, i = 1, 2, \dots, n\right\}$ is a random $2m$ -dependent sequence with mean zero and also $\text{Var}(W) = 1$. In addition, (LD) is satisfied with A_j and B_j as in (6.2). A simple triangle inequality gives that

$$(6.9) \leq |\mathbb{E}[h(W)] - \mathbb{E}[h(Z)]| \quad (6.12)$$

$$+ \left| \mathbb{E} \left[h \left(\frac{\sqrt{ni_2(\theta_0)} l'(\theta_0; \mathbf{X})}{\alpha} \right) - h(W) \right] \right|. \quad (6.13)$$

Since the assumptions of Lemma 2.3 are satisfied for W as in (6.11), one can directly use (2.15) in order to find an upper bound for (6.12). For (6.13), a first order Taylor expansion of $h \left(\frac{\sqrt{ni_2(\theta_0)} l'(\theta_0; \mathbf{X})}{\alpha} \right)$ about W yields

$$\begin{aligned} & h \left(\frac{\sqrt{ni_2(\theta_0)} l'(\theta_0; \mathbf{X})}{\alpha} \right) - h \left(\frac{l'(\theta_0; \mathbf{X})}{\sqrt{\text{Var}(l'(\theta_0; \mathbf{X}))}} \right) \\ &= \left(\frac{\sqrt{ni_2(\theta_0)} l'(\theta_0; \mathbf{X})}{\alpha} - \frac{l'(\theta_0; \mathbf{X})}{\sqrt{\text{Var}(l'(\theta_0; \mathbf{X}))}} \right) h'(t_1(\mathbf{X})), \end{aligned}$$

where $t_1(\mathbf{X})$ is between $\frac{\sqrt{ni_2(\theta_0)} l'(\theta_0; \mathbf{X})}{\alpha}$ and $\frac{l'(\theta_0; \mathbf{X})}{\sqrt{\text{Var}(l'(\theta_0; \mathbf{X}))}}$. Therefore,

$$\begin{aligned} (6.13) &\leq \|h'\| \left| \frac{\sqrt{ni_2(\theta_0)}}{\alpha} - \frac{1}{\sqrt{\text{Var}(l'(\theta_0; \mathbf{X}))}} \right| \mathbb{E} |l'(\theta_0; \mathbf{X})| \\ &\leq \|h'\| \left| \frac{\sqrt{ni_2(\theta_0) \text{Var}(l'(\theta_0; \mathbf{X}))}}{\alpha} - 1 \right|. \end{aligned} \quad (6.14)$$

Thus, for $h \in H_W$ as in (2.4), we conclude that

$$\begin{aligned} (6.9) &\leq \frac{2}{n^{\frac{3}{2}}} \left[\sum_{i=1}^n (\mathbb{E} |\xi_i \eta_i \tau_i|) + \sum_{i=1}^n (|\mathbb{E}(\xi_i \eta_i)| \mathbb{E} |\tau_i|) \right] + \frac{1}{n^{\frac{3}{2}}} \sum_{i=1}^n \mathbb{E} |\xi_i \eta_i^2| \\ &\quad + \left| \frac{\sqrt{ni_2(\theta_0) \text{Var}(l'(\theta_0; \mathbf{X}))}}{\alpha} - 1 \right|, \end{aligned} \quad (6.15)$$

with η_i and τ_i as in (2.14). The absolute expectations in (6.15) can be difficult to bound and we express the first three quantities of the above bound in terms of more easily

calculable terms. For the first term in (6.15), using Hölder's inequality

$$\begin{aligned}
\mathbb{E}|\xi_i \eta_i \tau_i| &= \mathbb{E} \left| \xi_i \sum_{j \in A_i} \xi_j \sum_{k \in B_i} \xi_k \right| \\
&\leq \sum_{j \in A_i} \sum_{k \in B_i} \mathbb{E} |\xi_i \xi_j \xi_k| \\
&\leq \sum_{j \in A_i} \sum_{k \in B_i} \left[\mathbb{E}(|\xi_i|^3) \mathbb{E}(|\xi_j|^3) \mathbb{E}(|\xi_k|^3) \right]^{\frac{1}{3}} \\
&\leq \sum_{j \in A_i} \sum_{k \in B_i} \left[\mathbb{E}((\xi_i)^4) \mathbb{E}((\xi_j)^4) \mathbb{E}((\xi_k)^4) \right]^{\frac{1}{4}}. \tag{6.16}
\end{aligned}$$

For the second term of the bound in (6.15), the Cauchy-Schwarz inequality yields

$$\begin{aligned}
|\mathbb{E}(\xi_i \eta_i)| \mathbb{E}|\tau_i| &= \left| \mathbb{E} \left(\xi_i \sum_{j \in A_i} \xi_j \right) \right| \left| \mathbb{E} \sum_{k \in B_i} \xi_k \right| \\
&\leq \sum_{j \in A_i} \mathbb{E} |\xi_i \xi_j| \sum_{k \in B_i} \mathbb{E} |\xi_k| \\
&\leq \sum_{j \in A_i} \sum_{k \in B_i} \left[\mathbb{E}((\xi_i)^2) \mathbb{E}((\xi_j)^2) \mathbb{E}((\xi_k)^2) \right]^{\frac{1}{2}}. \tag{6.17}
\end{aligned}$$

For the third term, we also employ Jensen's inequality to get that

$$\left(\sum_{i \in J} |a_i| \right)^z \leq J^{z-1} \sum_{i \in J} |a_i|^z, \quad \forall a_i \in \mathbb{R} \text{ and } z \in \mathbb{N}$$

and therefore

$$\begin{aligned}
\mathbb{E}|\xi_i \eta_i^2| &= \mathbb{E} \left| \xi_i \left(\sum_{j \in A_i} \xi_j \right)^2 \right| \\
&\leq |A_i| \mathbb{E} \left| \xi_i \sum_{j \in A_i} \xi_j^2 \right| \\
&\leq |A_i| \sum_{j \in A_i} \mathbb{E} |\xi_i \xi_j^2| \\
&\leq |A_i| \sum_{j \in A_i} \left[\mathbb{E}((\xi_i)^2) \mathbb{E}((\xi_j)^4) \right]^{\frac{1}{2}}. \tag{6.18}
\end{aligned}$$

The results in (6.16), (6.17) and (6.18) yield

$$\begin{aligned}
 (6.12) &\leq \frac{2}{n^{\frac{3}{2}}} \left[\sum_{i=1}^n (E|\xi_i \eta_i \tau_i|) + \sum_{i=1}^n (|E(\xi_i \eta_i)| E|\tau_i|) \right] + \frac{1}{n^{\frac{3}{2}}} \sum_{i=1}^n E|\xi_i \eta_i^2| \\
 &\leq \frac{2}{n^{\frac{3}{2}}} \sum_{i=1}^n \sum_{j \in A_i} \sum_{k \in B_i} \left[E((\xi_i)^4) E((\xi_j^4)) E((\xi_k^4)) \right]^{\frac{1}{4}} \\
 &\quad + \frac{2}{n^{\frac{3}{2}}} \sum_{i=1}^n \sum_{j \in A_i} \sum_{k \in B_i} \left[E((\xi_i)^2) E((\xi_j)^2) E((\xi_k)^2) \right]^{\frac{1}{2}} \\
 &\quad + \frac{1}{n^{\frac{3}{2}}} \sum_{i=1}^n |A_i| \sum_{j \in A_i} \left[E((\xi_i)^2) E((\xi_j^4)) \right]^{\frac{1}{2}}. \tag{6.19}
 \end{aligned}$$

The bound for (6.12) is now obviously a function only of $E(\xi_i^2)$ and $E(\xi_i^4)$.

Step 2: Bound for (6.10). As in Chapter 3 where the case of independent random variables was treated, the main tool used here is again Taylor expansions. However, comparing (A.D.1) with (R3) from Chapter 2, we do not need to use conditional expectations in this case as the third derivative of the log-likelihood function is assumed to be uniformly bounded in the whole parameter space and not just in a neighbourhood of it. This also means that no positive constant ε is used. For ease of presentation, let

$$\begin{aligned}
 \tilde{C}(\theta_0) &= \tilde{C}(h, \theta_0; \mathbf{X}) := h \left(\sqrt{n i_2(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) - h \left(\frac{\sqrt{n i_2(\theta_0)} l'(\theta_0; \mathbf{X})}{\alpha} \right) \\
 &= h \left(\frac{\sqrt{n i_2(\theta_0)} [l'(\theta_0; \mathbf{X}) + R_1(\theta_0; \mathbf{X}) + (\hat{\theta}_n(\mathbf{X}) - \theta_0) (l''(\theta_0; \mathbf{X}) + \alpha)]}{\alpha} \right) \\
 &\quad - h \left(\frac{\sqrt{n i_2(\theta_0)} l'(\theta_0; \mathbf{X})}{\alpha} \right),
 \end{aligned}$$

from the result in (6.8). Using now a first order Taylor expansion of $h \left(\frac{\sqrt{n i_2(\theta_0)} [l'(\theta_0; \mathbf{X}) + R_1(\theta_0; \mathbf{X}) + (\hat{\theta}_n(\mathbf{X}) - \theta_0) (l''(\theta_0; \mathbf{X}) + \alpha)]}{\alpha} \right)$ about $\frac{\sqrt{n i_2(\theta_0)} l'(\theta_0; \mathbf{X})}{\alpha}$ yields

$$\begin{aligned}
 (6.10) = |E[\tilde{C}(\theta_0)]| &\leq \frac{\sqrt{n i_2(\theta_0)}}{\alpha} \|h'\| \left(E \left[\frac{1}{2} (\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 |l^{(3)}(\theta^*; \mathbf{X})| \right] \right. \\
 &\quad \left. + E |(\hat{\theta}_n(\mathbf{X}) - \theta_0) (l''(\theta_0; \mathbf{X}) + \alpha)| \right)
 \end{aligned}$$

$$\leq \frac{\sqrt{ni_2(\theta_0)}}{\alpha} \|h'\| \left(\frac{S_d(n)}{2} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] + \sqrt{\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} \sqrt{\mathbb{E} \left[(l''(\theta_0; \mathbf{X}) + \alpha)^2 \right]} \right), \quad (6.20)$$

where for the last step we used Cauchy-Schwarz inequality and $S_d(n)$ is as in (6.3). We conclude that (6.14), (6.19) and (6.20) yield, for $h \in H_W$, the assertion of the theorem as expressed in (6.7). \square

6.2 Locally dependent normal random variables

To illustrate the general results, as an example assume that we have a sequence $\{S_1, S_2, \dots, S_n\}$ of random variables where for $k \in \mathbb{Z}^+$ and $\forall j \in \{1, 2, \dots, n\}$,

$$S_j = \sum_{i=(j-1)k}^{jk} X_i,$$

for $X_i, i = 0, 1, 2, \dots, nk$ i.i.d. random variables from the $N(\mu, \sigma^2)$ distribution with $\mu = \theta \in \mathbb{R}$ being the unknown parameter and σ^2 is known. It is easy to see that $\{S_i\}_{i=1,2,\dots,n}$ is a 1-dependent sequence of random variables since for $\delta \in \mathbb{Z} \setminus \{0\}$

$$\text{Cov}(S_i, S_{i+\delta}) = \begin{cases} \text{Var}(X_1) = \sigma^2, & \text{if } |\delta| = 1 \\ 0, & \text{if } |\delta| > 1. \end{cases}$$

Furthermore,

$$S_i \sim N((k+1)\theta, (k+1)\sigma^2) \quad (6.21)$$

as it is a sum of $k+1$ independent normally distributed random variables with mean θ and variance σ^2 .

It is standard, see (Casella and Berger, 2002, p.177), that if Y_1, Y_2 are two normally distributed random variables with $Y_1 \sim N(\mu_1, \sigma_1^2)$ and $Y_2 \sim N(\mu_2, \sigma_2^2)$, then

$$(Y_2|Y_1 = y_1) \sim N\left(\mu_2 + \frac{\sigma_2}{\sigma_1}\rho(y_1 - \mu_1), (1 - \rho^2)\sigma_2^2\right),$$

where $\rho = \frac{\text{Cov}(Y_1, Y_2)}{\sqrt{\text{Var}(Y_1)\text{Var}(Y_2)}}$ denotes the correlation between Y_1 and Y_2 . In our example,

$$\rho = \frac{\text{Cov}(S_{i-1}, S_i)}{\sqrt{\text{Var}(S_{i-1})\text{Var}(S_i)}} = \frac{\sigma^2}{(k+1)\sigma^2} = \frac{1}{k+1}, \forall i \in \{2, 3, \dots, n\}.$$

Therefore, for $i = 2, 3, \dots, n$

$$(S_i|S_{i-1} = s_{i-1}) \sim N\left((k+1)\theta + \frac{1}{k+1}(s_{i-1} - (k+1)\theta), \frac{k(k+2)}{k+1}\sigma^2\right). \quad (6.22)$$

The likelihood function for the parameter θ under $\mathbf{S} = (S_1, S_2, \dots, S_n)$ is

$$\begin{aligned} L(\theta; \mathbf{S}) &= f(S_1|\theta) \prod_{i=2}^n f(S_i|S_{i-1}; \theta) \\ &= \frac{1}{\sqrt{2\pi(k+1)\sigma^2}} \exp\left\{-\frac{(S_1 - (k+1)\theta)^2}{2(k+1)\sigma^2}\right\} \\ &\quad \times \prod_{i=2}^n \frac{\sqrt{k+1}}{\sqrt{2\pi k(k+2)\sigma^2}} \exp\left\{-\frac{(k+1)\left(S_i - \left((k+1)\theta + \frac{1}{k+1}(S_{i-1} - (k+1)\theta)\right)\right)^2}{2k(k+2)\sigma^2}\right\} \\ &= \frac{(k+1)^{\frac{n-1}{2}}}{\sqrt{2\pi(k+1)\sigma^2}(2\pi k(k+2)\sigma^2)^{\frac{n-1}{2}}} \exp\left\{-\frac{(S_1 - (k+1)\theta)^2}{2(k+1)\sigma^2}\right. \\ &\quad \left.- \frac{k+1}{2k(k+2)\sigma^2} \sum_{i=2}^n \left(S_i - \left((k+1)\theta + \frac{1}{k+1}(S_{i-1} - (k+1)\theta)\right)\right)^2\right\}. \end{aligned} \quad (6.23)$$

Having this closed-form expression for the likelihood allows us to derive the MLE under this local dependence structure. For

$$\begin{aligned}
l(\theta; \mathbf{S}) &:= \log(L(\theta; \mathbf{S})) \\
l'(\theta; \mathbf{S}) &:= \frac{d}{d\theta} l(\theta; \mathbf{S}) \\
l''(\theta; \mathbf{S}) &:= \frac{d^2}{d\theta^2} l(\theta; \mathbf{S}) \\
l^{(3)}(\theta; \mathbf{S}) &:= \frac{d^3}{d\theta^3} l(\theta; \mathbf{S}),
\end{aligned}$$

using (6.23) we obtain that

$$\begin{aligned}
l(\theta; \mathbf{S}) &= \frac{n-1}{2} \log(k+1) - \frac{1}{2} \log(2\pi(k+1)\sigma^2) - \frac{n-1}{2} \log(2\pi k(k+2)\sigma^2) \\
&\quad - \frac{(S_1 - (k+1)\theta)^2}{2(k+1)\sigma^2} \\
&\quad - \frac{k+1}{2k(k+2)\sigma^2} \sum_{i=2}^n \left(S_i - \left((k+1)\theta + \frac{1}{k+1}(S_{i-1} - (k+1)\theta) \right) \right)^2 \\
l'(\theta; \mathbf{S}) &= \frac{S_1 - (k+1)\theta}{\sigma^2} + \frac{k+1}{(k+2)\sigma^2} \sum_{i=2}^n \left\{ S_i - \left((k+1)\theta + \frac{1}{k+1}(S_{i-1} - (k+1)\theta) \right) \right\} \\
&= \frac{1}{(k+2)\sigma^2} \left\{ k \sum_{i=1}^n S_i + S_1 + S_n - (k+1)(nk+2)\theta \right\} \\
l''(\theta; \mathbf{S}) &= -\frac{(nk+2)(k+1)}{(k+2)\sigma^2} \\
l^{(3)}(\theta; \mathbf{S}) &= 0.
\end{aligned} \tag{6.24}$$

To find the solution of $l'(\theta, \mathbf{S}) = 0$ with respect to θ , using (6.24),

$$\begin{aligned}
l'(\theta; \mathbf{S}) = 0 &\Leftrightarrow k \sum_{i=1}^n S_i + S_1 + S_n - (k+1)(nk+2)\theta = 0 \\
&\Leftrightarrow (nk+2)(k+1)\theta = k \sum_{i=1}^n S_i + S_1 + S_n \\
&\Leftrightarrow \theta = \hat{\theta}_n(\mathbf{S}) = \frac{k \sum_{i=1}^n S_i + S_1 + S_n}{(nk+2)(k+1)}.
\end{aligned} \tag{6.25}$$

The solution is unique and also $l''(\theta; \mathbf{S}) < 0$. Hence, the unique MLE for θ is $\hat{\theta}_n(\mathbf{S})$ as in (6.25). The following Corollary gives the upper bound on the Wasserstein distance between the distribution of $\hat{\theta}_n(\mathbf{S})$ and the normal distribution.

Corollary 6.1. *Let S_1, S_2, \dots, S_n be a 1-dependent sequence of random variables with $S_i \sim N((k+1)\theta, (k+1)\sigma^2)$. The conditions (A.D.1)-(A.D.4) hold. For $Z \sim N(0, 1)$ and $i_2(\theta_0) = \frac{(k+1)^2}{(k+3)\sigma^2}$,*

$$\begin{aligned} & d_W \left(\sqrt{ni_2(\theta_0)} (\hat{\theta}_n(S) - \theta_0), Z \right) \\ & \leq 339(n-5) \left[\frac{k(k+1)(k+2)}{(nk^3 + (3n+2)k^2 + 10k + 2)} \right]^{\frac{3}{2}} \\ & \quad + \frac{(k+1)^{\frac{3}{2}}(k+2)^{\frac{3}{2}}}{(nk^3 + (3n+2)k^2 + 10k + 2)^{\frac{3}{2}}} \left\{ \left(1 + 3^{\frac{3}{4}} \right) \left(2\sqrt{k+2}(37k+2) + 4\sqrt{k}(61k+8) \right) \right. \\ & \quad \left. + \sqrt{3} \left(3\sqrt{k+2}(k+1) + \sqrt{k}(91k+18) \right) \right\} \\ & \quad + \left| \left(1 - \frac{2}{nk+2} \right) \left[\frac{k+3 + \frac{2}{n} + \frac{10}{nk} + \frac{2}{nk^2}}{k+3} \right]^{\frac{1}{2}} - 1 \right|. \end{aligned}$$

Remark 6.1. *The order of the bound with respect to the sample size is $\frac{1}{\sqrt{n}}$.*

Proof. We first check that the assumptions (A.D.1)-(A.D.4) are satisfied. The first assumption is satisfied from (6.24). Now, from (6.21) and (6.22) the derivatives of $\log f(s_1|\theta)$ and $\log f(s_i|s_{i-1}; \theta)$, $i = 2, 3, \dots, n$ with respect to θ are

$$\begin{aligned} \frac{d}{d\theta} \log f(s_1|\theta) &= \frac{s_1 - (k+1)\theta}{\sigma^2} \\ \frac{d}{d\theta} \log f(s_i|s_{i-1}; \theta) &= \frac{k+1}{(k+2)\sigma^2} \left(s_i - k\theta - \frac{s_{i-1}}{k+1} \right). \end{aligned} \quad (6.26)$$

Therefore,

$$\begin{aligned} E \left[\frac{d}{d\theta} \log f(S_1|\theta) \right] &= E \left[\frac{S_1 - (k+1)\theta}{\sigma^2} \right] = 0 \\ E \left[\frac{d}{d\theta} \log f(S_i|S_{i-1}; \theta) \right] &= E \left[\frac{k+1}{(k+2)\sigma^2} \left(S_i - k\theta - \frac{S_{i-1}}{k+1} \right) \right] = 0 \end{aligned}$$

and thus (A.D.2) holds. The assumption (A.D.3) is also satisfied since, using (6.21) and (6.25),

$$\mathbb{E} [\hat{\theta}_n(\mathbf{S})] = \frac{nk(k+1)\theta_0 + 2(k+1)\theta_0}{(nk+2)(k+1)} = \theta_0.$$

To show that (A.D.4) holds, we first calculate

$$\begin{aligned} \text{Var} [\hat{\theta}_n(\mathbf{S})] &= \frac{1}{(nk+2)^2(k+1)^2} \text{Var} \left(k \sum_{i=1}^n S_i + S_1 + S_n \right) \\ &= \frac{1}{(nk+2)^2(k+1)^2} \left\{ k^2 \text{Var} \left(\sum_{i=1}^n S_i \right) + \text{Var}(S_1) + \text{Var}(S_n) + 2k \text{Cov} \left(S_1, \sum_{i=1}^n S_i \right) \right. \\ &\quad \left. + 2k \text{Cov} \left(S_n, \sum_{i=1}^n S_i \right) \right\}. \end{aligned} \quad (6.27)$$

We know from (6.21) that $\text{Var}(S_i) = (k+1)\sigma^2, \forall i \in \{1, 2, \dots, n\}$. In addition, since $\{S_i\}_{i=1,2,\dots,n}$ is a 1-dependent sequence of random variables,

$$\begin{aligned} \text{Var} \left(\sum_{i=1}^n S_i \right) &= n \text{Var}(S_1) + 2(n-1) \text{Cov}(S_1, S_2) = n(k+1)\sigma^2 + 2(n-1)\sigma^2 \\ \text{Cov} \left(S_1, \sum_{i=1}^n S_i \right) &= \text{Var}(S_1) + \text{Cov}(S_1, S_2) = (k+2)\sigma^2. \end{aligned} \quad (6.28)$$

Applying the above results of (6.28) to (6.27) gives that

$$\text{Var} [\hat{\theta}_n(\mathbf{S})] = \frac{\sigma^2 (nk^3 + 3nk^2 + 2k^2 + 10k + 2)}{(nk+2)^2(k+1)^2}. \quad (6.29)$$

Therefore,

$$\begin{aligned} i_2(\theta_0) &= \lim_{n \rightarrow \infty} \frac{1}{n \text{Var} (\hat{\theta}_n(\mathbf{S}))} = \lim_{n \rightarrow \infty} \frac{(nk+2)^2(k+1)^2}{n\sigma^2(nk^3 + 3nk^2 + 2k^2 + 10k + 2)} \\ &= \lim_{n \rightarrow \infty} \frac{n^2(k+1)^2 \left(k^2 + \frac{4k}{n} + \frac{4}{n^2} \right)}{n^2\sigma^2 \left(k^3 + 3k^2 + \frac{2k^2}{n} + \frac{10k}{n} + \frac{2}{n} \right)} = \frac{(k+1)^2}{(k+3)\sigma^2} > 0, \end{aligned} \quad (6.30)$$

which shows that (A.D.4) is satisfied. To obtain α as defined in (6.6), we calculate the variance of the score function, which is

$$\begin{aligned}\text{Var}[l'(\theta_0; \mathbf{S})] &= \frac{1}{(k+2)^2 \sigma^4} \text{Var} \left(k \sum_{i=1}^n S_i + S_1 + S_n \right) \\ &= \frac{1}{(k+2)^2 \sigma^4} \left\{ k^2 \text{Var} \left(\sum_{i=1}^n S_i \right) + \text{Var}(S_1) + \text{Var}(S_n) + 2k \text{Cov} \left(S_1, \sum_{i=1}^n S_i \right) \right. \\ &\quad \left. + 2k \text{Cov} \left(S_n, \sum_{i=1}^n S_i \right) \right\}.\end{aligned}$$

Using (6.28), we conclude that

$$\begin{aligned}\text{Var}[l'(\theta_0; \mathbf{S})] &= \frac{1}{(k+2)^2 \sigma^2} [nk^2(k+1) + 2k^2(n-1) + 2(k+1) + 4k(k+2)] \\ &= \frac{1}{(k+2)^2 \sigma^2} [nk^3 + 3nk^2 + 2k^2 + 10k + 2].\end{aligned}\quad (6.31)$$

The above result and (6.29) yield

$$\alpha = \sqrt{\frac{\text{Var}[l'(\theta_0; \mathbf{S})]}{\text{Var}[\hat{\theta}_n(\mathbf{S})]}} = \sqrt{\frac{(nk+2)^2(k+1)^2}{\sigma^4(k+2)^2}} = \frac{(nk+2)(k+1)}{(k+2)\sigma^2}.\quad (6.32)$$

For $\xi_1 = \frac{d}{d\theta} \log f(S_1|\theta) \Big|_{\theta=\theta_0} \sqrt{\frac{n}{\text{Var}[l'(\theta_0; \mathbf{S})]}}$, $\xi_i = \frac{d}{d\theta} \log f(S_i|S_{i-1}; \theta) \Big|_{\theta=\theta_0} \sqrt{\frac{n}{\text{Var}[l'(\theta_0; \mathbf{S})]}}$, $i = 2, 3, \dots, n$, using (6.31) and (6.26), we get that

$$\xi_1 = \frac{\sqrt{n}(k+2)[S_1 - (k+1)\theta_0]}{\sigma \sqrt{nk^3 + (3n+2)k^2 + 10k + 2}}$$

and therefore

$$\begin{aligned}\text{E}(\xi_1^2) &= \frac{n(k+2)^2 \text{E}(S_1 - (k+1)\theta)^2}{(nk^3 + (3n+2)k^2 + 10k + 2)\sigma^2} = \frac{n(k+2)^2(k+1)}{nk^3 + (3n+2)k^2 + 10k + 2} \\ \text{E}(\xi_1^4) &= \frac{n^2(k+2)^4 \text{E}(S_1 - (k+1)\theta)^4}{(nk^3 + (3n+2)k^2 + 10k + 2)^2 \sigma^4} = \frac{3n^2(k+2)^4(k+1)^2}{(nk^3 + (3n+2)k^2 + 10k + 2)^2}.\end{aligned}\quad (6.33)$$

Furthermore, for $i = 2, 3, \dots, n$,

$$\xi_i = \frac{\sqrt{n}(k+1) \left[S_i - \left((k+1)\theta + \frac{1}{k+1} (S_{i-1} - (k+1)\theta) \right) \right]}{\sigma \sqrt{nk^3 + (3n+2)k^2 + 10k + 2}}$$

so that for $i = 2, 3, \dots, n$,

$$\begin{aligned} \mathbb{E}(\xi_i^2) &= \frac{n(k+1)^2 \mathbb{E} \left[S_i - \left((k+1)\theta + \frac{1}{k+1} (S_{i-1} - (k+1)\theta) \right) \right]^2}{\sigma^2 (nk^3 + (3n+2)k^2 + 10k + 2)} \\ &= \frac{nk(k+1)(k+2)}{nk^3 + (3n+2)k^2 + 10k + 2} \\ \mathbb{E}(\xi_i^4) &= \frac{n^2(k+1)^4 \mathbb{E} \left[S_i - \left((k+1)\theta + \frac{1}{k+1} (S_{i-1} - (k+1)\theta) \right) \right]^4}{\sigma^4 (nk^3 + (3n+2)k^2 + 10k + 2)^2} \\ &= \frac{3n^2k^2(k+1)^2(k+2)^2}{(nk^3 + (3n+2)k^2 + 10k + 2)^2}. \end{aligned} \quad (6.34)$$

We are now ready to calculate the first three terms of the general bound (6.7) denoted from now on by

$$\begin{aligned} Q_v = Q_v(k, n) &:= \frac{1}{n^{\frac{3}{2}}} \left\{ 2 \sum_{j \in A_v} \sum_{l \in B_v} \left[\mathbb{E}((\xi_v)^2) \mathbb{E}((\xi_j)^2) \mathbb{E}((\xi_l)^2) \right]^{\frac{1}{2}} \right. \\ &\quad \left. + 2 \sum_{j \in A_v} \sum_{l \in B_v} \left[\mathbb{E}((\xi_v)^4) \mathbb{E}((\xi_j)^4) \mathbb{E}((\xi_l)^4) \right]^{\frac{1}{4}} + |A_v| \sum_{j \in A_v} \left[\mathbb{E}((\xi_v)^2) \mathbb{E}((\xi_j)^4) \right]^{\frac{1}{2}} \right\}. \end{aligned}$$

We split our approach depending on whether 1 is an element of either A_i or B_i , for $i \in \{1, 2, \dots, n\}$.

Case 1: $i = 6, 7, \dots, n$. Using the results in (6.34) and since $|A_i| \leq 5$, $|B_i| \leq 9$, $\forall i \in \{6, 7, \dots, n\}$, we obtain that

$$\begin{aligned} Q_i &\leq \frac{1}{n^{\frac{3}{2}}} \left\{ 90 \left[\mathbb{E}((\xi_2)^4) \right]^{\frac{3}{4}} + 90 \left[\mathbb{E}((\xi_2)^2) \right]^{\frac{3}{2}} + 25 \left[\mathbb{E}((\xi_2)^2) \mathbb{E}((\xi_2)^4) \right]^{\frac{1}{2}} \right\} \\ &= \left[\frac{k(k+1)(k+2)}{(nk^3 + (3n+2)k^2 + 10k + 2)} \right]^{\frac{3}{2}} \left(90(3)^{\frac{3}{4}} + 90 + 25\sqrt{3} \right) \\ &< 339 \left[\frac{k(k+1)(k+2)}{(nk^3 + (3n+2)k^2 + 10k + 2)} \right]^{\frac{3}{2}}. \end{aligned} \quad (6.35)$$

Issues arise due to ξ_1 not having the same distribution as ξ_i for $i \in \{2, 3, \dots, n\}$. There are hence five more special cases corresponding to $i = 1, 2, \dots, 5$. We treat these cases separately.

Case 2: $i = 1$. For $A_1 = \{1, 2, 3\}$ and $B_1 = \{1, 2, \dots, 5\}$, the results in (6.33) and (6.34) yield

$$\begin{aligned}
 Q_1 &= \frac{1}{n^{\frac{3}{2}}} \left\{ 2 \left[\left[E((\xi_1)^2) \right]^{\frac{3}{2}} + 6E((\xi_1)^2) \left[E((\xi_2)^2) \right]^{\frac{1}{2}} + 8E((\xi_2)^2) \left[E((\xi_1)^2) \right]^{\frac{1}{2}} \right] \right. \\
 &\quad + 2 \left[\left[E((\xi_1)^4) \right]^{\frac{3}{4}} + 6 \left[E((\xi_1)^4) \right]^{\frac{1}{2}} \left[E((\xi_2)^4) \right]^{\frac{1}{4}} + 8 \left[E((\xi_1)^4) \right]^{\frac{1}{4}} \left[E((\xi_2)^4) \right]^{\frac{1}{2}} \right] \\
 &\quad \left. + 3 \left[\left[E((\xi_1)^2) \right]^{\frac{1}{2}} \left(\left[E((\xi_1)^4) \right]^{\frac{1}{2}} + 2 \left[E((\xi_2)^4) \right]^{\frac{1}{2}} \right) \right] \right\} \\
 &= \frac{(k+1)^{\frac{3}{2}}(k+2)^2}{(nk^3 + (3n+2)k^2 + 10k+2)^{\frac{3}{2}}} \left\{ 2 \left[(k+2) + 6\sqrt{k(k+2)} + 8k \right] \right. \\
 &\quad \left. + 2(3)^{\frac{3}{4}} \left[(k+2) + 6\sqrt{k(k+2)} + 8k \right] + 3\sqrt{3}(2k+2) \right\} \\
 &= \frac{2(k+1)^{\frac{3}{2}}(k+2)^2}{(nk^3 + (3n+2)k^2 + 10k+2)^{\frac{3}{2}}} \left\{ \left(9k+2+6\sqrt{k(k+2)} \right) \left(1+3^{\frac{3}{4}} \right) \right. \\
 &\quad \left. + 3\sqrt{3}(k+1) \right\}. \tag{6.36}
 \end{aligned}$$

Case 3: $i = 2$. For $A_2 = \{1, 2, 3, 4\}$ and $B_2 = \{1, 2, \dots, 6\}$, the results in (6.33) and (6.34) yield

$$\begin{aligned}
 Q_2 &= \frac{1}{n^{\frac{3}{2}}} \left\{ 2 \left[E((\xi_1)^2) \left[E((\xi_2)^2) \right]^{\frac{1}{2}} + 8E((\xi_2)^2) \left[E((\xi_1)^2) \right]^{\frac{1}{2}} + 15 \left[E((\xi_2)^2) \right]^{\frac{3}{2}} \right] \right. \\
 &\quad + 2 \left[\left[E((\xi_1)^4) \right]^{\frac{1}{2}} \left[E((\xi_2)^4) \right]^{\frac{1}{4}} + 8 \left[E((\xi_2)^4) \right]^{\frac{1}{2}} \left[E((\xi_1)^4) \right]^{\frac{1}{4}} + 15 \left[E((\xi_2)^4) \right]^{\frac{3}{4}} \right] \\
 &\quad \left. + 4 \left[\left[E((\xi_2)^2) \right]^{\frac{1}{2}} \left(\left[E((\xi_1)^4) \right]^{\frac{1}{2}} + 3 \left[E((\xi_2)^4) \right]^{\frac{1}{2}} \right) \right] \right\}
 \end{aligned}$$

$$\begin{aligned}
&= \frac{\sqrt{k}(k+1)^{\frac{3}{2}}(k+2)^{\frac{3}{2}}}{(nk^3 + (3n+2)k^2 + 10k + 2)^{\frac{3}{2}}} \left\{ 2 \left[k + 2 + 8\sqrt{k(k+2)} + 15k \right] \right. \\
&\quad \left. + 2(3)^{\frac{3}{4}} \left[k + 2 + 8\sqrt{k(k+2)} + 15k \right] + 4\sqrt{3}(4k+2) \right\} \\
&= \frac{4\sqrt{k}(k+1)^{\frac{3}{2}}(k+2)^{\frac{3}{2}}}{(nk^3 + (3n+2)k^2 + 10k + 2)^{\frac{3}{2}}} \left\{ \left(8k + 1 + 4\sqrt{k(k+2)} \right) \left(1 + 3^{\frac{3}{4}} \right) \right. \\
&\quad \left. + 2\sqrt{3}(2k+1) \right\}. \tag{6.37}
\end{aligned}$$

Case 4: $i = 3$. Following the same steps as in Case 3, now for $A_3 = \{1, 2, \dots, 5\}$ and $B_3 = \{1, 2, \dots, 7\}$, the results in (6.33) and (6.34) give that

$$\begin{aligned}
Q_3 &= \frac{\sqrt{k}(k+1)^{\frac{3}{2}}(k+2)^{\frac{3}{2}}}{(nk^3 + (3n+2)k^2 + 10k + 2)^{\frac{3}{2}}} \left\{ 2 \left(25k + 2 + 10\sqrt{k(k+2)} \right) \left(1 + 3^{\frac{3}{4}} \right) \right. \\
&\quad \left. + 5\sqrt{3}(5k+2) \right\}. \tag{6.38}
\end{aligned}$$

Case 5: $i = 4$. For $A_4 = \{2, 3, \dots, 6\}$, $B_4 = \{1, 2, \dots, 8\}$, the results in (6.33) and (6.34) yield

$$\begin{aligned}
Q_4 &= \frac{1}{n^{\frac{3}{2}}} \left\{ 2 \left[5\mathbb{E} \left((\xi_2)^2 \right) \left[\mathbb{E} \left((\xi_1)^2 \right) \right]^{\frac{1}{2}} + 35 \left[\mathbb{E} \left((\xi_2)^2 \right) \right]^{\frac{3}{2}} \right] \right. \\
&\quad + 2 \left[5 \left[\mathbb{E} \left((\xi_2)^4 \right) \right]^{\frac{1}{2}} \left[\mathbb{E} \left((\xi_1)^4 \right) \right]^{\frac{1}{4}} + 35 \left[\mathbb{E} \left((\xi_2)^4 \right) \right]^{\frac{3}{4}} \right] \\
&\quad \left. + 25 \left[\left[\mathbb{E} \left((\xi_2)^2 \right) \right]^{\frac{1}{2}} \left[\mathbb{E} \left((\xi_2)^4 \right) \right]^{\frac{1}{2}} \right] \right\} \\
&= \frac{5k(k+1)^{\frac{3}{2}}(k+2)^{\frac{3}{2}}}{(nk^3 + (3n+2)k^2 + 10k + 2)^{\frac{3}{2}}} \left\{ 2 \left[\sqrt{k+2} + 7\sqrt{k} \right] \left(1 + 3^{\frac{3}{4}} \right) + 5\sqrt{3k} \right\}. \tag{6.39}
\end{aligned}$$

Case 6: $i = 5$. Following the same steps as in Case 5, now for $A_5 = \{3, 4, \dots, 7\}$ and $B_5 = \{1, 2, \dots, 9\}$, the results in (6.33) and (6.34) yield

$$Q_5 = \frac{5k(k+1)^{\frac{3}{2}}(k+2)^{\frac{3}{2}}}{(nk^3 + (3n+2)k^2 + 10k + 2)^{\frac{3}{2}}} \left\{ 2 \left[\sqrt{k+2} + 8\sqrt{k} \right] \left(1 + 3^{\frac{3}{4}} \right) + 5\sqrt{3k} \right\}. \quad (6.40)$$

The sum of the results of (6.36), (6.37), (6.38), (6.39) and (6.40) with $(n-5)$ times the bound in (6.35) consists an upper bound for the first three terms of the general upper bound as expressed in (6.7). For the fourth term of the general upper bound, (6.30), (6.31) and (6.32) yield

$$\begin{aligned} \left| \frac{\sqrt{ni_2(\theta_0) \text{Var}[l'(\theta_0; \mathbf{X})]}}{\alpha} - 1 \right| &= \left| \frac{\left[\frac{n(k+1)^2(nk^3 + 3nk^2 + 2k^2 + 10k + 2)}{k+3} \right]^{\frac{1}{2}}}{(nk+2)(k+1)} - 1 \right| \\ &= \left| \frac{nk}{nk+2} \left[\frac{k+3 + \frac{2}{n} + \frac{10}{nk} + \frac{2}{nk^2}}{k+3} \right]^{\frac{1}{2}} - 1 \right| \\ &= \left| \left(1 - \frac{2}{nk+2} \right) \left[\frac{k+3 + \frac{2}{n} + \frac{10}{nk} + \frac{2}{nk^2}}{k+3} \right]^{\frac{1}{2}} - 1 \right|. \end{aligned} \quad (6.41)$$

The fifth term of the bound in (6.7) involves the calculation of $S_d(n)$, which is equal to zero from (6.24). Therefore, the fifth term of the general upper bound vanishes for this example. For the last term we have from (6.24) that

$$\mathbb{E}[l''(\theta_0; \mathbf{S})] = -\frac{(nk+2)(k+1)}{(k+2)\sigma^2} = -\alpha$$

and therefore

$$\sqrt{\mathbb{E}[(\hat{\theta}_n(\mathbf{S}) - \theta_0)^2] \mathbb{E}[(l''(\theta_0; \mathbf{S}) + \alpha)^2]} = \sqrt{\mathbb{E}[(\hat{\theta}_n(\mathbf{S}) - \theta_0)^2] \text{Var}[l''(\theta_0; \mathbf{S})]} = 0.$$

The results of Case 1 - Case 6 and (6.41) give the assertion of the corollary. \square

6.3 An alternative bound

In this section, under stricter regularity conditions we give the asymptotic normality result for the distribution of the MLE obtained from locally dependent random variables. Under these new conditions, an alternative bound on the Wasserstein distance between the exact distribution of the MLE and the normal distribution for locally dependent random variables is given.

We work under the following regularity conditions, taken from Bhat (1974):

(R.D.1) The range of x_k given x_1, x_2, \dots, x_{k-1} , for $k = 2, 3, \dots, n$, does not depend on θ and differentiation with respect to θ may be carried out under the integral sign, integration being with respect to x_k , up to the third order, for $\log f(x_k | x_{k-1}, \dots, x_1; \theta)$. These derivatives are also assumed to be continuous in $\theta \in \Theta$. In addition,

$$E_{k-1} \left(\frac{d}{d\theta} \log f(X_k | X_{k-1}, \dots, X_1; \theta) \right) \Big|_{\theta=\theta_0} = 0, \text{ almost surely}$$

and we denote by

$$i_k(\theta_0) := E_{k-1} \left[\left(\frac{d}{d\theta} \log f(X_k | X_{k-1}, \dots, X_1; \theta) \right)^2 \right] \Big|_{\theta=\theta_0}$$

where E_{k-1} denotes the conditional expectation given X_1, X_2, \dots, X_{k-1} .

(R.D.2) The third derivative $\frac{d^3}{d\theta^3} \log f(X_k | X_{k-1}, \dots, X_1; \theta)$ is bounded in probability uniformly in $\theta \in \Theta$, (X_1, X_2, \dots, X_k) .

(R.D.3) $\forall \theta_0 \in \Theta$, $\bar{i}_n(\theta_0) = \frac{1}{n} \sum_{j=1}^n i_j(\theta_0) \rightarrow \bar{i}(\theta_0) > 0$ as $n \rightarrow \infty$, where $\bar{i}(\theta_0)$ is a constant almost surely and θ_0 is the true value of θ .

(R.D.4) $\frac{1}{n^{1+\frac{\delta}{2}}} \sum_{k=2}^n E_{k-1} \left| \frac{\partial \log f(x_k | x_{k-1}, \dots, x_1; \theta)}{\partial \theta} \right|^{2+\delta} \rightarrow 0$ a.s. for some $\delta > 0$.

$$(R.D.5) \limsup_{n \rightarrow \infty} \left\{ \frac{1}{n^2} \sum_{k=1}^n \text{Var} \left(\frac{d^2}{d\theta^2} \log f(X_k | X_{k-1}, \dots, X_1; \theta) \right) \Big|_{\theta=\theta_0} \right\} < \infty.$$

The condition (R.D.3) implies that $\sum_{k=1}^n i_k(\theta_0) \rightarrow \infty$ as $n \rightarrow \infty$, which is essential for the CLT to hold for the dependent random variables $\frac{d}{d\theta} \log f(X_k | X_{k-1}, \dots, X_1; \theta) \Big|_{\theta=\theta_0}$. From now on, unless otherwise stated

$$\bar{i}(\theta_0) := \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n i_j(\theta_0). \quad (6.42)$$

Also, if $i_k(\theta_0)$ is an absolute constant $\forall k \in \{1, 2, \dots, n\}$ then (R.D.4) is not necessary as shown in Silvey (1961).

Assuming that the MLE exists and is unique and that the conditions (R.D.1)-(R.D.5) are satisfied, then Theorem 3 in Bhat (1974) gives the result for the asymptotic normality of the MLE in the case of dependent random variables. The proof can be found in Bhat (1974).

Theorem 6.2. (Bhat, 1974). *Let X_1, X_2, \dots, X_n be dependent and identically distributed random variables with probability density (or mass) functions $f(x_1 | \theta)$, $f(x_i | x_{i-1}, \dots, x_1; \theta)$, $i = 2, 3, \dots, n$ where $\theta \in \Theta \subset \mathbb{R}$, with Θ denoting the open parameter space. Assume that the regularity conditions (R.D.1)-(R.D.5) hold. Let $Z \sim N(0, 1)$. Then,*

$$\sqrt{n\bar{i}(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \xrightarrow[n \rightarrow \infty]{d} Z. \quad (6.43)$$

The following theorem gives an upper bound on the Wasserstein distance between the distribution of the MLE and its limiting, under (R.D.1)-(R.D.5), normal distribution in the case of m -dependent sequence of random variables.

Theorem 6.3. *Let $\{X_i, i = 1, 2, \dots, n\}$ be an m -dependent sequence of identically distributed random variables with probability density (or mass) functions $f(x_1 | \theta)$, $f(x_i | x_{i-1}, \dots, x_1; \theta)$, $i = 2, 3, \dots, n$, where $\theta \in \Theta \subset \mathbb{R}$ and $(X_1, X_2, \dots, X_n) \in S$, the support of the joint probability density (or mass) function. Assume that the regularity*

conditions (R.D.1)-(R.D.5) are satisfied and also that $\hat{\theta}_n(\mathbf{X})$ exists and is unique. Let $\bar{i}(\theta_0)$ as in (6.42), $\text{Var}[l'(\theta_0; \mathbf{X})] > 0$,

$$\xi_1 = \frac{d}{d\theta} \log f(X_1|\theta) \Big|_{\theta=\theta_0} \sqrt{\frac{n}{\text{Var}[l'(\theta_0; \mathbf{X})]}}$$

and for $i = 2, 3, \dots, n$,

$$\xi_i = \frac{d}{d\theta} \log f(X_i|X_{i-1}, \dots, X_{i-m}; \theta) \Big|_{\theta=\theta_0} \sqrt{\frac{n}{\text{Var}[l'(\theta_0; \mathbf{X})]}}.$$

With $S_d(n)$ as in (6.3) and for $Z \sim N(0, 1)$,

$$\begin{aligned} & d_W \left(\sqrt{n\bar{i}(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0), Z \right) \\ & \leq \frac{2}{n^{\frac{3}{2}}} \sum_{i=1}^n \sum_{j \in A_i} \sum_{k \in B_i} \left[E((\xi_i)^4) E((\xi_j)^4) E((\xi_k)^4) \right]^{\frac{1}{4}} \\ & \quad + \frac{2}{n^{\frac{3}{2}}} \sum_{i=1}^n \sum_{j \in A_i} \sum_{k \in B_i} \left[E((\xi_i)^2) E((\xi_j)^2) E((\xi_k)^2) \right]^{\frac{1}{2}} \\ & \quad + \frac{1}{n^{\frac{3}{2}}} \sum_{i=1}^n |A_i| \sum_{j \in A_i} \left[E((\xi_i)^2) E((\xi_j)^4) \right]^{\frac{1}{2}} \\ & \quad + \left| \frac{\sqrt{\text{Var}[l'(\theta_0; \mathbf{X})]}}{\sqrt{n\bar{i}(\theta_0)}} - 1 \right| + \frac{S_d(n)}{2\sqrt{n\bar{i}(\theta_0)}} E[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2] \\ & \quad + \frac{1}{\sqrt{n\bar{i}(\theta_0)}} \sqrt{E[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2]} \sqrt{E[n\bar{i}(\theta_0) + l''(\theta_0; \mathbf{X})]^2}. \end{aligned} \quad (6.44)$$

Remark 6.2. In order for the above bound to approach zero with rate \sqrt{n} as the sample size increases, we require that $\left| \frac{\sqrt{\text{Var}(l'(\theta_0; \mathbf{X}))}}{\sqrt{n\bar{i}(\theta_0)}} - 1 \right| = o\left(\frac{1}{\sqrt{n}}\right)$.

Proof. By the definition of the MLE and (R.D.1), $l'(\hat{\theta}_n(\mathbf{x}); \mathbf{x}) = 0$. A similar approach as the one in Chapter 3, based on a second order Taylor expansion, yields

$$\sqrt{n\bar{i}(\theta_0)} (\hat{\theta}_n(\mathbf{x}) - \theta_0) = \frac{l'(\theta_0; \mathbf{x}) + R_1(\theta_0; \mathbf{x}) + R_2(\theta_0; \mathbf{x})}{\sqrt{n\bar{i}(\theta_0)}}, \quad (6.45)$$

for $\bar{i}(\theta_0)$ as in (6.42). Also, $R_1(\theta_0; \mathbf{x})$ and $R_2(\theta_0; \mathbf{x})$ are as in (3.9) and (3.10), respectively. The triangle inequality gives that

$$\begin{aligned} & \left| \mathbb{E} \left[h \left(\sqrt{n\bar{i}(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) \right] - \mathbb{E}[h(Z)] \right| \\ & \leq \left| \mathbb{E} \left[h \left(\frac{l'(\theta_0; \mathbf{X})}{\sqrt{n\bar{i}(\theta_0)}} \right) \right] - \mathbb{E}[h(Z)] \right| \end{aligned} \quad (6.46)$$

$$+ \left| \mathbb{E} \left[h \left(\sqrt{n\bar{i}(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) - h \left(\frac{l'(\theta_0; \mathbf{X})}{\sqrt{n\bar{i}(\theta_0)}} \right) \right] \right|. \quad (6.47)$$

Step 1: Bound for (6.46). With W as in (6.11), the steps are the same as in the proof of Theorem 6.1 and lead to

$$(6.46) \leq |\mathbb{E}[h(W)] - \mathbb{E}[h(Z)]| \quad (6.48)$$

$$+ \left| \mathbb{E} \left[h \left(\frac{l'(\theta_0; \mathbf{X})}{\sqrt{n\bar{i}(\theta_0)}} \right) - h(W) \right] \right|. \quad (6.49)$$

The assumptions of Lemma 2.3 hold for (6.48) and thus we directly use (2.15) to bound this term. For (6.49) using a first order Taylor expansion of $h \left(\frac{l'(\theta_0; \mathbf{X})}{\sqrt{n\bar{i}(\theta_0)}} \right)$ about W and the Cauchy-Schwarz inequality, we have that

$$\begin{aligned} (6.49) & \leq \|h'\| \left| \frac{1}{\sqrt{n\bar{i}(\theta_0)}} - \frac{1}{\sqrt{\text{Var}(l'(\theta_0; \mathbf{X}))}} \right| \mathbb{E}|l'(\theta_0; \mathbf{X})| \\ & \leq \|h'\| \left| \frac{\sqrt{\text{Var}(l'(\theta_0; \mathbf{X}))}}{\sqrt{n\bar{i}(\theta_0)}} - 1 \right|. \end{aligned} \quad (6.50)$$

Step 2: Bound for (6.47). The approach is similar to the proof of Theorem 6.1. Let

$$\begin{aligned} C_d = C_d(h, \theta_0; \mathbf{X}) & := h \left(\sqrt{n\bar{i}(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) - h \left(\frac{l'(\theta_0; \mathbf{X})}{\sqrt{n\bar{i}(\theta_0)}} \right) \\ & = h \left(\frac{l'(\theta_0; \mathbf{X}) + R_1(\theta_0; \mathbf{X}) + R_2(\theta_0; \mathbf{X})}{\sqrt{n\bar{i}(\theta_0)}} \right) - h \left(\frac{l'(\theta_0; \mathbf{X})}{\sqrt{n\bar{i}(\theta_0)}} \right) \end{aligned}$$

using (6.45). For $t(\mathbf{X})$ between $\frac{l'(\theta_0; \mathbf{X})}{\sqrt{n\bar{i}(\theta_0)}}$ and $\frac{l'(\theta_0; \mathbf{X}) + R_1(\theta_0; \mathbf{x}) + R_2(\theta_0; \mathbf{X})}{\sqrt{n\bar{i}(\theta_0)}}$, a first order Taylor expansion of $h\left(\frac{l'(\theta_0; \mathbf{X}) + R_1(\theta_0; \mathbf{x}) + R_2(\theta_0; \mathbf{X})}{\sqrt{n\bar{i}(\theta_0)}}$ about $\frac{l'(\theta_0)}{\sqrt{n\bar{i}(\theta_0)}}$ yields

$$\begin{aligned} (6.47) = |E[C_d]| &\leq \left[E \left| \frac{R_1(\theta_0, \mathbf{X}) + R_2(\theta_0, \mathbf{X})}{\sqrt{n\bar{i}(\theta_0)}} h'(t(\mathbf{X})) \right| \right] \\ &\leq \frac{\|h'\|}{2\sqrt{n\bar{i}(\theta_0)}} E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \left| l^{(3)}(\theta^*; \mathbf{X}) \right| \right] \\ &\quad + \frac{\|h'\|}{\sqrt{n\bar{i}(\theta_0)}} E \left[|(\hat{\theta}_n(\mathbf{X}) - \theta_0) (n\bar{i}(\theta_0) + l''(\theta_0; \mathbf{X}))| \right]. \end{aligned}$$

Using now (R.D.2), the Cauchy-Schwarz inequality and (6.3) give

$$\begin{aligned} (6.47) &\leq \frac{\|h'\| S_d(n)}{2\sqrt{n\bar{i}(\theta_0)}} E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \\ &\quad + \frac{\|h'\|}{\sqrt{n\bar{i}(\theta_0)}} \sqrt{E \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} \sqrt{E [n\bar{i}(\theta_0) + l''(\theta_0; \mathbf{X})]^2}. \end{aligned} \quad (6.51)$$

We conclude that (6.50), (6.51) and the application of Lemma 2.3 to (6.48) gives for $h \in H_W$ the result of the theorem as expressed in (6.44). \square

Remark 6.3. *The bound in (6.44) is similar to the result in (6.7) and holds under the rather restrictive conditions (R.D.1)-(R.D.5). Notice the difference between $i_2(\theta_0)$ used in Theorem 6.1 with $\bar{i}(\theta_0)$ of Theorem 6.3 as defined in (R.D.3). The limit $i_2(\theta_0)$ requires the calculation of $\text{Var}(\hat{\theta}_n(\mathbf{X}))$, while $\bar{i}(\theta_0)$ is independent of the MLE. However, $\bar{i}(\theta_0)$ is not necessarily easy to calculate.*

6.4 Results for not analytically known MLE

The results of Theorems 6.1 and 6.3 can be used only in cases where a closed-form expression of the MLE is available. Thus, as in Section 3.5, a solution for situations where the MLE exists, but does not have a closed-form expression is desirable. Further assumptions are made in order to ensure both the validity of the procedure to be followed

and that a meaningful upper bound is obtained. Depending on these assumptions we split the results of this section into two categories. The first category requires the support of the distribution to be bounded, while for the second one we put the boundedness condition on the parameter space instead of the support of the distribution.

6.4.1 Bounded support

In addition to the regularity conditions (R.D.1)-(R.D.5), the results presented in this subsection hold under the following assumptions:

(F.A.1) The support of the distribution is bounded;

(F.A.2) The sample size, n , satisfies

$$n > \frac{[\|x\|S_d(n)]^{\frac{2}{3}}}{\bar{i}(\theta_0)}.$$

This ensures that $1 - \frac{\|x\|S_d(n)}{[n\bar{i}(\theta_0)]^{\frac{3}{2}}} > 0$ for $S_d(n)$ as in (6.3).

For ease of presentation, we denote by

$$\begin{aligned} D_1 &:= D_1(\theta_0, x, n) = 1 - \frac{\|x\|S_d(n)}{[n\bar{i}(\theta_0)]^{\frac{3}{2}}} \\ \gamma_1 &:= \gamma_1(\theta_0, x, n) = \frac{4\|x\|}{n^{\frac{3}{2}}} \sum_{i=1}^n \sum_{j \in A_i} \sum_{k \in B_i} \left[E((\xi_i)^4) E((\xi_j)^4) E((\xi_k)^4) \right]^{\frac{1}{4}} \\ &\quad + \frac{4\|x\|}{n^{\frac{3}{2}}} \sum_{i=1}^n \sum_{j \in A_i} \sum_{k \in B_i} \left[E((\xi_i)^2) E((\xi_j)^2) E((\xi_k)^2) \right]^{\frac{1}{2}} \\ &\quad + \frac{2\|x\|}{n^{\frac{3}{2}}} \sum_{i=1}^n |A_i| \sum_{j \in A_i} \left[E((\xi_i)^2) E((\xi_j)^4) \right]^{\frac{1}{2}} + 2\|x\| \left| \frac{\sqrt{\text{Var}(l'(\theta_0; \mathbf{X}))}}{\sqrt{n\bar{i}(\theta_0)}} - 1 \right|. \end{aligned} \quad (6.52)$$

The theorem that follows gives an upper bound for the MSE of the MLE, which is then used to bound the Wasserstein distance between the distribution of the MLE and the normal distribution in cases where the MLE is not analytically known and (F.A.1), (F.A.2) hold.

Theorem 6.4. Let $\{X_i, i = 1, 2, \dots, n\}$ be an m -dependent sequence of random variables with joint probability density (or mass) function $f(\mathbf{x}|\boldsymbol{\theta})$ as in (6.1). Assume that the regularity conditions (R.D.1)-(R.D.5), as well as the assumptions (F.A.1), (F.A.2) are satisfied. Assume that the MLE exists and is unique. Then $A_d = A_d(\boldsymbol{\theta}_0, n)$ is an upper bound for $\sqrt{\mathbb{E}[(\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0)^2]}$, where

$$A_d = \frac{1}{D_1} \left\{ \frac{\|\mathbf{x}\| \sqrt{\mathbb{E}[(l''(\boldsymbol{\theta}_0; \mathbf{X}) + n\bar{l}(\boldsymbol{\theta}_0))^2]}}{[n\bar{l}(\boldsymbol{\theta}_0)]^{\frac{3}{2}}} + \left[\frac{\|\mathbf{x}\|^2 \mathbb{E}[(l''(\boldsymbol{\theta}_0; \mathbf{X}) + n\bar{l}(\boldsymbol{\theta}_0))^2]}{[n\bar{l}(\boldsymbol{\theta}_0)]^3} + \frac{D_1(1 + \gamma_1)}{n\bar{l}(\boldsymbol{\theta}_0)} \right]^{\frac{1}{2}} \right\},$$

with D_1 and γ_1 as in (6.52). Assuming that $\text{Var}[l'(\boldsymbol{\theta}_0; \mathbf{X})] > 0$, then for

$$\xi_1 = \frac{d}{d\boldsymbol{\theta}} \log f(X_1|\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \sqrt{\frac{n}{\text{Var}[l'(\boldsymbol{\theta}_0; \mathbf{X})]}}$$

and

$$\xi_i = \frac{d}{d\boldsymbol{\theta}} \log f(X_i|X_{i-1}, \dots, X_{i-m}; \boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \sqrt{\frac{n}{\text{Var}[l'(\boldsymbol{\theta}_0; \mathbf{X})]}}, i = 2, 3, \dots, n,$$

we have that

$$\begin{aligned} & d_W \left(\sqrt{n\bar{l}(\boldsymbol{\theta}_0)}(\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0), Z \right) \\ & \leq \frac{2}{n^{\frac{3}{2}}} \sum_{i=1}^n \sum_{j \in A_i} \sum_{k \in B_i} \left[\mathbb{E}((\xi_i)^4) \mathbb{E}((\xi_j)^4) \mathbb{E}((\xi_k)^4) \right]^{\frac{1}{4}} \\ & \quad + \frac{2}{n^{\frac{3}{2}}} \sum_{i=1}^n \sum_{j \in A_i} \sum_{k \in B_i} \left[\mathbb{E}((\xi_i)^2) \mathbb{E}((\xi_j)^2) \mathbb{E}((\xi_k)^2) \right]^{\frac{1}{2}} \\ & \quad + \frac{1}{n^{\frac{3}{2}}} \sum_{i=1}^n |A_i| \sum_{j \in A_i} \left[\mathbb{E}((\xi_i)^2) \mathbb{E}((\xi_j)^4) \right]^{\frac{1}{2}} + \left| \frac{\sqrt{\text{Var}[l'(\boldsymbol{\theta}_0; \mathbf{X})]}}{\sqrt{n\bar{l}(\boldsymbol{\theta}_0)}} - 1 \right| \\ & \quad + \frac{A_d}{\sqrt{n\bar{l}(\boldsymbol{\theta}_0)}} \left(\sqrt{\mathbb{E}(l''(\boldsymbol{\theta}_0; \mathbf{X}) + n\bar{l}(\boldsymbol{\theta}_0))^2} + \frac{S_d(n)A_d}{2} \right). \end{aligned} \quad (6.53)$$

Remark 6.4. The order of A_d in terms of the sample size is $\frac{1}{\sqrt{n}}$. Assuming that $\left| \frac{\sqrt{\text{Var}(l'(\theta_0; \mathbf{X}))}}{\sqrt{n\bar{l}(\theta_0)}} - 1 \right| = \mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$, then the order of the final upper bound in (6.53) is also $\frac{1}{\sqrt{n}}$.

Proof. Using similar steps as in the proof of Theorem 6.3, but keeping the result in a general form and not specifying in terms of the Wasserstein distance, we obtain that for $h : \mathbb{R} \rightarrow \mathbb{R}$ being any absolutely continuous function,

$$\begin{aligned}
& \left| \mathbb{E} \left[h \left(\sqrt{n\bar{l}(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) \right] - \mathbb{E}[h(Z)] \right| \\
& \leq \frac{2\|h'\|}{n^{\frac{3}{2}}} \sum_{i=1}^n \sum_{j \in A_i} \sum_{k \in B_i} \left[\mathbb{E}((\xi_i)^4) \mathbb{E}((\xi_j)^4) \mathbb{E}((\xi_k)^4) \right]^{\frac{1}{4}} \\
& \quad + \frac{2\|h'\|}{n^{\frac{3}{2}}} \sum_{i=1}^n \sum_{j \in A_i} \sum_{k \in B_i} \left[\mathbb{E}((\xi_i)^2) \mathbb{E}((\xi_j)^2) \mathbb{E}((\xi_k)^2) \right]^{\frac{1}{2}} \\
& \quad + \frac{\|h'\|}{n^{\frac{3}{2}}} \sum_{i=1}^n |A_i| \sum_{j \in A_i} \left[\mathbb{E}((\xi_i)^2) \mathbb{E}((\xi_j)^4) \right]^{\frac{1}{2}} + \|h'\| \left| \frac{\sqrt{\text{Var}(l'(\theta_0; \mathbf{X}))}}{\sqrt{n\bar{l}(\theta_0)}} - 1 \right| \\
& \quad + \frac{\|h'\| S_d(n)}{2\sqrt{n\bar{l}(\theta_0)}} \mathbb{E}[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2] \\
& \quad + \frac{\|h'\|}{\sqrt{n\bar{l}(\theta_0)}} \sqrt{\mathbb{E}[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2]} \sqrt{\mathbb{E}[l''(\theta_0; \mathbf{X}) + n\bar{l}(\theta_0)]^2}. \tag{6.54}
\end{aligned}$$

The bound depends on $\hat{\theta}_n(\mathbf{X})$ only through the MSE. With γ_1 as in (6.52) we denote by

$$\begin{aligned}
D_{x^2} &= \gamma_1 + \frac{2\|x\|}{\sqrt{n\bar{l}(\theta_0)}} \sqrt{\mathbb{E}[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2]} \sqrt{\mathbb{E}[l''(\theta_0; \mathbf{X}) + n\bar{l}(\theta_0)]^2} \\
& \quad + \frac{\|x\| S_d(n)}{\sqrt{n\bar{l}(\theta_0)}} \mathbb{E}[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2]
\end{aligned}$$

the bound in (6.54) for the case of $h(x) = x^2$. Following the same approach as in (3.56),

$$\mathbb{E}[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2] \leq \frac{1}{n\bar{l}(\theta_0)} (D_{x^2} + 1) \tag{6.55}$$

and since $D_{x,2}$ is an expression of the mean squared error and its positive root, our next step is to solve the quadratic inequality in (6.55) with unknown

$$U_d := U_d(\theta_0) := \sqrt{\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]}. \quad (6.56)$$

We have that

$$\begin{aligned} U_d^2 &\leq \frac{\|x\| S_d(n)}{[n\bar{i}(\theta_0)]^{\frac{3}{2}}} U_d^2 + \frac{2\|x\| \sqrt{\mathbb{E} \left[(l''(\theta_0; \mathbf{X}) + n\bar{i}(\theta_0))^2 \right]}}{[n\bar{i}(\theta_0)]^{\frac{3}{2}}} U_d + \frac{1 + \gamma_1}{n\bar{i}(\theta_0)} \\ &\Leftrightarrow \left(1 - \frac{\|x\| S_d(n)}{[n\bar{i}(\theta_0)]^{\frac{3}{2}}} \right) U_d^2 - \frac{2\|x\| \sqrt{\mathbb{E} \left[(l''(\theta_0; \mathbf{X}) + n\bar{i}(\theta_0))^2 \right]}}{[n\bar{i}(\theta_0)]^{\frac{3}{2}}} U_d - \frac{1 + \gamma_1}{n\bar{i}(\theta_0)} \leq 0 \end{aligned} \quad (6.57)$$

and the two solutions of this inequality are

$$\begin{aligned} U_{d_{1,2}} = \frac{1}{2D_1} &\left\{ \frac{2\|x\| \sqrt{\mathbb{E} \left[(l''(\theta_0; \mathbf{X}) + n\bar{i}(\theta_0))^2 \right]}}{[n\bar{i}(\theta_0)]^{\frac{3}{2}}} \right. \\ &\left. \pm \left[4 \frac{\|x\|^2 \mathbb{E} \left[(l''(\theta_0; \mathbf{X}) + n\bar{i}(\theta_0))^2 \right]}{[n\bar{i}(\theta_0)]^3} + 4 \frac{D_1 (1 + \gamma_1)}{n\bar{i}(\theta_0)} \right]^{\frac{1}{2}} \right\}. \end{aligned} \quad (6.58)$$

To have meaningful results we need the quantity

$$4 \frac{\|x\|^2}{[n\bar{i}(\theta_0)]^3} \mathbb{E} \left[(l''(\theta_0; \mathbf{X}) + n\bar{i}(\theta_0))^2 \right] + 4 \frac{D_1 (1 + \gamma_1)}{n\bar{i}(\theta_0)},$$

which is in the square root above, to be non-negative. This is satisfied through (F.A.2) which ensures that D_1 is positive. It follows that the solution with the negative square root in (6.58) is negative, while the other with the positive square root is positive and the solution of (6.57) is

$$0 \leq U_d = \sqrt{\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} \leq A_d.$$

Applying this result to the bound in (6.44) gives the second result of the theorem. \square

6.4.2 Bounded parameter space

In a number of situations the support of the distribution is not bounded and (F.A.1) is not satisfied. This subsection gives another approach to the problem of bounding the MSE of the MLE without using (F.A.1). Apart from (R.D.1)-(R.D.5), we still make two further assumptions though, and these are

(Fur.Ass.1) The parameter space is a closed (or semiclosed) interval of the form $[a, b]$

(or $[a, b), (a, b]$), with $a < b \in \mathbb{R}$. From now on, $\kappa = b - a$.

(Fur.Ass.2) With κ as in (Fur.Ass.1) $\exists N \in \mathbb{N}$ such that $\forall n \geq N$

$$n > \frac{\kappa S_d(n)}{2\bar{l}(\theta_0)}.$$

Denote by

$$\begin{aligned} D_2 &:= D_2(\theta_0, x, n) = 1 - \frac{\kappa S_d(n)}{2n\bar{l}(\theta_0)} \\ \gamma_2 &:= \gamma_2(\theta_0, x, n) = \frac{2}{n^{\frac{3}{2}}} \sum_{i=1}^n \sum_{j \in A_i} \sum_{k \in B_i} \left[\mathbb{E}((\xi_i)^4) \mathbb{E}((\xi_j)^4) \mathbb{E}((\xi_k)^4) \right]^{\frac{1}{4}} \\ &\quad + \frac{2}{n^{\frac{3}{2}}} \sum_{i=1}^n \sum_{j \in A_i} \sum_{k \in B_i} \left[\mathbb{E}((\xi_i)^2) \mathbb{E}((\xi_j)^2) \mathbb{E}((\xi_k)^2) \right]^{\frac{1}{2}} \\ &\quad + \frac{1}{n^{\frac{3}{2}}} \sum_{i=1}^n |A_i| \sum_{j \in A_i} \left[\mathbb{E}((\xi_i)^2) \mathbb{E}((\xi_j)^4) \right]^{\frac{1}{2}} + \left| \frac{\sqrt{\text{Var}[l'(\theta_0; \mathbf{X})]}}{\sqrt{n\bar{l}(\theta_0)}} - 1 \right|. \end{aligned} \quad (6.59)$$

The following theorem gives a bound on the Wasserstein distance related to the MLE in cases where it is not known analytically and the support of the distribution is not necessarily bounded. The approach is again based on finding an upper bound for the MSE of the MLE, which we then apply to the expression of the general upper bound.

Theorem 6.5. Let $\{X_i, i = 1, 2, \dots, n\}$ be an m -dependent sequence of random variables with joint probability density (or mass) function $f(\mathbf{x}|\boldsymbol{\theta})$ as in (6.1). Assume that the regularity conditions (R.D.1)-(R.D.5), as well as the assumptions (Fur.Ass.1), (Fur.Ass.2) are satisfied. Assume that the MLE exists and is unique. Then $A_2 = A_2(\boldsymbol{\theta}_0, n)$ is an upper bound for $\sqrt{\mathbb{E}[(\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0)^2]}$, where

$$A_2 = \frac{1}{2D_2} \left\{ \frac{\kappa \sqrt{\mathbb{E}[(l''(\boldsymbol{\theta}_0; \mathbf{X}) + n\bar{l}(\boldsymbol{\theta}_0))^2]}}{n\bar{l}(\boldsymbol{\theta}_0)} + \left[\frac{\kappa^2 \mathbb{E}[(l''(\boldsymbol{\theta}_0; \mathbf{X}) + n\bar{l}(\boldsymbol{\theta}_0))^2]}{[n\bar{l}(\boldsymbol{\theta}_0)]^2} + 4 \frac{D_2 \kappa \left(\gamma_2 + \sqrt{\frac{2}{\pi}} \right)}{\sqrt{n\bar{l}(\boldsymbol{\theta}_0)}} \right]^{\frac{1}{2}} \right\}. \quad (6.60)$$

with D_2 and γ_2 as in (6.59) and κ as in (Fur.Ass.1). Then if $\text{Var}[l'(\boldsymbol{\theta}_0; \mathbf{X})] > 0$,

$$\xi_1 = \frac{d}{d\boldsymbol{\theta}} \log f(X_1|\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \sqrt{\frac{n}{\text{Var}[l'(\boldsymbol{\theta}_0; \mathbf{X})]}}$$

and

$$\xi_i = \frac{d}{d\boldsymbol{\theta}} \log f(X_i|X_{i-1}, \dots, X_{i-m}; \boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \sqrt{\frac{n}{\text{Var}[l'(\boldsymbol{\theta}_0; \mathbf{X})]}}, i = 2, 3, \dots, n,$$

we have that

$$\begin{aligned} d_W \left(\sqrt{n\bar{l}(\boldsymbol{\theta}_0)}(\hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \boldsymbol{\theta}_0), Z \right) &\leq \frac{2}{n^{\frac{3}{2}}} \sum_{i=1}^n \sum_{j \in A_i} \sum_{k \in B_i} \left[\mathbb{E}((\xi_i)^4) \mathbb{E}((\xi_j)^4) \mathbb{E}((\xi_k)^4) \right]^{\frac{1}{4}} \\ &+ \frac{2}{n^{\frac{3}{2}}} \sum_{i=1}^n \sum_{j \in A_i} \sum_{k \in B_i} \left[\mathbb{E}((\xi_i)^2) \mathbb{E}((\xi_j)^2) \mathbb{E}((\xi_k)^2) \right]^{\frac{1}{2}} \\ &+ \frac{1}{n^{\frac{3}{2}}} \sum_{i=1}^n |A_i| \sum_{j \in A_i} \left[\mathbb{E}((\xi_i)^2) \mathbb{E}((\xi_j)^4) \right]^{\frac{1}{2}} + \left| \frac{\sqrt{\text{Var}[l'(\boldsymbol{\theta}_0; \mathbf{X})]}}{\sqrt{n\bar{l}(\boldsymbol{\theta}_0)}} - 1 \right| \\ &+ \frac{A_2}{\sqrt{n\bar{l}(\boldsymbol{\theta}_0)}} \left(\sqrt{\mathbb{E}[(l''(\boldsymbol{\theta}_0; \mathbf{X}) + n\bar{l}(\boldsymbol{\theta}_0))^2]} + \frac{S_d(n)A_2}{2} \right). \end{aligned} \quad (6.61)$$

Proof. We use (Fur.Ass.1) to get that

$$\begin{aligned} \frac{|\hat{\theta}_n(\mathbf{X}) - \theta_0|}{\kappa} \leq 1 &\Rightarrow \frac{1}{\kappa^2} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \leq \frac{1}{\kappa} \mathbb{E} |\hat{\theta}_n(\mathbf{X}) - \theta_0| \\ &\Rightarrow \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \leq \kappa \mathbb{E} |\hat{\theta}_n(\mathbf{X}) - \theta_0|. \end{aligned} \quad (6.62)$$

For $h(x) = |x|$, (6.62) yields

$$\begin{aligned} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] &\leq \kappa \mathbb{E} |\hat{\theta}_n(\mathbf{X}) - \theta_0| \\ &= \frac{\kappa}{\sqrt{n\bar{i}(\theta_0)}} \left| \mathbb{E} \left[h \left(\sqrt{n\bar{i}(\theta_0)} (\hat{\theta}_n(\mathbf{X}) - \theta_0) \right) \right] - \mathbb{E}[h(Z)] + \mathbb{E}[h(Z)] \right| \\ &\leq \frac{\kappa}{\sqrt{n\bar{i}(\theta_0)}} \left(B_{|x|} + \sqrt{\frac{2}{\pi}} \right), \end{aligned} \quad (6.63)$$

where

$$\begin{aligned} B_{|x|} &= \gamma_2 + \frac{S_d(n)}{2\sqrt{n\bar{i}(\theta_0)}} \mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right] \\ &\quad + \frac{\sqrt{\mathbb{E} \left[(l'''(\theta_0; \mathbf{X}) + n\bar{i}(\theta_0))^2 \right]}}{\sqrt{n\bar{i}(\theta_0)}} \sqrt{\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} \end{aligned} \quad (6.64)$$

is the bound in (6.54) for $h(x) = |x|$. Obviously, in this case $\|h'\| = 1$. Thus, for U_d as in (6.56) we can easily see that (6.63) and (6.64) yield

$$\begin{aligned} U_d^2 &\leq \frac{\kappa S_d(n)}{2n\bar{i}(\theta_0)} U_d^2 + \frac{\kappa \sqrt{\mathbb{E} \left[(l'''(\theta_0; \mathbf{X}) + n\bar{i}(\theta_0))^2 \right]}}{n\bar{i}(\theta_0)} U_d + \frac{\kappa \left(\gamma_2 + \sqrt{\frac{2}{\pi}} \right)}{\sqrt{n\bar{i}(\theta_0)}} \\ &\Leftrightarrow \left(1 - \frac{\kappa S_d(n)}{2n\bar{i}(\theta_0)} \right) U_d^2 - \frac{\kappa \sqrt{\mathbb{E} \left[(l'''(\theta_0; \mathbf{X}) + n\bar{i}(\theta_0))^2 \right]}}{n\bar{i}(\theta_0)} U_d - \frac{\kappa \left(\gamma_2 + \sqrt{\frac{2}{\pi}} \right)}{\sqrt{n\bar{i}(\theta_0)}} \leq 0 \end{aligned} \quad (6.65)$$

with the solutions of this simple quadratic inequality found to be

$$U_{d_{1,2}} = \frac{1}{2D_2} \left\{ \frac{\kappa \sqrt{\mathbb{E} \left[(l''(\theta_0; \mathbf{X}) + n\bar{i}(\theta_0))^2 \right]}}{n\bar{i}(\theta_0)} \pm \left[\frac{\kappa^2 \mathbb{E} \left[(l''(\theta_0; \mathbf{X}) + n\bar{i}(\theta_0))^2 \right]}{[n\bar{i}(\theta_0)]^2} + 4 \frac{D_2 \kappa \left(\gamma_2 + \sqrt{\frac{2}{\pi}} \right)}{\sqrt{n\bar{i}(\theta_0)}} \right]^{\frac{1}{2}} \right\}. \quad (6.66)$$

For meaningful results, D_2 has to be positive which is ensured through (Fur.Ass.2). This leads to one of the solutions in (6.66) being negative and the other one being positive.

The positive solution of (6.65) is

$$0 \leq U_d = \sqrt{\mathbb{E} \left[(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2 \right]} \leq A_2,$$

with A_2 as in (6.60). Using this result in the expression of the upper bound in (6.44) gives the second result of the theorem as expressed in (6.61). \square

A prominent example where Theorem 6.5 may be useful is parameter estimation in MA(q) time series; see Subsection 7.2.1 for more details.

Chapter 7

Conclusion

This chapter summarises the work carried out in this thesis and places the research in a wider context. An outline of what we consider as open problems on which we plan to work in the future is also provided. We subdivide these future work ideas into two parts. The first one is concerned with further assessing the performance of our bounds through various, possibly complicated applications, while in the second part we go beyond the MLE and the derivation of new bounds for other asymptotic results is discussed.

7.1 Summary

In this thesis, we have derived under various settings upper bounds on the distributional distance between the exact, unknown distribution of the MLE and the normal distribution. These bounds are explicit and novel. Among other famous mathematical tools, such as Taylor expansions, our work has been partly based on an elegant probabilistic technique; Stein's method. The purpose and the motivation for this research project have been set out in the Introduction.

In Chapter 2, we gave insights to the CLT; the basic theory and properties related

to the MLE; as well as a concise introduction to Stein's method. The assumptions ensuring existence and uniqueness of the MLE along with sufficient regularity conditions for the asymptotic normality of a scalar MLE for i.i.d. random variables were given. After presenting and proving the asymptotic result for the distribution of the MLE, Stein's method was concisely reviewed with Lemma 2.2 being one of the main results we used in Chapters 3 and 5 while Lemma 2.3 was used in Chapter 6 in order to obtain bounds related to single-parameter distributions.

In Chapter 3, we introduced the bounded Wasserstein distance, which is the metric used throughout the chapter and in Chapter 5. A connection of this metric with the Kolmogorov distance allowed us to give insights to confidence intervals. The main result of the chapter is given in Theorem 3.1 with the derived bound being, in terms of the sample size, of order $\frac{1}{\sqrt{n}}$. After the general results, we turned our focus on single-parameter exponential families, with the bound applied on the exponential distribution presented in both its canonical and non-canonical form. We explained why the results in the canonical case are simpler and sharper and this is illustrated on simulated data. A slightly more complicated case was then treated based on the issue of the MLE having positive probability of being on the boundary of the parameter space for discrete distributions. The solution of this problem relies on a perturbation method which introduces a perturbed parameter and a perturbed MLE that are now interior to the parameter space. As probably expected, this procedure created error terms which were also bounded, leading to an upper bound for such cases of order still $\frac{1}{\sqrt{n}}$. The Poisson distribution with parameter space $[0, \infty)$ was the example we attacked. In the last part of Chapter 3, we focused on bounding the mean squared error, $E \left[\left(\hat{\theta}_n(\mathbf{X}) - \theta_0 \right)^2 \right]$, by a quantity of which the calculation does not require to know the MLE explicitly. This approach gave us the opportunity of presenting a new bound on the bounded Wasserstein distance between the distribution of the MLE and the normal distribution which holds even in cases where an analytic expression of the MLE is not

available. The Beta distribution served as an illustration of these results.

In Chapter 4, our interest shifted to the multi-parameter setting. We gave upper bounds related to the distribution of a vector MLE and the multivariate normal. We distinguished between two cases depending on whether the random vectors are assumed to be heterogeneous or not. Different regularity conditions were given in each case in order for the asymptotic normality of the MLE to hold. Keeping the dimensionality of the parameter constant, the bounds are of the same order as those in the single-parameter case. The general family of linear regression models was treated separately, with the results being simple and easy to calculate as can be seen from the example of the straight-line regression. The normal distribution with $\theta_0 = (\mu, \sigma^2)$ served as a test case for identically distributed random variables. The idea of deriving bounds for situations where the MLE does not have a closed-form expression has been expanded to the multi-parameter setting. An upper bound on $E \left[\sum_{j=1}^d (\hat{\theta}_n(\mathbf{X})_j - \theta_{0j})^2 \right]$ that does not depend on the MLE was derived and applied on the general bound, which gives an upper bound of which the calculation does not rely upon $\hat{\theta}_n(\mathbf{X})$. As an example we used again the Beta distribution but now with both shape parameters unknown.

Chapter 5 presents an exploration of the way that the Delta method can be used to obtain bounds on the distributional distance of interest, constrained to the scalar parameter case. The bounds given are of the same order and similar to those in Chapter 3. However, in some cases they have a simpler representation and are also sharper. A special interest was again put on single-parameter exponential family distributions, which made the comparison to the results in Chapter 3 straightforward. The generalised gamma distribution with the scale parameter considered unknown was treated.

In Chapter 6, the independence assumption between the random variables has been relaxed and we introduced a local dependence structure. The notion of m -dependence

was explained and employing the results from Stein's method as in Lemma 2.3, we achieved to obtain bounds on the Wasserstein distance between the distribution of the MLE and the normal distribution. An example from locally dependent normal random variables showed that the bound behaves well in terms of simplicity and order; it is still of order $\frac{1}{\sqrt{n}}$. Under further assumptions that require boundedness of the support of the distribution or of the parameter space, we obtained as in Chapter 3 an upper bound for cases where the MLE can not be explicitly derived.

7.2 Open problems for future exploration

Having developed the theory related to upper bounds on the distributional distance between the distribution of the MLE and the normal distribution, a next step is to cover a number of important applications and also develop similar results for different kinds of random variables—either estimators or not—which follow asymptotically a known distribution that is not necessarily the normal. For ease of presentation this section is split into two parts depending on whether we discuss further applications or theory.

7.2.1 Further applications

Time Series is the first broad area we would like to discuss. A time series can be seen as a sequence of data, which has traditionally been of great importance when it comes to statistical analysis and forecasting. Maximum likelihood estimation is among the two most used estimation methods—the other one is least squares—for the parameters in time series models. Under sufficient conditions (Theorem 10.8.1 from Brockwell and Davis (1991)), the asymptotic normality of the MLE holds and its proof for general $\text{ARMA}(p, q)$ processes is given in Brockwell and Davis (1991), p.390. We have carried out preliminary work on the $\text{MA}(q)$ process $\{X_t, t = 0, \pm 1, \pm 2, \dots\}$, where

the relationship

$$X_t = Z_t + \theta_1 Z_{t-1} + \dots + \theta_q Z_{t-q} \quad (7.1)$$

holds, with the process $\{Z_t\} \sim \text{WN}(0, \sigma^2)$. The notation $\text{WN}(0, \sigma^2)$ stands for white noise and the random variables Z_t have mean zero, variance σ^2 and also $\text{Cov}(Z_t, Z_{t+a}) = 0$, for $a \in \mathbb{Z} \setminus \{0\}$. The $\text{MA}(q)$ process gives q -dependent random variables and (7.1) defines the local dependence structure between the random variables X_t . More specifically, from Brockwell and Davis (1991) for $\theta_0 = 1$,

$$\text{Cov}(X_t, X_{t+a}) = \begin{cases} \sigma^2 \sum_{j=0}^{q-|a|} \theta_j \theta_{j+|a|}, & \text{if } |a| \leq q \\ 0, & \text{if } |a| > q. \end{cases} \quad (7.2)$$

For the $\text{MA}(1)$ process, assume that Z_t are normally distributed and denote by $\Sigma = \text{E}(\mathbf{X} \mathbf{X}^\top)$ the $n \times n$ covariance matrix of $\mathbf{X} = (X_1, X_2, \dots, X_n)$. The likelihood function of \mathbf{X} is

$$L(\theta; \mathbf{X}) = \frac{1}{(2\pi)^{\frac{n}{2}} [\det(\Sigma)]^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} \mathbf{X}^\top \Sigma^{-1} \mathbf{X} \right\}.$$

The covariance matrix Σ is obtained using (7.2) and

$$\Sigma = \sigma^2 \begin{pmatrix} 1 + \theta^2 & \theta & 0 & 0 & 0 & \dots & 0 \\ \theta & 1 + \theta^2 & \theta & 0 & 0 & \dots & 0 \\ 0 & \theta & 1 + \theta^2 & \theta & 0 & \dots & 0 \\ 0 & 0 & \theta & 1 + \theta^2 & \theta & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \theta & 1 + \theta^2 & \theta \\ 0 & 0 & 0 & 0 & 0 & \theta & 1 + \theta^2 \end{pmatrix}.$$

The MLE for θ exists but there is not an analytic expression available and numerical methods are used to maximize the likelihood. We could employ the results of Section

6.4 to find an upper bound on the distributional distance of interest under this time series framework of locally dependent random variables.

Another open problem is to obtain bounds related to the popular and rapidly expanding area of networks. Such a network is the Erdős-Renyi (Bernoulli) random graph, where an edge between two vertices is drawn with probability p , independently of the other edges. Suppose we have n vertices. The degree of a randomly picked vertex follows the Binomial($n - 1, p$) distribution with unknown parameter p . Thus, let $X_1, X_2, \dots, X_n \sim \text{Bin}(n - 1, p)$ be i.i.d. random variables for the degree of the vertices $1, 2, \dots, n$ respectively with $p \in [0, 1]$. The MLE for the parameter p in this case is $\hat{\theta}_n(\mathbf{X}) = \frac{1}{n(n-1)} \sum_{i=1}^n X_i$. Since we have a discrete distribution with compact parameter space, there is positive probability of the MLE being on the boundary of the parameter space. Thus, the perturbation approach as explained in Section 3.4 needs to be adapted before finding the required upper bound. Treating this case could lead to work on more complex examples, one of which is the geometric random graph. Two vertices are connected if and only if the distance between them is at most a radius r . The idea would be to estimate this radius and apply the theory for locally dependent random variables to the distribution of the estimator. In Raič (2004), relevant bounds on the distributional distance between the distribution of $W = \sum_{i=1}^n X_i$ and the normal have been obtained in the case of a dependence graph, with $X_i, i = 1, 2, \dots, n$ the degree of the vertex i . These results could be useful when it comes to bounding the relevant distributional distance for the MLE which can be of any form and not restricted to a sum of random variables.

The extension of the Delta method approach in order to treat multi-parameter distributions is also an interesting idea for future work. For this purpose, let $q : \Theta \rightarrow \mathbb{R}^d$, such as all the entries of the function q have continuous partial derivatives with respect to θ . In addition, we require that $q(\hat{\theta}_n(\mathbf{X})) = \frac{1}{n} \sum_{i=1}^n g(X_i)$, for some function $g : \mathbb{R} \rightarrow \mathbb{R}^d$. We use a multivariate generalisation of the Delta method on the result of Theorem 4.1 for

the asymptotic normality of the MLE. We denote by $K := [\nabla q(\boldsymbol{\theta}_0)][I(\boldsymbol{\theta}_0)]^{-1}[\nabla q(\boldsymbol{\theta}_0)]^\top$, where $\nabla q(\boldsymbol{\theta}_0)$ is just the $d \times d$ Jacobian matrix of the function q evaluated at $\boldsymbol{\theta}_0$. The Delta method gives

$$\sqrt{n} \left(q(\hat{\boldsymbol{\theta}}_n(\mathbf{X})) - q(\boldsymbol{\theta}_0) \right) \xrightarrow[n \rightarrow \infty]{d} K^{\frac{1}{2}} \mathbf{Z}.$$

Assuming that the Jacobian matrix of q is invertible, preliminary work has shown that the approach to be followed in order to find an upper bound for the quantity of interest is a combination of the steps explained in Chapters 4 and 5.

7.2.2 Beyond the classical framework of the MLE

The asymptotic theory for $\hat{\boldsymbol{\theta}}_n(\mathbf{X})$ allows to Taylor expand its bias and variance terms using expressions based on the sample size, n . For instance, according to Ferrari et al. (1996), expansions for the bias, $B(\hat{\boldsymbol{\theta}}_n(\mathbf{X}))$, and variance of the MLE to order n^{-2} can be written as

$$B(\hat{\boldsymbol{\theta}}_n(\mathbf{X})) = \frac{B_1(\hat{\boldsymbol{\theta}}_n(\mathbf{X}))}{n} + \frac{B_2(\hat{\boldsymbol{\theta}}_n(\mathbf{X}))}{n^2}, \quad \text{Var}(\hat{\boldsymbol{\theta}}_n(\mathbf{X})) = \frac{V_1(\hat{\boldsymbol{\theta}}_n(\mathbf{X}))}{n} + \frac{V_2(\hat{\boldsymbol{\theta}}_n(\mathbf{X}))}{n^2},$$

where closed formulas for the functions B_1 , B_2 , V_1 and V_2 can be easily derived. Using the above expansions, one can define various types of bias-corrected estimators, with the most straightforward being

$$\tilde{\boldsymbol{\theta}}_1(\mathbf{X}) = \hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \frac{B_1(\hat{\boldsymbol{\theta}}_n(\mathbf{X}))}{n}, \quad \tilde{\boldsymbol{\theta}}_2(\mathbf{X}) = \hat{\boldsymbol{\theta}}_n(\mathbf{X}) - \frac{B_1(\hat{\boldsymbol{\theta}}_n(\mathbf{X}))}{n} - \frac{B_2(\hat{\boldsymbol{\theta}}_n(\mathbf{X}))}{n^2}.$$

These modified estimators are bias-free to order $\frac{1}{n}$ as shown in Ferrari et al. (1996). An open problem raised directly from the results of the current thesis is to find upper bounds for the distributional distance related to these bias-corrected estimators. We would expect the constant of the bound to be better than the one linked with the MLE.

In addition, the order of the new bounds with respect to the sample size should also be examined; is it better than the one for the bound related to the MLE?

Another interesting idea for future work is to expand the approach explained in this thesis to a broad family of parametric estimators, called M -estimators. These are robust estimators defined as a root of an estimating function. This estimating function is usually the derivative of another statistical function, $\rho(x, \theta)$. An M -estimator is the minimum of $\rho(x, \theta)$. To set the scene in the case of independent random variables $X_i, i = 1, 2, \dots, n$, the M -estimator $\hat{\theta}_n^M(\mathbf{X})$ is the solution that minimises $\sum_{i=1}^n \rho(X_i, \theta)$. The choice of the function ρ differs—for example $\rho(x_i, \theta) = -L(\theta; x_i)$ in the case of the MLE—and when it is differentiable, the computation of $\hat{\theta}_n^M(\mathbf{X})$ comes from deriving the critical points of $\frac{d}{d\theta} \sum_{i=1}^n \rho(x_i, \theta)$. Under again some sufficient regularity conditions, it has been shown that the distribution of the M -estimators is asymptotically normal; see Huber (1964) for the first proof related to location parameters. The assessment of the quality of this asymptotic result remains an open problem. One could follow an approach based again on Taylor expansions of $0 = \sum_{i=1}^n \frac{d}{d\theta} \rho(x_i, \hat{\theta}_n^M(x))$ about the true value of the unknown parameter θ_0 . Then, Stein's method for normal approximation can be employed in order to find an upper bound on the distributional distance between the distribution of $W_M = \sum_{i=1}^n \frac{d}{d\theta} \rho(X_i, \theta_0)$ and the normal distribution. Alternative techniques can then be developed to find upper bounds for all the remainder terms.

A different, interesting direction is towards Bayesian Statistics. It has already been shown in the statistical literature (see for example van der Vaart (1998)) that under specific assumptions both the posterior measures and the Bayes point estimators are asymptotically normal. Specifically, according to the Bernstein-von Mises theorem as expressed in van der Vaart (1998), under remarkably weak conditions, the posterior distribution of $\sqrt{n}(\Theta_n - \theta_0)$ is under θ_0 asymptotically equal to $N\left(Y, \frac{1}{i(\theta_0)}\right)$. In this case, Y follows the $N\left(0, \frac{1}{i(\theta_0)}\right)$ distribution and Θ_n is the “random” parameter,

which in the Bayesian perspective is a random variable. Furthermore, under stronger regularity conditions we obtain that the distribution of Θ_n is approximately equal to $N\left(\hat{\theta}_n(\mathbf{X}), \frac{1}{ni(\theta_0)}\right)$, where $\hat{\theta}_n(\mathbf{X})$ is the MLE. This is again an asymptotic result of which the quality could be assessed through upper bounds on the distributional distance between the distribution of $\sqrt{n}(\Theta_n - \theta_0)$ and the normal. We would expect that new horizons in applications will be opened if such results in the area of Bayesian Statistics become available in the near future.

The log-likelihood ratio statistic is used in the so-called likelihood ratio test where the aim is to distinguish between two nested models. In the case of a scalar parameter, aiming to test whether $\theta = \theta_0$ or not, it is well known that

$$2(l(\hat{\theta}_n(\mathbf{X}); \mathbf{X}) - l(\theta_0; \mathbf{X})) \xrightarrow{d} \chi_1^2. \quad (7.3)$$

In the case of a vector parameter and composite hypothesis testing of whether $\theta \in \Theta_0$ or not, for Θ_0 a subset of the parameter space Θ , we get that

$$2(l(\hat{\theta}_2(\mathbf{X}); \mathbf{X}) - l(\hat{\theta}_1(\mathbf{X}); \mathbf{X})) \xrightarrow{d} \chi_d^2, \quad (7.4)$$

where $\hat{\theta}_1(\mathbf{X})$ and $\hat{\theta}_2(\mathbf{X})$ are the MLEs for $\theta \in \Theta_0$ and $\theta \in \Theta$, respectively. In addition, d is the difference in the dimensionality of Θ and Θ_0 . The result in (7.4) is known as Wilks's theorem and it was first proved in Wilks (1938). One could work on assessing the quality of the asymptotic results in (7.3) and (7.4). An idea of the process we have in mind is to split again the quantity of interest into two easier to handle terms. The first of these two terms could be upper bounded using Stein's method for χ^2 approximation; specifically the results in Gaunt et al. (2015) can be useful, where explicit bounds of order $\frac{1}{n}$ for smooth test functions are given. For the second term, one could exploit further expansions and known probability inequalities, such as the Markov and Cauchy-Schwarz inequalities.

References

- R. Arratia, L. Goldstein, and L. Gordon. Poisson Approximation and the Chen-Stein Method. *Statistical Science*, **5**:403–434, 1990.
- A. D. Barbour. Stein’s method for diffusion approximations. *Probability Theory and Related Fields*, **84**:297–322, 1990.
- A. D. Barbour and L. Holst. Some applications of the Stein-Chen method for proving Poisson convergence. *Advances in Applied Probability*, **21**:74–90, 1989.
- A. D. Barbour, L. Holst, and S. Janson. *Poisson Approximation*. Oxford University Press, 1992.
- A. D. Barbour, H. L. Gan, and A. Xia. Stein factors for negative binomial approximation in wasserstein distance. *Bernoulli*, **21**:1002–1013, 2015.
- R. H. Berk. Consistency and asymptotic normality of MLE’s for exponential models. *The Annals of Mathematical Statistics*, **43**:193–204, 1972.
- B. R. Bhat. On the Method of Maximum-Likelihood for Dependent Observations. *Journal of the Royal Statistical Society. Series B (Methodological)*, **36**:48–53, 1974.
- P. Billingsley. Statistical Methods in Markov Chains. *The Annals of Mathematical Statistics*, **32**:12–40, 1961.
- P. J. Brockwell and R. A. Davis. *Time Series: Theory and Methods*. Springer Series in Statistics. Springer-Verlag, New York, second edition, 1991.
- G. Casella and R. L. Berger. *Statistical Inference*. Brooks/Cole, Cengage Learning, Duxbury, Pacific Grove, second edition, 2002.
- S. Chatterjee and E. Meckes. Multivariate normal approximation using exchangeable pairs. *ALEA, Latin American Journal of Probability and Mathematical Statistics*, **4**: 257–283, 2008.
- S. Chatterjee, J. Fulman, and A. Röllin. Exponential approximation by exchangeable pairs and spectral graph theory. *ALEA, Latin American Journal of Probability and Mathematical Statistics*, **8**:1–27, 2011.
- L. H. Y. Chen. Poisson approximation for dependent trials. *Annals of Probability*, **3**: 534–545, 1975.
- L. H. Y. Chen, L. Goldstein, and Q. M. Shao. *Normal Approximation by Stein’s Method*. Springer-Verlag, Berlin Heidelberg, 2011.

- D. R. Cox and E. J. Snell. A General Definition of Residuals. *Journal of the Royal Statistical Society*, **30**:248–275, 1968.
- A. C. Davison. *Statistical Models*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2008.
- P. Diaconis. *An example for Stein's method*. Stanford University, Department of Statistics, Technical Report, 1989.
- L. H. Dicker and M. A. Erdogdu. Flexible results for quadratic forms with applications to variance components estimation. *Annals of Statistics*. To appear, 2016+.
- C. Döbler. Stein's method of exchangeable pairs for the Beta distribution and generalizations. *Electronic Journal of Probability*, **20**:1–34, 2015.
- W. Ehm. Binomial approximation to the Poisson binomial distribution. *Statistics & Probability Letters*, **11**:7–16, 1991.
- L. Fahrmeir and H. Kaufmann. Consistency and asymptotic normality of the maximum likelihood estimator in generalized linear models. *The Annals of Statistics*, **13**: 342–368, 1985.
- S. L. P. Ferrari, D. A. Botter, G. M. Cordeiro, and F. Cribari-Neto. Second and third order bias reduction for one-parameter family models. *Statistics and Probability Letters*, **30**:339–345, 1996.
- A. M. Fink and M. Jodeit. On Chebyshev's other inequality. *Lecture Notes-Monograph Series, Inequalities in Statistics and Probability*, **5**:115–120, 1984.
- R. A. Fisher. Theory of Statistical Estimation. *Mathematical Proceedings of the Cambridge Philosophical Society*, **22**:700–725, 1925.
- R. E. Gaunt. *Rates of Convergence of Variance-Gamma Approximations via Stein's Method*. PhD thesis, University of Oxford, 2013.
- R. E. Gaunt, A. Pickett, and G. Reinert. Chi-square approximation by Stein's method with application to Pearson's statistic. <http://arxiv.org/pdf/1507.01707v1.pdf>, 2015.
- C. J. Geyer. Asymptotics of exponential families. University lecture notes, <http://www.stat.umn.edu/geyer/8112/notes/expfam.pdf>, 2013.
- L. Goldstein and G. Reinert. Stein's method for the beta distribution and the Polya-Eggenberger urn. *Journal of Applied Probability*, **50**:1187–1205, 2013.
- F. Götze. On the rate of convergence in the multivariate CLT. *Annals of Probability*, **19**:724–739, 1991.
- L. Heinrich. A Method for the Derivation of Limit Theorems for Sums of m -dependent Random Variables. *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete*, **60**:501–515, 1982.
- B. Hoadley. Asymptotic Properties of Maximum Likelihood Estimators for the Independent Not Identically Distributed Case. *The Annals of Mathematical Statistics*, **42**: 1977–1991, 1971.

- P. J. Huber. Robust Estimation of a Location Parameter. *The Annals of Mathematical Statistics*, **35**:73–101, 1964.
- M. G. Kendall and A. Stuart. *The Advanced Theory of Statistics*, volume 1 Distribution Theory. Charles Griffin and Company Limited, London, third edition, 1969.
- N. M. Kiefer. Maximum likelihood estimation (mle). University Lecture, <https://courses.cit.cornell.edu/econ620/reviewm5.pdf>, 2008.
- S. Lauritzen. *Extremal Families and Systems of Sufficient Statistics*. Lecture Notes in Statistics, No.49. Springer-Verlag, Berlin-Heidelberg-New York, 1988.
- C. Ley and Y. Swan. Stein’s density approach and information inequalities. *Electronic Communications in Probability*, **18**:1–14, 2013a.
- C. Ley and Y. Swan. Local Pinsker inequalities via Stein’s discrete density approach. *IEEE Transactions on Information Theory*, **59**:5584–5591, 2013b.
- C. Ley, G. Reinert, and Y. Swan. Approximate computation of expectations: a canonical Stein operator. <http://arxiv.org/pdf/1408.2998.pdf>, 2014.
- C. Ley, G. Reinert, and Y. Swan. Distances between nested densities and a measure of the impact of the prior in Bayesian Statistics. <http://arxiv.org/pdf/1510.05826v1.pdf>, 2015.
- H. M. Luk. *Stein’s method for the gamma distribution and related statistical applications*. PhD thesis, University of Southern California, Los Angeles, USA, 1994.
- T. Makelainen, K. Schmidt, and G. P. H. Styan. On the existence and uniqueness of the maximum likelihood estimate of a vector-valued parameter in fixed-size samples. *The Annals of Statistics*, **9**:758–767, 1981.
- E. Meckes. On Stein’s method for multivariate normal approximation. *IMS Collections, High Dimensional Probability V: The Luminy Volume*, **5**:153–178, 2009.
- I. Nourdin and G. Peccati. Stein’s method on Wiener chaos. *Probability Theory and Related Fields*, **145**:75–118, 2009.
- I. Nourdin and G. Peccati. *Normal Approximations with Malliavin Calculus. From Stein’s method to universality*. Cambridge Tracts in Mathematics, No.192. Cambridge University Press, 2012.
- I. Nourdin, G. Peccati, and Y. Swan. Integration by parts and representation of information functionals. *Proceedings of the 2014 IEEE International Symposium on Information Theory (ISIT)*, pages 2217–2221, 2014.
- E. Peköz. Stein’s method for geometric approximation. *Journal of Applied Probability*, **33**:707–713, 1996.
- E. Peköz and A. Röllin. New rates for exponential approximation and the theorems of Rényi and Yaglom. *The Annals of Probability*, **39**:587–608, 2011.
- A. Pickett. *Rates of Convergence of χ^2 approximations via Stein’s method*. PhD thesis, University of Oxford, 2004.

- J. Pike and H. Ren. Stein's method and the Laplace distribution. *ALEA, Latin American Journal of Probability and Mathematical Statistics*, **11**:571–587, 2014.
- I. Pinelis and R. Molzon. Optimal-order bounds on the rate of convergence to normality in the multivariate delta method. *Electronic Journal of Statistics*, **10**:1001–1063, 2016.
- M. Raič. A Multivariate CLT for Decomposable Random Vectors with Finite Second Moments. *Journal of Theoretical Probability*, **17**:573–603, 2004.
- G. Reinert. A Weak Law of Large Numbers for Empirical Measures via Stein's Method. *The Annals of Probability*, **23**:334–354, 1995.
- G. Reinert. Couplings for normal approximations with Stein's method. *Microsurveys in Discrete Probability*. D. Aldous, J. Propp eds., *Dimacs series*. AMS, pages 193–207, 1998.
- G. Reinert and A. Röllin. Multivariate normal approximation with Stein's method of exchangeable pairs under a general linearity condition. *The Annals of Probability*, **37**:2150–2173, 2009.
- N. Ross. Fundamentals of Stein's method. *Probability Surveys*, **8**:210–293, 2011.
- I. Shevtsova. On the absolute constants in the Berry Esseen type inequalities. *Doklady Mathematics*, **89**:378–381, 2014.
- S. D. Silvey. A Note on Maximum-Likelihood in the Case of Dependent Random Variables. *Journal of the Royal Statistical Society. Series B (Methodological)*, **23**:444–452, 1961.
- E. C. Somayya. A Method for finding a Square Root of a 2×2 Matrix. *The Mathematics Education*, **XXXI**:52–54, 1997.
- C. Stein. A bound for the error in the normal approximation to the distribution of a sum of dependent random variables. In *Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability*, volume **2**, pages 586–602. Berkeley: University of California Press, 1972.
- C. Stein. *Approximate Computation of Expectations*. Lecture Notes-Monograph Series. Institute of Mathematical Statistics, Hayward, California, 1986.
- C. Stein, P. Diaconis, S. Holmes, and G. Reinert. Use of exchangeable pairs in the analysis of simulations. *IMS Lecture Notes-Monograph Series*, **46**:1–25, 2004.
- F. Tricomi and A. Erdélyi. The asymptotic expansion of a ratio of gamma functions. *Pacific Journal of Mathematics*, **1**:133–142, 1951.
- A. W. van der Vaart. *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 1998.
- S. S. Wilks. The Large-Sample Distribution of the Likelihood Ratio for Testing Composite Hypotheses. *The Annals of Mathematical Statistics*, **9**:60–62, 1938.