
What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 There are two major types of uncertainty one can model. *Aleatoric* uncertainty
2 captures noise inherent in the observations. On the other hand, *epistemic* uncer-
3 tainty accounts for uncertainty in the model – uncertainty which can be explained
4 away given enough data. Traditionally it has been difficult to model epistemic
5 uncertainty in computer vision, but with new Bayesian deep learning tools this
6 is now possible. We study the benefits of modeling epistemic vs. aleatoric un-
7 certainty in Bayesian deep learning models for vision tasks. For this we present
8 a Bayesian deep learning framework combining input-dependent aleatoric uncer-
9 tainty together with epistemic uncertainty. We study models under the framework
10 with per-pixel semantic segmentation and depth regression tasks. Further, our
11 explicit uncertainty formulation leads to new loss functions for these tasks, which
12 can be interpreted as learned attenuation. This makes the loss more robust to noisy
13 data, also giving new state-of-the-art results on segmentation and depth regression
14 benchmarks.

15 1 Introduction

16 Understanding what a model does not know is a critical part of many machine learning systems.
17 Today, deep learning algorithms are able to learn powerful representations which can map high di-
18 mensional data to an array of outputs. However these mappings are often taken blindly and assumed
19 to be accurate, which is not always the case. In two recent examples this has had disastrous con-
20 sequences. In May 2016 there was the first fatality from an assisted driving system, caused by the
21 perception system confusing the white side of a trailer for bright sky [1]. In a second recent ex-
22 ample, an image classification system erroneously identified two African Americans as gorillas [2],
23 raising concerns of racial discrimination. If both these algorithms were able to assign a high level
24 of uncertainty to their erroneous predictions, then the system may have been able to make better
25 decisions and likely avoid disaster.

26 Quantifying uncertainty in computer vision applications can be largely divided into regression set-
27 tings such as depth regression, and classification settings such as semantic segmentation. Existing
28 approaches to model uncertainty in such settings in computer vision include particle filtering and
29 conditional random fields [3, 4]. However many modern applications mandate the use of *deep learn-*
30 *ing* to achieve state-of-the-art performance [5], with most deep learning models not able to represent
31 uncertainty. Deep learning does not allow for uncertainty representation in regression settings for
32 example, and deep learning classification models often give normalised score vectors, which do not
33 necessarily capture model uncertainty. For both settings uncertainty can be captured with *Bayesian*
34 *deep learning* approaches – which offer a practical framework for understanding uncertainty with
35 deep learning models [6].

36 In Bayesian modeling, there are two main types of uncertainty one can model [7]. *Aleatoric* uncer-
37 tainty captures noise inherent in the observations. This could be for example sensor noise or motion
38 noise, resulting in uncertainty which cannot be reduced even if more data were to be collected. On

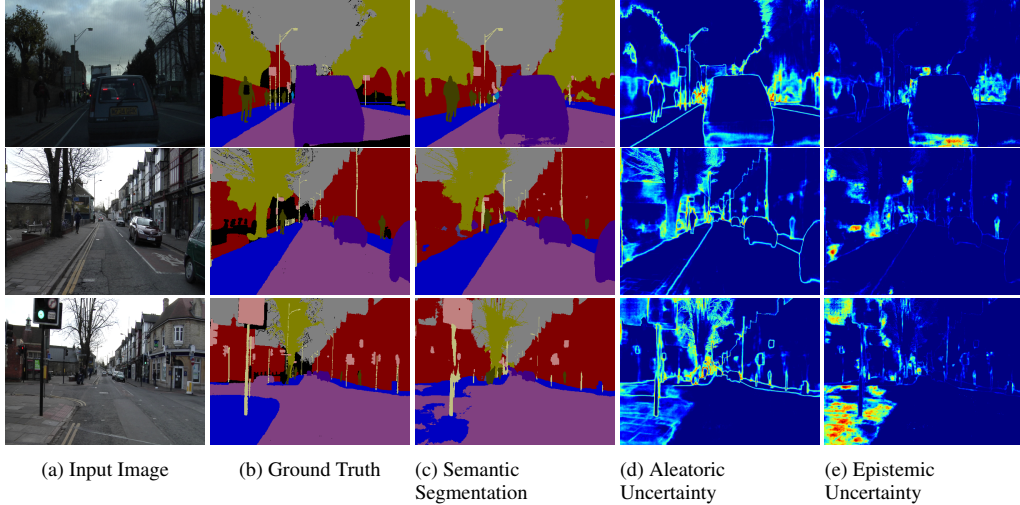


Figure 1: **Illustrating the difference between aleatoric and epistemic uncertainty** for semantic segmentation on the CamVid dataset. We observe aleatoric uncertainty captures object boundaries where labels are noisy. The bottom row shows a failure case of the segmentation model, when the model fails to segment the footpath, and the corresponding increased epistemic uncertainty.

39 the other hand, *epistemic* uncertainty accounts for uncertainty in the model parameters – uncertainty
 40 which captures our ignorance about which model generated our collected data. This uncertainty
 41 can be explained away given enough data, and is often referred to as *model uncertainty*. Aleatoric
 42 uncertainty can further be categorized into *homoscedastic* uncertainty, uncertainty which stays con-
 43 stant for different inputs, and *heteroscedastic* uncertainty. Heteroscedastic uncertainty depends on
 44 the inputs to the model, with some inputs potentially having more noisy outputs than others. Het-
 45 eroscedastic uncertainty is especially important for computer vision applications. For example, for
 46 depth regression, highly textured input images with strong vanishing lines are expected to result in
 47 confident predictions, whereas an input image of a featureless wall is expected to have very high
 48 uncertainty.

49 In this paper we make the observation that in many big data regimes (such as the ones common
 50 to deep learning with image data), it is most effective to model aleatoric uncertainty, uncertainty
 51 which cannot be explained away. This is in comparison to epistemic uncertainty which is mostly
 52 explained away with the large amounts of data often available in machine vision. We further show
 53 that modeling aleatoric uncertainty alone comes at a cost. Out-of-data examples, which can be
 54 identified with epistemic uncertainty, cannot be identified with aleatoric uncertainty alone.

55 For this we present a unified Bayesian deep learning framework which allows us to learn map-
 56 pings from input data to aleatoric uncertainty and compose these together with epistemic uncer-
 57 tainty approximations. We derive our framework for both regression and classification applications
 58 and present results for per-pixel depth regression and semantic segmentation tasks (see Figure 1 and
 59 the supplementary video for examples). We show how modeling aleatoric uncertainty in regression
 60 can be used to learn loss attenuation, and develop a complimentary approach for the classification
 61 case. This demonstrates the efficacy of our approach on difficult and large scale tasks.

62 The main contributions of this work are;

- 63 1. We capture an accurate understanding of aleatoric and epistemic uncertainties, in particular
 64 with a novel approach for classification,
- 65 2. We improve model performance by 1 – 3% over non-Bayesian baselines by reducing the
 66 effect of noisy data with the implied attenuation obtained from explicitly representing
 67 aleatoric uncertainty,
- 68 3. We study the trade-offs between modeling aleatoric or epistemic uncertainty by character-
 69 izing the properties of each uncertainty and comparing model performance and inference
 70 time.

71 **2 Related Work**

72 Existing approaches to Bayesian deep learning capture either epistemic uncertainty alone, or
 73 aleatoric uncertainty alone [6]. These uncertainties are formalised as probability distributions over
 74 either the model parameters, or model outputs, respectively. Epistemic uncertainty is modeled by
 75 placing a prior distribution over a model’s weights, and then trying to capture how much these
 76 weights vary given some data. Aleatoric uncertainty on the other hand is modeled by placing a dis-
 77 tribution over the output of the model. For example, in regression our outputs might be modeled as
 78 corrupted with Gaussian random noise. In this case we are interested in learning the noise’s variance
 79 as a function of different inputs (such noise can also be modeled with a constant value for all data
 80 points, but this is of less practical interest). These uncertainties, in the context of Bayesian deep
 81 learning, are explained in more detail in this section.

82 **2.1 Epistemic Uncertainty in Bayesian Deep Learning**

83 To capture epistemic uncertainty in a neural network (NN) we put a prior distribution over its
 84 weights, for example a Gaussian prior distribution: $\mathbf{W} \sim \mathcal{N}(0, I)$.

85 Such a model is referred to as a Bayesian neural network (BNN) [8–10]. Bayesian neural networks
 86 replace the deterministic network’s weight parameters with distributions over these parameters, and
 87 instead of optimising the network weights directly we average over all possible weights (referred
 88 to as *marginalisation*). Denoting the random output of the BNN as $\mathbf{f}^{\mathbf{W}}(\mathbf{x})$, we define the model
 89 likelihood $p(\mathbf{y}|\mathbf{f}^{\mathbf{W}}(\mathbf{x}))$. Given a dataset $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$, Bayesian inference
 90 is used to compute the posterior over the weights $p(\mathbf{W}|\mathbf{X}, \mathbf{Y})$. This posterior captures the set of
 91 plausible model parameters, given the data.

92 For regression tasks we often define our likelihood as a Gaussian with mean given by the model
 93 output: $p(\mathbf{y}|\mathbf{f}^{\mathbf{W}}(\mathbf{x})) = \mathcal{N}(\mathbf{f}^{\mathbf{W}}(\mathbf{x}), \sigma^2)$, with an observation noise scalar σ . For classification, on
 94 the other hand, we often squash the model output through a softmax function, and sample from the
 95 resulting probability vector: $p(\mathbf{y}|\mathbf{f}^{\mathbf{W}}(\mathbf{x})) = \text{Softmax}(\mathbf{f}^{\mathbf{W}}(\mathbf{x}))$.

96 BNNs are easy to formulate, but difficult to perform inference in. This is because the marginal
 97 probability $p(\mathbf{Y}|\mathbf{X})$, required to evaluate the posterior $p(\mathbf{W}|\mathbf{X}, \mathbf{Y}) = p(\mathbf{Y}|\mathbf{X}, \mathbf{W})p(\mathbf{W})/p(\mathbf{Y}|\mathbf{X})$,
 98 cannot be evaluated analytically. Different approximations exist [11–14]. In these approximate
 99 inference techniques, the posterior $p(\mathbf{W}|\mathbf{X}, \mathbf{Y})$ is fitted with a simple distribution $q_{\theta}^*(\mathbf{W})$, param-
 100 eterised by θ . This replaces the intractable problem of averaging over all weights in the BNN with an
 101 optimisation task, where we seek to optimise over the *parameters of the simple distribution* instead
 102 of optimising the original neural network’s parameters.

103 Dropout variational inference is a practical approach for approximate inference in large and complex
 104 models [14]. This inference is done by training a model with dropout before every weight layer,
 105 and by also performing dropout at test time to sample from the approximate posterior (stochastic
 106 forward passes, referred to as Monte Carlo dropout). More formally, this approach is equivalent
 107 to performing approximate variational inference where we find a simple distribution $q_{\theta}^*(\mathbf{W})$ in a
 108 tractable family which minimises the Kullback-Leibler (KL) divergence to the true model posterior
 109 $p(\mathbf{W}|\mathbf{X}, \mathbf{Y})$. Dropout can be interpreted as a variational Bayesian approximation, where the ap-
 110 proximating distribution is a mixture of two Gaussians with small variances and the mean of one of
 111 the Gaussians is fixed at zero. The minimisation objective is given by [15]:

$$\mathcal{L}(\theta, p) = -\frac{1}{N} \sum_{i=1}^N \log p(\mathbf{y}_i|\mathbf{f}^{\widehat{\mathbf{W}}_i}(\mathbf{x}_i)) + \frac{1-p}{2N} \|\theta\|^2 \quad (1)$$

112 with N data points, dropout probability p , samples $\widehat{\mathbf{W}}_i \sim q_{\theta}^*(\mathbf{W})$, and θ the set of the simple
 113 distribution’s parameters to be optimised (weight matrices in dropout’s case). In regression, for
 114 example, the negative log likelihood can be further simplified as

$$-\log p(\mathbf{y}_i|\mathbf{f}^{\widehat{\mathbf{W}}_i}(\mathbf{x}_i)) \propto \frac{1}{2\sigma^2} \|\mathbf{y}_i - \mathbf{f}^{\widehat{\mathbf{W}}_i}(\mathbf{x}_i)\|^2 + \log \sigma^2 \quad (2)$$

115 for a Gaussian likelihood, with σ the model’s observation noise parameter – capturing how much
 116 noise we have in the outputs.

117 Epistemic uncertainty in the weights can be reduced by observing more data. This uncertainty in-
 118 duces prediction uncertainty by marginalising over the (approximate) weights posterior distribution.

119 For classification this can be approximated using Monte Carlo integration as follows:

$$p(y = c|\mathbf{x}, \mathbf{X}, \mathbf{Y}) \approx \frac{1}{T} \sum_{t=1}^T \text{Softmax}(\mathbf{f}^{\widehat{\mathbf{W}}_t}(\mathbf{x})) \quad (3)$$

120 with T sampled masked model weights $\widehat{\mathbf{W}}_t \sim q_{\theta}^*(\mathbf{W})$, where $q_{\theta}(\mathbf{W})$ is the Dropout distribution
 121 [6]. The uncertainty of this probability vector \mathbf{p} can then be summarised using the entropy of the
 122 probability vector: $H(\mathbf{p}) = -\sum_{c=1}^C p_c \log p_c$. For regression this epistemic uncertainty is captured
 123 by the predictive variance, which can be approximated as:

$$\text{Var}(\mathbf{y}) \approx \sigma^2 + \frac{1}{T} \sum_{t=1}^T \mathbf{f}^{\widehat{\mathbf{W}}_t}(\mathbf{x})^T \mathbf{f}^{\widehat{\mathbf{W}}_t}(\mathbf{x}_t) - E(\mathbf{y})^T E(\mathbf{y}) \quad (4)$$

124 with predictions in this epistemic model done by approximating the predictive mean: $E(\mathbf{y}) \approx$
 125 $\frac{1}{T} \sum_{t=1}^T \mathbf{f}^{\widehat{\mathbf{W}}_t}(\mathbf{x})$. The first term in the predictive variance, σ^2 , corresponds to the amount of noise
 126 inherent in the data (which will be explained in more detail soon). The second part of the predictive
 127 variance measures how much the model is uncertain about its predictions – this term will vanish
 128 when we have zero parameter uncertainty (i.e. when all draws $\widehat{\mathbf{W}}_t$ take the same constant value).

129 2.2 Heteroscedastic Aleatoric Uncertainty

130 In the above we captured model uncertainty – uncertainty over the model parameters – by approxi-
 131 mating the distribution $p(\mathbf{W}|\mathbf{X}, \mathbf{Y})$. To capture aleatoric uncertainty in regression, we would have
 132 to tune the observation noise parameter σ .

133 Homoscedastic regression assumes constant observation noise σ for every input point \mathbf{x} . Het-
 134 eroscedastic regression, on the other hand, assumes that observation noise can vary with input \mathbf{x}
 135 [16, 17]. Heteroscedastic models are useful in cases where parts of the observation space might
 136 have higher noise levels than others. In non-Bayesian neural networks, this observation noise pa-
 137 rameter is often fixed as part of the model’s weight decay, and ignored. However, when made
 138 data-dependent, it can be learned as a function of the data:

$$\mathcal{L}_{\text{NN}}(\theta) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2\sigma(\mathbf{x}_i)^2} \|\mathbf{y}_i - \mathbf{f}(\mathbf{x}_i)\|^2 + \log \sigma(\mathbf{x}_i)^2 \quad (5)$$

139 with added weight decay parameterised by λ (and similarly for l_1 loss). Note that here, unlike
 140 the above, variational inference is *not* performed over the weights, but instead we perform MAP
 141 inference – finding a single value for the model parameters θ . This approach *does not* capture
 142 epistemic model uncertainty, as epistemic uncertainty is a property of the model and not of the data.

143 In the next section we will combine these two types of uncertainties together in a single model. We
 144 will see how heteroscedastic noise can be interpreted as model attenuation, and develop a compli-
 145 mentary approach for the classification case.

146 3 Combining Aleatoric and Epistemic Uncertainty in One Model

147 In the previous section we described existing Bayesian deep learning techniques. In this section we
 148 present novel contributions which extend this existing literature. We develop models that will allow
 149 us to study the effects of modeling either aleatoric uncertainty alone, epistemic uncertainty alone,
 150 or modeling both uncertainties together in a single model. This is followed by an observation that
 151 aleatoric uncertainty in regression tasks can be interpreted as learned loss attenuation – making the
 152 loss more robust to noisy data. We follow that by extending the ideas of heteroscedastic regression
 153 to classification tasks. This allows us to learn loss attenuation for classification tasks as well.

154 3.1 Combining Heteroscedastic Aleatoric Uncertainty and Epistemic Uncertainty

155 We wish to capture both epistemic and aleatoric uncertainty in a vision model. For this we turn the
 156 heteroscedastic NN in §2.2 into a Bayesian NN by placing a distribution over its weights, with our
 157 construction in this section developed specifically for the case of vision models¹.

158 We need to infer the posterior distribution for a BNN model \mathbf{f} mapping an input image, \mathbf{x} , to a unary
 159 output, $\hat{\mathbf{y}} \in \mathbb{R}$, and a measure of aleatoric uncertainty given by variance, σ^2 . We approximate the
 160 posterior over the BNN with a dropout variational distribution using the tools of §2.1. As before,

¹Although this construction can be generalised for any heteroscedastic NN architecture.

161 we draw model weights from the approximate posterior $\widehat{\mathbf{W}} \sim q(\mathbf{W})$ to obtain a model output, this
 162 time composed of both predictive mean as well as predictive variance:

$$[\hat{\mathbf{y}}, \hat{\sigma}^2] = \mathbf{f}^{\widehat{\mathbf{W}}}(\mathbf{x}) \quad (6)$$

163 where \mathbf{f} is a Bayesian convolutional neural network parametrised by model weights $\widehat{\mathbf{W}}$. We can use
 164 a single network to transform the input \mathbf{x} , with its head split to predict both $\hat{\mathbf{y}}$ as well as $\hat{\sigma}^2$.

165 We fix a Laplace likelihood to model our aleatoric uncertainty. This induces a minimisation objective
 166 given labeled output points x :

$$\mathcal{L}_{BNN}(\theta) = \frac{1}{D} \sum_i \|y_i - \hat{y}_i\| \hat{\sigma}_i^{-2} + \log \hat{\sigma}_i^2 \quad (7)$$

167 where D is the number of output pixels y_i corresponding to input image \mathbf{x} , indexed by i (addition-
 168 ally, the loss includes weight decay which is omitted for brevity). For example, we may set $D = 1$
 169 for image-level regression tasks, or D equal to the number of pixels for dense prediction tasks (pre-
 170 dicting a unary corresponding to each input image pixel). $\hat{\sigma}_i^2$ is the BNN output for the predicted
 171 variance for pixel i .

172 This loss consists of two components; the residual regression obtained with a stochastic sample
 173 through the model – making use of the uncertainty over the parameters – and an uncertainty regu-
 174 larization term. We do not need ‘uncertainty labels’ to learn uncertainty. Rather, we only need to
 175 supervise the learning of the regression task. We learn the variance, σ^2 , implicitly from the loss
 176 function. The second regularization term prevents the network from predicting infinite uncertainty
 177 (and therefore zero loss) for all data points.

178 In practice, we train the network to predict the log variance, $s_i := \log \hat{\sigma}_i^2$:

$$\mathcal{L}_{BNN}(\theta) = \frac{1}{D} \sum_i \|y_i - \hat{y}_i\| \exp(-s_i) + s_i. \quad (8)$$

179 This is because it is more numerically stable than regressing the variance, σ^2 , as the loss avoids a
 180 potential division by zero. The exponential mapping also allows us to regress unconstrained scalar
 181 values, where $\exp(-s_i)$ is resolved to the positive domain giving valid values for variance.

182 To summarize, the predictive uncertainty for pixel \mathbf{y} in this combined model can be approximated
 183 using:

$$\text{Var}(\mathbf{y}) \approx \frac{1}{T} \sum_{t=1}^T \hat{\mathbf{y}}_t^2 - \left(\frac{1}{T} \sum_{t=1}^T \hat{\mathbf{y}}_t \right)^2 + \frac{1}{T} \sum_{t=1}^T \hat{\sigma}_t^2 \quad (9)$$

184 with $\{\hat{\mathbf{y}}_t, \hat{\sigma}_t^2\}_{t=1}^T$ a set of T sampled outputs: $\hat{\mathbf{y}}_t, \hat{\sigma}_t^2 = \mathbf{f}^{\widehat{\mathbf{W}}_t}(\mathbf{x})$ for randomly masked weights
 185 $\widehat{\mathbf{W}}_t \sim q(\mathbf{W})$.

186 3.2 Heteroscedastic Uncertainty as Learned Loss Attenuation

187 We observe that allowing the network to predict uncertainty, allows it effectively to temper the
 188 residual loss by $\exp(-s_i)$, which depends on the data. This acts similarly to an intelligent robust
 189 regression function. It allows the network to adapt the residual’s weighting, and even allows the
 190 network to learn to attenuate the effect from erroneous labels. This makes the model more robust to
 191 noisy data: inputs for which the model learned to predict high uncertainty will have a smaller effect
 192 on the loss.

193 The model is discouraged from predicting high uncertainty for all points – in effect ignoring the
 194 data – through the $\log \sigma^2$ term. Large uncertainty increases the contribution of this term, and in turn
 195 penalizes the model: The model *can* learn to ignore the data – but is penalised for that. The model is
 196 also discouraged from predicting very low uncertainty for points with high residual error, as low σ^2
 197 will exaggerate the contribution of the residual and will penalize the model. It is important to stress
 198 that this learned attenuation is not an ad-hoc construction, but a consequence of the probabilistic
 199 interpretation of the model.

200 This learned loss attenuation property of heteroscedastic NNs in regression is a desirable effect for
 201 classification models as well. However, heteroscedastic NNs in classification are peculiar models
 202 because technically any classification task has input-dependent uncertainty. Nevertheless, the ideas
 203 above can be extended from regression heteroscedastic NNs to classification heteroscedastic NNs,
 204 discussed next.

205 Due to space limitations, we discuss heteroscedastic uncertainty in classification tasks in Appendix
206 A. We next assess the ideas above empirically.

207 4 Experiments

208 In this section we evaluate our methods with pixel-wise depth regression and semantic segmentation.
209 An analysis of these results is given in the following section. To show the robustness of our learned
210 loss attenuation – a side-effect of modeling uncertainty – we present results on an array of popular
211 datasets, CamVid, Make3D, and NYUv2 Depth, where we set new state-of-the-art benchmarks.

212 For the following experiments we use the DenseNet architecture [18] which has been adapted for
213 dense prediction tasks by [19]. We use our own independent implementation of the architecture
214 using TensorFlow [20] (which slightly outperforms the original authors’ implementation on CamVid
215 by 0.2%, see Table 1a). For all experiments we train with 224×224 crops of batch size 4, and then
216 fine-tune on full-size images with a batch size of 1. We train with RMS-Prop with a constant learning
217 rate of 0.001 and weight decay 10^{-4} .

218 We compare the results of the Bayesian neural network models outlined in §3. We model *epis-*
219 *temic* uncertainty using Monte Carlo dropout (§2.1). The DenseNet architecture places dropout
220 with $p = 0.2$ after each convolutional layer. Following [21], we use 50 Monte Carlo dropout sam-
221 ples. We model *aleatoric* uncertainty with MAP inference using loss functions (8) and (12 in the
222 appendix), for regression and classification respectively (§2.2). We model the benefit of combining
223 both epistemic uncertainty as well as aleatoric uncertainty using our developments presented in §3.

224 4.1 Semantic Segmentation

225 To demonstrate our method for semantic segmentation, we use two datasets, CamVid [22] and NYU
226 v2 [23]. CamVid is a road scene understanding dataset with 367 training images and 233 test images,
227 of day and dusk scenes, with 11 classes. We resize images to 360×480 pixels for training and
228 evaluation. In Table 1a we present results for our architecture. Our method sets a new state-of-the-art
229 on this dataset with mean intersection over union (IoU) score of 67.5%. We observe that modeling
230 both aleatoric and epistemic uncertainty improves over the baseline result. The implicit attenuation
231 obtained from the aleatoric loss provides a larger improvement than the epistemic uncertainty model.
232 However, the combination of both uncertainties improves performance even further. This shows that
233 for this application it is more important to model aleatoric uncertainty, suggesting that epistemic
234 uncertainty can be mostly explained away in this large data setting.

235 Secondly, NYUv2 [23] is a challenging indoor segmentation dataset with 40 different semantic
236 classes. It has 1449 images with resolution 640×480 from 464 different indoor scenes. Table 1b
237 shows our results. This dataset is much harder than CamVid because there is significantly less struc-
238 ture in indoor scenes compared to street scenes, and because of the increased number of semantic
239 classes. We use DeepLabLargeFOV [24] as our baseline model. We observe a similar result (qual-
240 itative results given in Figure 2); we improve baseline performance by giving the model flexibility
241 to estimate uncertainty and attenuate the loss. The effect is more pronounced, perhaps because the
242 dataset is more difficult.

243 4.2 Pixel-wise Depth Regression

244 We demonstrate the efficacy of our method for regression using two popular monocular depth regres-
245 sion datasets, Make3D [25] and NYUv2 Depth [23]. The Make3D dataset consists of 400 training
246 and 134 testing images, gathered using a 3-D laser scanner. We evaluate our method using the same
247 standard as [26], resizing images to 345×460 pixels and evaluating on pixels with depth less than
248 $70m$. NYUv2 Depth is taken from the same dataset used for classification above. It contains RGB-D
249 imagery from 464 different indoor scenes. We compare to previous approaches for Make3D in Table
250 2a and NYUv2 Depth in Table 2b, using standard metrics (for a description of these metrics please
251 see [27]).

252 These results show that aleatoric uncertainty is able to capture many aspects of this task which
253 are inherently difficult. For example, in the qualitative results in Figure 3 and 4 we observe that
254 aleatoric uncertainty is greater for large depths, reflective surfaces and occlusion boundaries in the
255 image. These are common failure modes of monocular depth algorithms [26]. On the other hand,
256 these qualitative results show that epistemic uncertainty captures difficulties due to lack of data. For
257 example, we observe larger uncertainty for objects which are rare in the training set such as humans
258 in the third example of Figure 3.

CamVid	IoU
SegNet [28]	46.4
FCN-8 [29]	57.0
DeepLab-LFOV [24]	61.6
Bayesian SegNet [21]	63.1
Dilation8 [30]	65.3
Dilation8 + FSO [31]	66.1
DenseNet [19]	66.9
<i>This work:</i>	
DenseNet (Our Implementation)	67.1
+ Aleatoric Uncertainty	67.4
+ Epistemic Uncertainty	67.2
+ Aleatoric & Epistemic	67.5

(a) CamVid dataset for road scene segmentation.

NYUv2 40-class	Accuracy	IoU
SegNet [28]	66.1	23.6
FCN-8 [29]	61.8	31.6
Bayesian SegNet [21]	68.0	32.4
Eigen and Fergus [32]	65.6	34.1
<i>This work:</i>		
DeepLabLargeFOV	70.1	36.5
+ Aleatoric Uncertainty	70.4	37.1
+ Epistemic Uncertainty	70.2	36.7
+ Aleatoric & Epistemic	70.6	37.3

(b) NYUv2 40-class dataset for indoor scenes.

Table 1: **Semantic segmentation performance.** Modeling both aleatoric and epistemic uncertainty gives a notable improvement in segmentation accuracy over state of the art baselines.

Make3D	rel	rms	log ₁₀
Karsch et al. [33]	0.355	9.20	0.127
Liu et al. [34]	0.335	9.49	0.137
Li et al. [35]	0.278	7.19	0.092
Laina et al. [26]	0.176	4.46	0.072
<i>This work:</i>			
DenseNet Baseline	0.167	3.92	0.064
+ Aleatoric Uncertainty	0.149	3.93	0.061
+ Epistemic Uncertainty	0.162	3.87	0.064
+ Aleatoric & Epistemic	0.149	4.08	0.063

(a) Make3D depth dataset [25].

NYU v2 Depth	rel	rms	log ₁₀	δ ₁	δ ₂	δ ₃
Karsch et al. [33]	0.374	1.12	0.134	-	-	-
Ladicky et al. [36]	-	-	-	54.2%	82.9%	91.4%
Liu et al. [34]	0.335	1.06	0.127	-	-	-
Li et al. [35]	0.232	0.821	0.094	62.1%	88.6%	96.8%
Eigen et al. [27]	0.215	0.907	-	61.1%	88.7%	97.1%
Eigen and Fergus [32]	0.158	0.641	-	76.9%	95.0%	98.8%
Laina et al. [26]	0.127	0.573	0.055	81.1%	95.3%	98.8%
<i>This work:</i>						
DenseNet Baseline	0.117	0.517	0.051	80.2%	95.1%	98.8%
+ Aleatoric Uncertainty	0.112	0.508	0.046	81.6%	95.8%	98.8%
+ Epistemic Uncertainty	0.114	0.512	0.049	81.1%	95.4%	98.8%
+ Aleatoric & Epistemic	0.110	0.506	0.045	81.7%	95.9%	98.9%

(b) NYUv2 depth dataset [23].

Table 2: **Monocular depth regression performance.** Comparison to previous approaches on depth regression dataset NYUv2 Depth. Modeling the combination of uncertainties improves accuracy.

259 In summary, we have demonstrated that our model can improve performance over non-Bayesian
 260 baselines by implicitly learning attenuation of systematic noise and difficult concepts. For example
 261 we observe high aleatoric uncertainty for distant objects and on object and occlusion boundaries.

262 5 Analysis: What Do Aleatoric and Epistemic Uncertainties Capture?

263 In §4 we showed that modeling aleatoric and epistemic uncertainties improves prediction perfor-
 264 mance, with the combination performing even better. In this section we wish to study the effec-
 265 tiveness of modeling aleatoric and epistemic uncertainty. In particular, we wish to quantify the
 266 performance of these uncertainty measurements and analyze what they capture.

267 5.1 Quality of Uncertainty Metric

268 Firstly, in Appendix C we show precision-recall curves for regression and classification models.
 269 They show how our model performance improves by removing pixels with uncertainty larger than
 270 various percentile thresholds. This illustrates two behaviors of aleatoric and epistemic uncertainty
 271 measures. Firstly, it shows that the uncertainty measurements are able to correlate well with accu-
 272 racy, because all curves are strictly decreasing functions. We observe that precision is lower when
 273 we have more points that the model is not certain about. Secondly, the curves for epistemic and
 274 aleatoric uncertainty models are very similar. This shows that each uncertainty ranks pixel con-
 275 fidence similarly to the other uncertainty, in the absence of the other uncertainty. This suggests
 276 that when only one uncertainty is explicitly modeled, it attempts to compensate for the lack of the
 277 alternative uncertainty when possible.

278 Secondly, in Appendix D we analyze the quality of our uncertainty measurement using calibration
 279 plots from our model on the test set. To form calibration plots for classification models, we dis-
 280 cretize our model’s predicted probabilities into a number of bins, for all classes and all pixels in the
 281 test set. We then plot the frequency of correctly predicted labels for each bin of probability values.
 282 Better performing uncertainty estimates should correlate more accurately with the line $y = x$ in
 283 the calibration plots. For regression models, we can form calibration plots by comparing the fre-
 284 quency of residuals lying within varying thresholds of the predicted distribution. Figure 6 shows the
 285 calibration of our classification and regression uncertainties.

Train dataset	Test dataset	RMS	Aleatoric variance	Epistemic variance
Make3D / 4	Make3D	5.76	0.506	7.73
Make3D / 2	Make3D	4.62	0.521	4.38
Make3D	Make3D	3.87	0.485	2.78
Make3D / 4	NYUv2	-	0.388	15.0
Make3D	NYUv2	-	0.461	4.87

(a) Regression

Train dataset	Test dataset	IoU	Aleatoric entropy	Epistemic logit variance ($\times 10^{-3}$)
CamVid / 4	CamVid	57.2	0.106	1.96
CamVid / 2	CamVid	62.9	0.156	1.66
CamVid	CamVid	67.5	0.111	1.36
CamVid / 4	NYUv2	-	0.247	10.9
CamVid	NYUv2	-	0.264	11.8

(b) Classification

Table 3: Accuracy and aleatoric and epistemic uncertainties for a range of different train and test dataset combinations. We show aleatoric and epistemic uncertainty as the mean value of all pixels in the test dataset. We compare reduced training set sizes (1, $\frac{1}{2}$, $\frac{1}{4}$) and unrelated test datasets. This shows that aleatoric uncertainty remains approximately constant, while epistemic uncertainty decreases the closer the test data is to the training distribution, demonstrating that epistemic uncertainty can be explained away with sufficient training data (but not for out-of-distribution data).

286 5.2 Uncertainty with Distance from Training Data

287 In this section we show two results:

- 288 1. Aleatoric uncertainty cannot be explained away with more data,
- 289 2. Aleatoric uncertainty does not increase for out-of-data examples (situations different from
- 290 training set), whereas epistemic uncertainty does.

291 In Table 3 we give accuracy and uncertainty for models trained on increasing sized subsets of
 292 datasets. This shows that epistemic uncertainty decreases as the training dataset gets larger. It also
 293 shows that aleatoric uncertainty remains relatively constant and cannot be explained away with more
 294 data. Testing the models with a different test set (bottom two lines) shows that epistemic uncertainty
 295 increases considerably on those test points which lie far from the training sets.

296 These results reinforce the case that epistemic uncertainty can be explained away with enough data,
 297 but is required to capture situations not encountered in the training set. This is particularly important
 298 for safety-critical systems, where epistemic uncertainty is required to detect situations which have
 299 never been seen by the model before.

300 5.3 Real-Time Application

301 Our model based on DenseNet [19] can process a 640×480 resolution image in $150ms$ on a NVIDIA
 302 Titan X GPU. The aleatoric uncertainty models add negligible compute. However, epistemic mod-
 303 els require expensive Monte Carlo dropout sampling. For models such as ResNet [4], this is pos-
 304 sible to achieve economically because only the last few layers contain dropout. Other models, like
 305 DenseNet, require the entire architecture to be sampled. This is difficult to parallelize due to GPU
 306 memory constraints, and often results in a $50 \times$ slow-down for 50 Monte Carlo samples.

307 6 Conclusions

308 We presented a novel Bayesian deep learning framework to learn a mapping to aleatoric uncertainty
 309 from the input data, which is composed on top of epistemic uncertainty models. We derived our
 310 framework for both regression and classification applications. We showed that it is important to
 311 model *aleatoric* uncertainty for:

- 312 • **Large data situations**, where epistemic uncertainty is explained away,
- 313 • **Real-time applications**, because we can form aleatoric models without expensive Monte
- 314 Carlo samples.

315 And *epistemic* uncertainty is important for:

- 316 • **Safety-critical applications**, because epistemic uncertainty is required to understand ex-
- 317 amples which are different from training data,
- 318 • **Small datasets** where the training data is sparse.

319 However aleatoric and epistemic uncertainty models are not mutually exclusive. We showed that
 320 the combination is able to achieve new state-of-the-art results on depth regression and semantic
 321 segmentation benchmarks.

322 The first paragraph in this paper posed two recent disasters which could have been averted by real-
 323 time Bayesian deep learning tools. Therefore, we leave finding a method for *real-time epistemic*
 324 *uncertainty* in deep learning as an important direction for future research.

325 **References**

326 [1] NHTSA. PE 16-007. Technical report, U.S. Department of Transportation, National Highway Traffic
327 Safety Administration, Jan 2017. Tesla Crash Preliminary Evaluation Report.

328 [2] Jessica Guynn. Google photos labeled black people 'gorillas'. *USA Today*, 2015.

329 [3] Andrew Blake, Rupert Curwen, and Andrew Zisserman. A framework for spatiotemporal control in the
330 tracking of visual contours. *International Journal of Computer Vision*, 11(2):127–145, 1993.

331 [4] Xuming He, Richard S Zemel, and Miguel Á Carreira-Perpiñán. Multiscale conditional random fields for
332 image labeling. In *Computer vision and pattern recognition, 2004. CVPR 2004. Proceedings of the 2004*
333 *IEEE computer society conference on*, volume 2, pages II–II. IEEE, 2004.

334 [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
335 In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778,
336 2016.

337 [6] Y. Gal. *Uncertainty in Deep Learning*. PhD thesis, University of Cambridge, 2016.

338 [7] Armen Der Kiureghian and Ove Ditlevsen. Aleatory or epistemic? does it matter? *Structural Safety*, 31
339 (2):105–112, 2009.

340 [8] John Denker and Yann LeCun. Transforming neural-net output levels to probability distributions. In
341 *Advances in Neural Information Processing Systems 3*. Citeseer, 1991.

342 [9] David JC MacKay. A practical Bayesian framework for backpropagation networks. *Neural Computation*,
343 4(3):448–472, 1992.

344 [10] Radford M Neal. *Bayesian learning for neural networks*. PhD thesis, University of Toronto, 1995.

345 [11] Alex Graves. Practical variational inference for neural networks. In *Advances in Neural Information*
346 *Processing Systems*, pages 2348–2356, 2011.

347 [12] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural
348 network. In *ICML*, 2015.

349 [13] José Miguel Hernández-Lobato, Yingzhen Li, Daniel Hernández-Lobato, Thang Bui, and Richard E
350 Turner. Black-box alpha divergence minimization. In *Proceedings of The 33rd International Confer-*
351 *ence on Machine Learning*, pages 1511–1520, 2016.

352 [14] Yarin Gal and Zoubin Ghahramani. Bayesian convolutional neural networks with Bernoulli approximate
353 variational inference. *ICLR workshop track*, 2016.

354 [15] Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An introduction to
355 variational methods for graphical models. *Machine learning*, 37(2):183–233, 1999.

356 [16] David A Nix and Andreas S Weigend. Estimating the mean and variance of the target probability dis-
357 tribution. In *Neural Networks, 1994. IEEE World Congress on Computational Intelligence., 1994 IEEE*
358 *International Conference On*, volume 1, pages 55–60. IEEE, 1994.

359 [17] Quoc V Le, Alex J Smola, and Stéphane Canu. Heteroscedastic Gaussian process regression. In *Proceed-*
360 *ings of the 22nd international conference on Machine learning*, pages 489–496. ACM, 2005.

361 [18] Gao Huang, Zhuang Liu, Kilian Q Weinberger, and Laurens van der Maaten. Densely connected convo-
362 lutional networks. *arXiv preprint arXiv:1608.06993*, 2016.

363 [19] Simon Jégou, Michal Drozdal, David Vazquez, Adriana Romero, and Yoshua Bengio. The one
364 hundred layers tiramisu: Fully convolutional densenets for semantic segmentation. *arXiv preprint*
365 *arXiv:1611.09326*, 2016.

366 [20] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin,
367 Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine
368 learning. In *Proceedings of the 12th USENIX Symposium on Operating Systems Design and Implementa-*
369 *tion (OSDI)*. Savannah, Georgia, USA, 2016.

370 [21] Alex Kendall, Vijay Badrinarayanan, and Roberto Cipolla. Bayesian segnet: Model uncertainty in deep
371 convolutional encoder-decoder architectures for scene understanding. *arXiv preprint arXiv:1511.02680*,
372 2015.

373 [22] Gabriel J Brostow, Julien Fauqueur, and Roberto Cipolla. Semantic object classes in video: A high-
374 definition ground truth database. *Pattern Recognition Letters*, 30(2):88–97, 2009.

375 [23] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support
376 inference from rgb-d images. In *European Conference on Computer Vision*, pages 746–760. Springer,
377 2012.

378 [24] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Se-
379 mantic image segmentation with deep convolutional nets and fully connected crfs. *arXiv preprint*
380 *arXiv:1412.7062*, 2014.

- 381 [25] Ashutosh Saxena, Min Sun, and Andrew Y Ng. Make3d: Learning 3d scene structure from a single still
382 image. *IEEE transactions on pattern analysis and machine intelligence*, 31(5):824–840, 2009.
- 383 [26] Iro Laina, Christian Rupprecht, Vasileios Belagiannis, Federico Tombari, and Nassir Navab. Deeper depth
384 prediction with fully convolutional residual networks. In *3D Vision (3DV), 2016 Fourth International
385 Conference on*, pages 239–248. IEEE, 2016.
- 386 [27] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a
387 multi-scale deep network. In *Advances in neural information processing systems*, pages 2366–2374,
388 2014.
- 389 [28] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder
390 architecture for scene segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
391 2017.
- 392 [29] Evan Shelhamer, Jonathon Long, and Trevor Darrell. Fully convolutional networks for semantic segmen-
393 tation. *IEEE transactions on pattern analysis and machine intelligence*, 2016.
- 394 [30] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint
395 arXiv:1511.07122*, 2015.
- 396 [31] Abhijit Kundu, Vibhav Vineet, and Vladlen Koltun. Feature space optimization for semantic video seg-
397 mentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages
398 3168–3175, 2016.
- 399 [32] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common
400 multi-scale convolutional architecture. In *Proceedings of the IEEE International Conference on Computer
401 Vision*, pages 2650–2658, 2015.
- 402 [33] Kevin Karsch, Ce Liu, and Sing Bing Kang. Depth extraction from video using non-parametric sampling.
403 In *European Conference on Computer Vision*, pages 775–788. Springer, 2012.
- 404 [34] Miaomiao Liu, Mathieu Salzmann, and Xuming He. Discrete-continuous depth estimation from a single
405 image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages
406 716–723, 2014.
- 407 [35] Bo Li, Chunhua Shen, Yuchao Dai, Anton van den Hengel, and Mingyi He. Depth and surface normal
408 estimation from monocular images using regression on deep features and hierarchical crfs. In *Proceedings
409 of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1119–1127, 2015.
- 410 [36] Lubor Ladicky, Jianbo Shi, and Marc Pollefeys. Pulling things out of perspective. In *Proceedings of the
411 IEEE Conference on Computer Vision and Pattern Recognition*, pages 89–96, 2014.