

## **Towards Algorithmic Analytics for Large-scale Datasets**

Danilo Bzdok<sup>1, 2, 3</sup>, Thomas E. Nichols<sup>4, 5</sup>, Stephen M. Smith<sup>4</sup>

<sup>1</sup> Department of Psychiatry, Psychotherapy and Psychosomatics, RWTH Aachen University, 52072 Aachen, Germany

<sup>2</sup> JARA, Translational Brain Medicine, Aachen, Germany

<sup>3</sup> Parietal Team, INRIA, Neurospin, bat 145, CEA Saclay, 91191 Gif-sur-Yvette, France

<sup>4</sup> Wellcome Trust Centre for Integrative Neuroimaging (WIN-FMRIB), University of Oxford, Oxford, UK

<sup>5</sup> Big Data Institute, University of Oxford, Oxford, UK

### **Abstract**

The traditional goals of quantitative analytics cherish simple, transparent models to generate explainable insights. Large-scale data acquisition, enabled for instance by brain scanning and genomic profiling with microarray-type techniques, has prompted a wave of statistical inventions and innovative applications. Modern analysis approaches 1) tame large variable arrays capitalizing on regularization and dimensionality-reduction strategies, 2) are increasingly backed up by empirical model validations rather than justified by mathematical proofs, 3) will compare against and build on open data and consortium repositories, as well as 4) often embrace more elaborate, less interpretable models in order to maximize prediction accuracy. Here we review these trends in learning from “big data” and illustrate examples from imaging neuroscience.

**Keywords:** reproducibility, open science, data science, machine learning, explainable AI, deep phenotyping

## Introduction

Tension is emerging in everyday data analysis in the biomedical sciences. Around the turn of the century, deployment of new measurement techniques, especially microarray-like techniques in genomics and brain scanning in neuroscience, have ignited data accumulation at massive scale <sup>1</sup>. As one consequence, the amount of health-related data is expected to double several times per year starting from 2020 (Nature editorial, 2016). Data-analytical methodology in turn has expanded more in the last two decades than any other point in history <sup>2,3</sup>. However, emerging opportunities to generate quantitative insight from accumulating data are adopted with hesitation in many empirical domains. Here, we portray the growing stack of algorithmic tools, illustrated with examples from the area of human neuroscience.

In many empirical sciences, classical statistics is still the dominant arsenal for deriving rigorous conclusions from data. A form of linear regression was already used by Gauss around 1795, and null-hypothesis testing emerged in the early 20<sup>th</sup> century to formally assess the significance of *tail-area statistics* <sup>2</sup>. *Closed-form* textbook formulae (for italic terms see glossary) were a necessity to avoid laborious paper-and-pencil calculations <sup>4</sup>. Electronic computations only became slowly available after World War II <sup>2</sup>. Hence, a new modeling approach was routinely validated by virtue of mathematical theory. Specifically, *consistency theorems* formally characterize how a particular analysis method behaves if the sample size increases indefinitely <sup>5,6</sup>. Much emphasis was put on straightforward linear models due to key advantages for understanding the relationships between carefully selected input variables. Simpler, less data-hungry models were also generally preferred because data acquisition was financially and logistically expensive for most of the last century. This is why experimental laboratory studies needed to be carefully planned in

advance, and research hypotheses had to be precisely defined beforehand <sup>7,8</sup>. From this traditional perspective of scientific investigation, re-analysing general-purpose data repositories would have been of little appeal. Moreover, it was only in the 90's that desktop computer software for many machine-learning algorithms or *Bayesian* modeling approaches became widely available <sup>2</sup>.

Impressive data aggregation in the early 21<sup>st</sup> century has given rise to a new kind of empirical research <sup>9,10</sup>. In many scientific domains, information has become cheaper, considerably more fine-grained and multi-faceted, as well as freely available. This shift of context opens the door to the principled exploration of already acquired “found” observational data. Increasingly, quantitative analysis tools are designed, evaluated, and deployed ad-hoc before complete formal analysis of their mathematical properties. Instead, empirical justifications are obtained from successful prediction performance in separate reference datasets <sup>4,11</sup>. More broadly, modern data analysis needs to negotiate trade-offs between statistical notions like effect *uncertainty* (e.g., ‘How sure are we about a detected effect?’) and computer-science notions like computational load and memory resources (e.g., ‘How expensive does the analysis become with an increasing number of input variables?’) <sup>9,12</sup>.

### **Global trends in empirical data analysis**

As a looming culture clash, university education in various empirical fields is still focused on classical methods from a time of scarce data and limited computation. Blossoming data resources however entail a need for exploiting analytical techniques suited

for today's data-rich setting. Getting back to our guiding example, imaging neuroscience has spawned increasingly wide and deep datasets over recent years. However, the adoption of analytical tools tailored for modern data is accelerating only recently<sup>13,14</sup>. The most commonly used methods from statistics and computing were not designed to solve the types of problems that data-rich scientists, including neuroscientists, are facing today. The present overview retraces this emerging transition from formally inspired modeling of a few hand-selected variables to learning complicated patterns from data with increasingly adaptive algorithms:

1) Many empirical sciences are now generating detailed phenotypical descriptions of organisms and phenomena like the brain. Investigators confronted with hundreds or thousands of quantitative measurements are also confronted with how to estimate statistical models with potentially hundreds or thousands of parameters. This new context questions the long-standing dogma that statistical analysis should strive to be maximally impartial. Instead, the unfolding analysis paradigm appears to ask 'What is the most useful a-priori knowledge that can inform and shape my quantitative analysis?'. The consequence is growing importance of *bias-inducing regularization* strategies and data transformations for *dimensionality reduction* like clustering and matrix decomposition.

2) Classical modeling tools have usually been trusted after their formal properties had been mathematically understood in detail. Increasing data availability is escalating the pace at which new quantitative methods are invented and re-purposed, even before studying their theoretical properties. As such, the modern quantitative investigator has special interest in asking 'How well does my obtained modeling solution hold up when

directly evaluated in other sampled observations?’ Hence, extracted candidate models are more often grounded in empirical *cross-validation* or *posterior predictive checks* to judge the model’s quality and practical usefulness.

3) Many empirical sciences have centered on careful planning and conducting of experiments in the laboratory. This predominance of acquiring expensive in-house datasets is re-balanced to ever-wider usage of openly accessible observational datasets. Investigators can increasingly ask ‘How does my research question play out in existing consortium data?’ or ‘how does my newly developed method scale to open population datasets?’. These opportunities can improve the reproducibility of scientific claims and the comparability of quantitative approaches.

4) Empirical sciences like imaging neuroscience capitalize on always more complex, sometimes untransparent modeling approaches. These investigators may want to ask ‘How well can a powerful pattern-learning algorithm forecast outcomes from the natural phenomenon under study?’. On the one hand, by maximising prediction accuracy on new data or settings, much harder problems can be tackled than before. On the other hand, uncompromising prediction studies may lose some interpretational grip on understanding the isolated role of each model parameter. The resulting interpretability trade-offs incur ethical and policy-related consequences for science, business, and government.

**Deeper phenotyping yields data with always more variables**

The brain sciences have recently been highlighted as the potentially most data-rich medical specialty (Nature Editorial 2016). In genomics and imaging genetics, jointly considering >1,000,000 single nucleotide polymorphisms (SNPs) typically exceeds the thousands of participants in the currently biggest human cohorts e.g., <sup>15</sup>. In imaging neuroscience, a brain scan with commonly available resolution offers measurements from ~100,000-500,000 locations ('high-p' scenario with many variables). Yet, sample sizes have reached hundreds or thousands of participants ('low-n' scenario with few observations) only over recent years. In this context, too few sample observations from the participants may be available to allow for rigorous statements about each separate input dimension like a specific gene or a particular brain location. At the extreme, detailed neuroanatomical studies with measurements at micrometer resolution may be available in one or only few participants <sup>16</sup>. Consequently, application of *dimensionality-reduction* techniques is becoming hard to avoid in various empirical sciences. The theme of reducing abundant multivariate information to the relevant essence is reflected in matrix decomposition techniques like principal component analysis (PCA), *canonical correlation analysis* (CCA), *partial least squares* (PLS), independent component analysis (ICA), and expanding tensor-decomposition techniques as well as clustering techniques like *k-means* <sup>17,18</sup>. In many workflows, it becomes an important pre-processing step to re-express the data in a simpler underlying form before applying the final data analysis model, such as linear regression <sup>19</sup>.

In traditionally small datasets with few variables, computing ordinary linear regression analysis, as a special instance of *maximum likelihood estimation* (MLE), readily provides parameter estimates, with almost optimal precision <sup>5</sup>. Standard linear regression corresponds well to the traditional goals of statistics that valued impartiality to what may be expected in the data - *unbiasedness*. Indeed, even datasets with ~30-40 variables were still

considered *high-dimensional* in the 80's and 90's <sup>1</sup>. Already in this setting, the formal theory backing up linear regression models starts to lose some of its optimality, although MLE successfully legitimized a series of classical statistics approaches still in pervasive use today <sup>1</sup>. Introducing an increasing number of input variables into a linear model usually leads to an increasing number of parameters to be fitted. Such larger models incur higher *variance* in estimating their model parameter values, e.g., due to ambiguities in the fitting of partially redundant variables that carry similar information about the target outcome <sup>19</sup>.

Even when adding more parameters to simple linear models, the expanded model capacity - with higher *degrees of freedom* - adds challenges to interpretation. With increasing number of model parameters, it becomes more difficult to clearly attribute the variance explained by each individual input variable <sup>1</sup>. Larger linear models are also more susceptible to picking up on idiosyncrasies and noise in data. In the *high-dimensional* scenario, hundreds or thousands of input variables (e.g., brain region volumes, functional connectivity strengths, or gene expression levels) can be submitted to model fitting. It becomes harder to tell, using classical goodness-of-fit tests, how well an obtained linear model actually encapsulates the data at hand. At the extreme, the number of measured variables exceeds the number of available samples or observations in many modern datasets <sup>20</sup>. Such data-rich settings hinder the reproducible identification of unique model parameter solutions. Such contexts render common linear-regression models *non-identifiable*, which makes the parameter value estimates difficult to interpret, and model performance on new data poor <sup>21</sup>.

The amount of information that can be gleaned from emerging *high-dimensional* datasets may sometimes remain low even if the sample size is increasing <sup>22,23</sup>. Probably no statistical approach performs well if thousands of input variables are truly individually

informative about the outcome to be predicted <sup>21</sup>. Specially, with high dimensions, pre-assuming a more parsimonious underlying representation in the relevant variables (e.g., identifying a smaller number of latent factors of variation) may be a pragmatic way to obtain useful and interpretable modeling solutions. Consequently, modern data often make it necessary to introduce some intentional *bias* into the data-analysis process. In addition to dimensionality-reduction techniques, a plethora of *regularization* strategies have flourished in response <sup>21</sup>. Often very simple extensions of traditional tools like linear regression can be effective in datasets with *high-dimensional* measurements. Such *penalized* linear models dedicated to many variables  $p$  are epitomized by the increasing adoptions of Ridge regression, Lasso, and Elastic Net <sup>21</sup>.

In particular, several classical analysis tools underwent a *sparsification* over the last 15-20 years: the introduced *biasing* assumption is that most input variables in the data are expected to be uninformative about the outcome. *Sparsified* model extensions were enthusiastically embraced by the machine-learning community <sup>21</sup>. Those often-*frequentist* modeling approaches try to learn from data which input variables can be ignored by encouraging a maximum of model parameters with exactly-zero values. This *penalized* modeling regime assumes that effective model estimation should be skewed towards finding only a subset of the input variables to be relevant to the research question <sup>21</sup>. Besides *sparsity*-inducing regression via Lasso and Elastic Net, these highly effective parsimony constraints recently motivated, for instance, *sparse PCA*, *sparse CCA*, or *sparse k-means*.

Instead, more *Bayesian*-minded analysts may prefer estimating the *uncertainty* of possible model parameter values and corresponding variable influence to be close to zero or not <sup>24-26</sup>. The investigator embracing *Bayesian* statistics intentionally *biases* model estimation by skewing parameter values towards existing knowledge expressed in prior distributions <sup>27</sup>.



For instance, such approaches can capitalize on hierarchical dependence structure in the data that exists between the quantitative measurements, such as to share statistical strengths between individual outcomes pooled from different time points of a longitudinal study. In *high dimensions*, many *Bayesian* approaches may however suffer as the prior probability distribution can take unexpected shapes, which may preclude the “true” parameter values from being recovered. As a manifestation of this so-called *curse of dimensionality*, imposing prior knowledge by guiding model estimation to expect certain ranges of model parameter values more than others may become ineffective <sup>28</sup>. Even if probabilistic parameter distributions could be obtained on each input dimension, it would be challenging for a domain expert to interpret every single parameter <sup>12</sup>. Algorithmically, even modern approximate methods for *Markov chain Monte Carlo sampling* widely used to infer *Bayesian* models are susceptible to the consequences of the *curse of dimensionality* <sup>29</sup>.

These side-effects of rich multi-variable phenotyping need to be tackled in an increasing number of modern neuroscience studies. Linear but flexible pattern-learning models have repeatedly yielded useful *dimensionality reductions* of *high-dimensional* subject descriptions and integration of different modalities of detailed measurements. In this spirit, CCA was used by Smith and colleagues to uncover population co-variation that links coupling measures of various brain networks and extensive phenotyping by a diversity of behavioral indicators <sup>30</sup>. Standard CCA can be viewed as reminiscent of classical statistics because this model is fitted based on *MLE* without deliberately imposing prior knowledge or *bias* that would guide parameter estimation. However, the same CCA method can be viewed to represent a proto-typical approach suited for modern datasets because of in-built *dimensionality reduction*, avoiding strong (*parametric*) assumptions about the distributions to be encountered in the data, the native ability to fuse two heterogeneous data modalities,

and acting in the *high-variance* regime due to the considerable *degrees of freedom*. CCA models have recently seen extensions and applications for *sparse penalization* <sup>31</sup>, *Bayesian modeling* <sup>32</sup>, and *deep learning* <sup>33</sup>. This *multivariate* method is complementary to so-called mass *univariate* approaches, which have been pervasively used in imaging neuroscience to study effects separately for each part of the brain <sup>34,35</sup>.

In a recent CCA application in imaging neuroscience <sup>30</sup>, one significant population mode of multi-modal co-variation was obtained based on one robust set of canonical correlations. This multivariate brain-behavior pattern extracted from rich phenotyping demonstrated a positive-negative axis: intelligence, memory and cognition tests and indices of life satisfaction on the positive end, and negative life-factor measures at the other end (Fig. 1). Additionally, the functional connectivity weights emphasized prominent modulation of the brain's "default mode network" (DMN). The doubly-*multivariate* CCA technique was recently re-purposed to revisit the idea that the DMN subserves some of the most human-defining cognitive processes by pooling neural information across the cortical landscape <sup>36</sup>. Profiting from multi-modal imaging data of 10,000 UK Biobank participants (Fig. 2), major nodes of the DMN could be shown to explain variance in how canonical brain networks communicate with each other. This population neuroscience study <sup>36</sup> thus provided robust indicators that the biological role of the DMN may emerge from propagating brain-wide information flow to orchestrate the cortical network repertoire, potentially mediated by the right and left temporoparietal junction of the DMN (cf., <sup>37</sup>). *Sparsity* to intentionally *bias* CCA estimation was recently used to provide a more complete understanding of the functional connectivity patterns of this major brain network during mind-wandering experience in humans <sup>38</sup>. Thus, imposing exactly-zero relevance weights, certain random-thought

behaviors among richly phenotyped experiences could be isolated to underlie functional connectivity signatures in the DMN.

### **Empirical model checks enabled by always larger sample sizes**

Classical statistics was conceived when datasets had modest sample size <sup>2,7</sup>. In the early 20<sup>th</sup> century, a primary concern was to gather information from scarce data points to achieve reasonable confidence in the model estimates for meaningful parameter interpretation. In this data-scarce context, a key theoretical property to judge the usefulness of a given modeling approach was *asymptotic consistency*. This formal guarantee of model performance has certified a host of long-standing statistical approaches. This criterion quantifies whether model estimation converges to the true variable relationships as the number of observation samples increases, mimicking for instance availability of brain scans from an infinite number of participants.

However, increasingly flexible algorithmic approaches, with data-hungry *deep neural-network algorithms* as an extreme case, can simply memorize much of the provided observation samples in certain settings. That is, modern complex models may enjoy *consistency guarantees*, but be highly prone to seriously *overfit* the provided data <sup>39</sup>. This adaptiveness of flexible modeling approaches can lead to spuriously high performance when evaluated on the observations used for model fitting (*in-sample performance*). This scenario may change the role of *consistency theorems* if the goal is to build models that perform well on observations to be sampled in the future, rather than the data sample at hand. Hence, in mathematical theory describing the convergence behavior of many machine-learning

models, *finite-sample theorems* are common where model performance is assessed as a function of the amount of observations available for model fitting, with its formal relation to the complexity of the chosen model and the *ground-truth* information density of the data rather than the raw number of input variables <sup>22,23</sup>.

With the increasing sample sizes of modern datasets, empirical evaluation procedures are becoming attractive to vouch model quality beyond those participants or observations used for model estimation, to new, independent data points. In fact, even a simple, inflexible model with few parameters can often *overfit* the available observations. This is because not every aspect of the measured data usually reflects the phenomenon of interest <sup>24</sup>. For instance, magnetoencephalographic brain measurements of neural activity responses can be influenced by passing trains hundreds of meters away, or by other electromagnetic fluctuations that happen to occur in the environment. Here, the model used may not have optimally fitted to the intended purpose in the participant sample at hand, even if the model enjoys the theoretical *consistency guarantee* to approach the true statistical relationship with unlimited amounts of data <sup>22</sup>. Moreover, at a given sample size of say 1,000 brain scans, it is possible that a purposefully *biased* model has already converged closer to the *ground-truth* solution based on the limited number of available observations than an *unbiased* model that is theoretically ensured to be correct in unlimited observations or participants.

As a consequence of staggering increase in sample size, modern quantitative analyses can increasingly be backed up by data-dependent optimality criteria rather than relying mostly on formal optimality guarantees. In the machine-learning community, *resampling* and *permutation* schemes are popular, typically *non-parametric* tools to glean further information from the data themselves. As an important example, *cross-validation*

procedures <sup>40,41</sup> repeatedly split the available observations to assess the discrepancy between potentially overly optimistic model performance on the *training data* used for model estimation and the previously unseen *test data*. This empirical model check can now be increasingly used to approximate the expected model performance in observations or participants yet-to-be-observed in the future <sup>2,19</sup>. An additional *validation data split* (internal to the training data) routinely serves for tuning any algorithm hyper-parameters to the data at hand, such as to optimize the strength of inducing zero parameters by *sparse regularization*.

Similarly, empirical *permutation* procedures allow *non-parametric* null-hypothesis testing based on *exchangeability* assumptions. This practical re-implementation of classical statistical inference is more general and flexible than what can typically be achieved by the common assumptions of ‘independent and identically distributed’ <sup>28</sup>. Additionally, *bagging* can improve prediction performance on new data based on data resampling and averaging hundreds of model solutions <sup>42</sup>. Moreover, *bootstrapping* can bestow population *uncertainty* intervals around almost any *frequentist* statistical approach, derived directly from the available observations themselves <sup>43</sup>. These empirical model checks are based on repeatedly resampling the data at hand, which yields more truthful results with more observations.

In a similar data-guided fashion, *Bayesian* approaches commonly re-adjust prior assumptions consecutively to enhance model estimation. The practical performance of each candidate model can be evaluated using *posterior predictive checks* that generate new data from candidate sets of posterior parameter distributions <sup>28</sup>. As *Bayesian* estimation conditions on the provided participants or observations, their fully specified probability intervals for each model parameter are valid for any sample size. These confidence bounds naturally tend to become always narrower as the amount of available observations grows.

Moreover, as the influence of the imposed prior knowledge gradually wanes with increasing sample size, the means of the inferred posterior parameter distributions ultimately converge with *frequentist* estimates of a particular parameter value (from *MLE*).

In imaging neuroscience, the recent surge in sample size led to re-evaluation of some established means to draw rigorous conclusions from brain measurements. In a sample of ~5,000 UK Biobank participants, Pearson correlation analyses between a behavioral phenotype and a brain imaging feature at  $r=0.1$  were found statistically significant for the most part (Fig. 3). This was even the case after correction for multiple comparisons <sup>44</sup>, which was anticipated long ago <sup>45</sup>. In this *univariate*-flavored approach, reporting effect estimates as interesting based on p-values alone may become insufficient if many observations are included in the analysis. This calls for systematic reporting of effect sizes (i.e., model parameter values) and other importance metrics such as prediction performance computed from *cross-validation* procedures <sup>8,46</sup>. Further, having larger participant samples, combined with more complicated *multivariate* analysis settings, has propelled the use of *non-parametric* null-hypothesis testing schemes based on more flexible *exchangeability* assumptions and data resampling schemes <sup>47,48</sup>. For instance, also in imaging neuroscience, statistical significance is increasingly drawn by generating a to-be-tested empirical null distribution directly from the data themselves. For instance, as by shuffling which brain scan is labeled as male versus female to make statements about statistically distinguishable sex differences in the brain <sup>49</sup>. The more participants' data are available in a neuroscience dataset, the more reliable conclusions from these resampling procedures can become <sup>50</sup>.

When will the sample sizes be sufficient for fitting and evaluating *deep-learning* approaches in the area of imaging neuroscience? Experts recently proposed a general rule of thumb <sup>51</sup>: In various application areas,  $n=5,000$  samples per category to be distinguished

were often necessary to achieve relevant model prediction performance. However, datasets with  $n > 10,000,000$  samples were repeatedly necessary to exceed human-level performance. The low sample( $n$ )-to-variables( $p$ ) ratio in today's neuroscience datasets may still hamper the potential of *multivariate deep-learning* techniques. Some current shortcomings on data availability in imaging neuroscience may be alleviated by *data augmentation* strategies, using *deep neural-network algorithms* whose parameters were already estimated on independent data, and other tricks <sup>51</sup>.

### **Open data become test bed and reference point**

The modus operandi in many empirical sciences is still to collect and analyse in-house data for publication in one paper. Various kinds of questions simply cannot be asked quantitatively using one small dataset, such as extracting links between a human's genetic blueprint and her vast diversity of behaviors <sup>52</sup>. Often genomics amasses data from participant samples collaboratively to chase small effects in multi-site consortia as a confederated research endeavor (e.g., Psychiatric Genomics Consortium). There are always more incentives and maturing practices to accumulate, curate, and distribute data for exploration, knowledge generation, and intervention <sup>53</sup>. This trend reverberates in various empirical research communities and is reinforced by data sharing mandates increasingly specified by funding agencies <sup>54</sup>.

Availability of rich open datasets enables using and intersecting data in unexpected ways, fuels continuous development of novel *multivariate* pattern-learning techniques, and renders new research questions actionable. A trusted community dataset can provide a

common test bed for those analysis methods as well as benchmarks to compare against a set of state-of-the-art methods in different processing pipelines. As an early tradition in machine learning, MNIST established itself as a community-wide dataset with 70,000 images of scanned hand-written digits '0' to '9'. New approaches are expected to beat the globally recorded status quo, and to compare against human performance in the task of number detection <sup>51</sup>. Kaggle-like competitions are also gaining momentum, where a data-analysis challenge is announced and a larger portion of an existing dataset is provided for model development ([www.kaggle.com](http://www.kaggle.com)). At the end of the competition, the modeling solutions from each team are evaluated and ranked on the private part of the dataset (examples from neuroimaging: <sup>55,56</sup>). Such prime example of healthy competition starts to show high efficacy to crowd-source novel analysis strategies to solve global challenges in biomedicine, as well as in business and government. Open data sharing is also an opportunity to dramatically reduce research costs. More broadly, in the future, new data-analysis methods will perhaps be validated empirically based on statistical performance on shared datasets across diverse existing studies and across various workflows <sup>4,11</sup>.

In imaging neuroscience, a majority of the large-scale data initiatives so far were retrospective collections of independently acquired data from different research centers <sup>57</sup>. Such data repositories can vary considerably in key properties, such as data quality and quality-control procedures. Across-site heterogeneity may explain why, counterintuitively, predictive model performance has been repeatedly reported to decrease as the available data increase <sup>58</sup>. As an ambitious attempt to create a large-scale neuroimaging dataset, the ENIGMA consortium launched in 2009 to centrally orchestrate research projects and recruitment of participating groups by providing analysis pipelines and quality control protocols. Several thousand participants were characterized with different imaging



modalities and genetic profiling. A smaller number of data initiatives realized prospectively planned collections with agreed-upon standards for data acquisition. Ensuing repositories offer higher data comparability due to strengthened efforts to, among many others, calibrated acquisition conditions, staff training, or traveling experts. The Human Connectome Project (HCP) was launched in 2009 <sup>59</sup> to promote insight into human brain connectivity by providing extensive multi-modal measurements of ~1,200 healthy adults (aged 22-35), including ~300 twin pairs. For each participant, the project gathered structural, functional, and diffusion MRI, genotyping data, as well as a variety of >400 demographic, behavioral, and lifestyle indicators. With genetic profiling and an extensive variety of phenotyping descriptors, UK Biobank Imaging is even more comprehensive. This data collection initiative set out in 2006 to gather genetic and environmental (e.g., nutrition, lifestyle, medications) data of ~500,000 volunteers (aged 40-69) and is currently the world's largest biomedical dataset. Its brain and body imaging extension was launched in 2014 (with the brain imaging gathering structural, functional, diffusion, and susceptibility-weighted MRI for ~100,000 participants by 2022) <sup>44</sup>.

Compared to more established application domains of machine learning, large datasets of human populations pose additional challenges. Extensive phenotypical profiling calls for a careful balance of trust between protecting each participant's privacy and providing rich open biomedical datasets to the larger research community. UK Biobank participants may be given the possibility to opt-out of sharing records and retroactively deny consent at any time. Industry can boost health-related big-data analytics by offering computing infrastructure as well as data gathering. However, possible conflicts of interest need to be taken into consideration. As public and media perception is very important,

transparent presentation, that is both enthusiastic about the benefits of a study but also completely honest, is crucial to avoid inaccurate negative messages being promulgated.

### **Powerful “black box” predictions supplement simple models**

As a core value of classical data analysis, insight is maximized by assuming linear additivity in how the input variables relate to each other and to the output prediction <sup>5,11</sup>. A traditional goal of statistics is to cleanly isolate the (*univariate*) effects of “special” variables on an outcome, such as a risk factor or a treatment response. All components of the model were supposed to be readily understandable by the investigator. The input variables were typically meticulously hand-picked and chosen to have meaningful units based on existing domain knowledge. This analysis paradigm of generating subject-matter understanding from “introspecting” isolated variable relationships has contributed tremendously to scientific progress in the 20<sup>th</sup> century <sup>2</sup>. However, this explainable modeling regime may also have exhausted the repertoire of natural phenomena that can be usefully described and understood by straightforward linear modeling (but see <sup>60</sup>).

In many empirical sciences, including imaging neuroscience, investigators started moving towards more complex modeling approaches and analysis pipelines. Expanding data resources is a prerequisite to estimating flexible, highly adaptive models that have a larger capacity to represent convoluted relationships between variables, such as hierarchical dependencies and higher-order non-linearity <sup>61</sup>. As more data become available, empirical scientists can now bring to bear more flexible models. In certain cases, the price one may have to pay is that some aspects of the estimated model remain partly opaque to human

intuition, pushing investigators to give up on uncompromised model transparency. As an early hint, *Bayesian hierarchical modeling* can gain traction on complicated datasets that handle nested data settings (e.g., brain scans from participants in different cities) with many more model parameters than input variables. These extensions of classical linear models allow for integrating disparate information sources, sharing statistical strengths between variability sources, de-escalating concerns of class imbalance and selection bias, and estimating full *uncertainty* distributions <sup>28</sup>. As a side effect of increased model complexity, however, not every single parameter value of such a *Bayesian hierarchical model* may merit equal attention for scientific interpretation.

As a continuation of this theme, in adaptive machine-learning algorithms, and especially in *deep neural-network algorithms*, much emphasis is put on the output of a model. This change of focus is why *identifiability* may receive lesser attention in certain studies, although model interpretability was key in classical statistics. Take for example a neurosurgeon who wants to remove brain tissue without impairing language. By relying on a linear model, she predicts outcome in a language task from neural activity measured across the cortex <sup>34</sup>. If quite different model parameter solutions yield an identical prediction accuracy, the model is not *identifiable*. This fitted linear model cannot be physiologically interpreted as a brain map indicating where tissue resection is safer to preserve language capacity. Different candidate models would have other parameters with small absolute values that can suggest diverging brain locations to be less implicated in language processes, which hampers the classical goal towards mechanistic explanation.

Instead, the predictive analyst would typically neglect such arbitrary values of estimated model parameters and prioritize successful prediction of language performance. This neuroscience example illustrates that parameter interpretation is often more

challenging in *multivariate* models optimized for prediction performance <sup>2</sup>. The difficulty in explaining the role of individual input variables is even bigger in current *deep neural-network architectures* <sup>62</sup>. Here, the output predictions can result from highly nonlinear processing cascades from the input variables. There may be little hope to exhaustively understand every single one of the thousands or millions of model parameters in some of today's machine-learning models (<sup>63</sup>; blogpost by Frank Harrell at <http://www.fharrell.com/post/medml/>). Although some remedies have been recently proposed <sup>64,65</sup>, these largely provide understandable simplifications of or linear approximations to the actual non-linear prediction function.

Consequently, for increasingly popular, powerful prediction models, classical (*parametric*) inference may become more challenging to obtain statistically significant p-values. In complicated non-linear models it is partly infeasible to assess the “trueness” of an effect of individual input variables, as an exclusive path for scientific knowledge creation <sup>8,66,67</sup>. These developments do not belittle the importance of working theories in guiding the cumulative construction of scientific understanding. However, there are certain hard problems in empirical research where estimating complex “black box” models may be one of the very few viable solutions. This is probably the case in weather forecasting, perhaps also in some areas of neuroscience. Assessing which aspects of a phenomenon have been successfully captured or inadvertently ignored in an estimated model will probably more often rely on predictive simulation of new data or querying the obtained model for predictions on unseen participants or observations in the 21<sup>th</sup> century. As a side effect of optimizing uncompromised prediction performance, some neuroscience applications may move away from the goal of causal discovery, or even move away from cumulative creation of scientific knowledge <sup>68</sup>. Aiming for crude prediction performance estranges the investigator from asking ‘why?’ - the reason behind statistical relationships between certain

input variables. Today, there is still no commonly agreed upon framework for causal inference.

Instead of carving out new biological mechanisms in nature, prediction metrics can capture how well an estimated complex “black box” model can “imitate” or “reproduce” the studied phenomenon. The *Bayesian-frequentist* debate, which ignited much controversy in 20<sup>th</sup> century statistics (e.g., <sup>69</sup>), may give way to a new antagonistic discourse. One candidate dilemma is classical statistical inference in interpretable models, versus prediction accuracy of complicated natural phenomena, with particularly flexible analysis approaches to quantitatively describe particularly complex systems, such as the human brain in health and disease. In many settings, empirical scientists may have to prioritize ‘providing insight’ (i.e., classical statistical inference targeted at single input variables) against ‘accurately modeling the world’ (i.e., model prediction outputs) <sup>70</sup>. The implied domain interpretability trade-offs will have important consequences for ethical considerations and policy making <sup>71</sup>.

In imaging neuroscience, the prediction-inference antagonism has surfaced as whether or not the prediction accuracies of increasingly used *multivariate* pattern-learning approaches and machine-learning algorithms should undergo post-hoc statistical significance testing (cf. <sup>46,58,72</sup>). This discussion appears to highlight a culture clash between different data-analysis communities seldom in contact before <sup>11</sup>. Neuroimaging and other empirical academic fields have been dominated by a decade-long legacy of initial linear-regression-type estimation and subsequent statistical null-hypothesis testing. Since its inception, machine learning, however, has put a premium on prediction performance as “hard currency” <sup>11,19</sup>. This is especially the case given the immediate practical relevance for various data-intensive industries, such as recommendation systems or micro-targeted customer advertisement <sup>73</sup>. In general, input variables that do enhance prediction accuracy

do not always declare to be statistically significant <sup>13,24,57</sup>. Conversely, variables that are assessed to be statistically significant can be useless for the goal of prediction in new data in certain cases (cf. <sup>67</sup>). To recapitulate these diverging modeling notions of importance, validating a built machine-learning model based on the metric of successful *out-of-sample* prediction is based on extracting patterns in the training data and evaluating how well these identified relationships extrapolate to independent observations drawn from the same distribution. In contrast, classical null-hypothesis significance testing pursues a different analytical goal in asking the question whether an obtained prediction accuracy exceeds two standard deviations of happening by chance under some null hypothesis.

## Conclusion

Historically, innovation and changing practices in quantitative analytics have been shaped by trends in application domains. The recent advent of massive data in neuroscience and biomedicine is ushering towards larger revisions in everyday analytical practices: 1) Unconstrained linear regression models have been a workhorse in 20<sup>th</sup> century empirical research. However, in the 21<sup>st</sup> century, analysis that biases model estimation by parameter regularization or involves dimensionality-reducing transformations may become ubiquitous as extensive datasets become more available. 2) Many classical models have routinely been backed up by consistency theorems, emulating infinite sample size to characterize the model's quality for converging to a good parameter solution. Assessing model quality based on the variance explained in the fitted data will probably be increasingly supplemented by empirical validation procedures such as prediction accuracy in untouched data. 3) Many

empirical sciences like neuroscience may transition from the predominance of a-priori planned, gathered, and published experimental data to conducting more re-analyses of freely available data resources with deep and wide phenotyping. Traditionally, scientific value was seen in unique private datasets. Now, creative modeling strategies may become key to making the most of mushrooming open datasets. 4) Powerful predictive models may not easily lend themselves to exhaustive understanding of all extracted variable-variable relationships. However, the democratization and feasibility of advanced pattern-learning algorithms may enable quantifying always more sophisticated phenomena in nature. Hence, an important dilemma arises between classical relevance claims about single model parameters and successful “black box” predictions of complicated natural phenomena. These megatrends in data analytics hopefully propel the scientific description of particularly complex systems, epitomized by the human brain in health and disease.

## Glossary

Technical term (examples from neuroscience)	Core intuition
<b>Bagging</b> ( <sup>74</sup> )	‘Wisdom of crowds’ strategy to enhance predictive performance by averaging several outcome predictions from models that have been fitted to resampled versions of the same dataset.
<b>Bias vs. variance</b> ( <sup>75</sup> )	The <i>bias-variance trade-off</i> calibrates between losing information (bad fit to data at hand) or succumbing to noise (bad extrapolation to new data). A model with high <i>bias</i> tends to ignore relevant patterns in the data - <i>underfitting</i> due to low <i>effective degrees of freedom</i> . A model with high <i>variance</i> tends to extract arbitrary patterns in the data - <i>overfitting</i> due to high <i>effective degrees of freedom</i> .
<b>Bayesian vs. frequentist modeling</b> ( <sup>76-78</sup> )	<i>Bayesian</i> modeling assumes the model parameters to be random and the data to be fixed; vice versa for <i>frequentist</i> modeling. Consequently, <i>Bayesian</i> analysis provides certainty distributions for each model parameter value, while <i>frequentist</i> analysis yields a single best-guess value for each model parameter. <i>Bayesian</i> model posterior distributions are conditioned on the data at hand, while <i>frequentist</i> model estimation implicitly averages across other data one could have observed.
<b>Canonical correlation analysis</b> ( <sup>79</sup> )	( <i>Multivariate</i> ) pattern-discovery approach that extends the idea of principal component analysis to two variable sets or two datasets. Mutual dependences are extracted as correlated linear combinations of these two data matrices.
<b>Closed-form solutions</b>	A mathematical formula that solves a problem “in one shot” by a circumscribed set of computing operations (i.e., non-iterative), always yielding the same result in the same amount of time.
<b>Consistency theorems or asymptotic guarantees</b>	A widely-used class of mathematical proofs that study the properties of a given modeling approach by taking the number of data points to infinity. It is without asymptotic consistency guarantees, that <i>finite-sample theorems</i> describe properties of modeling approaches as a function of the number of available data points.
<b>Cross-validation</b> ( <sup>80,81</sup> )	A ( <i>non-parametric</i> ) sequential resampling procedure used as the gold standard to practically quantify the performance of predictive models to extrapolate discovered patterns to future data. First, <i>model estimation</i> is carried out by fitting the parameter values to the training data ( <i>in-sample</i> ). Second, if model hyper-parameters need to be set, <i>model selection</i> can be carried out on another independent data split - <i>validation data</i> - to automatically tune towards a winning hyper-parameter combination. Third, <i>model evaluation</i> then quantifies the pattern generalization based on predictive performance in independent <i>hold-out data</i> ( <i>out-of-sample</i> ). The overall process is repeated for different splits of the available data (usually 5 or 10 times). <i>Underfitting</i> yields bad in- and bad out-of-sample generalization performance. <i>Overfitting</i> yields excellent in- and bad out-of-sample prediction accuracy.
<b>Curse of dimensionality (“High dimensions”)</b> ( <sup>82</sup> )	If data have abundant input variables, relative to the available number of data points, each such input dimension is populated with and represented by less data points. Hence, even classical ( <i>unregularized</i> ) linear regression can be over-parameterized. In this setting, common models have trouble finding patterns existing in the data and goodness-of-fit metrics (computed <i>in-sample</i> ) may become impotent. <i>Variance</i> in model parameter estimation escalates and thus data <i>overfitting</i> becomes a core challenge.
<b>Data augmentation</b>	A heterogeneous group of ad-hoc engineering tricks to repeatedly duplicate and modify the original data points while trying to keep their characteristics realistic. The increased effective sample size can allow for estimating more robust model parameters.
<b>Deep neural-network algorithms (“Deep learning”)</b> ( <sup>83-85</sup> )	A growing class of pattern-learning algorithms that perform prediction based on a non-linear, hierarchical, multi-layer neural-network model. Deep neural-network algorithms are able to fit parameters of a particularly high number of nested non-linear processing layers and have extreme freedom in fitting patterns in data.



<b>Degrees of freedom</b>	The number of separate pieces of information to be estimated from data. In the setting of classical linear regression, the degrees of freedom typically refer to the number of independent data points $n$ minus the number of fitted model parameters $p$ to estimate residual errors. This conception is starting to struggle or is difficult to compute for many modern adaptive modeling approaches.
<b>Dimensionality reduction</b> ( <sup>86,87</sup> )	Breaking down the number of input variables to a (much) smaller number of quintessential summary variables. Examples include clustering approaches, like <i>k-means</i> , to partition an array of input variables into typically few non-overlapping variable groups, and matrix decomposition approaches, like PCA and ICA, to extract new continuous representations spanning across input variables that may have partial overlap with each other.
<b>Exchangeability</b> ( <sup>49</sup> )	A characteristic of the data that is a more general form of the independent-and-identically-distributed (i.i.d.) assumption. For example, null hypothesis testing may be used to try to reject that males and females have the same average height. Here, exchangeability may be imposed by shuffling which height measurement belongs to male or female participants to assess whether the summary statistics differ given otherwise identical joint distributions among the input variables.
<b>Ground truth</b> ( <sup>88</sup> )	The true pattern in nature to be approximated by using quantitative modeling of empirical measurements.
<b>Hierarchical multi-level regression</b> ( <sup>27,89</sup> )	Extension of classical linear regression, where the model parameters are also themselves modeled. Linear interactions are introduced by data-level regression parameters being <i>regularized</i> in groups towards upper-level model parameters to “borrow statistical strength”, such as between study sites. Can be carried out in the <i>frequentist</i> regime, but lends itself particularly well to <i>Bayesian</i> modeling. Linear hierarchical regression can more readily fit models with more parameters $p$ than observations $n$ .
<b>Identifiability</b> ( <sup>34</sup> )	Whether the parameter values of a given model can be unambiguously estimated and thus meaningfully interpreted. The combination of data scenario and model properties may result in <i>identifiability</i> (highly valued in classical statistics) or <i>non-identifiability</i> (often a smaller concern in predictive machine-learning applications). <i>Non-identifiable</i> model parameters can result in very different fitted values despite identical prediction performance of the overall model.
<b>K-means clustering</b> ( <sup>90</sup> )	A popular clustering algorithm that partitions the $p$ input variables into $k$ non-overlapping groups.
<b>Markov chain Monte Carlo (MCMC) sampling</b> ( <sup>91</sup> )	An iterative sampling procedure for numerical approximation of challenging posterior integrals, such as those often arising in <i>Bayesian</i> statistics. Each random ‘draw’ yields one candidate set of model parameters that are jointly plausible as to how the data could have come about.
<b>Maximum likelihood estimation (MLE)</b> ( <sup>35</sup> )	Formal ( <i>parametric</i> ) framework on how to find one good set of model parameter values (assumed to be fixed) that maximize the plausibility of how the data may have come about given a pre-specified model. Ordinary linear regression and other classical approaches are special cases. MLE enjoys strong <i>asymptotic guarantees</i> , but can incur problems as the number of input variables $p$ increase.
<b>Multivariate vs. univariate modeling</b> ( <sup>34,72,92,93</sup> )	Technically, <i>univariate</i> methods consider one variable at a time, whereas <i>multivariate</i> methods consider several, possibly many, variables at a time. In neuroscience applications, “univariate” analysis has often been taken to refer to estimating effects for a single brain location, particular brain connection, or specific gene at a time. Instead, “multivariate” analysis would jointly assess patterns in many such biological measurements.
<b>Parametric vs. non-parametric</b> ( <sup>47</sup> )	A <i>parametric</i> approach explicitly assumes structure or a particular form of how the input variables relate to the output. A <i>non-parametric</i> approach tries to fully model the data themselves, for instance by avoiding assuming Gaussian normality in the data.
<b>Partial least squares</b> ( <sup>17,18</sup> )	(Multivariate) pattern-discovery approach aimed at decomposition of covariation, similar to <i>canonical correlation analysis</i> . While partial least squares operates on the un-normalized covariation, <i>canonical correlation analysis</i> acts on the data in a scale/unit-invariant fashion.

<b>Permutation procedures</b> ( <sup>47,49</sup> )	A computation-intensive group of typically <i>non-parametric</i> resampling procedures, which can control error rates (including correction for multiple comparisons), while making few theoretical assumptions. For instance, such procedures enable computation of empirical distributions under some null hypothesis for significance testing in a much wider range of analysis scenarios.
<b>Posterior predictive checks</b>	In <i>Bayesian</i> modeling, generating new data (typically outcome predictions) from model parameter sets sampled from <i>MCMC</i> chains to assess discrepancies between an obtained probabilistic model and the actual data at hand.
<b>Regularization / penalization / shrinkage / sparsification</b> ( <sup>94,95</sup> )	<i>Bias</i> is introduced on purpose in model estimation, for instance to address <i>the curse of dimensionality</i> . As one widespread example, <i>sparse</i> modeling via <i>L1</i> term characteristically drives towards <i>variable selection</i> by encouraging exactly-zero model parameter values (cf. Lasso regression). As another widely-used example, <i>L2</i> terms intentionally skew model parameter values to be closer to zero (cf. Ridge regression).
<b>Tail-area statistics</b> ( <sup>96</sup> )	Instead of some aggregate statistic (e.g., mean, median, mode), interest lies in the shape of a data distribution; especially its extremes with low probability. Special interest was placed outside of the 95% interval assessing whether or not an observation exceeds two standard deviations as in null-hypothesis significance testing.
<b>Variability</b> ( <sup>97</sup> )	A property of the data. E.g., how does the volume of the amygdala really differ between individuals? The <i>variability</i> of a parameter estimate does not go to zero as the sample size approaches infinity, in contrast to <i>uncertainty</i> .
<b>Uncertainty</b> ( <sup>97</sup> )	A property of the modeling approach. E.g., how sure are we about the modeling estimate of amygdala volumes? The <i>uncertainty</i> of an estimated parameter value goes to zero as the number of data points <i>n</i> increases indefinitely, in contrast to <i>variability</i> . <i>Frequentist</i> standard deviations or error bars may mix aspects of <i>variability</i> and <i>uncertainty</i> , in contrast to <i>Bayesian</i> posterior density intervals.

## References

- 1 Efron, B. *Large-scale inference: empirical Bayes methods for estimation, testing, and prediction*. Vol. 1 (Cambridge University Press, 2012).
- 2 Efron, B. & Hastie, T. *Computer-Age Statistical Inference*. (Cambridge University Press, 2016).
- 3 Jordan, M. I. On statistics, computation and scalability. *Bernoulli* **19**, 1378-1390 (2013).
- 4 Donoho, D. 50 Years of Data Science. *Journal of Computational and Graphical Statistics* **26**, 745-766 (2017).
- 5 Casella, G. & Berger, R. L. *Statistical inference*. Vol. 2 (Duxbury Pacific Grove, CA, 2002).
- 6 Efron, B. & Tibshirani, R. J. Statistical data analysis in the computer age. *Science* **253**, 390-395, doi:10.1126/science.253.5018.390 (1991).
- 7 Nuzzo, R. Scientific method: statistical errors. *Nature* **506**, 150-152, doi:10.1038/506150a (2014).
- 8 Wasserstein, R. L. & Lazar, N. A. The ASA's statement on p-values: context, process, and purpose. *Am Stat* **70**, 129-133 (2016).
- 9 Blei, D. M. & Smyth, P. Science and data science. *Proceedings of the National Academy of Sciences* **114**, 8689-8692 (2017).
- 10 Halevy, A., Norvig, P. & Pereira, F. The unreasonable effectiveness of data. *Intelligent Systems, IEEE* **24**, 8-12 (2009).
- 11 Breiman, L. Statistical Modeling: The Two Cultures. *Statistical Science* **16**, 199-231 (2001).
- 12 Jordan, M. I. *et al. Frontiers in Massive Data Analysis*. (The National Academies Press, 2013).
- 13 Bzdok, D. & Yeo, B. T. T. Inference in the age of big data: Future perspectives on neuroscience. *NeuroImage* **155**, 549-564 (2017).
- 14 Smith, S. M. & Nichols, T. E. Statistical challenges in "Big Data" human neuroimaging. *Neuron* **97**, 263-268 (2018).
- 15 Elliott, L. T. *et al.* Genome-wide association studies of brain imaging phenotypes in UK Biobank. *Nature* **562**, 210 (2018).
- 16 Amunts, K. *et al.* BigBrain: an ultrahigh-resolution 3D human brain model. *Science* **340**, 1472-1475, doi:10.1126/science.1235381 (2013).
- 17 McIntosh, A. R. & Mišić, B. Multivariate statistical analyses for neuroimaging data. *Annual review of psychology* **64**, 499-525 (2013).
- 18 McIntosh, A., Bookstein, F., Haxby, J. V. & Grady, C. Spatial pattern analysis of functional brain images using partial least squares. *NeuroImage* **3**, 143-157 (1996).
- 19 Hastie, T., Tibshirani, R. & Friedman, J. *The Elements of Statistical Learning*. (Springer Series in Statistics, 2001).
- 20 Giraud, C. *Introduction to high-dimensional statistics*. (CRC Press, 2014).
- 21 Hastie, T., Tibshirani, R. & Wainwright, M. *Statistical Learning with Sparsity: The Lasso and Generalizations*. (CRC Press, 2015).
- 22 Mohri, M., Talwalkar, A. & Rostamizadeh, A. *Foundations of machine learning (adaptive computation and machine learning series)*. (Mit Press Cambridge, MA, 2012).
- 23 Shalev-Shwartz, S. & Ben-David, S. *Understanding machine learning: From theory to algorithms*. (Cambridge University Press, 2014).
- 24 McElreath, R. (Chapman & Hall/CRC, Boca Raton, FL, USA, 2015).
- 25 Kruschke, J. K. *Doing Bayesian Data Analysis*. (Elsevier, 2011).
- 26 Wipf, D. P. & Nagarajan, S. S. in *Advances in neural information processing systems*. 1625-1632.
- 27 Chen, G. *et al.* Handling Multiplicity in Neuroimaging through Bayesian Lenses with Multilevel Modeling. *Neuroinformatics*, 1-31 (2018).
- 28 Gelman, A., Carlin, J. B., Stern, H. S. & Rubin, D. B. *Bayesian data analysis*. Vol. 2 (Chapman & Hall/CRC Boca Raton, FL, USA, 2014).
- 29 MacKay, D. J. C. *Information theory, inference and learning algorithms*. (Cambridge university press, 2003).

- 30 Smith, S. M. *et al.* A positive-negative mode of population covariation links brain  
connectivity, demographics and behavior. *Nature neuroscience* **18**, 1565-1567,  
doi:10.1038/nn.4125 (2015).
- 31 Witten, D. M., Tibshirani, R. & Hastie, T. A penalized matrix decomposition, with applications  
to sparse principal components and canonical correlation analysis. *Biostatistics*, kxp008  
(2009).
- 32 Virtanen, S., Klami, A. & Kaski, S. in *Proceedings of the 28th International Conference on  
International Conference on Machine Learning*. 457-464 (Omnipress).
- 33 Andrew, G., Arora, R., Bilmes, J. & Livescu, K. in *International Conference on Machine  
Learning*. 1247-1255.
- 34 Haufe, S. *et al.* On the interpretation of weight vectors of linear models in multivariate  
neuroimaging. *NeuroImage* **87**, 96-110 (2014).
- 35 Friston, K. J. *et al.* Statistical parametric maps in functional imaging: a general linear  
approach. *Human brain mapping* **2**, 189-210 (1994).
- 36 Kernbach, J. M. *et al.* Subspecialization within default mode nodes characterized in 10,000  
UK Biobank participants. *Proceedings of the National Academy of Sciences of the United  
States of America* **115**, 12295-12300 (2018).
- 37 Bzdok, D. *et al.* Characterization of the temporo-parietal junction by combining data-driven  
parcellation, complementary connectivity analyses, and functional decoding. *NeuroImage* **81**,  
381-392, doi:10.1016/j.neuroimage.2013.05.046 (2013).
- 38 Wang, H.-T. *et al.* Dimensions of experience: exploring the heterogeneity of the wandering  
mind. *Psychological science* **29**, 56-71 (2018).
- 39 Zhang, C., Bengio, S., Hardt, M., Recht, B. & Vinyals, O. Understanding deep learning requires  
rethinking generalization. *arXiv preprint arXiv:1611.03530* (2016).
- 40 Stone, M. Cross-validated choice and assessment of statistical predictions. *Journal of the  
royal statistical society. Series B (Methodological)*, 111-147 (1974).
- 41 Geisser, S. The predictive sample reuse method with applications. *Journal of the American  
statistical Association* **70**, 320-328 (1975).
- 42 Breiman, L. Bagging predictors. *Machine learning* **24**, 123-140 (1996).
- 43 Efron, B. & Tibshirani, R. J. *An introduction to the bootstrap*. (CRC press, 1994).
- 44 Miller, K. L. *et al.* Multimodal population brain imaging in the UK Biobank prospective  
epidemiological study. *Nature neuroscience* (2016).
- 45 Berkson, J. Some difficulties of interpretation encountered in the application of the chi-  
square test. *Journal of the American Statistical Association* **33**, 526-536 (1938).
- 46 Bzdok, D. Classical statistics and statistical learning in imaging neuroscience. *Frontiers in  
neuroscience* **11**, 543 (2017).
- 47 Nichols, T. E. & Holmes, A. P. Nonparametric permutation tests for functional neuroimaging:  
a primer with examples. *Human brain mapping* **15**, 1-25 (2002).
- 48 Winkler, A. M. *et al.* Non-parametric combination and related permutation tests for  
neuroimaging. *Human brain mapping* **37**, 1486-1511 (2016).
- 49 Ge, T., Yeo, B. T. T. & Winkler, A. A brief overview of permutation testing with examples.  
*Organization for Human Brain Mapping*, [https://www.ohbmbbrainmappingblog.com/blog/a-  
brief-overview-of-permutation-testing-with-examples](https://www.ohbmbbrainmappingblog.com/blog/a-brief-overview-of-permutation-testing-with-examples) (2018).
- 50 Varoquaux, G. Cross-validation failure: small sample sizes lead to large error bars.  
*NeuroImage* (2017).
- 51 Goodfellow, I. J., Bengio, Y. & Courville, A. *Deep learning*. (MIT Press, 2016).
- 52 Medland, S. E., Jahanshad, N., Neale, B. M. & Thompson, P. M. Whole-genome analyses of  
whole-brain data: working within an expanded search space. *Nature neuroscience* **17**, 791  
(2014).
- 53 Leonelli, S. *Data-centric biology: a philosophical study*. (University of Chicago Press, 2016).
- 54 Poldrack, R. A. & Gorgolewski, K. J. Making big data open: data sharing in neuroimaging.  
*Nature neuroscience* **17**, 1510-1517, doi:10.1038/nn.3818 (2014).

- 55 Bron, E. E. *et al.* Standardized evaluation of algorithms for computer-aided diagnosis of dementia based on structural MRI: the CADDementia challenge. *NeuroImage* **111**, 562-579 (2015).
- 56 Sarica, A., Cerasa, A., Quattrone, A. & Calhoun, V. Editorial on special issue: Machine learning on MCI. *Journal of neuroscience methods* **302**, 1 (2018).
- 57 Arbabshirani, M. R., Plis, S., Sui, J. & Calhoun, V. D. Single subject prediction of brain disorders in neuroimaging: Promises and pitfalls. *NeuroImage* **145**, 137-165 (2017).
- 58 Woo, C.-W., Chang, L. J., Lindquist, M. A. & Wager, T. D. Building better biomarkers: brain models in translational neuroimaging. *Nature neuroscience* **20**, 365-377, doi:10.1038/nn.4478 (2017).
- 59 Van Essen, D. C. *et al.* The Human Connectome Project: a data acquisition perspective. *NeuroImage* **62**, 2222-2231 (2012).
- 60 Petkova, E. *et al.* Statistical analysis plan for stage 1 EMBARC (Establishing Moderators and Biosignatures of Antidepressant Response for Clinical Care) study. *Contemporary clinical trials communications* **6**, 22-30 (2017).
- 61 Ghahramani, Z. Probabilistic machine learning and artificial intelligence. *Nature* **521**, 452-459, doi:10.1038/nature14541 (2015).
- 62 LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436-444, doi:10.1038/nature14539 (2015).
- 63 Shmueli, G. To explain or to predict? *Statistical science*, 289-310 (2010).
- 64 Lundberg, S. M. & Lee, S.-I. in *Advances in Neural Information Processing Systems*. 4765-4774.
- 65 Chen, J., Song, L., Wainwright, M. J. & Jordan, M. I. Learning to Explain: An Information-Theoretic Perspective on Model Interpretation. *arXiv preprint arXiv:1802.07814* (2018).
- 66 Szucs, D. & Ioannidis, J. When null hypothesis significance testing is unsuitable for research: a reassessment. *Frontiers in human neuroscience* **11**, 390 (2017).
- 67 Bzdok, D. & Ioannidis, J. P. A. Exploration, inference and prediction in neuroscience and biomedicine. *Trend in Neurosciences* **42**, 251-262 (2019).
- 68 Pearl, J. & Mackenzie, D. *The book of why: the new science of cause and effect*. (Basic Books, 2018).
- 69 Efron, B. Why Isn't Everyone a Bayesian? *The American Statistician* **40**, 1-5, doi:10.2307/2683105 (1986).
- 70 Norvig, P. On chomsky and the two cultures of statistical learning. *Author Homepage* (2011).
- 71 O'Neil, C. Weapons of Math Destruction. *How Big Data Increases Inequality and Threatens Democracy*, New York: Crown (2016).
- 72 Haynes, J.-D. A primer on pattern-based approaches to fMRI: Principles, pitfalls, and perspectives. *Neuron* **87**, 257-270 (2015).
- 73 Henke, N. *et al.* The age of analytics: Competing in a data-driven world. *Technical report, McKinsey Global Institute*. (2016).
- 74 Hoyos-Idrobo, A., Varoquaux, G., Schwartz, Y. & Thirion, B. FReM—scalable and stable decoding with fast regularized ensemble of models. *NeuroImage* **180**, 160-172 (2018).
- 75 Yarkoni, T. & Westfall, J. Choosing prediction over explanation in psychology: Lessons from machine learning. *Perspectives on Psychological Science* (2016).
- 76 Friston, K. J. *et al.* Classical and Bayesian inference in neuroimaging: applications. *NeuroImage* **16**, 484-512, doi:10.1006/nimg.2002.1091 (2002).
- 77 Friston, K. J. *et al.* Classical and Bayesian inference in neuroimaging: theory. *NeuroImage* **16**, 465-483, doi:10.1006/nimg.2002.1090 (2002).
- 78 Körding, K. P. & Wolpert, D. M. Bayesian integration in sensorimotor learning. *Nature* **427**, 244 (2004).
- 79 Friston, K. J., Liddle, P. F., Frith, C. D., Hirsch, S. R. & Frackowiak, R. S. J. The left medial temporal region and schizophrenia. *Brain : a journal of neurology* **115**, 367-382 (1992).
- 80 Varoquaux, G. *et al.* Assessing and tuning brain decoders: cross-validation, caveats, and guidelines. *NeuroImage* **145**, 166-179 (2017).

- 81 Pereira, F., Mitchell, T. & Botvinick, M. Machine learning classifiers and fMRI: a tutorial  
overview. *NeuroImage* **45**, 199-209, doi:10.1016/j.neuroimage.2008.11.007 (2009).
- 82 Allen, E. A., Erhardt, E. B. & Calhoun, V. D. Data visualization in the neurosciences:  
overcoming the curse of dimensionality. *Neuron* **74**, 603-608 (2012).
- 83 Marblestone, A. H., Wayne, G. & Kording, K. P. Toward an integration of deep learning and  
neuroscience. *Frontiers in computational neuroscience* **10**, 94 (2016).
- 84 Plis, S. M. *et al.* Deep learning for neuroimaging: a validation study. *Frontiers in neuroscience*  
**8** (2014).
- 85 Hassabis, D., Kumaran, D., Summerfield, C. & Botvinick, M. Neuroscience-inspired artificial  
intelligence. *Neuron* **95**, 245-258 (2017).
- 86 Doria, V. *et al.* Emergence of resting state networks in the preterm human brain. *Proceedings  
of the National Academy of Sciences of the United States of America* **107**, 20015-20020,  
doi:10.1073/pnas.1007921107 (2010).
- 87 Sui, J. *et al.* A CCA+ ICA based model for multi-task brain imaging data fusion and its  
application to schizophrenia. *NeuroImage* **51**, 123-134 (2010).
- 88 Jonas, E. & Kording, K. P. Could a neuroscientist understand a microprocessor? *PLoS  
computational biology* **13**, e1005268 (2017).
- 89 Dai, T., Guo, Y. & Initiative, A. s. D. N. Predicting individual brain functional connectivity using  
a Bayesian hierarchical model. *NeuroImage* **147**, 772-787 (2017).
- 90 Eickhoff, S. B., Thirion, B., Varoquaux, G. & Bzdok, D. Connectivity-based parcellation:  
Critique and implications. *Human brain mapping*, doi:10.1002/hbm.22933 (2015).
- 91 Woolrich, M. W. Bayesian inference in FMRI. *NeuroImage* **62**, 801-810 (2012).
- 92 Haxby, J. V. *et al.* Distributed and overlapping representations of faces and objects in ventral  
temporal cortex. *Science* **293**, 2425-2430 (2001).
- 93 Kriegeskorte, N., Goebel, R. & Bandettini, P. Information-based functional brain mapping.  
*PNAS* **103**, 3863-3868, doi:10.1073/pnas.0600244103 (2006).
- 94 Rasmussen, P. M., Hansen, L. K., Madsen, K. H., Churchill, N. W. & Strother, S. C. Model  
sparsity and brain pattern interpretation of classification models in neuroimaging. *Pattern  
Recognition* **45**, 2085-2100 (2012).
- 95 Baldassarre, L., Pontil, M. & Mourão-Miranda, J. Sparsity is better with stability: combining  
accuracy and stability for model selection in brain decoding. *Frontiers in neuroscience* **11**, 62  
(2017).
- 96 Woo, C. W., Krishnan, A. & Wager, T. D. Cluster-extent based thresholding in fMRI analyses:  
pitfalls and recommendations. *NeuroImage* **91**, 412-419,  
doi:10.1016/j.neuroimage.2013.12.058 (2014).
- 97 Faisal, A. A., Selen, L. P. & Wolpert, D. M. Noise in the nervous system. *Nature reviews.  
Neuroscience* **9**, 292-303, doi:10.1038/nrn2258 (2008).

## Figure 1

**Significant population mode that relates patterns of inter-network connectivity to patterns within deep behavioral phenotyping.** In the ~500 participant release of the Human Connectome Project, canonical correlation analysis lends itself particularly well to uncovering multi-modal correspondences between brain and behavior. Intrinsic network coupling fluctuations between 200 nodes were demonstrated to bear rich relationships with >100 cognitive assessments, demographic profiles, and life-factor indicators. A functional connectivity fingerprint emerged with rich profiling of behavioral associations that varied along a global positive-negative axis with high intelligence, memory and cognition performances on the one end, and negative lifestyle measures and events on the other end. The brain regions exhibiting strongest contributions to coherent connectivity changes were reminiscent of the default mode network, which is implicated in episodic memory and semantic capacity, mental scene construction, and complex social reasoning such as taking other people's perspective. Reprinted with permission from <sup>30</sup>.

## Figure 2

**Strongest population mode that links intra-network connectivity patterns and inter-network connectivity patterns.** In ~10,000 UK Biobank participants, canonical correlation analysis was used to identify robust correspondences between functional connectivity shifts inside a major brain network (*top*), the default mode network, and functional connectivity shifts between a set of major brain networks (*bottom*). This large-scale analysis made apparent that specific subregions inside the default mode network, namely, the right and left anterior temporoparietal junction, could play a dominant role in the process of global network reconfiguration in humans. Reprinted with permission from <sup>36</sup>.

## Figure 3

**Relevance of population associations between six brain-imaging modalities and thousands of behavioral phenotypes.** For ~21,000 UK Biobank participants, this Manhattan plot depicts results from ~15 million cross-subject association tests (*each color indicates a different neuroimaging modality*). The horizontal dashed lines indicate significance after correction for multiple statistical comparisons based on Bonferroni's more stringent method (Bonf, *upper line*) or more modern false discovery rate (FDR, *lower line*). Even after accounting for family-wise error, ~28,000 (Bonferroni) or ~180,000 (FDR) brain-behavior associations remained statistically significant at the population level. These results demonstrate the rich relationships between different brain tissue measurements and extensive phenotyping of many thousand individuals. Figure computed analogous to previous study on the UK Biobank <sup>44</sup>.