

Bias/Variance is not the same as Approximation/Estimation

Anonymous authors

Paper under double-blind review

Abstract

We study the relation between two classical results: the bias-variance decomposition, and the approximation-estimation decomposition. Both are important conceptual tools in Machine Learning, helping us describe the nature of model fitting. It is commonly stated that the two decompositions “are closely related”, or are “similar in spirit”. However, sometimes it is said they “are equivalent” (spoiler: no, they’re not). In reality, they have subtle connections, cutting across learning theory and classical statistics, that (very surprisingly) have not been previously observed. In this paper we uncover these connections, and build a bridge between the two.

1 Introduction

Geman et al. (1992) introduced the bias-variance decomposition to the Machine Learning community, and Vapnik & Chervonenkis (1974) introduced the approximation-estimation decomposition, founding the field of statistical learning theory. Both decompositions help us understand model fitting: involving model size, and some kind of tradeoff. The terms are sometimes used interchangeably. And yet, they are different things. Given their fundamental nature and similar purposes, it is surprising that explicit connections are not widely known—perhaps due to differing notations and conventions of their respective communities. Our goal is to uncover these connections and build a bridge between these two seminal results.

The approximation-estimation decomposition refers to models drawn from some function class \mathcal{F} , and considers their *excess risk*—that is, the risk above that of the Bayes model—breaking it into two components:

$$\text{excess risk} = \text{approximation error} + \text{estimation error}. \quad (1)$$

We might choose to increase the size of our function class, perhaps by adding more parameters to our model. In this situation it is commonly understood that the approximation error will decrease, and the estimation error will increase (Von Luxburg & Schölkopf, 2011), beyond a certain point resulting in over-fitting of the model. In contrast to the abstract notion of a “function class”, the *bias-variance* decomposition considers the risk of real, trained models, in *expectation over possible training sets*. It breaks the expected risk into two components:

$$\text{expected risk} = \text{bias} + \text{variance}. \quad (2)$$

As we increase model size: the bias tends to decrease, and the variance tends to increase, again determining the degree of over-fitting. Recently, it has become apparent that this trade-off is not always simple, e.g. with over-parameterised models; but, the decomposition still holds even if a simple tradeoff does not.

We therefore have two decompositions: both referring to model size, with some kind of trade-off between their terms, and with bearing on the nature of over-fitting. It is easy, and common, to conflate these. From online discussion forums, to the lecture notes of esteemed institutions and well-cited research articles, one can observe innocent statements such as “the two are closely related”, but also the more extreme (and incorrect) “the approximation-estimation tradeoff is also known as the bias-variance tradeoff”.

To the best of our knowledge, this is the first work to discuss their connection explicitly. We consider a range of loss functions: including Bregman divergences and the 0/1 loss, identifying the properties of each decomposition, studying where they coincide and where they do not.

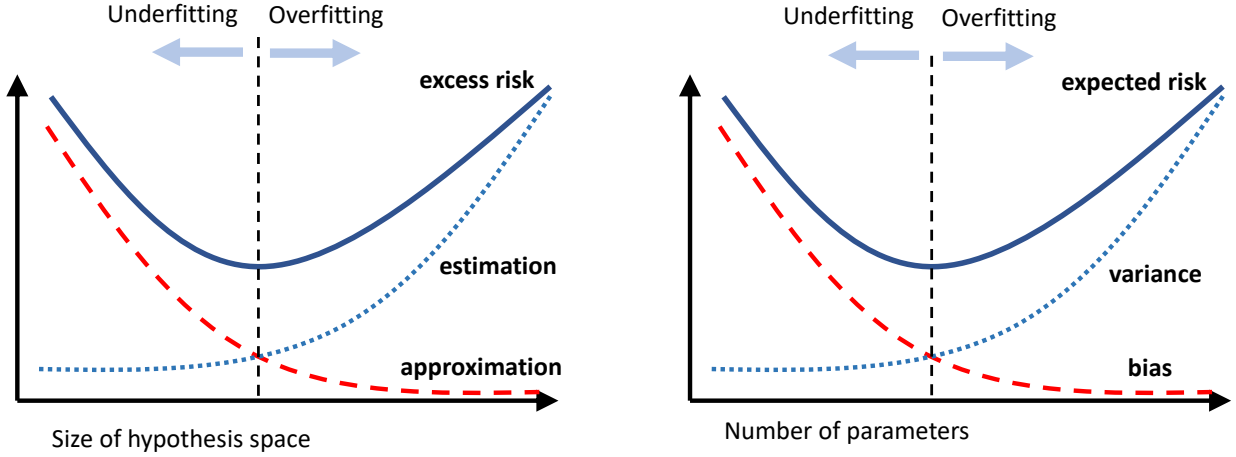


Figure 1: Two diagrams (same on left/right is intentional) illustrating how the approximation/estimation and bias/variance tradeoffs are commonly described, and easily confused.

2 Background

We introduce notation and review the two decompositions.

2.1 Preliminaries

Consider a standard supervised learning setup, where the task is to map from an input $\mathbf{x} \in \mathcal{X}$ to an output $y \in \mathcal{Y}$, and we assume there exists an unknown joint probability distribution $P(\mathbf{x}, y)$. This is achieved by learning the parameters of a model $f : \mathbf{x} \rightarrow y$, which can be seen as selecting an f from a set $\mathcal{F} \subset \mathcal{F}_{all}$, a subset of the space of all measurable functions. The discrepancy of $f(\mathbf{x})$ from the true y is quantified with a loss function $\ell(y, f(\mathbf{x}))$, which may or may not be symmetric. Using this, we define the *risk* of a given f ,

$$R(f) = \mathbb{E}_{\mathbf{x}y}[\ell(y, f(\mathbf{x}))] = \int \ell(y, f(\mathbf{x})) dP(\mathbf{x}, y). \quad (3)$$

The *Bayes* model g is that, drawn from \mathcal{F}_{all} , which minimizes this quantity, i.e.

$$g := \arg \inf_{f \in \mathcal{F}_{all}} R(f). \quad (4)$$

Given that we picked a family $\mathcal{F} \subset \mathcal{F}_{all}$, this may not be achievable. The best-in-family model f^* is defined

$$f^* := \arg \inf_{f \in \mathcal{F}} R(f). \quad (5)$$

Both of these are defined in terms of the true distribution $P(\mathbf{x}, y)$. In practice, we have only a sample $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ where each point is drawn i.i.d. from $P(\mathbf{x}, y)$. We write the *empirical risk* as:

$$R_{emp}(f; D) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i)), \quad (6)$$

and a model in \mathcal{F} that minimizes this, known as an *empirical risk minimizer*, is defined:

$$\hat{f}_{erm} := \arg \inf_{f \in \mathcal{F}} R_{emp}(f; D). \quad (7)$$

We use a ‘hat’ notation to emphasize that the ERM is dependent on D , the random training data sample. The model \hat{f}_{erm} is non-unique, but the best performance we could hope to find with our limited training data. For some models/losses an ERM is achievable—e.g. the closed-form solution for linear models under squared loss. However in general, training a model \hat{f} will not necessarily result in an ERM. We can now cover the specifics for the two decompositions.

2.2 The Approximation-Estimation decomposition

The approximation-estimation decomposition is a seminal observation from the 1970s work of Vapnik and Chervonenkis, reviewed in Vapnik (1999). An excellent historical account can be found in Bottou (2013). The result deals with the *excess* risk $R(\hat{f}_{erm}) - R(g)$, i.e. the risk of \hat{f}_{erm} above that of the Bayes model, g . The approximation-estimation decomposition, applicable for any loss ℓ , breaks this into two terms:

$$\underbrace{R(\hat{f}_{erm}) - R(g)}_{\text{excess risk}} = \underbrace{R(\hat{f}_{erm}) - R(f^*)}_{\text{estimation error}} + \underbrace{R(f^*) - R(g)}_{\text{approximation error}}. \quad (8)$$

The approximation error is the additional risk due to using a restricted family \mathcal{F} , rather than the space of all functions \mathcal{F}_{all} . This is a systematic quantity, not dependent on any particular data sample. The estimation error is the additional risk due to our finite training data, when trying to identify $f^* \in \mathcal{F}$. This is a random variable, dependent on the particular data sample used to obtain \hat{f}_{erm} . There is a natural tradeoff (see Figure 1) as we change the size of \mathcal{F} , keeping data size fixed. As we increase $|\mathcal{F}|$, approximation error will likely decrease (potentially to zero, if $g \in \mathcal{F}$), but estimation error will increase, as it becomes harder to identify f^* in the larger space. The reason behind this is, in effect, the classical *multiple hypothesis testing* problem—we cannot reliably distinguish many hypotheses when our dataset is small. *Sample complexity bounds* address this by bounding estimation error in terms of n and $|\mathcal{F}|$ (Von Luxburg & Schölkopf, 2011).

Bottou & Bousquet (2007) proposed an extension of Equation 8, which turns out important for our work. For many learning scenarios, it is infeasible to find a global minimum on the training data (i.e. \hat{f}_{erm}), and we can only have a sub-optimal model \hat{f} . An additional risk component then emerges, and the excess risk of \hat{f} decomposes into a sum of optimisation error, estimation error, and approximation error:

$$\underbrace{R(\hat{f}) - R(g)}_{\text{excess risk of } \hat{f}} = \underbrace{R(\hat{f}) - R(\hat{f}_{erm})}_{\text{optimisation error}} + \underbrace{R(\hat{f}_{erm}) - R(f^*)}_{\text{estimation error}} + \underbrace{R(f^*) - R(g)}_{\text{approximation error}}. \quad (9)$$

These three terms describe the learning process in abstract form: accounting respectively for the choice of learning algorithm, the quality/amount of data, and the capacity of the model family.

2.3 The Bias-Variance decomposition

A bias-variance decomposition involves the *expected* risk of a model \hat{f} , where the expectation \mathbb{E}_D is over *all possible* training sets D of a fixed size n . For the squared loss $\ell(y, \hat{f}) = (y - \hat{f})^2$, Geman et al. (1992) showed:

$$\underbrace{\mathbb{E}_D [R(\hat{f})]}_{\text{expected risk}} = \underbrace{\mathbb{E}_{xy} [(y - \mathbb{E}_{y|x}[y])^2]}_{\text{noise}} + \underbrace{\mathbb{E}_x \left[\left(\mathbb{E}_D [\hat{f}(x)] - \mathbb{E}_{y|x}[y] \right)^2 \right]}_{\text{bias}} + \underbrace{\mathbb{E}_x \left[\mathbb{E}_D [(\hat{f}(x) - \mathbb{E}_D [\hat{f}(x)])^2] \right]}_{\text{variance}}. \quad (10)$$

The bias is a systematic component, independent of any particular training set D , and commonly regarded as measuring the ‘strength’ of a model. The variance measures the sensitivity of \hat{f} to changes in the training sample, independent of the true label y . The noise is a constant, independent of any model parameters. There is again a perceived tradeoff with these terms. With many models, as the size of the (un-regularised) model increases: bias *tends* to decrease, and variance *tends* to increase. However, the tradeoff can be more complex (e.g. with over-parameterized models) and the exact dynamics are an open research issue.

It is a common misconception that this holds only for squared loss. In fact, the same form holds for any Bregman divergence (Bregman, 1967; Pfau, 2013), e.g. KL-divergence, or Poisson regression. If we define ℓ as an arbitrary Bregman divergence, parameterised by a strictly convex generator function ϕ , then:

$$\underbrace{\mathbb{E}_D [\mathbb{E}_{xy} [\ell(y, \hat{f}(x))]]}_{\text{expected risk}} = \underbrace{\mathbb{E}_{xy} [\ell(y, \bar{y})]}_{\text{noise}} + \underbrace{\mathbb{E}_x [\ell(\bar{y}, \hat{f}_\phi(x))]}_{\text{bias}} + \underbrace{\mathbb{E}_x [\mathbb{E}_D [\ell(\hat{f}_\phi(x), \hat{f}(x))]]}_{\text{variance}}. \quad (11)$$

where $\bar{y} = \mathbb{E}_{y|x}[y]$. If we use $\phi(f) = f^2$, $\ell(y, f) = (y - f)^2$, and $\hat{f}_\phi = \mathbb{E}_D[\hat{f}]$, we have Equation 10. In general, \hat{f}_ϕ is not always the expectation $\mathbb{E}_D[\hat{f}]$, instead a measure of centrality derived from the loss, referred to as a *centroid* (Nielsen & Nock, 2009). As this will be important in the coming sections, we define it formally.

Definition 1 (Centroid for a Bregman divergence) Assume a Bregman divergence $B_\phi(y, \hat{f})$ defined by a strictly convex generator function ϕ . For a model distribution induced by a random variable D , the centroid is that closest (as defined by the particular Bregman divergence) on average to all others.

$$\mathring{f}_\phi := \arg \min_z \mathbb{E}_D [B_\phi(z, \hat{f})] = [\nabla \phi]^{-1} \left(\mathbb{E}_D [\nabla \phi(\hat{f})] \right). \quad (12)$$

For example with a KL-divergence, the centroid¹ is a normalized geometric mean. Of course, the notion of a centroid can be defined for *any* loss of interest: $\mathring{f} := \arg \min_z \mathbb{E}_D [\ell(z, \hat{f})]$. To distinguish a Bregman centroid from an arbitrary centroid, we use notation \mathring{f}_ϕ versus \mathring{f} . Further examples of Bregman divergences (and their centroids) can be found in the table below. The bias/variance terms take *different functional forms* for each particular loss, instantiated by different Bregman generators. This has a consequence for nomenclature: the term defined in Geman et al. (1992) is sometimes referred to as “*squared bias*”, but the square turns out to be an artefact from using squared loss, not present in other decompositions, hence we use simply ‘bias’.

Table 1: Examples of losses (and their centroids) which admit a bias-variance decomposition.

Name	Range	Loss $\ell(y, \hat{f})$	Centroid \mathring{f}_ϕ
Squared	$y \in \mathbb{R}$	$(y - \hat{f})^2$	$\mathbb{E}_D[\hat{f}]$
Poisson	$y \in \{0, 1, 2, \dots\}$	$y \ln \frac{y}{\hat{f}} - (y - \hat{f})$	$\exp(\mathbb{E}_D[\ln \hat{f}])$
KL-divergence	$\mathbf{y} \in \mathbb{R}^k, s.t. \sum_c y_c = 1$	$D_{KL}(\mathbf{y} \parallel \hat{f})$	$Z^{-1} \exp(\mathbb{E}_D[\ln \hat{f}])$
Ikatura-Saito	$y \in [0, \infty)$	$\frac{y}{\hat{f}} - \ln \frac{y}{\hat{f}} - 1$	$1/\mathbb{E}_D[\hat{f}^{-1}]$

Note also that the KL-divergence example also implies a decomposition for the ubiquitous cross-entropy loss, since the two differ only by a constant. It is furthermore interesting to note that, whilst the generalised decomposition above only appeared in the ML community in Pfau (2013), the result seems to be known much earlier in classical statistics, e.g. mentioned in passing by Hastie & Tibshirani (1986, Eq. 19).

The bias-variance decomposition does not hold for all losses. The approximation-estimation decomposition, Equation 9, applies for *any* loss function. This is not the case for bias-variance. The variance term in Equation 11 is *independent of the true label, y* . This is an elemental property: a *requirement* for something to be called a ‘bias-variance’ decomposition. As such, Equation 11 does not hold for all losses, e.g. with the 0/1 mis-classification loss, the residual term that we might call ‘variance’ is necessarily dependent on the label (Wood et al., 2023). The necessary and sufficient conditions for such a decomposition are an open question.

2.4 Summary

The two decompositions are *conceptual* tools to describe the nature of model fitting. They are by no means perfect reflections of the process, most especially in the context of over-parameterized models (Nagarajan & Kolter, 2019; Zhang et al., 2021). However, it is *extremely* common to see papers making the incorrect assumption/claim that the two are equivalent, or that one is a special case of the other. Our purpose with this work is simply to correct these false assumptions, identifying *precisely* how the two connect.

¹Note that this is a minimization over the first (left-hand) argument, so it is technically a *left* centroid. The right centroid can be similarly defined, turning out to be simply $\mathbb{E}_D[\hat{f}]$ for any valid ϕ (Nielsen & Nock, 2009).

3 Bias/Variance is not the same as Approximation/Estimation

By now it should be evident that these decompositions are related, but are not quite the same thing. They were conceived and built upon by different sub-communities—one in classical statistics, the other in statistical learning theory—with differing notations and purposes. We now build a bridge between the two.

The estimation error involves $R(\hat{f}_{erm})$, making it a random variable dependent on D . We take the expectation with respect to D , and separate it into two components.

Theorem 1 (Decomposing the Expected Estimation Error) *For an arbitrary loss ℓ , and corresponding risk $\mathbb{E}_{\mathbf{x}y}[\ell(y, f)]$ the (expected) estimation error decomposes as so:*

$$\underbrace{\mathbb{E}_D [R(\hat{f}_{erm}) - R(f^*)]}_{\mathcal{E}_{est}} = \underbrace{\mathbb{E}_D [R(\hat{f}_{erm}) - R(\hat{f})]}_{\mathcal{E}_{est(v)}} + \underbrace{R(\hat{f}) - R(f^*)}_{\mathcal{E}_{est(b)}}. \quad (13)$$

The *estimation bias*, $\mathcal{E}_{est(b)}$, is independent of any particular training sample D . The *estimation variance*, $\mathcal{E}_{est(v)}$, measures the *random* variations of \hat{f}_{erm} around the centroid, with respect to changes in D . If the risk of \hat{f}_{erm} varies greatly with changes in D , this will be large. The proof of this is trivial: adding/subtracting the risk of the centroid \hat{f} . Whilst this seems arbitrary (we could add/subtract other terms) it is the insertion of this particular $R(\hat{f})$ that will allow us to see the relation between the decompositions.

Theorem 2 (Bias-Variance in terms of Approximation-Estimation) *For a loss $\ell(y, f) = B_\phi(y, f)$, the following decomposition of the bias and variance applies.*

$$\underbrace{\mathbb{E}_{\mathbf{x}} [\ell(\bar{y}, \hat{f}_\phi(\mathbf{x}))]}_{\text{bias}} = \underbrace{R(f^*) - R(g)}_{\text{approximation error}} + \underbrace{R(\hat{f}_\phi) - R(f^*)}_{\text{estimation bias}} \quad (14)$$

$$\underbrace{\mathbb{E}_{\mathbf{x}} [\mathbb{E}_D [\ell(\hat{f}_\phi(\mathbf{x}), \hat{f}(\mathbf{x}))]]}_{\text{variance}} = \underbrace{\mathbb{E}_D [R(\hat{f}) - R(\hat{f}_{erm})]}_{\text{optimisation error}} + \underbrace{\mathbb{E}_D [R(\hat{f}_{erm}) - R(\hat{f}_\phi)]}_{\text{estimation variance}} \quad (15)$$

This confirms the premise of our paper. Bias is *not* approximation error, and variance is *not* estimation error. It is not even the case that one is a special case of the other, as is sometimes stated. The true relation is more subtle. The approximation error is in fact just *one component of the bias*, and, the estimation error *contributes to both bias and variance*. The theorem above is illustrated in [Figure 2](#).

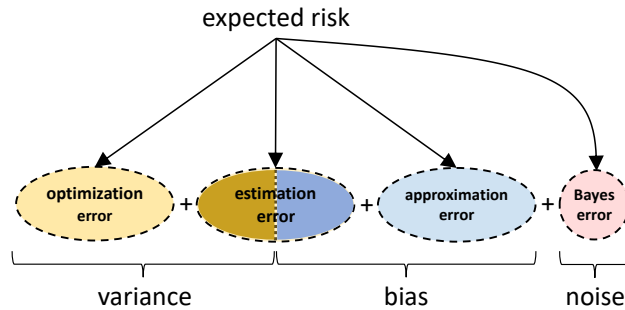


Figure 2: Illustration of Theorem 2. The bias is only partly determined by approximation error (i.e. choice/size of \mathcal{F}), while the rest comes from the expected estimation error (i.e. choice of data). Similarly, variation in data accounts for only part of the variance, and the remainder is due to the optimisation error (i.e. choice of learning algorithm).

4 Discussion

A simplistic description of bias and variance might say they are the error ‘*due to the model*’ (bias) and the error ‘*due to the data*’ (variance). [Theorem 2](#) shows there is more nuance to understand. We now discuss the subtleties and implications of these results.

4.1 The bias is a flawed proxy for model capacity.

It is common to assume the bias is an indication of how simple/complex a model is—expected to be lower if the model has higher ‘capacity’. But what is model ‘capacity’? If we take it to be the ability to minimize the population risk, then the *ultimate* measure of model capacity is the *approximation error*. We see from [Equation 14](#) that the bias contains exactly this, but also the additional *estimation bias* term, which gives it some surprising dynamics.

For a squared loss with a linear model, [Equation 14](#) has been noted² before ([Hastie et al., 2017](#), Eq 7.14). Our result generalises it to a broader range of losses and arbitrary *non-linear* models. For tractability, their analyses were restricted to linear models. However tractable, they were unable to observe a critical fact—that in the general case, estimation bias $R(\hat{f}) - R(f^*)$, can take *negative values*, i.e.

$$\text{bias} = \underbrace{\left[\begin{array}{c} \text{approximation} \\ \text{error} \end{array} \right]}_{\text{always } \geq 0} + \underbrace{\left[\begin{array}{c} \text{estimation} \\ \text{bias} \end{array} \right]}_{\text{can be negative}} \quad (16)$$

To understand how this can be, we must accept the somewhat non-intuitive idea that the centroid can be outside the hypothesis class \mathcal{F} . It turns out this is trivially possible, even with a simple regression stump.

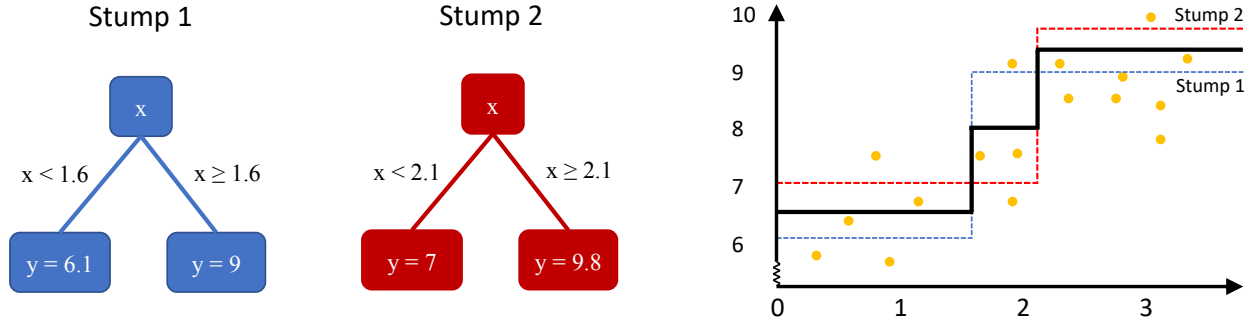


Figure 3: Two regression stumps (red/blue lines), and their centroid (black line, the arithmetic mean). Notice the centroid is *outside* the hypothesis class, i.e. it cannot be represented as a binary stump, as it has four turning points. As a result, the centroid fits the data better than any $f \in \mathcal{F}$, and $\mathcal{E}_{est(b)}$ is negative.

The possibility of negative values here has significant implications. There are two ways in which bias can be zero. If \mathcal{F} contains the Bayes model g , then we might have $\mathcal{E}_{app} = \mathcal{E}_{est(b)} = 0$. But, there is another way. For some $\epsilon > 0$, we might have $\mathcal{E}_{app} = \epsilon$, and $\mathcal{E}_{est(b)} = -\epsilon$. In this case, the model family does *not* have sufficient capacity, since $\mathcal{E}_{app} > 0$. And yet, the bias is zero. Hence, the bias is a flawed proxy for the true model capacity.

To illustrate the point, we show experiments on a synthetic regression problem. Details in [Appendix C](#).

²Hastie et al. described the first term on the right as “the error between the best-fitting linear approximation and the true function”. This is exactly the definition of approximation error for the linear model. We provide a proof of this relation and further discussion in the appendix.

Figure 4 shows results increasing the depth of a decision tree. The left panel shows excess risk, and the bias/variance components. We observe the classical bias/variance tradeoff, including overfitting, as the depth increases beyond a certain point. It is notable that *the bias decreases to zero*, after depth 6. *Does this imply the model is ‘unbiased’, in the sense that it has sufficient capacity to capture the full data distribution?*

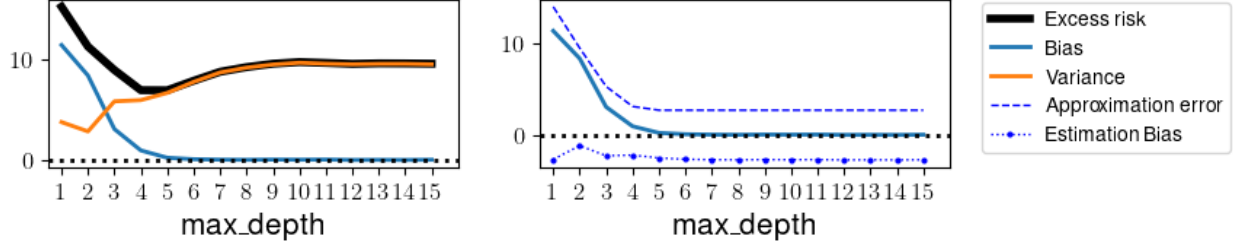


Figure 4: Risk components as we increase the depth of a regression tree.

The answer is no. A decomposition of the bias into two components (right panel) shows that the approximation error is non-zero, i.e. the best possible model *cannot* achieve zero testing error. The cause of the bias going to zero is that the estimation bias is negative, hence the bias is not a good proxy for the true model capacity. Very similar results are obtained with a k -nearest neighbour regression (Figure 5), where increasing complexity is obtained by *decreasing* the value of k .

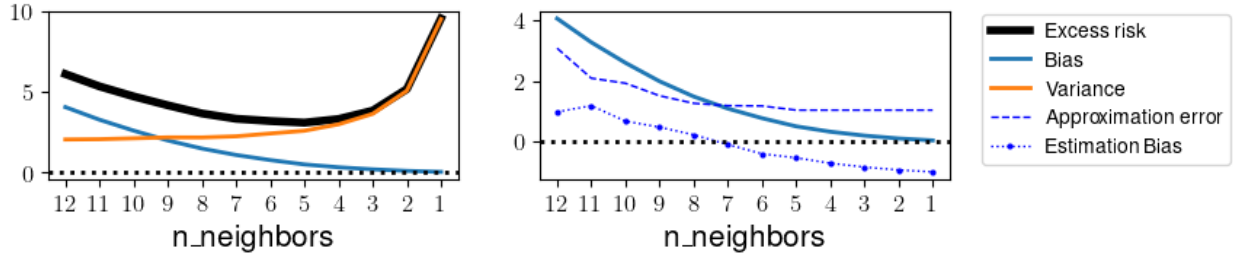


Figure 5: Risk components as we decrease the number of neighbours in a k-nn.

We can formally characterise this phenomenon, by studying the geometry of the hypothesis class \mathcal{F} . In particular, if the set \mathcal{F} is *dual-convex* (Amari, 2008) with respect to the generator ϕ , then $\hat{f} \in \mathcal{F}$, and hence estimation bias is guaranteed to be non-negative.

Theorem 3 (Sufficient condition for a non-negative estimation bias.) *If the hypothesis class \mathcal{F} is dual-convex, then the estimation bias is non-negative.*

A simple example of a non-dual convex set is the class of regression stumps evaluated by squared loss, where $\hat{f} = \mathbb{E}_D[\hat{f}]$, illustrated in Figure 3. A simple example of a convex set is the class of Generalized Linear Models evaluated by their corresponding deviance measure.

Theorem 4 (GLMs have non-negative estimation bias.) *For a Bregman divergence with generator ϕ , define \mathcal{F} as the set of all GLMs with inverse link $[\nabla\phi]^{-1}$ and natural parameters $\theta \in \mathbb{R}^d$. Then, the estimation bias is non-negative.*

An example of this would be a logistic regression, $\hat{f}(\mathbf{x}) = [\nabla\phi]^{-1}(\hat{\theta}^T \mathbf{x}) = 1/(1 + \exp(-\hat{\theta}^T \mathbf{x}))$, which results from $\phi(f) = f \ln f + (1 - f) \ln(1 - f)$ and the binary KL as the deviance.

4.2 New insights into the bias/variance trade-off.

In the age of deep learning, the relevance (and even existence) of a *trade-off* between bias and variance has been debated, with voices both against (Neal et al., 2018; Dar et al., 2021) and in favour (Witten, 2020). [Theorem 1](#) places a constraint between $\mathcal{E}_{est(b)}$ and $\mathcal{E}_{est(v)}$, the estimation bias and the estimation variance. When the estimation bias is negative (e.g. [Figure 3](#)), it obviously *reduces* the bias. However, it simultaneously *increases* the variance, since $\mathcal{E}_{est(v)}$ has an imposed lower bound, satisfying the constraint. Therefore, for every single reduction in bias attributable to a negative estimation bias, *the same quantity will be lost* in the increased variance. This is an *unavoidable* tradeoff. Obviously other components of the bias/variance may mask this behaviour, making the tradeoff less visible.

4.3 The estimation variance plays a role in double descent.

In many models, an increasing degree of over-parameterisation has been associated with a ‘peaking’ trend in the variance ([Nakkiran, 2019](#); [Yang et al., 2020](#)), ultimately causing a *double descent* in the risk.

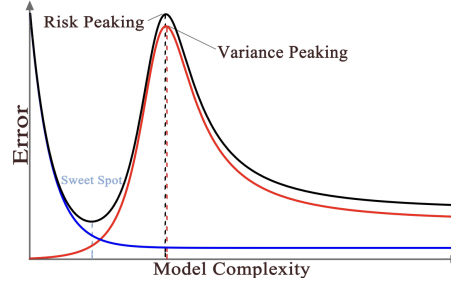


Figure 6: Illustration of double descent, caused by a ‘peaking’ variance (red line) and monotonically decreasing bias (blue line). Image credit [Yang et al. \(2020\)](#).

Such models often fit their training data perfectly ([Belkin et al., 2019](#); [Zhang et al., 2021](#)), i.e. they *interpolate* the data. If we consider this in the context of [Equation 15](#), we see that:

$$\text{variance} = \begin{bmatrix} \text{estimation} \\ \text{variance} \end{bmatrix} + \underbrace{\begin{bmatrix} \text{optimisation} \\ \text{error} \end{bmatrix}}_{\approx 0 \text{ for interpolating models}}.$$

i.e., the optimization error is close to zero. This observed ‘peaking’ variance must therefore be primarily due to the *estimation variance*, $\mathcal{E}_{est(v)}$. Furthermore, very deep models are likely to be able to fit any function, i.e. their approximation error is zero. In these scenarios, the *only terms* remaining in the expected risk are $\mathcal{E}_{est(b)}$ and $\mathcal{E}_{est(v)}$. *Why* such models can push training error to zero, even on random labels ([Zhang et al., 2021](#)), and still generalise well, remains an open question for modern machine learning. Overall, we believe this warrants further study in the context of deep models.

5 Conclusions

We analysed the precise connections between two seminal results: the bias-variance decomposition, and the approximation-estimation decomposition. Perhaps the most surprising aspect of this work was that it had not been explored before—two such foundational ideas, not previously connected. On a literature review, we found numerous sources incorrectly stating the two were equivalent, or related as a special case / general case. This is false. The true relation, given by [Theorem 2](#), is more intricate, and yielded interesting novel observations that we detailed through [section 4](#). The centroid of a model distribution, \hat{f} , turned out to be a key mathematical object in bridging the two decompositions. We conjecture that further study of this object, and its role in generalisation, may yield yet deeper and interesting insights.

References

- Shun-ichi Amari. Information geometry and its applications: Convex function and dually flat manifold. In *LIX Fall Colloquium on Emerging Trends in Visual Computing*, pp. 75–102. Springer, 2008.
- Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proc. National Academy of Sciences*, 116(32):15849–15854, 2019.
- Léon Bottou. In *Hindsight: Doklady Akademii Nauk SSSR*, 181(4), 1968, pp. 3–5. Springer Berlin Heidelberg, Berlin, Heidelberg, 2013. ISBN 978-3-642-41136-6. doi: 10.1007/978-3-642-41136-6_1. URL https://doi.org/10.1007/978-3-642-41136-6_1.
- Léon Bottou and Olivier Bousquet. The tradeoffs of large scale learning. *Advances in neural information processing systems*, 20, 2007.
- Lev M Bregman. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Comp. Mathematics*, 7(3):200–217, 1967.
- Yehuda Dar, Vidya Muthukumar, and Richard G Baraniuk. A farewell to the bias-variance tradeoff? an overview of the theory of overparameterized machine learning. *arXiv preprint arXiv:2109.02355*, 2021.
- Stuart Geman, Elie Bienenstock, and René Doursat. Neural networks and the bias/variance dilemma. *Neural computation*, 4(1):1–58, 1992.
- Trevor Hastie and Robert Tibshirani. Generalized Additive Models. *Statistical Science*, 1(3):297 – 310, 1986. doi: 10.1214/ss/1177013604. URL <https://doi.org/10.1214/ss/1177013604>.
- Trevor Hastie, Robert Tibshirani, and Jerome H Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, volume 12th printing, January 13th 2017. Springer, 2017.
- Vaishnavh Nagarajan and J Zico Kolter. Uniform convergence may be unable to explain generalization in deep learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- Preetum Nakkiran. More data can hurt for linear regression: Sample-wise double descent. *arXiv preprint arXiv:1912.07242*, 2019.
- Brady Neal, Sarthak Mittal, Aristide Baratin, Vinayak Tantia, Matthew Scicluna, Simon Lacoste-Julien, and Ioannis Mitliagkas. A modern take on the bias-variance tradeoff in neural networks. *arXiv preprint arXiv:1810.08591*, 2018.
- Frank Nielsen and Richard Nock. Sided and symmetrized bregman centroids. *IEEE transactions on Information Theory*, 55(6):2882–2904, 2009.
- David Pfau. A Generalized Bias-Variance Decomposition for Bregman Divergences. Technical report, Columbia University, 2013.
- Vladimir Vapnik. *The nature of statistical learning theory*. Springer science & business media, 1999.
- Vladimir Vapnik and Alexey Chervonenkis. *Theory of pattern recognition*. Nauka, Moscow, 1974.
- Ulrike Von Luxburg and Bernhard Schölkopf. Statistical learning theory: Models, concepts, and results. In *Handbook of the History of Logic*, volume 10, pp. 651–706. Elsevier, 2011.
- Daniela Witten. Twitter thread: *The Bias-Variance Trade-Off & "DOUBLE DESCENT"*, 2020. URL https://x.com/daniela_witten/status/1292293102103748609. Posted 3.54am, 9th August, 2020.
- Danny Wood, Tingting Mu, Andrew Webb, Henry Reeve, Mikel Lujan, and Gavin Brown. A unified theory of diversity in ensemble learning. *arXiv preprint arXiv:2301.03962*, 2023.
- Zitong Yang, Yaodong Yu, Chong You, Jacob Steinhardt, and Yi Ma. Rethinking bias-variance trade-off for generalization of neural networks. In *International Conf. on Machine Learning*, 2020.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.

Appendix

A Proofs of Theorems

A.1 Proof of Theorem 1 (Decomposing the Expected Estimation Error).

Allow terms to cancel on the right hand side, leading to the left hand side.

A.2 Proof of Theorem 2 (Bias-Variance in terms of Approximation-Estimation).

We wish to prove the following statements:

$$\underbrace{\mathbb{E}_{\mathbf{x}} \left[\ell(\bar{y}, \hat{f}_{\phi}(\mathbf{x})) \right]}_{\text{bias}} = \underbrace{R(f^*) - R(g)}_{\text{approximation error}} + \underbrace{R(\hat{f}_{\phi}) - R(f^*)}_{\text{estimation bias}} \quad (17)$$

$$\underbrace{\mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_D \left[\ell(\hat{f}_{\phi}(\mathbf{x}), \hat{f}(\mathbf{x})) \right] \right]}_{\text{variance}} = \underbrace{\mathbb{E}_D \left[R(\hat{f}) - R(\hat{f}_{\text{erm}}) \right]}_{\text{optimisation error}} + \underbrace{\mathbb{E}_D \left[R(\hat{f}_{\text{erm}}) - R(\hat{f}_{\phi}) \right]}_{\text{estimation variance}} \quad (18)$$

We first note the definition of the Bayes model, $g := \arg \inf_z \mathbb{E}_{\mathbf{x}y} [\ell(y, z)] = \mathbb{E}_{y|\mathbf{x}}[y]$. Note that this ‘right-hand’ centroid is the same regardless of the Bregman divergence (Nielsen & Nock, 2009). Now, to show Equation 17, we note that the $R(f^*)$ terms cancel, hence we just need to prove:

$$\mathbb{E}_{\mathbf{x}} \left[\ell(\bar{y}, \hat{f}_{\phi}(\mathbf{x})) \right] = R(\hat{f}_{\phi}) - R(g). \quad (19)$$

The proof below builds on the *Bregman 3-point property* (Nielsen & Nock, 2009).

Definition (Bregman three-point identity) *The Bregman three-point property states, for any p, q, r ,*

$$B_{\phi}(p, r) = B_{\phi}(p, q) + B_{\phi}(q, r) + \langle p - q, \nabla \phi(q) - \nabla \phi(r) \rangle \quad (20)$$

We then have the following, where we apply the three-point property to \bar{y}, \hat{f}_{ϕ} , with \bar{y} as the mid-point.

$$B_{\phi}(y, \hat{f}_{\phi}) = B_{\phi}(y, \bar{y}) + B_{\phi}(\bar{y}, \hat{f}_{\phi}) + \langle y - \bar{y}, \nabla \phi(\bar{y}) - \nabla \phi(\hat{f}_{\phi}) \rangle \quad (21)$$

Take the expected value w/r $p(y|\mathbf{x})$ and the inner product term vanishes, since $\bar{y} = \mathbb{E}_{y|\mathbf{x}}[y]$. Rearranging terms and further taking expectation w/r \mathbf{x} , we recover:

$$R(\hat{f}_{\phi}) - R(\bar{y}) = \mathbb{E}_{\mathbf{x}} \left[B_{\phi}(\bar{y}, \hat{f}_{\phi}) \right] \quad (22)$$

which is the desired result, proving Equation 17.

To show Equation 18, we follow a similar pattern. Take the 3-point property for y, \hat{f} with \hat{f}_{ϕ} as the mid-point.

$$B_{\phi}(y, \hat{f}) = B_{\phi}(y, \hat{f}_{\phi}) + B_{\phi}(\hat{f}_{\phi}, \hat{f}) + \langle y - \hat{f}_{\phi}, \nabla \phi(\hat{f}_{\phi}) - \nabla \phi(\hat{f}) \rangle \quad (23)$$

Take the expected value w/r D and the inner product term vanishes, since $\nabla \phi(\hat{f}_{\phi}) = \mathbb{E}_D [\nabla \phi(\hat{f})]$.

Rearranging terms and further taking expectation over $p(\mathbf{x})$, we recover:

$$\mathbb{E}_D \left[R(\hat{f}) - R(\hat{f}_{\phi}) \right] = \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_D \left[B_{\phi}(\hat{f}_{\phi}, \hat{f}) \right] \right] \quad (24)$$

which is the desired result, completing the theorem. \blacksquare

Special case of Theorem 2 for squared loss. The following presents the special case of squared loss, included for didactic purposes due to its ubiquity and links to the results for linear models. We wish to prove the following statements:

$$\mathbb{E}_{\mathbf{x}} \left[(\mathbb{E}_D[\hat{f}(\mathbf{x})] - \mathbb{E}_{y|\mathbf{x}}[y])^2 \right] = \mathcal{E}_{app} + \mathcal{E}_{est(b)} \quad (25)$$

$$\mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_D[(\hat{f}(\mathbf{x}) - \mathbb{E}_D[\hat{f}(\mathbf{x})])^2] \right] = \mathcal{E}_{opt} + \mathcal{E}_{est(v)} \quad (26)$$

$$\mathbb{E}_{\mathbf{x}y} [(y - \mathbb{E}_{y|\mathbf{x}}[y])^2] = R(g) \quad (27)$$

To show equation 27, we simply note the definition of the Bayes model g for the special case of squared loss. This is, $g := \arg \inf_f R(f) = \arg \inf_f \mathbb{E}_{\mathbf{x}y} [(y - f)^2] = \mathbb{E}_{y|\mathbf{x}}[y]$. Plugging this into $R(g)$ is the desired result.

To show Equation 25 we note, as an intermediate step, that:

$$\mathcal{E}_{app} + \mathcal{E}_{est(b)} = \left(R(f^*) - R(g) \right) + \left(R(\mathbb{E}_D[\hat{f}]) - R(f^*) \right) = R(\mathbb{E}_D[\hat{f}]) - R(g). \quad (28)$$

We then have the following, again using the definition of g .

$$\begin{aligned} R(\mathbb{E}_D[\hat{f}]) - R(g) &= \mathbb{E}_{\mathbf{x}y} \left[(\mathbb{E}_D[\hat{f}] - y)^2 \right] - \mathbb{E}_{\mathbf{x}y} \left[(y - \mathbb{E}_{y|\mathbf{x}}[y])^2 \right] \\ &= \mathbb{E}_{\mathbf{x}y} \left[\left(\mathbb{E}_D[\hat{f}] \right)^2 - 2y\mathbb{E}_D[\hat{f}] - \mathbb{E}_{y|\mathbf{x}}[y]^2 + 2y\mathbb{E}_{y|\mathbf{x}}[y] \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[\left(\mathbb{E}_D[\hat{f}] \right)^2 - 2\mathbb{E}_{y|\mathbf{x}}[y]\mathbb{E}_D[\hat{f}] - \mathbb{E}_{y|\mathbf{x}}[y]^2 + 2\mathbb{E}_{y|\mathbf{x}}[y]^2 \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[\left(\mathbb{E}_D[\hat{f}] \right)^2 - 2\mathbb{E}_{y|\mathbf{x}}[y]\mathbb{E}_D[\hat{f}] + \mathbb{E}_{y|\mathbf{x}}[y]^2 \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[(\mathbb{E}_D[\hat{f}] - \mathbb{E}_{y|\mathbf{x}}[y])^2 \right], \end{aligned}$$

which is the bias, and the desired result.

To show Equation 26, we follow a similar pattern. From definitions:

$$\mathcal{E}_{opt} + \mathcal{E}_{est(v)} = \mathbb{E}_D \left[R(\hat{f}) - R(\hat{f}_{erm}) \right] + \mathbb{E}_D \left[R(\hat{f}_{erm}) - R(\mathbb{E}_D[\hat{f}]) \right] = \mathbb{E}_D \left[R(\hat{f}) - R(\mathbb{E}_D[\hat{f}]) \right]. \quad (29)$$

We then have the following.

$$\begin{aligned} \mathbb{E}_D \left[R(\hat{f}) - R(\mathbb{E}_D[\hat{f}]) \right] &= \mathbb{E}_D \left[\mathbb{E}_{\mathbf{x}y} \left[(\hat{f} - y)^2 \right] - \mathbb{E}_{\mathbf{x}y} \left[(\mathbb{E}_D[\hat{f}] - y)^2 \right] \right] \\ &= \mathbb{E}_D \left[\mathbb{E}_{\mathbf{x}y} \left[\hat{f}^2 - 2y\hat{f} - \mathbb{E}_D[\hat{f}]^2 + 2y\mathbb{E}_D[\hat{f}] \right] \right] \\ &= \mathbb{E}_{\mathbf{x}y} \left[\mathbb{E}_D[\hat{f}^2] - 2y\mathbb{E}_D[\hat{f}] - \mathbb{E}_D[\hat{f}]^2 + 2y\mathbb{E}_D[\hat{f}] \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_D[\hat{f}^2] - \mathbb{E}_D[\hat{f}]^2 \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_D[(\hat{f} - \mathbb{E}_D[\hat{f}])^2] \right] \end{aligned}$$

where the final step is the standard definition of variance, giving the desired result. ■

A.3 Proof of Theorem 3 (Sufficient condition for a non-negative estimation bias).

To prove Theorem 3, we demonstrate that under a certain condition, $\mathring{f} \in \mathcal{F}$, which implies $R(\mathring{f}) \geq R(f^*)$, and therefore $R(\mathring{f}) - R(f^*) \geq 0$. We use the following definition, due to Amari (2008).

Definition 2 (Dual convex set) Let ϕ be a strictly convex function. A set \mathcal{F} is dually convex with respect to ϕ iff, for any pair of points $f, g \in \mathcal{F}$ and for all $\lambda \in [0, 1]$

$$\lambda \nabla \phi(f) + (1 - \lambda) \nabla \phi(g) \in \mathcal{F}$$

i.e. the set \mathcal{F} is dual-convex iff it is convex in its dual coordinate representation.

An arbitrary set C is convex iff for any random variable X defined over elements of C , its expectation is also in C , i.e. $\mathbb{E}[X] \in C$.

Therefore, for a dual convex set \mathcal{F} , we have that the point $\mathbb{E}_D[\nabla \phi(f)] \in \mathcal{F}$. The primal coordinate representation of this point, $\nabla \phi^{-1}(\mathbb{E}_D[\nabla \phi(f)])$, is also a member of \mathcal{F} , i.e. $\mathring{f} \in \mathcal{F}$, proving the theorem. ■

A.4 Proof of Theorem 4 (GLMs have non-negative estimation bias).

We demonstrate that $\mathcal{E}_{est(b)} \geq 0$ if \hat{f} is a GLM of a particular form. We give two proofs: a direct one and one that makes use of Theorem 3.

Direct proof. The estimation bias is defined:

$$\mathcal{E}_{est(b)} = R(\mathring{f}) - R(f^*). \quad (30)$$

This involves the definition of the centroid, which for a Bregman divergence is,

$$\mathring{f}_\phi := [\nabla \phi]^{-1} \left(\mathbb{E}_D \left[\nabla \phi(\hat{f}) \right] \right). \quad (31)$$

Given a Bregman divergence with convex generator ϕ , define \mathcal{F} as the class of GLMs with inverse link $[\nabla \phi]^{-1}$, parameterised by $\theta \in \mathbb{R}^d$. In this case, each $\hat{f} \in \mathcal{F}$ takes the form:

$$\hat{f} := [\nabla \phi]^{-1} (\theta^T \mathbf{x}), \quad (32)$$

where θ are the natural parameters. Substituting this into the centroid gives us,

$$\begin{aligned} \mathring{f}_\phi &= [\nabla \phi]^{-1} \left(\mathbb{E}_D \left[\nabla \phi \left([\nabla \phi]^{-1} (\theta^T \mathbf{x}) \right) \right] \right), \\ &= [\nabla \phi]^{-1} \left(\mathbb{E}_D [\theta]^T \mathbf{x} \right). \end{aligned} \quad (33)$$

Since $\mathbb{E}_D[\theta]$ is within the convex hull of the distribution of θ induced by D , the centroid is the same form of GLM as \hat{f} , and therefore $\mathring{f}_\phi \in \mathcal{F}$. Then, since by definition f^* is the risk minimizer in \mathcal{F} , we must have that $R(\mathring{f}_\phi) \geq R(f^*)$, and therefore Equation 30 is non-negative. ■

Proof using Theorem 3. To show that the estimation bias is non-negative, it suffices to show that the class of GLMs of a particular form is *dually-convex*. We verify that the property of dual-convexity holds. Define \mathcal{F} = GLMs with inverse link $\nabla \phi^{-1}$

By definition, if $f \in \mathcal{F}$, it is parameterised by a vector θ as follows: $f = \nabla \phi^{-1}(\theta x)$.

Let $h \in \mathcal{F}$ be a GLM corresponding to the convex combination of two arbitrary GLMs in their dual coordinates, i.e. $h = \nabla\phi^{-1}(\lambda\nabla\phi(f) + (1-\lambda)\nabla\phi(g))$, with $\lambda \in [0, 1]$, and with f and g two GLMs $f = \nabla\phi^{-1}(\theta x)$ and $g = \nabla\phi^{-1}(\xi x)$. Then,

$$\begin{aligned} h &= \nabla\phi^{-1}(\lambda\nabla\phi(f) + (1-\lambda)\nabla\phi(g)) \\ &= \nabla\phi^{-1}(\lambda\nabla\phi(\nabla\phi^{-1}(\theta x)) + (1-\lambda)\nabla\phi(\nabla\phi^{-1}(\xi x))) \\ &= \nabla\phi^{-1}(\lambda\theta x + (1-\lambda)\xi x) \\ &= \nabla\phi^{-1}((\lambda\theta + (1-\lambda)\xi)x) \end{aligned}$$

which is again a GLM in \mathcal{F} .

B Discussion of related work by Hastie et al, 2017

We detail related observations by [Hastie et al. \(2017\)](#), who assume a linear model with squared loss, i.e. $\ell(y, \hat{f}) = (y - \hat{\theta}^T \mathbf{x})^2$. The optimal parameters are $\theta_* := \arg \min_{\theta} \mathbb{E}_{\mathbf{x}y}[(y - \hat{\theta}^T \mathbf{x})^2] = \mathbb{E}_{\mathbf{x}}[\mathbf{x}\mathbf{x}^T]^{-1}\mathbb{E}_{\mathbf{x}}[\mathbf{x}\bar{y}]$, and the Bayes model is $\bar{y} = \arg \min_z \mathbb{E}_{y|\mathbf{x}}[(y - z)^2] = \mathbb{E}_{y|\mathbf{x}}[y]$. In this case, the bias-variance decomposition is,

$$\underbrace{\mathbb{E}_D \left[\mathbb{E}_{\mathbf{x}y} \left[(y - \hat{\theta}^T \mathbf{x})^2 \right] \right]}_{\text{expected risk}} = \underbrace{\mathbb{E}_{\mathbf{x}y} \left[(y - \bar{y})^2 \right]}_{\text{noise}} + \underbrace{\mathbb{E}_{\mathbf{x}} \left[(\bar{y} - \mathbb{E}_D[\hat{\theta}^T \mathbf{x}])^2 \right]}_{\text{bias}} + \underbrace{\mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_D \left[(\hat{\theta}^T \mathbf{x} - \mathbb{E}_D[\hat{\theta}^T \mathbf{x}])^2 \right] \right]}_{\text{variance}}. \quad (34)$$

[Hastie et al. \(2017, Eq 7.14\)](#) show that the bias decomposes more finely:

$$\mathbb{E}_{\mathbf{x}} \left[(\bar{y} - \mathbb{E}_D[\hat{\theta}^T \mathbf{x}])^2 \right] = \mathbb{E}_{\mathbf{x}} \left[(\bar{y} - \theta_*^T \mathbf{x})^2 \right] + \mathbb{E}_{\mathbf{x}} \left[(\theta_*^T \mathbf{x} - \mathbb{E}_D[\hat{\theta}^T \mathbf{x}])^2 \right]. \quad (35)$$

Hastie et al. describe this expression as:

“The first term on the right-hand side is the average squared model bias, the error between the best-fitting linear approximation and the true function. The second term is the average squared estimation bias, the error between the average estimate [...] and the best-fitting linear approximation.”

The first point to note is that Hastie et al. deal only with a decomposition for squared loss. When referring to bias, they use nomenclature ‘squared bias’, whereas in fact the square is an artefact of using squared loss, and is not present in the general case.

When they refer to the error between “the best-fitting linear approximation and the true function.”, we note that this is precisely a description of the *approximation error* for a linear model. This is no coincidence. To see the connection precisely, we note the following property of the approximation error.

Theorem 5 *For a loss ℓ , if a bias-variance decomposition holds, then the approximation error $R(f^*) - R(g)$ simplifies as follows.*

$$\begin{aligned} \mathcal{E}_{app} &= R(f^*) - R(g) \\ &= \mathbb{E}_D \mathbb{E}_{\mathbf{x}y} \ell(g, f^*) \end{aligned} \quad (36)$$

i.e. the difference-of-risks is equal to the divergence between the two models themselves.

Proof sketch. Use the 3-point theorem in exactly the same manner as in the proof of [Theorem 2](#), i.e. between y, f^* with \bar{y} as the mid-point, then take expectation over $y|\mathbf{x}$. ■

For the case of squared loss, this yields the term from Hastie et al, shown above in [Equation 35](#), i.e.,

$$R(\theta_*^T \mathbf{x}) - R(g) = \mathbb{E}_{\mathbf{x}} \left[(\bar{y} - \theta_*^T \mathbf{x})^2 \right]. \quad (37)$$

Overall, this shows that [Hastie et al. \(2017, Eq 7.14\)](#) is the special case of our [Equation 14](#) for squared loss/linear models.

Estimation bias: The second term on the right of equation [35](#) is described as the error between the expected model and the best-fitting linear approximation. This is equivalent to the standard definition of estimation bias, but for the specific case of a linear model and squared loss, i.e.

$$\underbrace{R(\mathbb{E}_D[\hat{\theta}^T \mathbf{x}]) - R(\theta_*^T \mathbf{x})}_{\text{estimation bias}} = \mathbb{E}_{\mathbf{x}} \left[(\theta_*^T \mathbf{x} - \mathbb{E}_D[\hat{\theta}^T \mathbf{x}])^2 \right]. \quad (38)$$

However, the squared loss seems to be unique in that [Equation 38](#) holds. This is a consequence of taking an expectation over \mathbf{x} , and the properties of the OLS solution. In the general Bregman case we have an inequality:

$$\underbrace{R(\hat{f}) - R(f^*)}_{\text{estimation bias}} \neq \mathbb{E}_{\mathbf{x}} [B_{\phi}(f^*, \hat{f})]. \quad (39)$$

Hastie et al. observed that in unregularized linear models, the estimation bias will³ be zero. With ridge regression, $\theta_* := \mathbb{E}_{\mathbf{x}}[\mathbf{x}\mathbf{x}^T + \lambda]^{-1} \mathbb{E}_{\mathbf{x}}[\mathbf{x}\bar{y}]$, and thus the estimation bias will be non-negative for $\lambda > 0$. Since the OLS solution is closed-form, $\mathcal{E}_{opt} = 0$, and the estimation variance is simply $\text{Var}(\hat{\theta}^T \mathbf{x})$.

[Hastie et al. \(2017, Figure 7.2\)](#) also briefly alludes to the idea of estimation *variance*—from this we assume that Hastie *et al.* were well aware of these terms in the context of squared loss / linear models. However, the *difference-of-risks* formulation that we use generalises these ideas to any model family, and any Bregman divergence.

C Experimental details

We summarise our methodology to generate the illustrative experiments shown in the paper. **For the purposes of anonymous submission, an outline is provided below. For the final submission we will supply all code for reproducible research.**

We use a synthetic 1-d problem: $x \in [0, 15]$, and the true label is $y = x + 5 \sin(2x) + \epsilon$, where ϵ is Gaussian noise with zero mean and $\sigma = 3$. Training data is $n = 100$ points, illustrated below.

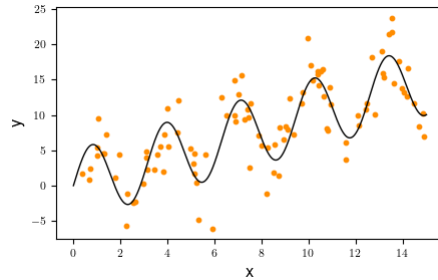


Figure 7: Synthetic problem for experiments.

Since this is a regression problem, $\ell(y, f) = (y - f)^2$, and $\hat{f} := \mathbb{E}_D[\hat{f}]$.

The hypothesis class \mathcal{F} is defined as the set of all trained models obtained over $T = 1000$ independently sampled datasets, each of size $n = 100$. The best-in-class model is the minimum across the T trials:

$$f^* := \arg \min_D \hat{R}(\hat{f}_D) \quad (40)$$

where the risk \hat{R} is approximated by sample of uniformly sampled points at a resolution of 0.001, giving a total of $n = 15,000$ test points. To simplify analysis, we assume $\hat{f} = \hat{f}_{erm}$.

³Assuming the Gauss-Markov conditions hold.