# STAT 479: Machine Learning
# Lecture Notes

Sebastian Raschka
Department of Statistics
University of Wisconsin–Madison

http://stat.wisc.edu/~sraschka/teaching/stat479-fs2019/

Fall 2019

## Contents

# STAT 479: Machine Learning
# Lecture Notes

Sebastian Raschka
Department of Statistics
University of Wisconsin–Madison

http://stat.wisc.edu/~sraschka/teaching/stat479-fs2019/

Fall 2019

## 8 Model Evaluation 1: Overfitting and Underfitting

### 8.1 Overview

- In this lecture, we discuss some of the basic terms and machine learning fundamentals that are relevant for model evaluation, namely, *bias and variance*, and *overfitting and underfitting*.
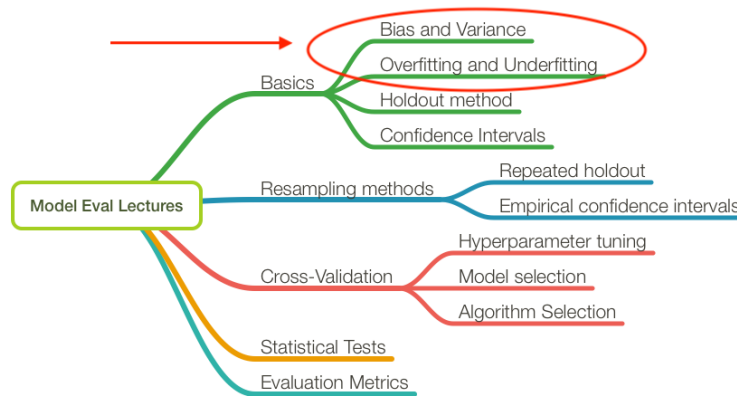


**Figure 1:** Overview of topics being covered in this lecture in the context of topics related to model evaluation that we will cover at a later point in time.

### 8.2 Overfitting and Underfitting

- The overall goal in machine learning is to obtain a model/hypothesis that generalizes well to new, unseen data.

- In other words, we want a model that *generalizes* well to unseen data, which we can measure, for example, by using an independent test set – while it sounds like this should be very straightforward, there are some pitfalls which we will discuss in the next lecture.

- Some of the evaluation metrics we can use to measure the performance on the test set are the prediction accuracy and misclassification error in the context of classification models – we say that a good model has a "high generalization accuracy" or "low generalization error" (or, simply "good generalization performance").

- The assumptions we generally make are the following:

  - i.i.d. assumption: inputs are independent, and training and test examples are identically distributed (drawn from the same probability distribution).
  - For some random model that has not been fit to the training set, we expect both the training and test error to be equal.
  - The training error or accuracy provides an (optimistically) biased estimate of the generalization performance.

Now, *overfitting* and *underfitting* are two terms that we can use to diagnose a machine learning model based on the training and test set performance. I.e., a model that suffers from underfitting does *not* perform well on the test *and* training set. In contrast, a model that overfits (e.g., from fitting the noise in the training dataset) can be usually recognized by a high training set accuracy, but low test set accuracy. Intuitively, as a rule of thumb, the larger the hypothesis space a model has access to, the higher the risk of overfitting.

A more technical term for the size of the hypothesis space is the so-called *capacity*. There are different measures for specific models and datasets that can be used to calculate the capacity of the model such as the VC dimension[1][2] (however, topics from learning theory such as VC dimension are beyond the scope of this course).
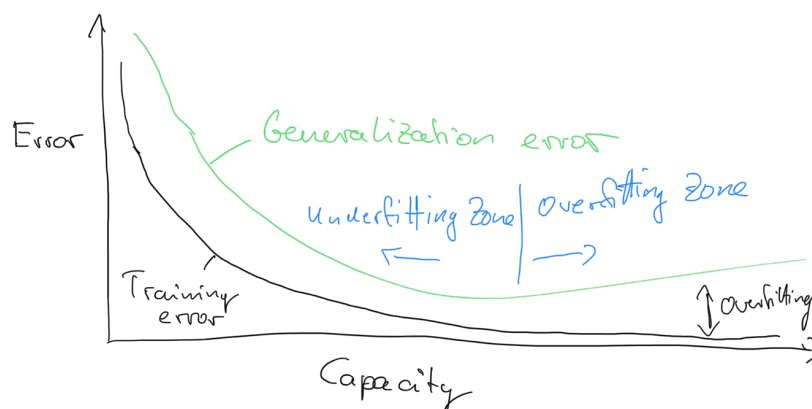


**Figure 2:** Illustration of overfitting and underfitting in relation to the training and test error.

---

[1]VC dimension stands for Vapnik-Chervonenkis dimension.
[2]Vladimir N Vapnik and A Ya Chervonenkis. "On the uniform convergence of relative frequencies of events to their probabilities". In: *Measures of complexity*. Springer, 2015, pp. 11–30.

## 8.3  Bias and Variance

Often, researchers use the terms *bias* and *variance* or "bias-variance tradeoff" to describe the performance of a model – i.e., you may stumble upon talks, books, or articles where people say that a model has a high variance or high bias. So, what does that mean? In general, we might say that "high variance" is proportional to overfitting, and "high bias" is proportional to underfitting. However, in this lecture, we are going to define these terms more precisely.
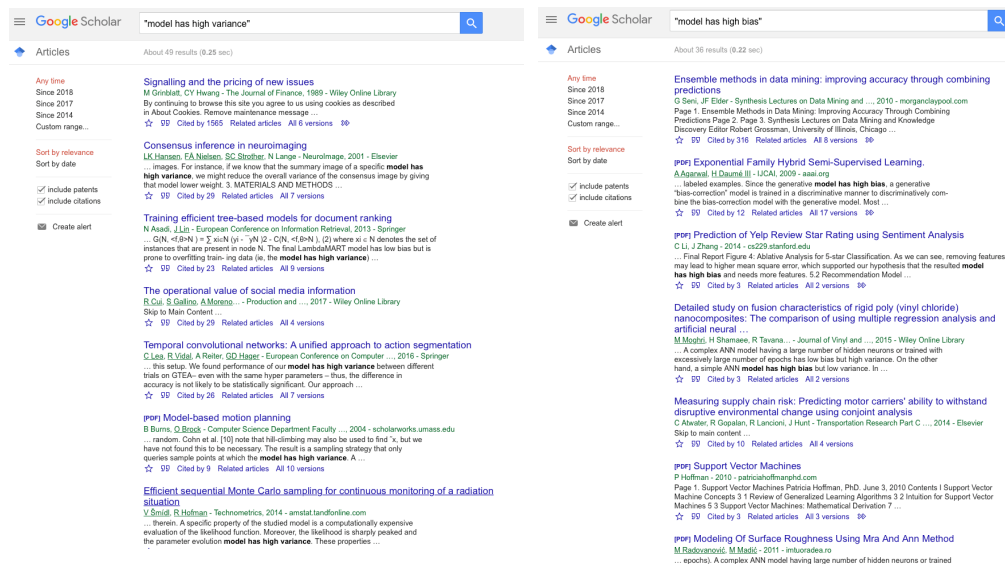


**Figure 3:** Results from searching the terms "model has high variance" and "model has high bias" on GoogleScholar

Note that the so-called Bias-Variance decomposition we are talking about in this lecture was initially formulated for regression losses (i.e., mean squared error[3]); however, we are also going to look into formulations for the 0-1 loss that we use to measure the misclassification error (or accuracy).

- **Why are we doing this?** The Decomposition of the loss into bias and variance help us understand learning algorithms, concepts are correlated to underfitting and overfitting.

- Thinking back of the ensemble lecture, the bias-variance decomposition and tradeoff help explain why ensemble methods might perform better than single models (i.e., why bagging reduces the variance, and why a boosting model has a lower bias than individual weak learners like decision tree stumps).

---

[3]In particular, in statistics, we evaluate the goodness of an estimator in relation to the true parameter or function
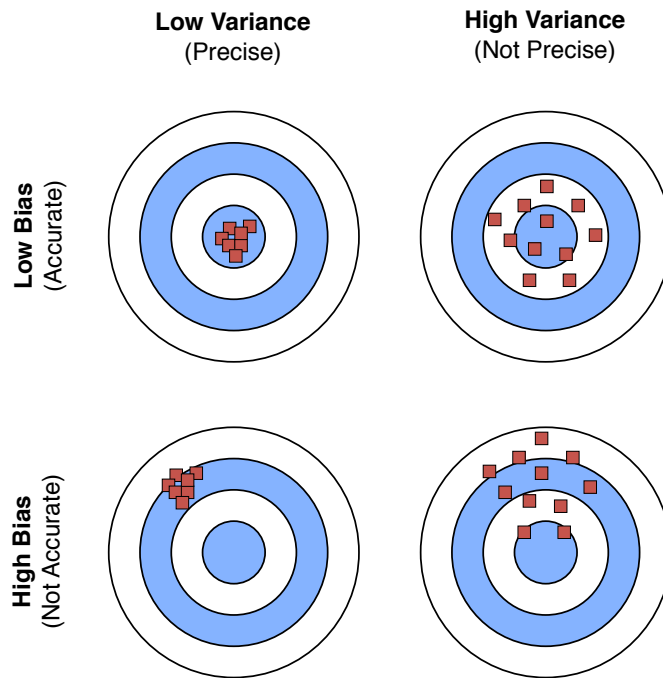
**Figure 4:** Bias-variance intuition.

To use the more formal terms for bias and variance, assume we have a point estimator $\hat{\theta}$ of some parameter or function $\theta$. Then, the bias is commonly defined as the difference between the expected value of the estimator and the parameter that we want to estimate:

$$\text{Bias} = E[\hat{\theta}] - \theta. \tag{1}$$

If the bias is larger than zero, we also say that the estimator is positively biased, if the bias is smaller than zero, the estimator is negatively biased, and if the bias is exactly zero, the estimator is unbiased. Similarly, we define the variance as the difference between the expected value of the squared estimator minus the squared expectation of the estimator:

$$\text{Var}(\hat{\theta}) = E\left[\hat{\theta}^2\right] - \left(E\left[\hat{\theta}\right]\right)^2. \tag{2}$$

Note that in the context of this lecture, it will be more convenient to write the variance in its alternative form:

$$\text{Var}(\hat{\theta}) = E[(E[\hat{\theta}] - \hat{\theta})^2]. \tag{3}$$
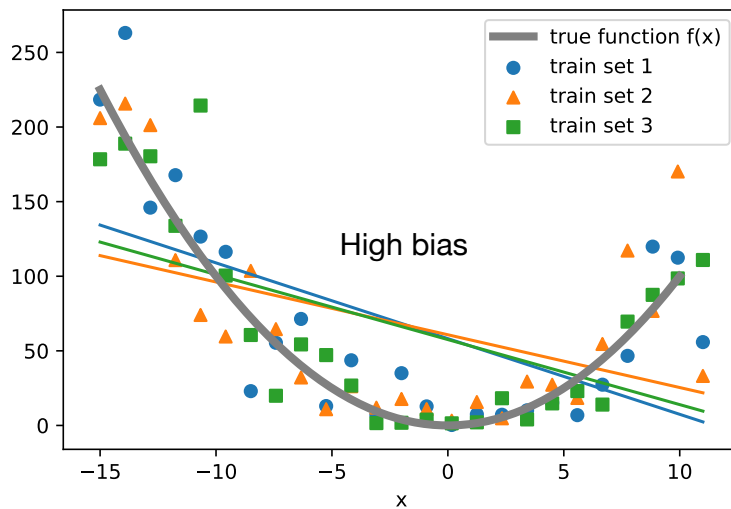
**Figure 5:** Suppose that there is an unknown target function or "true function" $f(x)$ which we want to approximate. Now, suppose we have different training sets drawn from an unknown distribution defined as "true function + noise." This plot shows different linear regression models, each fit to a different training set. None of these hypotheses approximate the true function well, except at two points (around x=-10 and x=6). In general, we can say that the bias is large because the difference between the true value and the predicted value, on average (here, average means "expectation of the training sets" not "expectation over examples in the training set"), is large (for most points).
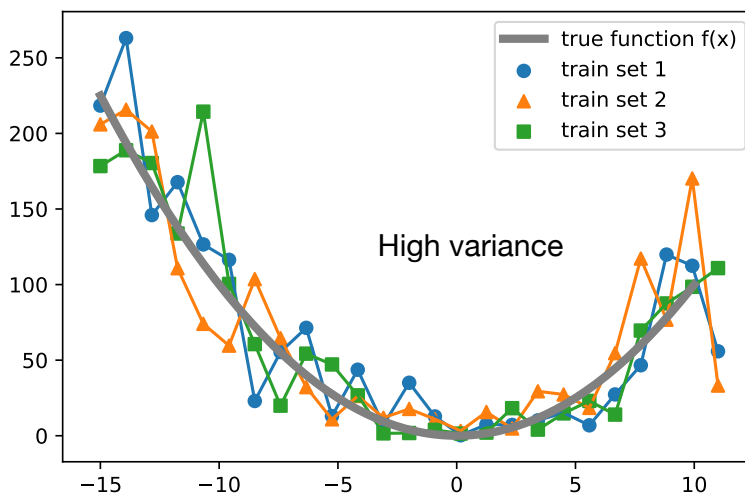


**Figure 6:** Suppose that there is an unknown target function or "true function" which we want to approximate. Now, suppose we have different training sets drawn from an unknown distribution defined as "true function + noise." This plot shows different unpruned decision tree models, each fit to a different training set. Note that these hypotheses fit the training data very closely. However, if we would consider the expectation over training sets, the average hypothesis would fit the true function perfectly (given that the noise is unbiased and has an expected value of 0). As we can see, the variance is very high, since on average, a prediction differs a lot from the expectation value of the prediction.

## 8.4   Bias-Variance Decomposition of the Squared Loss

We can decompose a loss function such as the squared loss into three terms, a variance, bias, and a noise term (and the same is true for the decomposition of the 0-1 loss later). However, for simplicity, we will ignore the noise term in this lecture (some of the literature referenced in later sections include the noise term if you are eager to learn more about variance-bias decomposition.).

Before we introduce the bias-variance decomposition of the 0-1 loss for classification, let us start with the decomposition of the squared loss as an easy warm-up exercise to get familiar with the overall concept.

The previous section already listed the common formal definitions of bias and variance, however, let us define them again for convenience:

$$\text{Bias}(\hat{\theta}) = E[\hat{\theta}] - \theta, \quad \text{Var}(\hat{\theta}) = E[(E[\hat{\theta}] - \hat{\theta})^2]. \tag{4}$$

Recall that in the context of these machine learning lecture (notes), we defined

- the true or target function as $y = f(x)$,
- the predicted target value as $\hat{y} = \hat{f}(x) = h(x)$,
- and the squared loss as $S = (y - \hat{y})^2$. (I use $S$ here because it will be easier to tell it apart from the $E$, which we use for the *expectation* in this lecture.)

**Note that unless noted otherwise, the expectation is over training sets!**

To get started with the squared error loss decomposition into bias and variance, let use do some algebraic manipulation, i.e., adding and subtracting the expected value of $\hat{y}$ and then expanding the expression using the quadratic formula $(a + b)^2 = a^2 + b^2 + 2ab$:

$$S = (y - \hat{y})^2 \tag{5}$$
$$(y - \hat{y})^2 = (y - E[\hat{y}] + E[\hat{y}] - \hat{y})^2 \tag{6}$$
$$= (y - E[\hat{y}])^2 + (E[\hat{y}] - \hat{y})^2 + 2(y - E[\hat{y}])(E[\hat{y}] - \hat{y}). \tag{7}$$

Next, we just use the expectation on both sides, and we are already done:

$$E[S] = E[(y - \hat{y})^2] \tag{8}$$
$$E[(y - \hat{y})^2] = (y - E[\hat{y}])^2 + E[(E[\hat{y}] - \hat{y})^2] \tag{9}$$
$$= [\text{Bias}]^2 + \text{Variance} \tag{10}$$

You may wonder what happened to the "$2ab$" term $(2(y - E[\hat{y}])(E[\hat{y}] - \hat{y}))$ when we used the expectation. It turns that it evaluates to zero and hence vanishes from the equation, which can be shown as follows:

$$E[2(y - E[\hat{y}])(E[\hat{y}] - \hat{y})] = 2E[(y - E[\hat{y}])(E[\hat{y}] - \hat{y})] \tag{11}$$
$$= 2(y - E[\hat{y}])E[(E[\hat{y}] - \hat{y})] \tag{12}$$
$$= 2(y - E[\hat{y}])(E[E[\hat{y}]] - E[\hat{y}]) \tag{13}$$
$$= 2(y - E[\hat{y}])(E[\hat{y}] - E[\hat{y}]) \tag{14}$$
$$= 0. \tag{15}$$

So, this is the canonical decomposition of the squared error loss into bias and variance. The next section will discuss some approaches that have been made to decompose the 0-1 loss that we commonly use for classification accuracy or error.
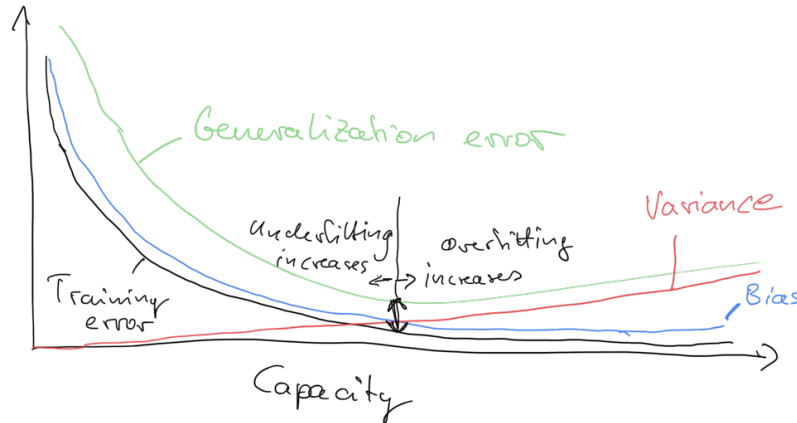


**Figure 7:** A sketch of variance and bias in relation to the training error and generalization error – how high variance related to overfitting, and how large bias relates to underfitting.

## 8.5   Bias-Variance Decomposition of the 0-1 Loss

Note that decomposing the 0-1 loss into bias and variance components is not as straight-forward as for the squared error loss. To quote Pedro Domingos, a well-known machine learning researcher and professor at University of Washington: "several authors have proposed bias-variance decompositions related to zero-one loss (Kong & Dietterich, 1995; Breiman, 1996b; Kohavi & Wolpert, 1996; Tibshirani, 1996; Friedman, 1997). However, each of these decompositions has significant shortcomings."[4]. In fact, the paper this quote was taken from may offer the most intuitive and general formulation at this point. However, we will first, for simplicity, go over Kong & Dietterich formulation[5] of the 0-1 loss decomposition, which is the same as Domingos's but excluding the noise term (for simplicity).

The table below summarizes the relevant terms we used for the squared loss in relation to the 0-1 loss. Recall that the 0-1 loss, $L$, is 0 if a class label is predicted correctly, and it is 1, otherwise. The main prediction for the squared error loss is simply the average over the predictions $E[\hat{y}]$ (the expectation is over training sets), for the 0-1 loss Kong & Dietterich and Domingos defined it as the mode. I.e., if a model predicts the label one more than 50% of the time (considering all possible training sets), then the main prediction is 1, and 0 otherwise.

|                         | Squared Loss          | 0-1 Loss           |
| ----------------------- | --------------------- | ------------------ |
| Single loss             | $(y - \hat{y})^2$     | $L(y, \hat{y})$    |
| Expected loss           | $E[(y - \hat{y})^2]$  | $E[L(y, \hat{y})]$ |
| Main prediction $E[\hat{y}]$ | mean (average)   | mode               |
| Bias$^2$                | $(y - E[\hat{y}])^2$  | $L(y, E[\hat{y}])$ |
| Variance                | $E[(E[\hat{y}] - \hat{y})^2]$ | $E[L(\hat{y}, E[\hat{y}])]$ |

---

[4]Pedro Domingos. "A unified bias-variance decomposition". In: *Proceedings of 17th International Conference on Machine Learning*. 2000, pp. 231–238.

[5]Thomas G Dietterich and Eun Bae Kong. *Machine learning bias, statistical bias, and statistical variance of decision tree algorithms*. Tech. rep. Technical report, Department of Computer Science, Oregon State University, 1995.

Hence, as result from using the mode to define the main prediction of the 0-1 loss, the bias is 1 if the main prediction does not agree with the true label $y$, and 0 otherwise:

$$Bias = \begin{cases} 1 \text{ if } y \neq E[\hat{y}], \\ 0 \text{ otherwise.} \end{cases} \qquad (16)$$

The variance of the 0-1 loss is defined as the probability that the predicted label does not match the main prediction:

$$Variance = P(\hat{y} \neq E[\hat{y}]). \qquad (17)$$

Next, let us take a look at what happens to the loss if the bias is 0. Given the general definition of the loss, loss = bias + variance, if the bias is 0, then we define the loss as the variance:

$$Loss = 0 + Variance = Loss = P(\hat{y} \neq y) = Variance = P(\hat{y} \neq E[\hat{y}]). \qquad (18)$$

In other words, if a model has zero bias, it's loss is entirely defined by the variance, which is intuitive if we think of variance in the context of being proportional overfitting.

The more surprising scenario is if the bias is equal to 1. If the bias is equal to 1, as explained by Pedro Domingos, the increasing the variance can decrease the loss, which is an interesting observation. This can be seen by first rewriting the 0-1 loss function as

$$Loss = P(\hat{y} \neq y) = 1 - P(\hat{y} = y). \qquad (19)$$

(Note that we have not done anything new, yet.) Now, if we look at the previous equation of the bias, if the bias is 1, we have $y \neq E[\hat{y}]$. If $y$ is not equal to the main prediction and $y$ is is equal to $\hat{y}$, then $\hat{y}$ cannot be equal to the main prediction. Using the "inverse" ("1 minus"), we can then write the loss as

$$Loss = P(\hat{y} \neq y) = 1 - P(\hat{y} = y) = 1 - P(\hat{y} \neq E[\hat{y}]). \qquad (20)$$

Since the bias is 1, the loss is hence defined as "loss = bias - variance" if the bias is 1 (or "loss = 1 - variance"). This might be quite unintuitive at first, but the explanations Kong, Dietterich, and Domingos offer was that if a model has a very high bias such that it main prediction is always wrong, increasing the variance can be beneficial, since increasing the variance would push the decision boundary, which might lead to some correct predictions just by chance then. In other words, for scenarios with high bias, increasing the variance can improve (decrease) the loss!

## 8.6   Conclusion

In this lecture, we decomposed the squared error loss into variance and bias terms and discussed how these components relate to overfitting and underfitting. Then, we referred to a bias-variance decomposition that Kong & Dietterich defined for the 0-1 loss. Pedro Domingos later generalized this further, including the noise term, which we did not discuss in this lecture. However, interested students are encouraged to read the original paper:

- Domingos, P. (2000). A unified bias-variance decomposition. In *Proceedings of 17th International Conference on Machine Learning* (pp. 231-238).[6]

---

[6]https://homes.cs.washington.edu/~pedrod/bvd.pdf

Now, we should be more familiar with the terms bias and variance (or, as statistics students, consider this as a refresher) and how it relates to overfitting and underfitting if we say that a model has a high variance or high bias, respectively.

In the next lecture, we will take a closer look at the holdout method for model evaluation (and estimating the generalization performance). Also, we will discuss several methods for constructing confidence intervals.