

# **Stat 442 Lecture Notes**

**Spring Semester, 2025**

*Last Updated:* February 7, 2025

---

# 1 Wednesday 2/5: Welcome, what is statistical learning, basics of supervised learning

## 1.1 What is statistical learning?

Classical statistics is built on the scientific method. It assumes that we start with a hypothesis, and then we go out and we collect data that can be used to test this hypothesis. Importantly, the statistical methods that we will use to test this hypothesis are pre-specified; we choose them before we ever look at our data.

For example, maybe we come up with the hypothesis that college students who sleep more than 7 hours per night have higher average GPAs than those who sleep less than 7 hours per night. We then go out and take a random sample of students, and ask them all if they sleep more than 7 hours per night, and also ask for their GPA. We then use a t-test to test if there is a statistically significant difference in GPA between students who sleep more or less than 7 hours per night.

Statistical learning is different. While there is no single, satisfying definition of statistical learning, it refers generally to a collection of tools used to make sense of complex data. Unlike classical statistics, where a hypothesis is defined before data collection, statistical learning often involves discovering patterns in data without an initial hypothesis.

In the context of our motivating example, suppose that we collect a large survey of students at a college. We are generally interested in what factors impact the GPA of a college student. We search through the data, try different models, and end up realizing that “hours of sleep per night” seems to be an important factor in determining GPA: more important than, say, major, whether or not you are a varsity athlete, etc. However, we also realize that the relationship between hours of sleep per night and GPA does not look linear. We decide to turn “hours of sleep per night” into a binary variable: yes/no do you sleep more than 7 hours per night. We have ended up with the same “model” as before, but this was a much more exploratory process. It started with a general task of *prediction* (as opposed to a specific model or hypothesis), and used *variable selection* and *feature engineering* to settle on a final model. I would consider this to be an example of statistical learning.

In recent decades, as datasets have become larger and computers have become more powerful, there have been a proliferation of complex methods for analyzing data that are often labeled as machine learning methods. In general, however, there is no clear dividing line between statistics and machine learning. Linear regression and logistic regression, which you all studied in detail in Stat 346, are certainly examples of machine learning methods: they help us learn from data. In fact, they are both methods for *supervised learning*, which is one of two primary types of statistical learning.<sup>1</sup>

- In *supervised learning*, we start with a response variable of interest and a set of potential predictor variables (also known as explanatory variables). Our general goal is to use the predictor variables to explain or predict the response variable.
- In *unsupervised learning*, we just start with a collection of variables: there is no specific response variable. Our goal is to find structure or patterns in the data.

In this course, we will spend around 9 weeks on supervised learning, and only 2 weeks at the end on unsupervised learning. However, I do not want to give you the impression that unsupervised learning is less important! Unsupervised learning is extremely important, and is in fact the back-bone of recent generative AI technology. We start with supervised learning because it has a more natural connection to classical statistics and builds more naturally on your previous work. Hopefully, some of the concepts that you learn during our supervised learning unit will help you be equipped to delve more deeply into unsupervised learning in the future if you choose.

There is one more type of learning that has become really important in recent years, partly due to the proliferation of generative AI. This is the field of *semi-supervised learning*, in which we have a (potentially small) set of *labeled data* (data with predictors + response) and a (potentially much larger) set of *unlabeled data* (data that is missing the response). While we will not cover this in lecture, it would make a great *final project topic*. Throughout the semester, I will try to flag these topics as they come up!

---

<sup>1</sup>I would be remiss if I did not mention *reinforcement learning*, which is a very important topic in machine learning but does not relate as directly to drawing insights from a given dataset.

## 1.2 Supervised learning

### 1.2.1 Introduction

As mentioned, you all already know at least two methods for supervised learning: linear regression and logistic regression. In general, these represent one from each of two classes of supervised learning methods.

- Regression: we use this term for any method that is used to predict a quantitative response variable.
- Classification: we use this term for any method that is used to predict a categorical response variable.

As a reminder, a quantitative variable can be either continuous or discrete, but it must be a number. A categorical variable can be either ordinal (the categories have orders) or unordered. In this class, we will not cover any methods for considering ordinal response variables: if this topic interests you, it is another *potential final project topic*.

Linear regression and logistic regression are basic examples of supervised learning methods. In the past few decades, much more complex algorithms have been developed, and have exploded in popularity. There are a few reasons for this proliferation in complex methods.

- Datasets have gotten bigger. As we will learn in this class, complex models are often only appropriate when applied to enough data. Thus, it is natural that the “big data” era has gone hand-in-hand with the development of complex statistical learning methods.
- Computers have gotten better, faster, and cheaper, which has enabled the fitting of complex models.
- Free, open-source software programs like R and python have allowed researchers to develop algorithms and then easily share them with non-experts. This has allowed many methods to become popular.

Given the huge number of statistical learning methods, we could spend each class this semester learning a new method, such that at the end you are left with a huge cookbook of methods to choose from for your future data analysis needs. In fact, this class will be a little bit like this. But, there is a problem with this approach of treating statistical learning like a cookbook. New machine learning methods get published every day, and the state-of-the-art is constantly changing. If your goal is to learn a list of methods or algorithms, you will find yourself continually behind and out-of-date! Thus, the goal of this class is not really to teach you a list of methods. Instead, the goal is to teach you the fundamental concepts and overarching themes that go into the development of the methods, such that you can compare methods, recognize their limitations, and learn new methods on your own in the future. This goal helps explain two features of this course:

- We will spend a fair amount of time on older, less state-of-the-art methods. We will use these methods to illustrate general principles and themes, even if they may not be the methods that you will use in practice in your future.
- You will all be doing a *teaching presentation* this semester. You will be responsible for doing independent reading to learn about a method or a concept. You will need to synthesize what you learn, and give a 20 minute presentation to the class explaining the method. This skill mirrors the way you will all interact with machine learning in the future: you will be continually learning and synthesizing new methods.

### 1.2.2 Notation for datasets and random variables

The backbone of statistical learning is a matrix of predictors  $\mathbf{X} \in \mathbb{R}^{n \times p}$  and a vector of length  $n$  that stores the responses  $\mathbf{y}$ . We let  $n$  denote the number of *observations* (rows of our dataset) and let  $p$  be the number of *variables* (columns of our dataset, usually not including the response): this notation is quite standard in statistics, but isn't always particularly precise.

Your textbook (ISL) denotes pieces of the dataset  $\mathbf{X}$  and  $\mathbf{y}$  as follows. We observe  $x_i$  for  $i = 1, \dots, n$ , where  $x_i = (x_{i1}, \dots, x_{ip})^\top$  is a  $p$ -vector for each observation. We use  $\mathbf{x}_1, \dots, \mathbf{x}_p$  to denote the columns of  $\mathbf{X}$ ; each is an  $n$ -vector representing its own variable. Your textbook says that it will use bold lowercase letters anytime the vector has length  $n$  and use unbold lowercase letters otherwise: I am sure that I will start messing this up soon, but I will try to be consistent.

The other important distinction is between a random variable and its observed realization. We will use capital, unbold letters to denote random variables. In our example from Section 1.1, we might let  $X_1$  denote the number of hours that a randomly selected student sleeps per night and let  $Y$  denote their GPA. Once we observe a particular  $n$  students for our dataset, we let their sleep values be  $\mathbf{x}_1 = (x_{11}, x_{12}, \dots, x_{1n})^\top$  and we let their GPA values be  $\mathbf{y} = (y_1, \dots, y_n)^\top$ .

Consider a dataset that contains information on  $n$  pets. The response variable is the weight of the pet  $\mathbf{y}$ , and we have one single predictor variable: the type of pet. This is a categorical variable that takes on values {cat, dog, hamster, rabbit}. I could represent  $\mathbf{X}$  in this case as a matrix with one column; where the column stores the actual values {dog, dog,

cat, dog, hamster, cat, dog, ...}. But, as you all know from Stat 346, it is likely going to be convenient for fitting to instead let  $\mathbf{X}$  be a matrix with four columns, that store binary variables indicating “is this a dog”, “is this a cat”, etc. This simple example shows us how notation can get complicated: do we say that  $p = 4$  or  $p = 1$  for this dataset? Practically speaking, the  $p = 4$  is likely more relevant! But it can depend!

We can also get into problems when counting  $n$  sometimes. Suppose we take measurements on 920 students, who are spread out between 8 schools. The schools are randomly assigned to either have required yoga PE class, or to have traditional PE class. While we have 920 students, the students within a given school are not *independent* of one another: there may be factors other than the yoga class that are making their test scores more similar to one another. On the other hand, the schools underwent random assignment and can be safely assumed to be independent of one another. Thus, is the sample size  $n$  the 920 students or the 8 schools? It turns out that for this *clustered* or *hierarchical* data, we have a notion of *effective sample size* that is somewhere between 920 and 8, and depends on how correlated the students are within the schools.

We should also mention that data does not always come to you in a form where it looks like  $\mathbf{X} \in \mathbb{R}^{n \times p}$  and  $\mathbf{y} \in \mathbb{R}^n$ . Consider the task of image classification. You get a set of 500 images, which are stored as  $400 \times 400$  pixel images, where each pixel records a numerical value between 0 and 255 for red, green, and blue. Each image is associated with a label in  $\{\text{cat}, \text{dog}\}$ , and you observe  $y_1, \dots, y_{500}$ . In this case, the RGB values in the pixels in the images are your predictor values, but they don't come in a simple  $n \times p$  matrix. We need to reshape the data so that each image  $i$  has an associated vector  $x_i \in \mathbb{R}^{400 \times 400 \times 3}$  that stores all three color values for all  $400 \times 400$  pixels. It will also likely be convenient to code  $y_i$  such that a 1 denotes a dog and a 0 denotes a cat, such that  $y$  becomes numerical. When there are more than 2 classes, such as  $\{\text{cat}, \text{dog}, \text{fish}\}$ , we will need more than one column to store our response variables numerically.

Who knew that counting  $n$  and  $p$  could be so complicated! In this class, just think of it as the number of rows and columns (minus the response columns) of the data, and remember these subtleties in case they ever come up.

In our supervised learning unit, we will start with *regression*. Thus, for the rest of today, we assume that  $Y$  is numerical.

## 1.3 Regression

In a regression problem, we assume that our response variable  $Y$  is related to a set of predictors  $X = (X_1, \dots, X_p)$  via some model that can be written as

$$Y = f(X) + \epsilon. \quad (1)$$

The function  $f()$  is a fixed but unknown function that represents the systematic information that  $X$  provides about  $Y$ , and  $\epsilon$  is a random error term, which should satisfy the condition that  $E[\epsilon] = 0$  and that  $\epsilon$  is independent of  $X$ . These conditions ensure that  $f()$  is capturing all parts of  $Y$  that can be explained by  $X$ , and  $\epsilon$  captures the parts of  $Y$  that cannot be explained by  $X$ . The task in regression is to come up with a good approximation for the unknown function  $f()$ . More specifically, we want to find a function  $\hat{f}()$  such that,  $Y \approx \hat{f}(X)$ .

We talked really briefly about two methods that you might use to find a good  $\hat{f}$ , that sort of represent opposite ends of the spectrum in terms of how “wiggly” they are. We talked about:

- Linear regression, which you have all seen before.
- KNN, which many of you have not seen before. We will go into a lot more detail on KNN next class.

### 1.3.1 Two possible goals

There are two reasons why we might want to come up with a function  $\hat{f}()$  such that,  $Y \approx \hat{f}(X)$ .

1. We might want to be able to come up with predictions  $\hat{Y} = \hat{f}(X)$  for future realizations of data where  $Y$  itself is not observed.
2. We might want to understand which predictors in  $X$  are most associated with  $Y$  in order to draw a scientific conclusion, or at least take a step towards drawing a scientific conclusion.

When beginning a regression problem, it is important to clearly define which of these two reasons is more important to you. This will aid you in choosing the best strategy for approaching the problem!

As we will see throughout the semester, these two goals can sometimes be at odds with one another! If we only care about prediction, we do not need our function  $\hat{f}()$  to be interpretable: we can use an extremely complex model (known as a “black box”), and we will be happy as long as  $\hat{Y} \approx Y$ . On the other hand, if we want to draw scientific

conclusions, then we might wish to use a simple, interpretable model for  $\hat{f}()$  that lends itself to rigorous inference. We should also note that it is not always an either-or situation. Sometimes, we want to make good predictions and also understand which variables are contributing significantly to the predictions. Managing the tradeoff, or lack thereof, between these two goals will be one of our fundamental themes for the semester.

We will pick up with everything else next class!

## 2 Monday, Feb 10: The bias variance tradeoff, as illustrated with KNN and linear regression

Like last class, we will focus on a really simple regression setting. We observe a single response  $\mathbf{y} \in \mathbb{R}^n$  and a single predictor  $\mathbf{x} \in \mathbb{R}^n$ . We assume that our observations are i.i.d. realizations of random variables  $(X, Y)$ , where

$$Y = f(X) + \epsilon.$$

We assume that  $E[\epsilon] = 0$  and  $\epsilon \perp\!\!\!\perp X$ , but the function  $f()$  is unknown. Our goal is to find a function  $\hat{f}$  such that  $Y \approx \hat{f}(X)$ . In this unit, we will talk about different algorithms for finding such a function  $\hat{f}$ !

We will assume that we are looking for  $\hat{f}$  that makes the squared error loss:

$$\left(Y - \hat{f}(X)\right)^2$$

small on average. Today we are going to talk a lot about this goal. And we are going to talk about how well KNN and linear regression each achieve this goal.

For today, we are using squared error loss everywhere. But please remember that this is not the only way to measure model accuracy! The exact math of the bias-variance decomposition that we will discuss today holds only for squared error loss, but similar concepts apply for other loss functions.

Here is the agenda for today:

1. Training error vs. expected test error.
2. The bias-variance decomposition of the expected test error.
3. What is the ideal function for minimizing the expected test error?
4. How do KNN and Linear Regression each approximate this ideal function?
5. R Demo of KNN and Linear regression and the bias variance tradeoff.

We might move a little bit fast, but HW1 is going to give you a chance to practice all of these concepts.

### 2.1 Training error vs. test error

If we use our dataset  $\mathbf{y} \in \mathbb{R}^n$  and  $\mathbf{x} \in \mathbb{R}^n$  to *train* a function  $\hat{f}$ , then our *training error* is:

$$MSE_{\text{train}} = \frac{1}{n} \sum_{i=1}^n \left(y_i - \hat{f}(x_i)\right)^2.$$

For a fixed function  $\hat{f}$  that we already trained, if we observe some test points  $x_i^{\text{test}}, y_i^{\text{test}}$  for  $i = 1, \dots, n_{\text{test}}$ , we can compute our test error as

$$MSE_{\text{test}} = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \left(y_i^{\text{test}} - \hat{f}(x_i^{\text{test}})\right)^2,$$

where the important thing is that the points  $x_i^{\text{test}}, y_i^{\text{test}}$  were not used to compute the function  $\hat{f}$ .

While in practice we usually do compute  $MSE_{\text{test}}$  over some fixed test set that has size  $n_{\text{test}}$ , we don't want our notion of model performance to be specific to the test set that we happened to see. We are likely studying  $MSE_{\text{test}}$  to estimate how big our error will be on a random new datapoint. In this case, we might be computing this for the particular function  $\hat{f}$  that we already computed. In this case, our expected prediction error conditional on our particular function  $\hat{f}$  that we already computed is:

$$E_{X^{\text{test}}, Y^{\text{test}}} \left[ \left(Y^{\text{test}} - \hat{f}(X^{\text{test}})\right)^2 \right],$$

where the function  $\hat{f}$  is not treated as random.

If we want to evaluate the potential of an algorithm or a procedure to make good predictions, then we want to take into account the randomness in our training set. We want to acknowledge that we could have seen a different training set that would have given us a different function  $\hat{f}$ . With this in mind, in the next section we will define a more complex notion of expected prediction error. This more complex notion will let us come up with a very beautiful decomposition!

## 2.2 The bias-variance tradeoff

Suppose that we train our model  $\hat{f}$  on a dataset  $\mathbf{X}^{\text{train}}, \mathbf{y}^{\text{train}}$ . Since this dataset is random, the function  $\hat{f}()$  is also random. Now, suppose that we would like to know about the expected prediction error for a new datapoint with  $X = x^{\text{test}}$  and unknown  $Y = y^{\text{test}}$ . The prediction error itself is  $(y^{\text{test}} - \hat{f}(x^{\text{test}}))^2$ , but to find the *expected* prediction error we need to take the expected value over multiple sources of randomness.

For our purposes, let's treat  $x^{\text{test}}$  as fixed. We want to find the average prediction error that we would obtain if we repeatedly (1) sampled training sets of size  $n$  from the population, (2) refit the function  $\hat{f}$  using this training data, and (3) evaluated the squared prediction error for a datapoint drawn with  $X = x^{\text{test}}$ . So, we want to evaluate:

$$E_{\mathbf{X}^{\text{train}}, \mathbf{y}^{\text{train}}, Y | X = x^{\text{test}}} \left[ (Y - \hat{f}(X))^2 \mid X = x^{\text{test}} \right] = E_{\mathbf{X}^{\text{train}}, \mathbf{y}^{\text{train}}, Y | X = x^{\text{test}}} \left[ (Y - \hat{f}(x^{\text{test}}))^2 \mid X = x^{\text{test}} \right].$$

For brevity, I am not going to keep writing the subscripts in the expected values. I am going to assume that we are treating  $x^{\text{test}}$  as fixed but taking the expected value over all other randomness. But remember that it is always good to know what exactly you are taking the expected value over!

To break this quantity down, we are going to do a clever trick that involves adding and subtracting 0 twice. This is a very common proof technique! After we do this, we expand our big squared quantity by thinking of it as  $(a + b + c)^2$ , and we apply linearity of expectation to put each term in its own expected value.

$$\begin{aligned} E \left[ (Y - \hat{f}(x^{\text{test}}))^2 \right] &= E \left[ (Y - \textcolor{red}{f}(x^{\text{test}}) + \textcolor{red}{f}(x^{\text{test}}) - E[\hat{f}(x^{\text{test}})] + E[\hat{f}(x^{\text{test}})] - \hat{f}(x^{\text{test}}))^2 \right] \\ &= E \left[ (Y - f(x^{\text{test}}))^2 \right] + E \left[ (f(x^{\text{test}}) - E[\hat{f}(x^{\text{test}})])^2 \right] + E \left[ (E[\hat{f}(x^{\text{test}})] - \hat{f}(x^{\text{test}}))^2 \right] \\ &\quad + 2E \left[ (Y - f(x^{\text{test}})) (f(x^{\text{test}}) - E[\hat{f}(x^{\text{test}})]) \right] \\ &\quad + 2E \left[ (Y - f(x^{\text{test}})) (E[\hat{f}(x^{\text{test}})] - \hat{f}(x^{\text{test}})) \right] \\ &\quad + 2E \left[ (f(x^{\text{test}}) - E[\hat{f}(x^{\text{test}})]) (E[\hat{f}(x^{\text{test}})] - \hat{f}(x^{\text{test}})) \right]. \end{aligned}$$

We will now argue that each of the cross terms is 0. We start with the orange term, and note that:

$$\begin{aligned} E \left[ (Y - f(x^{\text{test}})) (f(x^{\text{test}}) - E[\hat{f}(x^{\text{test}})]) \right] &= (f(x^{\text{test}}) - E[\hat{f}(x^{\text{test}})]) E[Y - f(x^{\text{test}})] \\ &= (f(x^{\text{test}}) - E[\hat{f}(x^{\text{test}})]) (E[Y] - f(x^{\text{test}})) \\ &= (f(x^{\text{test}}) - E[\hat{f}(x^{\text{test}})]) (f(x^{\text{test}}) - f(x^{\text{test}})) = 0. \end{aligned}$$

The first two equalities follow because the quantities  $f(x^{\text{test}})$  and  $E[\hat{f}(x^{\text{test}})]$  are not random, since  $x^{\text{test}}$  is a constant. The last equality follows since, because we are conditioning on  $X = x^{\text{test}}$ ,  $Y = f(x^{\text{test}}) + \epsilon$ , and so by our assumption that  $E[\epsilon] = 0$ ,  $E[Y] = f(x^{\text{test}})$ .

We now consider the purple term. We note that  $Y - f(x^{\text{test}})$  is independent of the training dataset, since  $f()$  is non-random and  $Y$  is a new realization. On the other hand,  $(E[\hat{f}(x^{\text{test}})] - \hat{f}(x^{\text{test}}))$  is a function only of the training data, since  $x^{\text{test}}$  is non-random. Thus, this term has the form  $E[ab]$  where  $a$  and  $b$  are independent. So we can write

it as:

$$\begin{aligned} E \left[ (Y - f(x^{\text{test}})) (E[\hat{f}(x^{\text{test}})] - \hat{f}(x^{\text{test}})) \right] &= E[Y - f(x^{\text{test}})] E \left[ (E[\hat{f}(x^{\text{test}})] - \hat{f}(x^{\text{test}})) \right] \\ &= E[Y - f(x^{\text{test}})] (E[\hat{f}(x^{\text{test}})] - E[\hat{f}(x^{\text{test}})]) = 0. \end{aligned}$$

Finally, we consider the green term.

$$\begin{aligned} E \left[ (f(x^{\text{test}}) - E[\hat{f}(x^{\text{test}})]) (E[\hat{f}(x^{\text{test}})] - \hat{f}(x^{\text{test}})) \right] &= (f(x^{\text{test}}) - E[\hat{f}(x^{\text{test}})]) E \left[ (E[\hat{f}(x^{\text{test}})] - \hat{f}(x^{\text{test}})) \right] \\ &= (f(x^{\text{test}}) - E[\hat{f}(x^{\text{test}})]) (E[\hat{f}(x^{\text{test}})] - E[\hat{f}(x^{\text{test}})]) = 0. \end{aligned}$$

Now that we know that all three cross terms are 0, we know that

$$E \left[ (Y - \hat{f}(x^{\text{test}}))^2 \right] = E \left[ (Y - f(x^{\text{test}}))^2 \right] + E \left[ (f(x^{\text{test}}) - E[\hat{f}(x^{\text{test}})])^2 \right] + E \left[ (E[\hat{f}(x^{\text{test}})] - \hat{f}(x^{\text{test}}))^2 \right].$$

Our next goal is to understand each of these terms!

1. The first term can be written as:

$$E \left[ (Y - f(x^{\text{test}}))^2 \right] = E \left[ (f(x^{\text{test}}) + \epsilon - f(x^{\text{test}}))^2 \right] = E[\epsilon^2] = \text{Var}(\epsilon).$$

This is our irreducible error term! It comes from the inherent noise in our data that cannot be explained by  $X$ .

2. In the second term,  $f(x^{\text{test}})$  and  $E[\hat{f}(x^{\text{test}})]$  are not random, so we do not need the expected value! This just becomes

$$(f(x^{\text{test}}) - E[\hat{f}(x^{\text{test}})])^2,$$

which we call the squared bias. Is  $\hat{f}(x^{\text{test}})$  equal to  $f(x^{\text{test}})$  on average? If so, then  $\hat{f}$  is unbiased and this term will disappear.

3. Finally, in the third term:

$$E \left[ (E[\hat{f}(x^{\text{test}})] - \hat{f}(x^{\text{test}}))^2 \right] = \text{Var} \left[ (E[\hat{f}(x^{\text{test}})] - \hat{f}(x^{\text{test}})) \right] + E \left[ (E[\hat{f}(x^{\text{test}})] - \hat{f}(x^{\text{test}}))^2 \right].$$

First, note that  $\text{Var} \left[ (E[\hat{f}(x^{\text{test}})] - \hat{f}(x^{\text{test}})) \right] = \text{Var}[\hat{f}(x^{\text{test}})]$ , because  $E[\hat{f}(x^{\text{test}})]$  is not random. Then, note that  $E \left[ (E[\hat{f}(x^{\text{test}})] - \hat{f}(x^{\text{test}})) \right] = E[\hat{f}(x^{\text{test}})] - E[\hat{f}(x^{\text{test}})] = 0$ . So, this third term is simply the variance (with respect to the training set) of  $\hat{f}(x^{\text{test}})$ .

This was a lot of work, but we can finally say that, at a given point  $x^{\text{test}}$ ,

$$E \left[ (Y - \hat{f}(x^{\text{test}}))^2 \right] = \text{Var}(\epsilon) + \text{Bias}(\hat{f}(x^{\text{test}}))^2 + \text{Var}(\hat{f}(x^{\text{test}})).$$

This is the bias variance decomposition, which will be important throughout the semester!

The  $\text{Var}(\epsilon)$  term is our irreducible error. No matter how good we are at estimating  $f$ , this is just the amount of noise in our data, so we will always have this error. On the other hand, we can reduce the bias and the variance terms if we pick a better statistical learning algorithm for estimating  $f$ .

Without giving too much away: very simple models tend to have high bias but low variance. Very complex (wiggly) models have low bias (they are never constrained!) but very high variance (they overfit to the training data). To minimize our expected prediction error, we are always looking for functions that hit the sweet spot of complexity!

## 2.3 The ideal function $\hat{f}$

Let's briefly forget the training data! If our goal was to minimize

$$E_{X,Y} \left[ (Y - \hat{f}(X))^2 \right]$$

and we had all of the resources in the world, what would we choose for  $\hat{f}$ ?

Well, under the law of total expectation, we can write this as:

$$E_X \left[ E_{Y|X} \left[ (Y - \hat{f}(X))^2 \mid X \right] \right].$$

In the inner expected value,  $X$  is no longer random, and solving for the point-wise minimum is easy. We have that:

$$\operatorname{argmin}_c E_{Y|X=x} [(Y - c)^2 \mid X = x] = E[Y \mid X = x].$$

This comes from the fact that an expected squared error is always minimized at a mean. You will prove this on HW1, and you have likely seen it before.

If  $E[Y \mid X = x]$  minimizes the point-wise expected prediction error at every  $x$ , then the function  $\hat{f}(x) = E[Y \mid X = x]$  is the function that minimizes the overall expected prediction error. Of course, we do not know the joint distribution of  $X$  and  $Y$ , and so we cannot simply set  $\hat{f}(x)$  to be this conditional expectation. But we can develop methods that attempt to approximate  $\hat{f}(x) = E[Y \mid X = x]$  under various sets of assumptions! As we will see today, both KNN and linear regression try to approximate  $E[Y \mid X]$ .

## 2.4 How close do KNN and Linear Regression get to this ideal function?

Let's review a really simple case. We have a single numerical response variable  $Y$  and a single numerical predictor variable  $X$ . We observe 100 datapoints, so  $n = 100$  and  $p = 1$ . Our goal is to use our dataset  $\mathbf{x} = (x_1, \dots, x_n)$ ,  $\mathbf{y} = (y_1, \dots, y_n)$  to come up with a function  $\hat{f}$  such that, on new, unseen data,  $Y \approx \hat{f}(X)$ . Today, we will consider two methods for coming up with  $\hat{f}$ .

### 2.4.1 Simple linear regression

We restrict our attention to  $\hat{f}$  that have the form  $b_0 + b_1x$  for  $b_0, b_1 \in \mathbb{R}$ . This makes our problem of finding the best  $\hat{f}$  easier, because we limited the complexity of the model class that we are considering.

Based on our training data, we let:

$$\hat{\beta}_0, \hat{\beta}_1 = \operatorname{argmin}_{b_0, b_1} \sum_{i=1}^n (y_i - b_0 - b_1x_i)^2. \quad (2)$$

We then let  $\hat{f}(X) = \hat{\beta}_0 + \hat{\beta}_1X$ . Under the assumption that  $E[Y] = \beta_0 + \beta_1X$  for some true, unknown  $\beta_0$  and  $\beta_1$ , this solution directly approximates the best possible function  $E[Y \mid X]$ .

However, if this assumption does *not* hold, then the linear model is too limiting. The result of limiting our model class so much is that we might have a lot of bias. If our linearity assumption does not hold, then no matter how much data we have  $E[\hat{f}(x)]$  just cannot be that close to  $f$ . This is why simple models can have a lot of bias— we did not give them enough complexity to be able to capture the true function!

You are all experts in linear regression! So we will not spend too much more space in the notes on it.

### 2.4.2 K-nearest neighbors

KNN is at the opposite end of the spectrum from linear regression. The premise is quite simple. KNN lets

$$\hat{f}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i, \quad (3)$$

where  $N_k(x)$  is a function that returns the  $k$  points in the training set that are closest to the input point  $x$  according to some distance metric (for us, Euclidean distance). So, avoiding equations: KNN makes predictions by finding the  $k$  training observations with  $x_i$  closest to  $x$ , and predicts that the response for  $x$  will be the average of the responses for those  $k$  points.

The hyper-parameter  $k$  is chosen by the user. When  $k = n$ , KNN just predicts  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$  for all observations, regardless of  $x$ . As  $k$  decreases, the functions returned by KNN get increasingly wiggly and flexible. When  $k = 1$ , the KNN function is a step function that can be arbitrarily wiggly, and that always has 0 training error.

Note that KNN directly approximates  $E[Y \mid X = x]$  *locally*, making no assumptions on the structure of  $E[Y \mid X]$ . When  $k$  is small, the estimation of  $E[Y \mid X = x]$  uses only data that are really close to  $x$ . This lets our prediction function get arbitrarily wiggly— we will not run into an issue of bias. However, with small  $k$ , there is a lot of variance



in each individual prediction, since it is made using very few datapoints. When  $k$  is large, there is less variance in the predictions, but we do not let our predictions be as wiggly, and so we introduce more bias. Thus, KNN is a great place to see the bias-variance tradeoff at work.

A few notes on KNN:

- KNN is non-parametric; we cannot write down the function  $\hat{f}$  without storing basically the entire dataset. The number of effective parameters grows with the sample size  $n$ .
- Linear regression in theory takes some time to train, but it is nearly instantaneous to make new predictions. KNN involves no training step, but it could be computationally expensive to obtain new predictions.

## 2.5 Time for an R demo!

The RMarkdown document that we go over in class will be posted on GLOW!