

# LECTURE 1: INTRODUCTION AND REGRESSION

Text references: Chapter 1 in Shalizi (James et al. Chapter 2)

---

## Introduction to Regression Analysis

Regression analysis is used to answer questions about how one variable depends on the level of one or more other variables. For example, does diet correlate with cholesterol level, and does this relationship depend on other factors, such as age, smoking status, and level of exercise?

Given data  $(X_1, Y_1), \dots, (X_n, Y_n)$  we have two goals:

---

---

---

---

---

---

---

Follow-up questions:

---

---

1

2

3

4

5

---

## 6 Regression and Prediction

7 Let's start with the simplest regression setting where we want to predict  
8 a real-valued variable  $Y \in \mathbb{R}$  from nothing at all. For example, suppose  
9 I'm taking a particular cholesterol medication, and I want to predict my  
10 improvement in blood cholesterol level.

11 **Q:** What is the optimal point prediction for  $Y$ ?

12

13

14

15

16

17

18

19

1 **Q:** In practice, how would we make this point prediction?

2 \_\_\_\_\_

3 \_\_\_\_\_

4 \_\_\_\_\_

5 \_\_\_\_\_

## 6 **The Regression Function**

7 Prediction with nothing to distinguish one random variable from the next  
8 has limited applicability. More often we want to predict  $Y$  based on an-  
9 other variable  $X$ , called the *predictor*, *covariate*, or *input*. In this context  $Y$  is  
10 now often called the *response*, *outcome*, or *output*.

11 Our prediction for  $Y$  will be a function of  $X$ , denoted  $r(X)$ . For example,  
12 what will be the amount of improvement in blood cholesterol level as a  
13 function of the dose  $X$  of a particular medication?

1 **Q:** What is the “best” prediction of  $Y$  given  $X$ ?

2 \_\_\_\_\_

3 \_\_\_\_\_

4 \_\_\_\_\_

5 \_\_\_\_\_

6 \_\_\_\_\_

7 We call  $r(x)$  the *regression function*, it's what we use to predict  $Y$  from  $X$ .

8 *Our statistical model:*

9 \_\_\_\_\_

10 \_\_\_\_\_

11 \_\_\_\_\_

12 \_\_\_\_\_

13 \_\_\_\_\_

14 \_\_\_\_\_

15 *Check:* What happens when we take expectation conditional on  $X$  on both

16 sides?

17 \_\_\_\_\_

18 \_\_\_\_\_

19 \_\_\_\_\_

20 \_\_\_\_\_

1 *Note (or disclaimer):*

2 \_\_\_\_\_

3 \_\_\_\_\_

4 \_\_\_\_\_

5 \_\_\_\_\_

6 \_\_\_\_\_

7 \_\_\_\_\_

## 8 **How to estimate the regression function?**

9 From data  $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ , we could estimate the regression  
10 function via

$$11 \quad \hat{r}(x) = \frac{1}{\#\{i : X_i = x\}} \sum_{i: X_i = x} Y_i \quad (1)$$

12 *But* this can work only when  $X$  takes discrete values; otherwise for  $X$  con-  
13 tinuous, at any  $x$ , the probability of getting a sample at any particular value  
14 of  $x$  is zero. Back to cholesterol example: what happens if we didn't ob-  
15 serve any person in our study who took exactly  $x = 31.5$  mg weekly of the  
16 drug? Important concepts: *interpolation, extrapolation, and smoothing...*

17 **Linear regression (revisited).** If we're willing to use a *linear* relationship,  
18 we can find the best fitting linear function, that is

19 \_\_\_\_\_

20 \_\_\_\_\_

1 Taking derivatives, we get

2 \_\_\_\_\_

3 \_\_\_\_\_

4 \_\_\_\_\_

5 \_\_\_\_\_

6 \_\_\_\_\_

7 **Q:** How do we estimate the linear oracle from data? (Mention two ap-  
8 proaches)

9 Method 1:

10 \_\_\_\_\_

11 \_\_\_\_\_

12 \_\_\_\_\_

13 \_\_\_\_\_

14 \_\_\_\_\_

15 \_\_\_\_\_

16 Method 2:

17 \_\_\_\_\_

18 \_\_\_\_\_

19 \_\_\_\_\_

20 \_\_\_\_\_

1 Let's look at this slightly differently: to make a prediction at an arbitrary  
2 point  $x$  (center  $X$  and  $Y$  to simplify the notation), we use

3 \_\_\_\_\_

4 \_\_\_\_\_

5 \_\_\_\_\_

6 \_\_\_\_\_

7 \_\_\_\_\_

## 8 **Linear Smoothers**

9 Define “linear smoothers” and go back to previous estimators.

10 \_\_\_\_\_

11 \_\_\_\_\_

12 \_\_\_\_\_

13 \_\_\_\_\_

14 \_\_\_\_\_

15 \_\_\_\_\_

16 \_\_\_\_\_

17 \_\_\_\_\_

1 Another common, more flexible linear smoother is the *k*-nearest-neighbors  
2 regression:

3 \_\_\_\_\_

4 \_\_\_\_\_

5 \_\_\_\_\_

6 \_\_\_\_\_

7 To use this in practice, we're going to need to choose *k*, the number of  
8 nearest neighbors.

9 **Q: What are the tradeoffs at either end?** [See Figures 1 and 2]

10 \_\_\_\_\_

11 \_\_\_\_\_

12 \_\_\_\_\_

13 \_\_\_\_\_

14 \_\_\_\_\_

15 \_\_\_\_\_

16 \_\_\_\_\_

17 \_\_\_\_\_

18 \_\_\_\_\_

19 \_\_\_\_\_



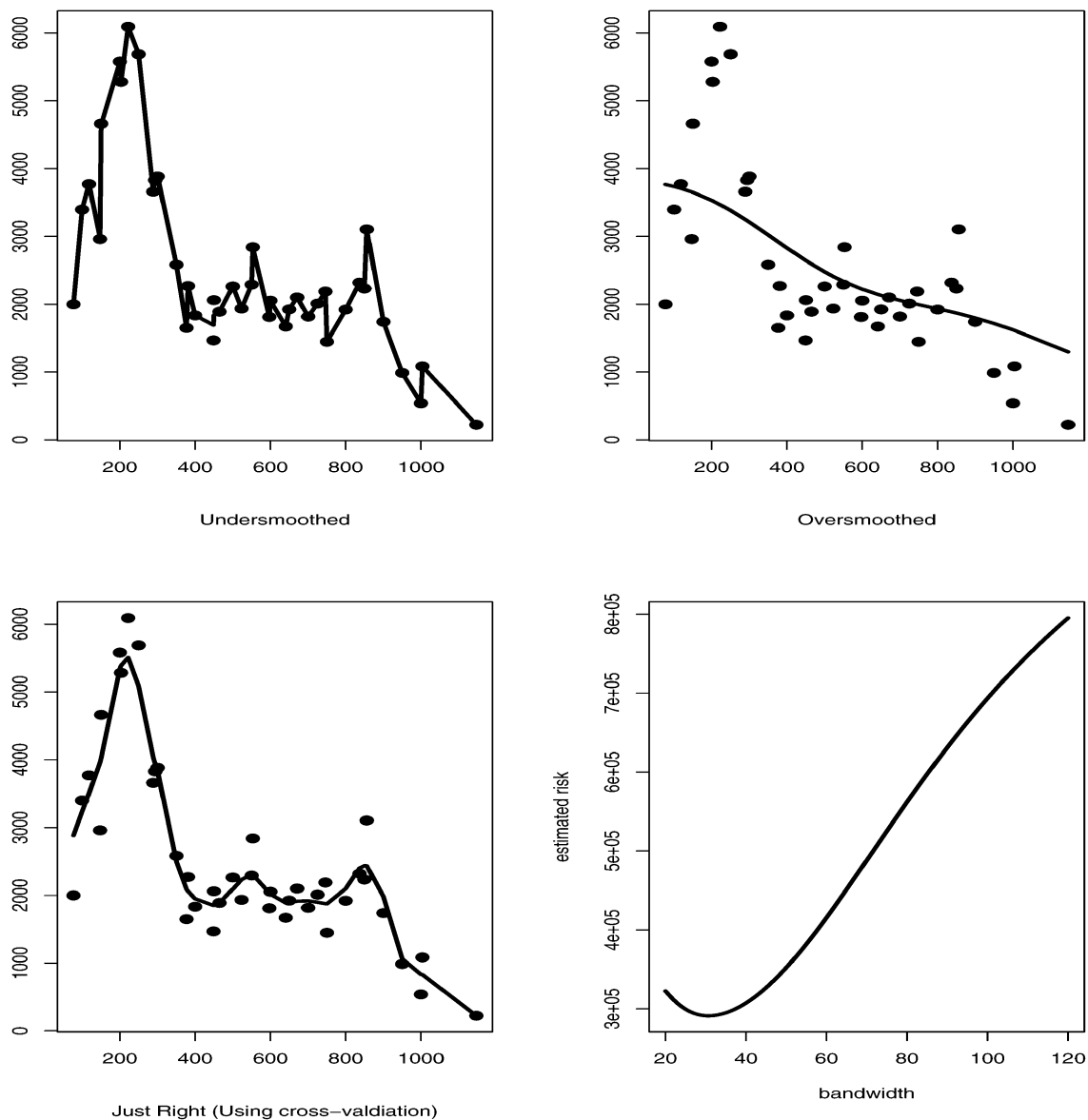
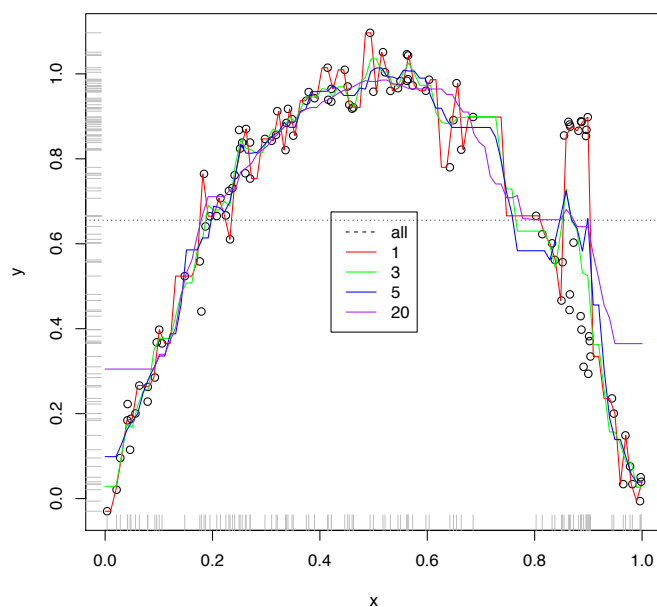


FIGURE 20.8. Regression analysis of the CMB data. The first fit is undersmoothed, the second is oversmoothed, and the third is based on cross-validation. The last panel shows the estimated risk versus the bandwidth of the smoother. The data are from BOOMERaNG, Maxima, and DASI.

Figure 1: From *All of Statistics* by Larry Wasserman

## 1.5. LINEAR SMOOTHERS

42



```
library(FNN)
plot.seq <- matrix(seq(from = 0, to = 1, length.out = 100), byrow = TRUE)
lines(plot.seq, knn.reg(train = all.x, test = plot.seq, y = all.y, k = 1)$pred,
      col = "red")
lines(plot.seq, knn.reg(train = all.x, test = plot.seq, y = all.y, k = 3)$pred,
      col = "green")
lines(plot.seq, knn.reg(train = all.x, test = plot.seq, y = all.y, k = 5)$pred,
      col = "blue")
lines(plot.seq, knn.reg(train = all.x, test = plot.seq, y = all.y, k = 20)$pred,
      col = "purple")
legend("center", legend = c("all", "1", "3", "5", "20"), lty = c("dashed", rep("solid",
  4)), col = c("black", "red", "green", "blue", "purple"))
```

FIGURE 1.5: Points from Figure 1.1 with horizontal dashed line at the mean and the  $k$ -nearest-neighbor regression curves for the indicated  $k$ . Increasing  $k$  smoothes out the regression line, pulling it towards the mean. — The code is repetitive; can you define functions to simplify it?

17:09 Monday 30<sup>th</sup> January, 2017

Figure 2: From Shalizi

## 1 Training and test errors

2 **Q:** How are we going to quantify this (overfitting vs underfitting)?

3 Let's call  $(X_1, Y_1), \dots (X_n, Y_n)$ , the sample of data that we used to fit  $\hat{r}$ , our  
4 **training sample**. What's wrong with looking at how well we do in fitting  
5 the training points themselves, i.e., the **training error** or *expected training*  
6 *error*, defined as

---

---

---

---

---

---

---

---

---

---

15 Now, suppose that we an independent **test sample**  $(X'_1, Y'_1), (X'_2, Y'_2), \dots (X'_m, Y'_m)$   
16 (following the same distribution as our training sample). We could then  
17 look at the **expected test error**, defined as

---

---

---

---

Note that the expectation here is taken over *all* that is random (both training and test samples). This really does capture what we want, and has the right behavior with  $k$ !

Exercise: Summarize how underfitting and overfitting relate to training and test errors. [See R Demo 1]

(Later in Lecture 3 we are going to see that underfitting and overfitting are related to two quantities called *estimation bias* and *estimation variance*, or just bias and variance.)

## **Smoother Linear Smoothers**

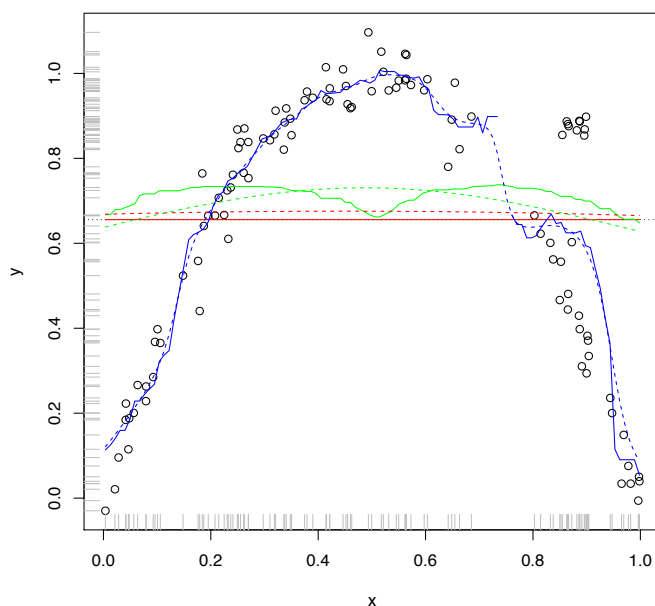
Let's look back at the  $k$ -nearest-neighbors regression, the weights  $w(x, x_i)$  are either  $1/k$  or 0, depending on  $x_i$ . This will often produce jagged looking

fits. Why? Are these weights smooth over  $x$ ?

How about choosing a smoother way of averaging? For example, a *kernel regression* with a *Gaussian kernel*:

The parameter  $h$  here is called the **bandwidth**, like  $k$  in  $k$ -nearest-neighbors regression, it controls the level of adaptivity/flexibility of the fit. One reason to prefer kernel regression with Gaussian kernels is that it can produce smoother (less jagged) looking fits than  $k$ -nearest-neighbors or kernel regression with box kernels [See Figure 3]

We are also going to see later in Lecture 4 that these smoother (nonparametric) estimators converge faster toward the true regression function.



```
lines(ksmooth(all.x, all.y, "box", bandwidth = 2), col = "red")
lines(ksmooth(all.x, all.y, "box", bandwidth = 1), col = "green")
lines(ksmooth(all.x, all.y, "box", bandwidth = 0.1), col = "blue")
lines(ksmooth(all.x, all.y, "normal", bandwidth = 2), col = "red", lty = "dashed")
lines(ksmooth(all.x, all.y, "normal", bandwidth = 1), col = "green", lty = "dashed")
lines(ksmooth(all.x, all.y, "normal", bandwidth = 0.1), col = "blue", lty = "dashed")
```

FIGURE 1.6: Data from Figure 1.1 together with kernel regression lines. Solid colored lines are box-kernel estimates, dashed colored lines Gaussian-kernel estimates. Red,  $h = 2$ ; green,  $h = 1$ ; blue,  $h = 0.1$  (per the definition of bandwidth in the `ksmooth` function). Note the abrupt jump around  $x = 0.75$  in the box-kernel/ $h = 0.1$  (solid green) line — with a small bandwidth the box kernel is unable to interpolate smoothly across the break in the training data, while the Gaussian kernel can.

17:09 Monday 30<sup>th</sup> January, 2017

Figure 3: From Shalizi

## 1 R Demo 1

### 2 (a) Load data.

3 Download the file [nonlin.Rdata](#) from Canvas, and load it into your R ses-  
4 sion with `load("nonlin.Rdata")`. You can type `ls()` to see the R ob-  
5 jects that have been loaded into memory.

6 The matrices `xtrain` and `ytrain`, are each  $100 \times 50$ , containing 50 training  
7 data sets of  $x$  and  $y$  points along its columns. That is, the first column of  
8 `xtrain` and the first column of `ytrain` make up a training data set of 100  
9  $x$ - $y$  pairs.

10 For the next bit, we will restrict our attention to just the first training set,  
11 i.e., the first columns of `xtrain` and `ytrain`. Plot the  $x$  points versus the  
12  $y$  points for these data to get an idea of the trend.

1 **(b) Linear regression on the first training set.**

2 Fit a linear regression on the first training set using the function `lm()`. Plot  
3 the estimated regression function on top of the training points.

4 **(c) k-nn regression on the first training set.**

5 Similarly to Part (b), now fit a k-nn regression to the training sample.  
6 Install the R library `FNN`. Use the function `knn.reg()` to fit a k-nn re-  
7 gression with 3 different values of the number of nearest neighborhoods:  
8  $k = 3, 15, 45$ . For each value of  $k$ , plot the estimated regression function  
9 from kernel regression on top of the training points.

10 **(d) Training and test errors as a function of  $k$  for the first data set.**

11 Examine training and test errors for the linear regression and for k-nn with  
12  $k = 1, \dots, 60$ . What happens as  $k$  decreases/increases? Why?

13 **(e) Compute training and test errors averaged over 50 data sets.**

14 Now repeat Part (d) for all 50 pairs of training and test sets. Average these  
15 results over the 50 data sets. Plot the results. Compare with the results  
16 from Part (d). What do you see? Why?