

Predicting Quantitative Features: Regression

36-350, Data Mining

6 October 2008

READING: sections 6.1–6.3 and 11.1 in *Principles of Data Mining*.

OPTIONAL READING: chapter 1 of Berk.

We’ve already looked at some examples of predictive modeling in the form of classification and factor analysis, but now we’ll get into it more seriously, and it will largely occupy us for the rest of the course. The place we’ll begin is with predictive quantitative features, i.e., regression. The math here is not so bad, and it connects more smoothly with your previous statistics courses; having learned some lessons here we’ll come back to classification, and then go to more complicated kinds of prediction.

1 Guessing the Value of a Random Variable

We have a quantitative, numerical feature, which we’ll imaginatively call Y . We’ll suppose that it’s a random variable, and try to predict it by guessing a single value for it. (Other kinds of predictions are possible — we might guess whether Y will fall within certain limits, or the probability that it does so, or even the whole probability distribution of Y . But some lessons we’ll learn here will apply to these other kinds of predictions as well.) What is the best value to guess? Or, more formally, what is the **optimal point forecast** for Y ?

To answer this question, we need to pick a function to be optimized, which should measure how good the guesses are — or equivalent how bad they are, how big an error is involved. A reasonable start point is the **mean squared error**:

$$\text{MSE}(a) \equiv \mathbf{E} \left[(Y - a)^2 \right] \quad (1)$$

So we’d like to find the value r where $\text{MSE}(a)$ is smallest.

$$\text{MSE}(a) = \mathbf{E} \left[(Y - a)^2 \right] \quad (2)$$

$$= (\mathbf{E}[Y - a])^2 + \text{Var}[Y - a] \quad (3)$$

$$= (\mathbf{E}[Y - a])^2 + \text{Var}[Y] \quad (4)$$

$$= (\mathbf{E}[Y] - a)^2 + \text{Var}[Y] \quad (5)$$

$$\frac{d\text{MSE}}{da} = 2(\mathbf{E}[Y] - a) + 0 \quad (6)$$

$$2(\mathbf{E}[Y] - r) = 0 \quad (7)$$

$$r = \mathbf{E}[Y] \quad (8)$$

So, if we gauge the quality of our prediction by mean-squared error, the best prediction to make is the expected value.

1.1 Estimating the Expected Value

Of course, to make the prediction $\mathbf{E}[Y]$ we would have to know the expected value of Y . Typically, we do not. However, if we have sampled values, y_1, y_2, \dots, y_n , we can estimate the expectation from the sample mean:

$$\hat{r} \equiv \frac{1}{n} \sum_{i=1}^n y_i \quad (9)$$

If the samples are IID, then the law of large numbers tells us that

$$\hat{r} \rightarrow \mathbf{E}[Y] = r \quad (10)$$

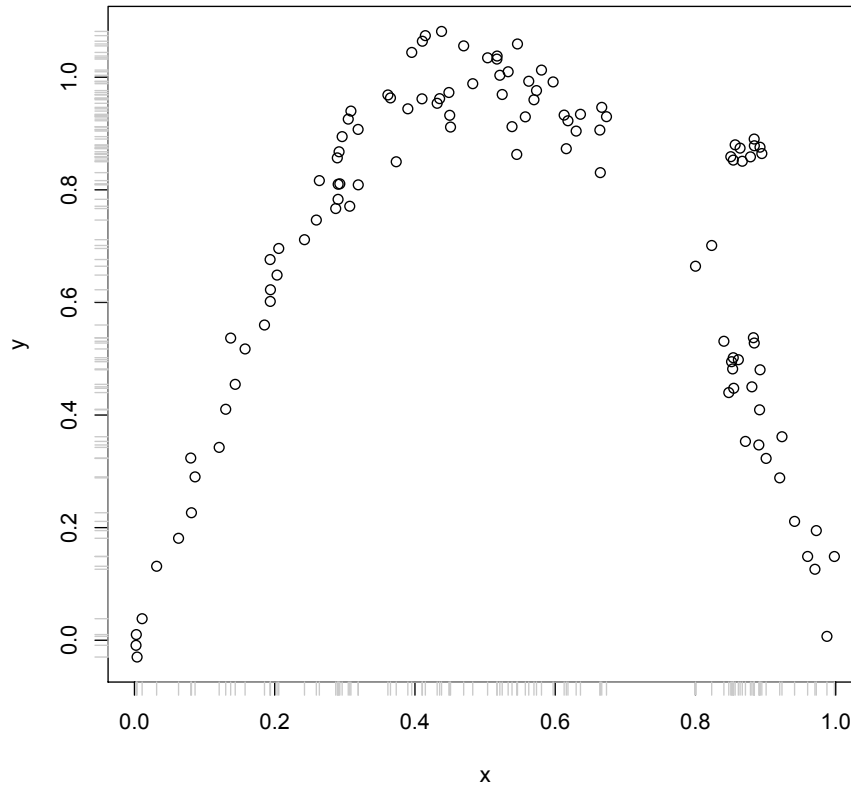
and the central limit theorem tells us something about how fast the convergence is (namely the squared error will typically be about $\text{Var}[Y]/n$).

Of course the assumption that the y_i come from IID samples is a strong one, but we can assert pretty much the same thing if they're just uncorrelated with a common expected value. Even if they are correlated, but the correlations decay fast enough, all that changes is the rate of convergence. So “sit, wait, and average” is a pretty reliable way of estimating the expectation value.

2 The Regression Function

Of course, it's not very useful to predict *just one number* for a feature. Typically, we have lots of features in our data, and we believe that there is *some* relationship between them. For example, suppose that we have data on two features, X and Y , which might look like Figure 1. The feature Y is what we are trying to predict, a.k.a. the **dependent variable** or **output** or **response**, and X is the **predictor** or **independent variable** or **covariate** or **input**. Y might be something like the profitability of a customer and X their credit rating, or, if you want a less mercenary example, Y could be some measure of improvement in blood cholesterol and X the dose taken of a drug. Typically we won't have just one input feature X but rather many of them, but that gets harder to draw and doesn't change the points of principle.

Figure 2 shows the same data as Figure 1, only with the sample mean added on. This clearly tells us something about the data, but also it seems like we should be able to do better — to reduce the average error — by using X , rather than by ignoring it.



```
plot(all.x,all.y,xlab="x",ylab="y")
axis(1,at=all.x,labels=FALSE,col="grey")
axis(2,at=all.y,labels=FALSE,col="grey")
```

Figure 1: Scatterplot of the example data. (These are made up.) The `axis` commands add horizontal and vertical ticks to the axes to mark the location of the data (in grey so they're less strong than the main tick-marks). This isn't necessary but is often helpful. The data are in the `example.dat` file.

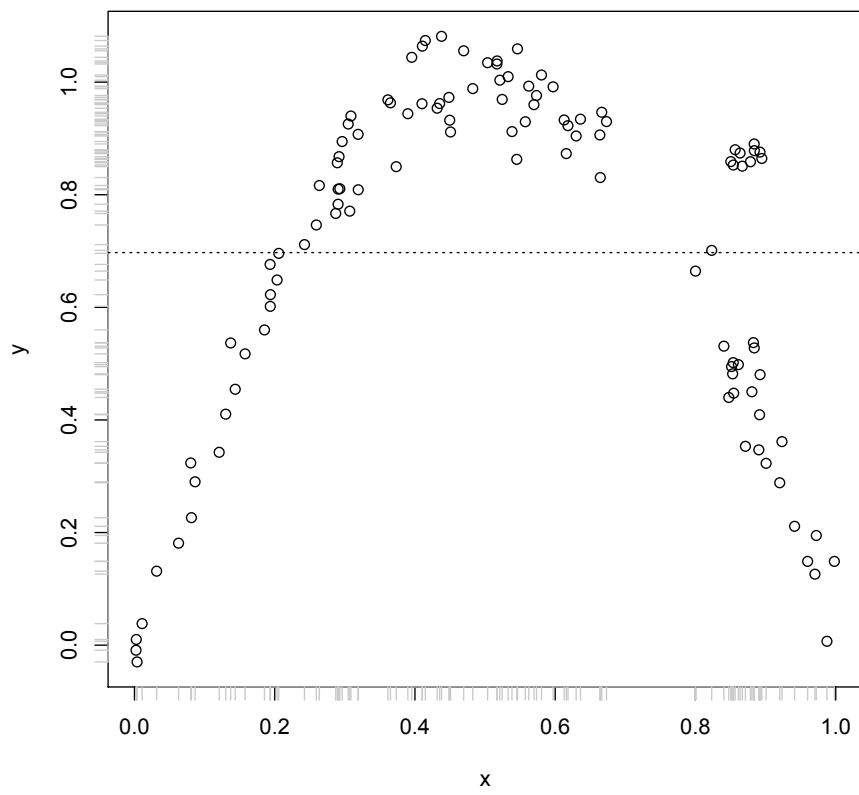


Figure 2: Data from Figure 1, with a horizontal line showing the sample mean of Y .

Let's say that we want our prediction to be a *function* of X , namely $f(X)$. What should that function be, if we still use mean squared error? We can work this out by using the law of total expectation, i.e., the fact that $\mathbf{E}[U] = \mathbf{E}[\mathbf{E}[U|V]]$ for any random variables U and V .

$$\text{MSE}(f(X)) = \mathbf{E}[(Y - f(X))^2] \tag{11}$$

$$= \mathbf{E}[\mathbf{E}[(Y - f(X))^2|X]] \tag{12}$$

$$= \mathbf{E}[\text{Var}[Y|X] + (\mathbf{E}[Y - f(X)|X])^2] \tag{13}$$

When we want to minimize this, the first term inside the expectation doesn't depend on our prediction, and the second term looks just like our previous optimization only with all expectations conditional on X , so for our optimal function $r(x)$ we get

$$r(x) = \mathbf{E}[Y|X = x] \tag{14}$$

In other words, the (mean-squared) optimal *conditional* prediction is just the conditional expected value. The function $r(x)$ is called the **regression function**. This is what we would like to know when we want to predict Y .

2.1 Some Disclaimers

It's important to be clear on what is and is not being assumed here. Talking about X as the “independent variable” and Y as the “dependent” one suggests a causal model, which we might write

$$Y \leftarrow r(X) + \epsilon \tag{15}$$

where the direction of the arrow, \leftarrow , indicates the flow from causes to effects, and ϵ is some noise variable. If the gods of inference are very, very kind, then ϵ would have a fixed distribution, independent of X , and we could without loss of generality take it to have mean zero. (“Without loss of generality” because if it has a non-zero mean, we can incorporate that into $r(X)$ as an additive constant.) This is the kind of thing we saw with the factor model. However, *no* such assumption is required to get Eq. 14. It works when predicting effects from causes, or the other way around when predicting (or “retrodicting”) causes from effects, or indeed when there is no causal relationship whatsoever between X and Y . It *is* always true that

$$Y|X = r(X) + \eta(X) \tag{16}$$

where $\eta(X)$ is a noise variable with mean zero, but as the notation indicates the distribution of the noise generally depends on X .

It's also important to be clear that when we find the regression function is a constant, $r(x) = r_0$ for all x , that this does not mean that X and Y are independent. If they are independent, then the regression function is a constant, but turning this around is the logical fallacy of “affirming the consequent”.¹

¹As in combining the fact that all human beings are featherless bipeds, and the observation

3 Estimating the Regression Function

We want to find the regression function $r(x) = \mathbf{E}[Y|X = x]$. Suppose that we have a big set of pairs $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$. Then this is a *supervised* learning problem — we know the label (value of Y) for each of our n training points. How can we estimate the regression function from these training examples?

If X takes on only a finite set of values, then a simple strategy is to use the conditional sample means:

$$\hat{r}(x) = \frac{1}{\#\{i : x_i = x\}} \sum_{i: x_i=x} y_i \quad (17)$$

By the same kind of law-of-large-numbers reasoning, we can be confident that $\hat{r}(x) \rightarrow \mathbf{E}[Y|X = x]$.

Unfortunately, this *only* works if X has only a finite set of values. If X is continuous, then in general the probability of our getting a sample at any *particular* value is zero, is the probability of getting *multiple* samples at *exactly* the same value of x . This is a basic issue with estimating any kind of function from data — the function will always be **undersampled**, and we need to fill in between the values we see. We also need to somehow take into account the fact that each y_i is a *sample* from the conditional distribution of $Y|X = x_i$, and so is not generally equal to $\mathbf{E}[Y|X = x_i]$. So any kind of function estimation is going to involve interpolation, extrapolation, and smoothing.

Different methods of estimating the regression function — different regression methods, for short — involve different choices about how we interpolate, extrapolate and smooth. This involves our making a choice about how to approximate $r(x)$ by a limited class of functions which we know (or at least hope) we can estimate. There is no guarantee that our choice leads to a *good* approximation in the case at hand, though it is sometimes possible to say that the approximation error will shrink as we get more and more data. This is an extremely important topic and deserves an extended discussion, coming next.

3.1 The Bias-Variance Tradeoff

Suppose that the true regression function is $r(x)$, but we use the function \hat{r} to make our predictions. Let's look at the mean squared error at $X = x$ in a slightly different way than before, which will make it clearer what happens when we can't use r to make predictions. We'll begin by expanding $(Y - \hat{r}(x))^2$, since the MSE at x is just the expectation of this.

$$\begin{aligned} (Y - \hat{r}(x))^2 & \quad (18) \\ &= (Y - r(x) + r(x) - \hat{r}(x))^2 \\ &= (Y - r(x))^2 + 2(Y - r(x))(r(x) - \hat{r}(x)) + (r(x) - \hat{r}(x))^2 \quad (19) \end{aligned}$$

that a cooked turkey is a featherless biped, to conclude that cooked turkeys are human beings. An econometrician stops there; an econometrician who wants to be famous writes a best-selling book about how this proves that Thanksgiving is really about cannibalism.

We saw above (Eq. 16) that $Y - r(x) = \eta$, a random variable which has expectation zero (and is uncorrelated with x). When we take the expectation of Eq. 19, nothing happens to the last term (since it doesn't involve any random quantities); the middle term goes to zero (because $\mathbf{E}[Y - r(x)] = \mathbf{E}[\eta] = 0$), and the first term becomes the variance of η . This depends on x , in general, so let's call it σ_x^2 . We have

$$\text{MSE}(\widehat{r}(x)) = \sigma_x^2 + ((r(x) - \widehat{r}(x))^2) \quad (20)$$

The σ_x^2 term doesn't depend on our prediction function, just on how hard it is, intrinsically, to predict Y at $X = x$. The second term, though, is the extra error we get from not knowing r . (Unsurprisingly, not knowing r cannot *improve* our predictions.) This is our first **bias-variance decomposition**: the total MSE at x is decomposed into a (squared) bias $r(x) - \widehat{r}(x)$, the amount by which our predictions are *systematically* off, and a variance σ_x^2 , the unpredictable, "statistical" fluctuation around even the best prediction.

All of the above assumes that \widehat{r} is a single fixed function. In practice, of course, \widehat{r} is something we estimate from earlier data. But if those data are random, the exact regression function we get is random too; let's call this random function \widehat{R}_n , where the subscript reminds us of the finite amount of data we used to estimate it. What we have analyzed is really $\text{MSE}(\widehat{R}_n(x) | \widehat{R}_n = \widehat{r})$, the mean squared error *conditional on* a particular estimated regression function. What can we say about the prediction error of the *method*, averaging over all the possible training data sets?

$$\text{MSE}(\widehat{R}_n(x)) = \mathbf{E} \left[(Y - \widehat{R}_n(X))^2 | X = x \right] \quad (21)$$

$$= \mathbf{E} \left[\mathbf{E} \left[(Y - \widehat{R}_n(X))^2 | X = x, \widehat{R}_n = \widehat{r} \right] | X = x \right] \quad (22)$$

$$= \mathbf{E} \left[\sigma_x^2 + (r(x) - \widehat{R}_n(x))^2 | X = x \right] \quad (23)$$

$$= \sigma_x^2 + \mathbf{E} \left[(r(x) - \widehat{R}_n(x))^2 | X = x \right] \quad (24)$$

$$= \sigma_x^2 + \mathbf{E} \left[(r(x) - \mathbf{E}[\widehat{R}_n(x)] + \mathbf{E}[\widehat{R}_n(x)] - \widehat{R}_n(x))^2 \right] \quad (25)$$

$$= \sigma_x^2 + (r(x) - \mathbf{E}[\widehat{R}_n(x)])^2 + \text{Var}[\widehat{R}_n(x)] \quad (26)$$

This is our second bias-variance decomposition — I pulled the same trick as before, adding and subtract a mean inside the square. The first term is just the variance of the process; we've seen that before and isn't, for the moment, of any concern. The second term is the bias in using \widehat{R}_n to estimate r — the **approximation bias** or **approximation error**. The third term, though, is the variance in our *estimate* of the regression function. Even if we have an unbiased *method* ($r(x) = \mathbf{E}[\widehat{R}_n(x)]$), if there is a lot of variance in our estimates, we can expect to make large errors.

The approximation bias has to depend on the true regression function. For example, if $\mathbf{E}[\widehat{R}_n(x)] = 42 + 37x$, the error of approximation will be zero if

$r(x) = 42 + 37x$, but it will be larger and x -dependent if $r(x) = 0$. However, there are flexible methods of estimation which will have small approximation biases for *all* r in a broad range of regression functions. The catch is that, at least past a certain point, decreasing the approximation bias can only come through increasing the estimation variance. This is the **bias-variance trade-off**. However, nothing says that the trade-off has to be one-for-one. Sometimes we can lower the total error by *introducing* some bias, since it gets rid of more variance than it adds approximation error. The next section gives an example.

In general, both the approximation bias and the estimation variance depend on n . A method is **consistent**² when both of these go to zero as $n \rightarrow 0$ — that is, if we recover the true regression function as we get more and more data.³ Again, consistency depends on how well the method matches the actual data-generating process, not just on the method, and again, there is a bias-variance trade-off. There can be multiple consistent methods for the same problem, and their biases and variances don't have to go to zero at the same *rates*.

3.2 The Bias-Variance Trade-Off in Action

Let's take an extreme example: we could decide to approximate $r(x)$ by a constant r_0 . The implicit smoothing here is very strong, but sometimes appropriate. For instance, it's appropriate when $r(x)$ really is a constant! Then trying to estimate any additional structure in the regression function is just so much wasted effort. Alternately, if $r(x)$ is *nearly* constant, we may still be better off approximating it as one. For instance, suppose the true $r(x) = r_0 + a \sin(\nu x)$, where $a \ll 1$ and $\nu \gg 1$ (Figure 3 shows an example). With limited data, we can actually get better predictions by estimating a constant regression function than one with the correct functional form.

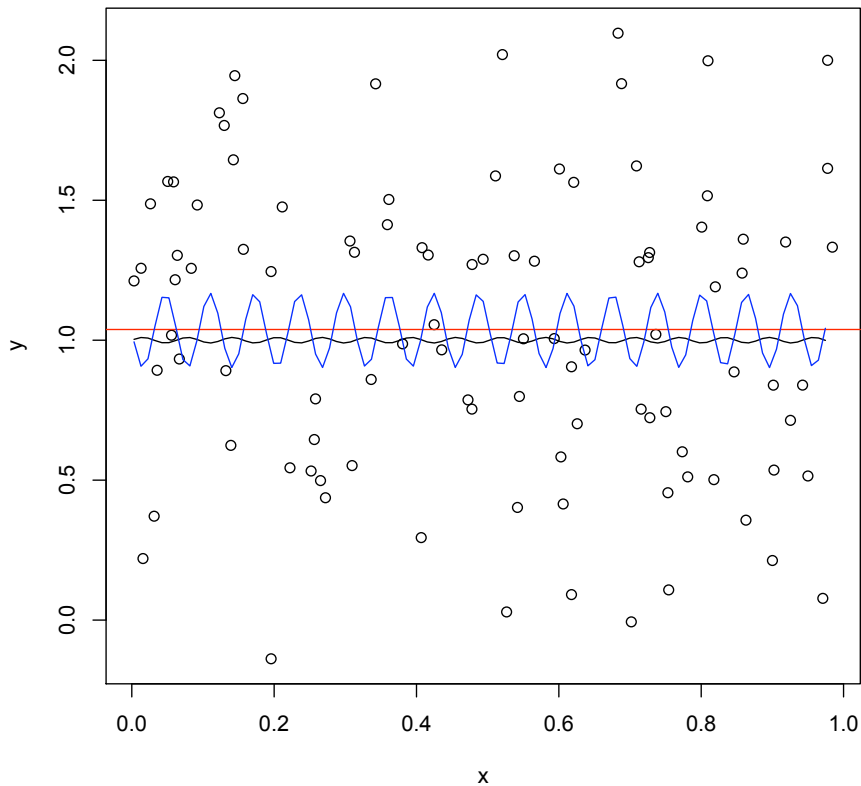
3.3 Ordinary Least Squares Linear Regression as Smoothing

Let's revisit ordinary least-squares linear regression from this point of view. Let's assume that the independent variable X is one-dimensional, and that both X and Y are centered (i.e. have mean zero) — neither of these assumptions is really necessary, but they reduce the book-keeping.

We *choose* to approximate $r(x)$ by $\alpha + \beta x$, and ask for the best values a, b of those constants. These will be the ones which minimize the mean-squared

²To be precise, **consistent for r** , or **consistent for conditional expectations**. More generally, an estimator of any property of the data, or of the whole distribution, is consistent if it converges on the truth.

³You might worry about this claim, especially if you've taken more probability theory — aren't we just saying something about average performance of the \widehat{R} , rather than any *particular* estimated regression function? But notice that if the estimation variance goes to zero, then by Chebyshev's inequality each $\widehat{R}_n(x)$ comes arbitrarily close to $\mathbf{E}[\widehat{R}_n(x)]$ with arbitrarily high probability. If the approximation bias goes to zero, therefore, the estimated regression functions converge *in probability* on the true regression function, not just *in mean*.



```

ugly.func = function(x) {1 + 0.01*sin(100*x)}
r = runif(100); y = ugly.func(r) + rnorm(length(r),0,0.5)
plot(r,y,xlab="x",ylab="y"); curve(ugly.func,add=TRUE)
abline(h=mean(y),col="red")
sine.fit = lm(y ~ 1+ sin(100*r))
curve(sine.fit$coefficients[1]+sine.fit$coefficients[2]*sin(100*x),
      col="blue",add=TRUE)

```

Figure 3: A rapidly-varying but nearly-constant regression function; $Y = 1 + 0.01 \sin 100x + \epsilon$, with $\epsilon \sim \mathcal{N}(0, 0.1)$. (The x values are uniformly distributed between 0 and 1.) Red: constant line at the sample mean. Blue: estimated function of the same form as the true regression function, i.e., $r_0 + a \sin 100x$. If the data set is small enough, the constant actually generalizes better — the bias of using the wrong functional form is smaller than the additional variance from the extra degrees of freedom. Here, the RMS error of the constant on new data is 0.50, while that of the estimated sine function is 0.51 — using the right function actually hurts us!

error.

$$MSE(\alpha, \beta) = \mathbf{E} \left[(Y - \alpha - \beta X)^2 \right] \quad (27)$$

$$= \mathbf{E} \left[(Y - \alpha - \beta X)^2 | X \right] \quad (28)$$

$$= \mathbf{E} \left[\text{Var} [Y|X] + (\mathbf{E} [Y - \alpha - \beta X | X])^2 \right] \quad (29)$$

$$= \mathbf{E} [\text{Var} [Y|X]] + \mathbf{E} \left[(\mathbf{E} [Y - \alpha - \beta X | X])^2 \right] \quad (30)$$

The first term doesn't depend on α or β , so we can drop it for purposes of optimization. Taking derivatives, and then brining them inside the expectations,

$$\frac{\partial MSE}{\partial \alpha} = \mathbf{E} [2(Y - \alpha - \beta X)(-1)] \quad (31)$$

$$\mathbf{E} [Y - a - bX] = 0 \quad (32)$$

$$a = \mathbf{E} [Y] - b\mathbf{E} [X] = 0 \quad (33)$$

using the fact that X and Y are centered; and,

$$\frac{\partial MSE}{\partial \beta} = \mathbf{E} [2(Y - \alpha - \beta X)(-X)] \quad (34)$$

$$\mathbf{E} [XY] - b\mathbf{E} [X^2] = 0 \quad (35)$$

$$b = \frac{\text{Cov} [X, Y]}{\text{Var} [X]} \quad (36)$$

again using the centering of X and Y . That is, the mean-squared optimal linear prediction is

$$r(x) = x \frac{\text{Cov} [X, Y]}{\text{Var} [X]} \quad (37)$$

Now, if we try to estimate this from data, there are two approaches. One is to replace the true population values of the covariance and the variance with their sample values, respectively

$$\frac{1}{n} \sum_i y_i x_i \quad (38)$$

and

$$\frac{1}{n} \sum_i x_i^2 \quad (39)$$

(again, assuming centering). The other is to minimize the residual sum of squares,

$$RSS(\alpha, \beta) \equiv \sum_i (y_i - \alpha - \beta x_i)^2 \quad (40)$$

You may or may not find it surprising that both approaches lead to the same answer:

$$\hat{a} = 0 \tag{41}$$

$$\hat{b} = \frac{\sum_i y_i x_i}{\sum_i x_i^2} \tag{42}$$

Provided that $\text{Var}[X] > 0$, this will converge with IID samples, so we have a consistent estimator.⁴

We are now in a position to see how the least-squares linear regression model is really a smoothing of the data. Let's write the estimated regression function explicitly in terms of the training data points.

$$\hat{r}(x) = \hat{b}x \tag{43}$$

$$= x \frac{\sum_i y_i x_i}{\sum_i x_i^2} \tag{44}$$

$$= \sum_i y_i \frac{x_i}{\sum_j x_j^2} x \tag{45}$$

$$= \sum_i y_i \frac{x_i}{ns_X^2} x \tag{46}$$

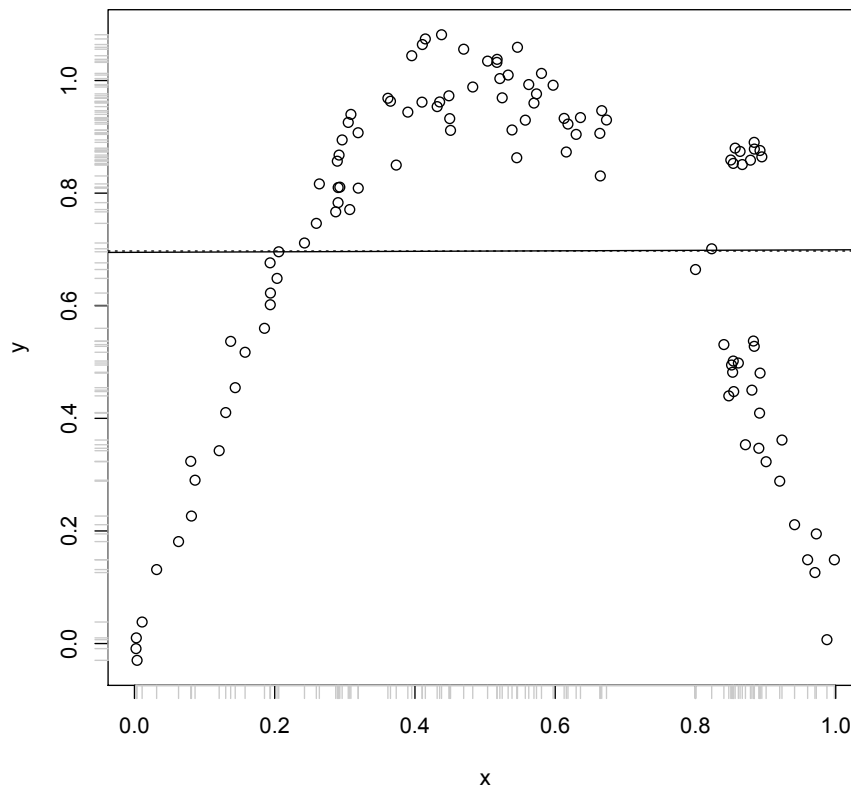
where s_X^2 is the sample variance of X . In words, our prediction is a weighted average of the observed values y_i of the dependent variable, where the weights are proportional to how far x_i is from the center, relative to the variance, and proportional to the magnitude of x . If x_i is on the same side of the center as x , it gets a positive weight, and if it's on the opposite side it gets a negative weight.

Figure 4 shows the data from Figure 1 with the least-squares regression line added. It will not escape your notice that this is very, very slightly different from the constant regression function; the coefficient on X is 0.004438. Visually, the problem is that there should be a positive slope in the left-hand half of the data, and a negative slope in the right, but the slopes are the densities are balanced so that the best *single* slope is zero.⁵

Mathematically, the problem arises from the somewhat peculiar way in which least-squares linear regression smoothes the data. As I said, the weight of a data point depends on how far it is from the *center* of the data, not how far it is from the *point at which we are trying to predict*. This works when $r(x)$ really is a straight line, but otherwise — e.g., here — it's a recipe for trouble. However, it does suggest that if we could somehow just tweak the way we smooth the data, we could do better than linear regression.

⁴Eq. 41 may look funny, but remember that we're assuming X and Y have been centered. Centering doesn't change the slope of the least-squares line but does change the intercept; if we go back to the un-centered variables the intercept becomes $\bar{Y} - \hat{b}\bar{X}$, where the bar denotes the sample mean.

⁵The standard test of whether this coefficient is zero is about as far from rejecting the null hypothesis as you will ever see, $p = 0.965$. You should remember this the next time you look at regression output.



```

abline(h=mean(all.y),lty=2)
fit.all = lm(all.y~all.x)
abline(a=fit.all$coefficients[1],b=fit.all$coefficients[2],col="blue")

```

Figure 4: Data from Figure 1, with a horizontal line at the mean (dotted) and the ordinary least squares regression line (solid). If you zoom in online you will see that there really are two lines there.

4 Linear Smoothers

The sample mean and the linear regression line are both special cases of **linear smoothers**, which are estimates of the regression function with the following form:

$$\hat{r}(x) = \sum_i y_i \hat{w}(x_i, x) \quad (47)$$

The sample mean is the special case where $\hat{w}(x_i, x) = 1/n$, regardless of what x_i and x are.

Ordinary linear regression is the special case where $\hat{w}(x_i, x) = (x_i/n s_X^2)x$.

Both of these, as remarked, ignore how far x_i is from x .

4.1 k -Nearest-Neighbor Regression

At the other extreme, we could do **nearest-neighbor regression**:

$$\hat{w}(x_i, x) = \begin{cases} 1 & x_i \text{ nearest neighbor of } x \\ 0 & \text{otherwise} \end{cases} \quad (48)$$

This is very sensitive to the distance between x_i and x . If $r(x)$ does not change too rapidly, and X is pretty thoroughly sampled, then the nearest neighbor of x among the x_i is probably close to x , so that $r(x_i)$ is probably close to $r(x)$. However, $y_i = r(x_i) + \text{noise}$, so nearest-neighbor regression will include the noise into its prediction. We might instead do k -nearest neighbor regression,

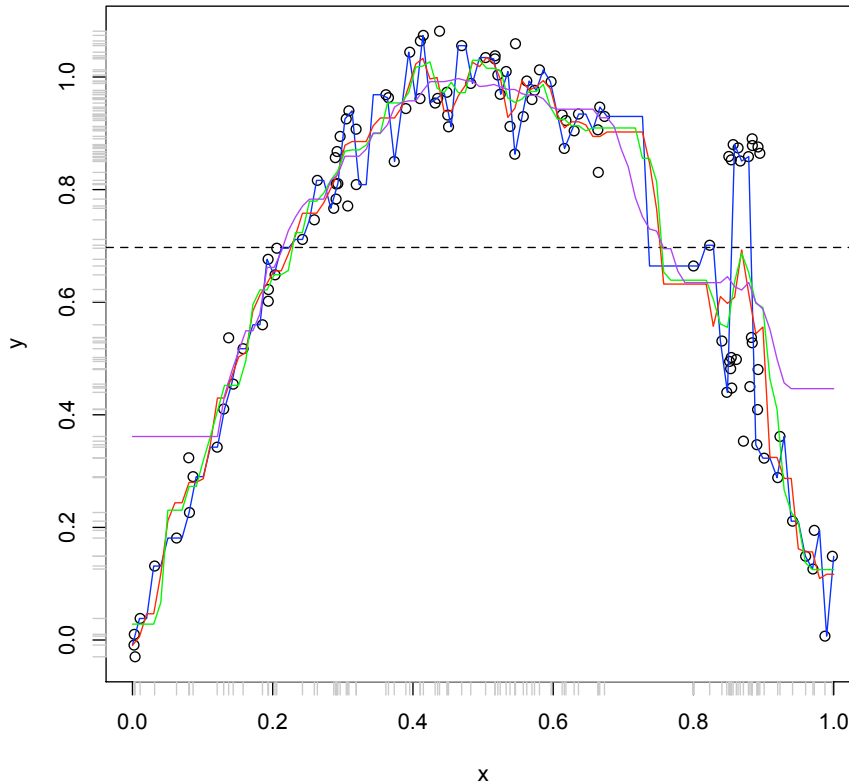
$$\hat{w}(x_i, x) = \begin{cases} 1/k & x_i \text{ one of the } k \text{ nearest neighbors of } x \\ 0 & \text{otherwise} \end{cases} \quad (49)$$

Again, with enough samples all the k nearest neighbors of x are probably close to x , so their regression functions there are going to be close to the regression function at x . But because we average their values of y_i , the noise terms should tend to cancel each other out. As we increase k , we get smoother functions — in the limit $k = n$ and we just get back the constant. Figure 5 illustrates this for our running example data.⁶

To use k -nearest-neighbors regression, we need to pick k somehow. This means we need to decide *how much* smoothing to do, and this is not trivial. We will return to this point.

Because k -nearest-neighbors averages over only a fixed number of neighbors, each of which is a noisy sample, it always has some noise in its prediction, and is generally not consistent. This may not matter very much with moderately-large data (especially once we have a good way of picking k). However, it is sometimes useful to let k systematically grow with n , but not too fast, so as to avoid just doing a global average; say $k \propto \sqrt{n}$. Such schemes *can* be consistent.

⁶The code uses the k -nearest neighbor function provided by the package `knnflex` (available from CRAN). This requires one to pre-compute a matrix of the distances between all the points of interest, i.e., training data and testing data (using `knn.dist`); the `knn.predict` function then needs to be told which rows of that matrix come from training data and which from testing data. See `help(knnflex.predict)` for more, including examples.



```

library(knnflex)
all.dist = knn.dist(c(all.x,seq(from=0,to=1,length.out=100)))
all.nn1.predict = knn.predict(1:110,111:210,all.y,all.dist,k=1)
abline(h=mean(all.y),lty=2)
lines(seq(from=0,to=1,length.out=100),all.nn1.predict,col="blue")
all.nn3.predict = knn.predict(1:110,111:210,all.y,all.dist,k=3)
lines(seq(from=0,to=1,length.out=100),all.nn3.predict,col="red")
all.nn5.predict = knn.predict(1:110,111:210,all.y,all.dist,k=5)
lines(seq(from=0,to=1,length.out=100),all.nn5.predict,col="green")
all.nn20.predict = knn.predict(1:110,111:210,all.y,all.dist,k=20)
lines(seq(from=0,to=1,length.out=100),all.nn20.predict,col="purple")

```

Figure 5: Data points from Figure 1 with horizontal dashed line at the mean and the k -nearest-neighbor regression curves for $k = 1$ (blue), $k = 3$ (red), $k = 5$ (green) and $k = 20$ (purple). Note how increasing k smooths out the regression line, and pulls it back towards the mean. ($k = 100$ would give us back the dashed horizontal line.)

4.2 Kernel Smoothers

Changing k in a k -nearest-neighbors regression lets us change how much smoothing we're doing on our data, but it's a bit awkward to express this in terms of a number of data points. It feels like it would be more natural to talk about a range in the independent variable over which we smooth or average. Another problem with k -NN regression is that each testing point is predicted using information from only a few of the training data points, unlike linear regression or the sample mean, which always uses all the training data. If we could somehow use all the training data, but in a location-sensitive way, that would be nice.

There are several ways to do this, as we'll see, but a particularly useful one is to use a **kernel smoother**, a.k.a. **kernel regression** or **Nadaraya-Watson regression**. To begin with, we need to pick a **kernel function**⁷ $K(x_i, x)$ which satisfies the following properties:

1. $K(x_i, x) \geq 0$
2. $K(x_i, x)$ depends only on the distance $x_i - x$, not the individual arguments
3. $\int xK(0, x)dx = 0$
4. $0 < \int x^2K(0, x)dx < \infty$

These conditions together (especially the last one) imply that $K(x_i, x) \rightarrow 0$ as $|x_i - x| \rightarrow \infty$. Two examples of such functions are the density of the $\text{Unif}(-h/2, h/2)$ distribution, and the density of the standard Gaussian $\mathcal{N}(0, \sqrt{h})$ distribution. Here h can be any positive number, and is called the **bandwidth**.

The Nadaraya-Watson estimate of the regression function is

$$\hat{r}(x) = \sum_i y_i \frac{K(x_i, x)}{\sum_j K(x_j, x)} \quad (50)$$

i.e., in terms of Eq. 47,

$$\hat{w}(x_i, x) = \frac{K(x_i, x)}{\sum_j K(x_j, x)} \quad (51)$$

(Notice that here, as in k -NN regression, the sum of the weights is always 1. Why?)⁸

What does this achieve? Well, $K(x_i, x)$ is large if x_i is close to x , so this will place a lot of weight on the training data points close to the point where we are trying to predict. More distant training points will have smaller weights,

⁷There are many other mathematical objects which are *also* called “kernels”. Some of these meanings are related, but not all of them. (Cf. “normal”.)

⁸What do we do if $K(x_i, x)$ is zero for some x_i ? Nothing; they just get zero weight in the average. What do we do if *all* the $K(x_i, x)$ are zero? Different people adopt different conventions; popular ones are to return the global, unweighted mean of the y_i , to do some sort of interpolation from regions where the weights are defined, and to throw up our hands and refuse to make any predictions (computationally, return **NA**).

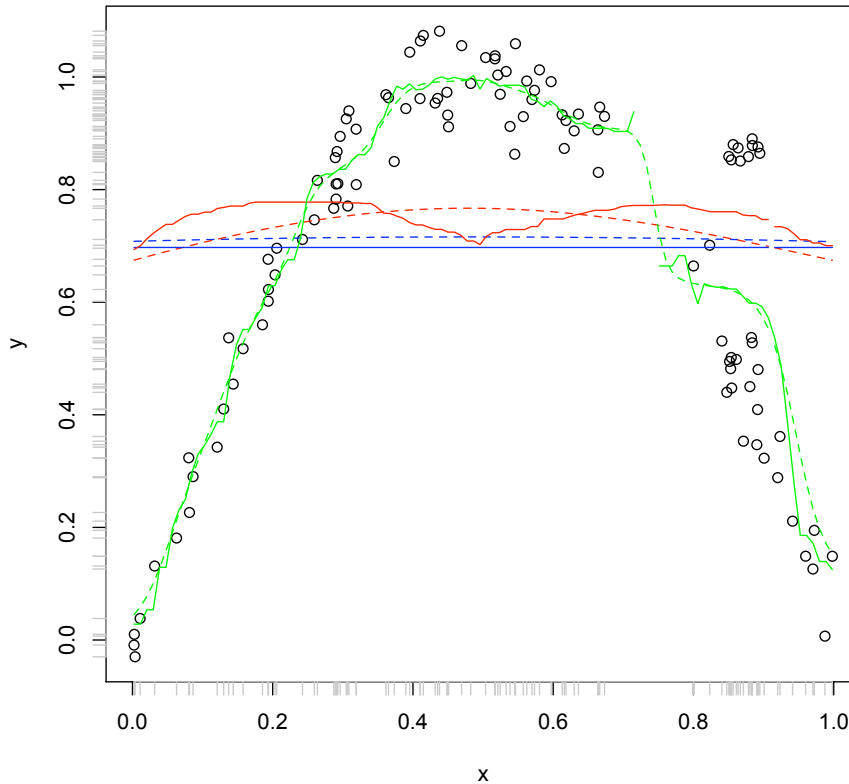
falling off towards zero. If we try to predict at a point x which is very far from any of the training data points, the value of $K(x_i, x)$ will be small for all x_i , but it will typically be *much, much smaller* for all the x_i which are not the nearest neighbor of x , so $\hat{w}(x_i, x) \approx 1$ for the nearest neighbor and ≈ 0 for all the others.⁹ That is, far from the training data, our predictions will tend towards nearest neighbors, rather than going off to $\pm\infty$, as linear regression's predictions do. Whether this is good or bad of course depends on the true $r(x)$ — and how often we have to predict what will happen very far from the training data.

Figure 6 shows our running example data, together with kernel regression estimates formed by combining the uniform-density, or **box**, and Gaussian kernels with different bandwidths. The box kernel simply takes a region of width h around the point x and averages the training data points it finds there. The Gaussian kernel gives reasonably large weights to points within h of x , smaller ones to points within $2h$, tiny ones to points within $3h$, and so on, shrinking like $e^{-(x-x_i)^2/2h}$. As promised, the bandwidth h controls the degree of smoothing. As $h \rightarrow \infty$, we revert to taking the global mean. As $h \rightarrow 0$, we tend to get spikier functions — with the Gaussian kernel at least it tends towards the nearest-neighbor regression.

If we want to use kernel regression, we need to choose both which kernel to use, and the bandwidth to use with it. Experience, like Figure 6, suggests that the bandwidth usually matters a lot more than the kernel. This puts us back to roughly where we were with k -NN regression, needing to control the degree of smoothing, without knowing how smooth $r(x)$ really is. Similarly again, with a fixed bandwidth h , kernel regression is generally not consistent. However, if $h \rightarrow 0$ as $n \rightarrow \infty$, but doesn't shrink *too* fast, then we can get consistency.

Next time, we'll look more at linear regression and some extensions, and then come back to nearest-neighbor and kernel regression, and say something about how to handle things like the blob of data points around $(0.9, 0.9)$ in the scatter-plot.

⁹Take a Gaussian kernel in one dimension, for instance, so $K(x_i, x) \propto e^{-(x_i-x)^2/2h^2}$. Say x_i is the nearest neighbor, and $|x_i - x| = L$, with $L \gg h$. So $K(x_i, x) \propto e^{-L^2/2h^2}$, a small number. But now for any other x_j , $K(x_j, x) \propto e^{-L^2/2h^2} e^{-(x_j-x_i)L/2h^2} e^{-(x_j-x_i)^2/2h^2} \ll e^{-L^2/2h^2}$. — This assumes that we're using a kernel like the Gaussian, which never quite goes to zero, unlike the box kernel.



```

plot(all.x,all.y,xlab="x",ylab="y")
axis(1,at=all.x,labels=FALSE)
axis(2,at=all.y,labels=FALSE)
lines(ksmooth(all.x, all.y, "normal", bandwidth=2),col="blue",lty=2)
lines(ksmooth(all.x, all.y, "normal", bandwidth=1),col="red",lty=2)
lines(ksmooth(all.x, all.y, "normal", bandwidth=0.1),col="green",lty=2)
lines(ksmooth(all.x, all.y, "box", bandwidth=2),col="blue")
lines(ksmooth(all.x, all.y, "box", bandwidth=1),col="red")
lines(ksmooth(all.x, all.y, "box", bandwidth=0.1),col="green")

```

Figure 6: Data from Figure 1 together with kernel regression lines. Solid colored lines are box-kernel estimates, dashed colored lines Gaussian-kernel estimates. Blue, $h = 2$; red, $h = 1$; green, $h = 0.5$; purple, $h = 0.1$ (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 purple) 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.

Exercises

These are for you to think through, not to hand in.

1. Suppose we use the **mean absolute error** instead of the mean squared error:

$$\text{MAE}(a) = \mathbf{E}[|Y - a|] \tag{52}$$

Is this also minimized by taking $a = \mathbf{E}[Y]$? If not, what value \tilde{r} minimizes the MAE? Should we use MSE or MAE to measure error?

2. Derive Eqs. 41 and 42 by minimizing Eq. 40.
3. What does it mean for Gaussian kernel regression to approach nearest-neighbor regression as $h \rightarrow 0$? Why does it do so? Is this true for all kinds of kernel regression?