Chapter 4

Using Nonparametric Smoothing in Regression

Having spent long enough running down linear regression, and thought through evaluating predictive models, it is time to turn to constructive alternatives, which are (also) based on smoothing.

Recall the basic kind of smoothing we are interested in: we have a response variable Y, some input variables which we bind up into a vector X, and a collection of data values, $(x_1, y_1), (x_2, y_2), \dots, x_n, y_n)$. By "smoothing", I mean that predictions are going to be weighted averages of the observed responses in the training data:

$$\widehat{r}(x) = \sum_{i=1}^{n} y_i w(x, x_i, b)$$
(4.1)

Most smoothing methods have a control setting, here written h, that says how much to smooth. With k nearest neighbors, for instance, the weights are 1/k if x_i is one of the k-nearest points to x, and w = 0 otherwise, so large k means that each prediction is an average over many training points. Similarly with kernel regression, where the degree of smoothing is controlled by the bandwidth.

Why do we want to do this? How do we pick how much smoothing to do?

4.1 How Much Should We Smooth?

When we smooth very little $(h \rightarrow 0)$, then we can match very small, fine-grained or sharp aspects of the true regression function, if there are such. That is, less smoothing leads to less bias. At the same time, less smoothing means that each of our predictions is going to be an average over (in effect) fewer observations, making the prediction noisier. Smoothing less increases the variance of our estimate. Since

$$(total error) = (noise) + (bias)^{2} + (variance)$$
(4.2)



curve(2*x^4,from=0,to=1,lty=2,xlab="Smoothing",ylab="Generalization error")
curve(0.12+x-x,lty=3,add=TRUE)
curve(1/(10*x),lty=4,add=TRUE)
curve(0.12+2*x^4+1/(10*x),add=TRUE)

Figure 4.1: Decomposition of the generalization error of smoothing: the total error (solid) equals process noise (dotted) plus approximation error from smoothing (=squared bias, dashed) and estimation variance (dot-and-dash). The numerical values here are arbitrary, but the functional forms (squared bias $\propto h^4$, variance $\propto n^{-1}h^{-1}$) are representative of kernel regression (Eq. 4.12.

(Eq. 1.26), if we plot the different components of error as a function of h, we typically get something that looks like Figure 4.1. Because changing the amount of smoothing has opposite effects on the bias and the variance, there is an optimal amount of smoothing, where we can't reduce one source of error without increasing the other. We therefore want to find that optimal amount of smoothing, which is where cross-validation comes in.

You should note, at this point, that the optimal amount of smoothing depends on the real regression curve, our smoothing method, *and* how much data we have. This is because the variance contribution generally shrinks as we get more data.¹ If we get more data, we go from Figure 4.1 to Figure 4.2. The minimum of the over-all error curve has shifted to the left, and we should smooth less.

Strictly speaking, **parameters** are properties of the data-generating process alone, so the optimal amount of smoothing is not really a parameter. If you do think of it as a parameter, you have the problem of why the "true" value changes as you get more data. It's better thought of as a setting or control variable in the smoothing method, to be adjusted as convenient.

¹Sometimes bias changes as well. Noise does not (why?).





Figure 4.2: Consequences of adding more data to the components of error: noise (dotted) and bias (dashed) are unchanged, but the new variance curve (dotted and dashed, black) is to the left of the old (greyed), so the new over-all error curve (solid black) is lower, and has its minimum at a smaller amount of smoothing than the old (solid grey).

4.2 Adapting to Unknown Roughness

Figure 4.3, which graphs two functions, f and g. Both are "smooth" functions in the mathematical sense². We could Taylor-expand both functions to approximate their values anywhere, just from knowing enough derivatives at one point x_0 .³ If instead of knowing the derivatives at x_0 we have the values of the functions at a sequence of points $x_1, x_2, \ldots x_n$, we could use interpolation to fill out the rest of the curve. Quantitatively, however, f(x) is less smooth than g(x) — it changes much more rapidly, with many reversals of direction. For the same degree of accuracy in the interpolation $f(\cdot)$ needs more, and more closely spaced, training points x_i than does $g(\cdot)$.

Now suppose that we don't get to actually get to see f(x) and g(x), but rather just $f(x)+\epsilon$ and $g(x)+\eta$, where ϵ and η are noise. (To keep things simple I'll assume they're constant-variance, IID Gaussian noises, say with $\sigma = 0.15$.) The data now look something like Figure 4.4. Can we recover the curves?

As remarked in Chapter 1, if we had many measurements at the same x, then we could find the expectation value by averaging: the regression function r(x) = E[Y|X = x], so with multiple observations at x_i , the mean of the corresponding y_i would (by the law of large numbers) converge on r(x). Generally, however, we have at most one measurement per value of x, so simple averaging won't work. Even if we just confine ourselves to the x_i where we have observations, the mean-squared error will always be σ^2 , the noise variance. However, our estimate would be unbiased.

Smoothing methods try to use multiple measurements at points x_i which are *near* the point of interest x. If the regression function is smooth, as we're assuming it is, $r(x_i)$ will be close to r(x). Remember that the mean-squared error is the sum of bias (squared) and variance. Averaging values at $x_i \neq x$ is going to introduce bias, but averaging independent terms together also reduces variance. If smoothing gets rid of more variance than it adds bias, we come out ahead.

Here's a little math to see it. Let's assume that we can do a first-order Taylor expansion (Figure C.1), so

$$r(x_i) \approx r(x) + (x_i - x)r'(x) \tag{4.3}$$

and

$$y_i \approx r(x) + (x_i - x)r'(x) + \epsilon_i \tag{4.4}$$

Now we average: to keep the notation simple, abbreviate the weight $w(x_i, x, h)$ by

²They are " C^{∞} ": not just continuous, but with continuous derivatives to all orders. ³See App. C for a refresher on Taylor expansions and Taylor series.

^{19:55} Thursday 22nd January, 2015



par(mfcol=c(2,1))
true.f <- function(x) {sin(x)*cos(20*x)}
true.g <- function(x) {log(x+1)}
curve(true.f(x),from=0,to=3,xlab="x",ylab=expression(f(x)))
curve(true.g(x),from=0,to=3,xlab="x",ylab=expression(g(x)))</pre>

Figure 4.3: Two curves for the running example. Above, f(x); below, g(x). (As it happens, $f(x) = \sin x \cos 20x$, and $g(x) = \log x + 1$, but that doesn't really matter.)



curve(true.g(x),col="grey",add=TRUE)

Figure 4.4: The curves of Fig. 4.3 (in grey), plus IID Gaussian noise with mean 0 and standard deviation 0.15. The two curves are sampled at the same x values, but with different noise realizations.

just w;.

$$\widehat{r}(x) = \sum_{i=1}^{n} y_i w_i$$
(4.5)

$$= \sum_{i=1}^{n} (r(x) + (x_i - x)r'(x) + \epsilon_i)w_i$$
(4.6)

$$= r(x) + \sum_{i=1}^{n} w_i \epsilon_i + r'(x) \sum_{i=1}^{n} w_i (x_i - x)$$
(4.7)

$$\widehat{r}(x) - r(x) = \sum_{i=1}^{n} w_i \epsilon_i + r'(x) \sum_{i=1}^{n} w_i(x_i - x)$$
(4.8)

$$\mathbf{E}\left[(\hat{r}(x) - r(x))^{2}\right] = \sigma^{2} \sum_{i=1}^{n} w_{i}^{2} + \mathbf{E}\left[\left(r'(x) \sum_{i=1}^{n} w_{i}(x_{i} - x)\right)^{2}\right]$$
(4.9)

(Remember that: $\sum w_i = 1$; $\mathbf{E}[\epsilon_i] = 0$; ϵ is uncorrelated with everything; and $\operatorname{Var}\left[\epsilon_{i}\right] = \sigma^{2}$.)

The first term on the final right-hand side is an estimation variance, which will tend to shrink as n grows. (If we just do a simple global mean, $w_i = 1/n$ for all i, so we get σ^2/n , just like in baby stats.) The second term, an expectation, on the other hand, is bias, which grows as x_i gets further from x, and as the magnitudes of the derivatives grow, i.e., its growth varies with how smooth or wiggly the regression function is. For smoothing to work, w_i had better shrink as $x_i - x$ and r'(x) grow.⁴ Finally, all else being equal, w_i , should also shrink with n, so that the over-all size of the sum shrinks as we get more data.

To illustrate, let's try to estimate f(1.6) and g(1.6) from the noisy observations. We'll try a simple approach, just averaging all values of $f(x_i) + \epsilon_i$ and $g(x_i) + \eta_i$ for $1.5 < x_i < 1.7$ with equal weights. For f, this gives 0.46, while f(1.6) = 0.83. For g, this gives 0.91, with g(1.6) = 0.95. (See figure 4.5.) The same window size creates a much larger bias with the rougher, more rapidly changing f than with the smoother, more slowly changing g. Varying the size of the averaging window will change the amount of error, and it will change it in different ways for the two functions.

If one does a more careful second-order Taylor expansion like that leading to Eq. 4.9, specifically for kernel regression, one can show that the bias at x is

$$\mathbf{E}\left[\hat{r}(x) - r(x)|X_1 = x_1, \dots, X_n = x_n\right] = h^2 \left[\frac{1}{2}r''(x) + \frac{r'(x)f'(x)}{f(x)}\right]\sigma_K^2 + o(h^2)$$
(4.10)

where f is the density of x, and $\sigma_K^2 = \int u^2 K(u) du$, the variance of the probability density corresponding to the kernel⁵. The r'' term just comes from the second-order part of the Taylor expansion. To see where the r'f' term comes from, imagine first

⁴The higher derivatives of r also matter, since we should really keep more than just the first term in the Taylor expansion. The details get messy, but Eq. 4.12 below gives the upshot for kernel smoothing. ⁵If you are not familiar with the "order" symbols O and *o*, see Appendix B.



par(mfcol=c(2,1))
x.focus <- 1.6; x.lo <- x.focus-0.1; x.hi <- x.focus+0.1
colors=ifelse((x<x.hi)&(x>x.lo),"black","grey")
plot(x,yf,xlab="x",ylab=expression(f(x)+epsilon),col=colors)
curve(true.f(x),col="grey",add=TRUE)
points(x.focus,mean(yf[(x<x.hi)&(x>x.lo)]),pch=18,cex=2)
plot(x,yg,xlab="x",ylab=expression(g(x)+eta),col=colors)
curve(true.g(x),col="grey",add=TRUE)
points(x.focus,mean(yg[(x<x.hi)&(x>x.lo)]),pch=18,cex=2)

Figure 4.5: Relationship between smoothing and function roughness. In both panels we estimate the value of the regression function at x = 1.6 by averaging observations where $1.5 < x_i < 1.7$ (black points, others are "ghosted" in grey). The location of the average in shown by the large black diamond. This works poorly for the rough function f(x) in the upper panel (the bias is large), but much better for the smoother function in the lower panel (the bias is small).



Figure 4.6: Error of estimating f(1.6) (solid line) and g(1.6) (dashed) from averaging observed values at 1.6 - h < x < 1.6 + h, for different radii h. The grey is σ , the standard deviation of the noise — how can the estimation error be smaller than that?

that x is a mode of the distribution, so f'(x) = 0. As h shrinks, only training points where X_i is very close to x will have any weight in $\hat{r}(x)$, and their distribution will be roughly symmetric around x (at least once h is sufficiently small). So, at mode, $\mathbf{E}\left[w(X_i, x, h)(X_i - x)\hat{r}(x)\right] \approx 0$. Away from a mode, there will tend to be more training points on one side or the other of x, depending on the sign of f'(x), and this induces a bias. The tricky part of the analysis is concluding that the bias has exactly the form given above.⁶

One can also work out the variance of the kernel regression estimate,

$$\operatorname{Var}\left[\hat{r}(x)|X_{1} = x_{1}, \dots, X_{n} = x_{n}\right] = \frac{\sigma^{2}(x)R(K)}{nhf(x)} + o((nh)^{-1})$$
(4.11)

where $R(K) = \int K^2(u)du$. Roughly speaking, the width of the region where the kernel puts non-trivial weight is about h, so there will be about nhf(x) training points available to estimate $\hat{r}(x)$. Each of these has a y_i value, equal to r(x) plus noise of variance $\sigma^2(x)$. The final factor of R(K) accounts for the average weight.

Putting the bias together with the variance, we get an expression for the mean squared error of the kernel regression at *x*:

$$MSE(x) = \sigma^{2}(x) + h^{4} \left[\frac{1}{2} r''(x) + \frac{r'(x)f'(x)}{f(x)} \right]^{2} (\sigma_{K}^{2})^{2} + \frac{\sigma^{2}(x)R(K)}{nhf(x)} + o(h^{4}) + o(1/nh)$$
(4.12)

Eq. 4.12 tells us that, in principle, there is a single optimal choice of bandwidth h, an optimal degree of smoothing. We could find it by taking Eq. 4.12, differentiating with respect to the bandwidth, and setting everything to zero (neglecting the *o* terms):

$$0 = 4h^{3} \left[\frac{1}{2} r''(x) + \frac{r'(x)f'(x)}{f(x)} \right]^{2} (\sigma_{K}^{2})^{2} - \frac{\sigma^{2}(x)R(K)}{nh^{2}f(x)}$$
(4.13)

$$b = \left(n\frac{4f(x)(\sigma_K^2)^2 \left[\frac{1}{2}r''(x) + \frac{r'(x)f'(x)}{f(x)}\right]^2}{\sigma^2(x)R(K)}\right)^{-1/5}$$
(4.14)

Of course, this expression for the optimal h involves the unknown derivatives r'(x) and r''(x), plus the unknown density f(x) and its unknown derivative f'(x). But if we knew the derivative of the regression function, we would basically know the function itself (just integrate), so we seem to be in a vicious circle, where we need to know the function before we can learn it.⁷

One way of expressing this is to talk about how well a smoothing procedure *would* work, if an Oracle were to tell us the derivatives, or (to cut to the chase) the optimal bandwidth h_{opt} . Since most of us do not have access to such oracles, we

⁶Exercise 1 sketches the demonstration for the special case of the uniform ("boxcar") kernel.

⁷You may be wondering why I keep talking about *the* optimal bandwidth, when Eq. 4.14 makes it seem that the bandwidth should vary with *x*. One can go through pretty much the same sort of analysis in terms of the *expected* values of the derivatives, and the qualitative conclusions will be the same, but the notational overhead is even worse. Alternatively, there are techniques for variable-bandwidth smoothing.

need to *estimate* h_{opt} . Once we have this estimate, \hat{h} , then we get out weights and our predictions, and so a certain mean-squared error. Basically, our MSE will be the Oracle's MSE, plus an extra term which depends on how far \hat{h} is to h_{opt} , and how sensitive the smoother is to the choice of bandwidth.

What would be really nice would be an **adaptive** procedure, one where our actual MSE, using \hat{h} , approaches the Oracle's MSE, which it gets from h_{opt} . This would mean that, in effect, we are *figuring out* how rough the underlying regression function is, and so how much smoothing to do, rather than having to guess or be told. An adaptive procedure, if we can find one, is a partial⁸ substitute for prior knowledge.

4.2.1 Bandwidth Selection by Cross-Validation

The most straight-forward way to pick a bandwidth, and one which generally manages to be adaptive, is in fact cross-validation; *k*-fold CV is usually somewhat better than leave-one-out, but the latter often works acceptably too. The usual procedure is to come up with an initial grid of candidate bandwidths, and then use cross-validation to estimate how well each one of them would generalize. The one with the lowest error under cross-validation is then used to fit the regression curve to the whole data⁹.

Code Example 3 shows how it would work in R, with a one predictor variable, borrowing the npreg function from the np library (Hayfield and Racine, 2008).¹⁰

The return value has three parts. The first is the actual best bandwidth. The second is a vector which gives the cross-validated mean-squared mean-squared errors of all the different bandwidths in the vector bandwidths. The third component is an array which gives the MSE for each bandwidth on each fold. It can be useful to know things like whether the difference between the CV score of the best bandwidth and the runner-up is bigger than their fold-to-fold variability.

Figure 4.7 plots the CV estimate of the (root) mean-squared error versus bandwidth for our two curves. Figure 4.8 shows the data, the actual regression functions and the estimated curves with the CV-selected bandwidths. This illustrates *why* picking the bandwidth by cross-validation works: the curve of CV error against bandwidth is actually a pretty good approximation to the true curve of generalization error (which would look like Figure 4.1), so optimizing the CV error is close to optimizing the generalization error.

Notice, by the way, in Figure 4.7, that the rougher curve is more sensitive to the choice of bandwidth, and that the smoother curve always has a lower mean-squared error. Also notice that, at the minimum, one of the cross-validation estimates of generalization error is smaller than the true system noise level; this shows that cross-validation doesn't completely correct for optimism¹¹.

⁸Only partial, because we'd *always* do better if the Oracle would just tell us h_{opt} .

⁹Since the optimal bandwidth is $\propto n^{-1/5}$, and the training sets in cross-validation are smaller than the whole data set, one might adjust the bandwidth proportionally. However, if *n* is small enough that this makes a big difference, the sheer noise in bandwidth estimation usually overwhelms this.

¹⁰The package has methods for automatically selecting bandwidth by cross-validation — see §4.6 below.

¹¹Tibshirani and Tibshirani (2009) gives a fairly straightforward way to adjust the estimate of the generalization error for the selected model or bandwidth, but that doesn't influence the choice of the best bandwidth.

```
cv_bws_npreg <- function(x,y,bandwidths=(1:50)/50,nfolds=10) {</pre>
  require(np)
  n <- length(x)</pre>
  stopifnot(n> 1, length(y) == n)
  stopifnot(length(bandwidths) > 1)
  stopifnot(nfolds > 0, nfolds==trunc(nfolds))
  fold_MSEs <- matrix(0,nrow=nfolds,ncol=length(bandwidths))</pre>
  colnames(fold_MSEs) = bandwidths
  case.folds <- sample(rep(1:nfolds,length.out=n))</pre>
  for (fold in 1:nfolds) {
    train.rows = which(case.folds!=fold)
    x.train = x[train.rows]
    y.train = y[train.rows]
    x.test = x[-train.rows]
    y.test = y[-train.rows]
    for (bw in bandwidths) {
      fit <- npreg(txdat=x.train,tydat=y.train,</pre>
                    exdat=x.test,eydat=y.test,bws=bw)
      fold_MSEs[fold,paste(bw)] <- fit$MSE</pre>
    }
  }
  CV_MSEs = colMeans(fold_MSEs)
  best.bw = bandwidths[which.min(CV_MSEs)]
  return(list(best.bw=best.bw,CV_MSEs=CV_MSEs,fold_MSEs=fold_MSEs))
}
```

Code Example 3: Cross-validation for univariate kernel regression (see online for comments). The colnames trick: component names have to be character strings; other data types will be coerced into characters when we assign them to be names. Later, when we want to refer to a bandwidth column by its name, we wrap the name in another coercing function, such as paste. — The vector of default bandwidths is only for illustration; it should not be used blindly on data, or in homework.





Figure 4.7: Cross-validated estimate of the (root) mean-squard error as a function of the bandwidth (solid curve, f(x) data; dashed, g(x) data; grey line, true noise σ). Notice that the rougher curve is more sensitive to the choice of bandwidth, and that the smoother curve is more predictable at every choice of bandwidth. CV selects bandwidths of 0.020 for f and 0.055 for g.



x.ord=order(x)
par(mfcol=c(2,1))
plot(x,yf,xlab="x",ylab=expression(f(x)+epsilon))
fhat <- npreg(bws=fbws\$best.bw,txdat=x,tydat=yf)
lines(x[x.ord],fitted(fhat)[x.ord],lwd=4)
curve(true.f(x),col="grey",add=TRUE,lwd=2)
plot(x,yg,xlab="x",ylab=expression(g(x)+eta))
ghat <- npreg(bws=fbws\$best.bw,txdat=x,tydat=yg)
lines(x[x.ord],fitted(ghat)[x.ord],lwd=4)
curve(true.g(x),col="grey",add=TRUE,lwd=2)</pre>

Figure 4.8: Data from the running examples (circles), true regression functions (grey) and kernel estimates of regression functions with CV-selected bandwidths (black). R NOTES: The x values aren't sorted, so we need to put them in order before drawing lines connecting the fitted values; then we need to put the fitted values in the same order. Alternately, we could have used predict on the sorted values, as in §4.3.

We still need to come up with an initial set of candidate bandwidths. For reasons which will drop out of the math in Chapter 16, it's often reasonable to start around $1.06s_X/n^{1/5}$, where s_X is the sample standard deviation of X. However, it is hard to be very precise about this, and good results often require some honest trial and error.

4.2.2 Convergence of Kernel Smoothing and Bandwidth Scaling

Go back to Eq. 4.12 for the mean squared error of kernel regression. As we said, it involves some unknown constants, but we can bury them inside big-O order symbols, which also absorb the little-*o* remainder terms:

$$MSE(h) = \sigma^{2}(x) + O(h^{4}) + O(1/nh)$$
(4.15)

The $\sigma^2(x)$ term is going to be there no matter what, so let's look at the excess risk over and above the intrinsic noise:

$$MSE(h) - \sigma^{2}(x) = O(h^{4}) + O(1/nh)$$
(4.16)

That is, the (squared) bias from the kernel's only approximately getting the curve is proportional to the fourth power of the bandwidth, but the variance is inversely proportional to the product of sample size and bandwidth. If we kept h constant and just let $n \to \infty$, we'd get rid of the variance, but we'd be left with the bias. To get the MSE to go to zero, we need to let the bandwidth h change with n — call it h_n . Specifically, suppose $h_n \to 0$ as $n \to \infty$, but $nh_n \to \infty$. Then, by Eq. 4.16, the risk (generalization error) of kernel smoothing is approaching that of the ideal predictor.

What is the best bandwidth? We saw in Eq. 4.14 that it is (up to constants)

$$h_{\rm opt} = O(n^{-1/5}) \tag{4.17}$$

If we put this bandwidth into Eq. 4.16, we get

$$MSE(h) - \sigma^{2}(x) = O\left(\left(n^{-1/5}\right)^{4}\right) + O\left(n^{-1}\left(n^{-1/5}\right)^{-1}\right) = O\left(n^{-4/5}\right) + O\left(n^{-4/5}\right) = O\left(n^{-4/5}\right)$$
(4.18)

That is, the excess prediction error of kernel smoothing over and above the system noise goes to zero as $1/n^{0.8}$. Notice, by the way, that the contributions of bias and variance to the generalization error are both of the same order, $n^{-0.8}$.

Is this fast or small? We can compare it to what would happen with a parametric model, say with parameter θ . (For linear regression, θ would be the vector of slopes and the intercept.) The optimal value of the parameter, θ_0 , minimizes the mean-squared error. At θ_0 , the parametric model has MSE

$$MSE(\theta_0) = \sigma^2(x) + b(x, \theta_0) \tag{4.19}$$

where b is the bias of the parametric model; this is zero when the parametric model is true¹². Since θ_0 is unknown and must be estimated, one typically has $\hat{\theta} - \theta_0 =$

¹²When the model is wrong, the optimal parameter value θ_0 is often called the **pseudo-truth**.

 $O(1/\sqrt{n})$. Because the error is minimized at θ_0 , the first derivatives of *MSE* at θ_0 are 0. Doing a second-order Taylor expansion of the parametric model contributes an error $O((\hat{\theta} - \theta_0)^2)$, so altogether

$$MSE(\hat{\theta}) - \sigma^2(x) = b(x,\theta_0) + O(1/n)$$
(4.20)

This means parametric models converge more quickly $(n^{-1} \text{ goes to zero faster than } n^{-0.8})$, but they typically converge to the wrong answer $(b^2 > 0)$. Kernel smoothing converges more slowly, but always converges to the right answer¹³.

This doesn't change much if we use cross-validation. Writing $\widehat{h_{CV}}$ for the bandwidth picked by cross-validation, it turns out (Simonoff, 1996, ch. 5) that

$$\frac{\widehat{h_{CV}} - h_{\text{opt}}}{h_{\text{opt}}} - 1 = O(n^{-1/10})$$
(4.21)

Given this, one concludes (Exercise 2) that the MSE of using $\widehat{h_{CV}}$ is also $O(n^{-4/5})$.

4.2.3 Summary on Kernel Smoothing in 1D

Suppose that X and Y are both one-dimensional, and the true regression function $r(x) = \mathbb{E}[Y|X = x]$ is continuous and has first and second derivatives¹⁴. Suppose that the noise around the true regression function is uncorrelated between different observations. Then the bias of kernel smoothing, when the kernel has bandwidth h, is $O(h^2)$, and the variance, after n samples, is O(1/nh). The optimal bandwidth is $O(n^{-1/5})$, and the excess mean squared error of using this bandwidth is $O(n^{-4/5})$. If the bandwidth is selected by cross-validation, the excess risk is still $O(n^{-4/5})$.

4.3 Kernel Regression with Multiple Inputs

For the most part, when I've been writing out kernel regression I have been treating the input variable x as a scalar. There's no reason to insist on this, however; it could equally well be a vector. If we want to enforce that in the notation, say by writing $\vec{x} = (x^1, x^2, \dots, x^d)$, then the kernel regression of y on \vec{x} would just be

$$\hat{r}(\vec{x}) = \sum_{i=1}^{n} y_i \frac{K(\vec{x} - \vec{x_i})}{\sum_{j=1}^{n} K(\vec{x} - \vec{x_j})}$$
(4.22)

In fact, if we want to predict a vector, we'd just substitute $\vec{y_i}$ for y_i above.

¹³It is natural to wonder if one couldn't do better than kernel smoothing's $O(n^{-4/5})$ while still having no asymptotic bias. Resolving this is very difficult, but the answer turns out to be "no" in the following sense (Wasserman, 2006). Any curve-fitting method which can learn arbitrary smooth regression functions will have some curves where it cannot converge any faster than $O(n^{-4/5})$. (In the jargon, that is the **minimax rate**.) Methods which converge faster than this for some kinds of curves have to converge more slowly for others. So this is the best rate we can hope for on truly unknown curves.

¹⁴Or can be approximated to arbitrarily closely by such functions.

To make this work, we need kernel functions for vectors. For scalars, I said that any probability density function would work so long as it had mean zero, and a finite, strictly positive (not 0 or ∞) variance. The same conditions carry over: any distribution over vectors can be used as a multivariate kernel, provided it has mean zero, and the variance matrix is finite and strictly positive¹⁵. In practice, the overwhelmingly most common and practical choice is to use **product kernels**¹⁶.

A product kernel simply uses a different kernel for each component, and then multiplies them together:

$$K(\vec{x} - \vec{x_i}) = K_1(x^1 - x_i^1)K_2(x^2 - x_i^2)\dots K_d(x^d - x_i^d)$$
(4.23)

Now we just need to pick a bandwidth for each kernel, which in general should not be equal $- \operatorname{say} \vec{h} = (h_1, h_2, \dots h_d)$. Instead of having a one-dimensional error curve, as in Figure 4.1 or 4.2, we will have a *d*-dimensional error surface, but we can still use cross-validation to find the vector of bandwidths that generalizes best. We generally can't, unfortunately, break the problem up into somehow picking the best bandwidth for each variable without considering the others. This makes it slower to select good bandwidths in multivariate problems, but still often feasible.

(We can actually turn the need to select bandwidths together to our advantage. If one or more of the variables are irrelevant to our prediction given the others, crossvalidation will tend to give them the maximum possible bandwidth, and smooth away their influence. In Chapter 16, we'll look at formal tests based on this idea.)

Kernel regression will recover almost any regression function. This is true even when the true regression function involves lots of interactions among the input variables, perhaps in complicated forms that would be very hard to express in linear regression. For instance, Figure 4.9 shows a contour plot of a reasonably complicated regression surface, at least if one were to write it as polynomials in x^1 and x^2 , which would be the usual approach. Figure 4.11 shows the estimate we get with a product of Gaussian kernels and only 1000 noisy data points. It's not perfect, of course (in particular the estimated contours aren't as perfectly smooth and round as the true ones), but the important thing is that we got this without having to know, and describe in Cartesian coordinates, the type of shape we were looking for. Kernel smoothing *discovered* the right general form.

There are limits to these abilities of kernel smoothers; the biggest one is that they require more and more data as the number of predictor variables increases. We will see later (Chapter 9) exactly how much data is required, generalizing the kind of analysis done §4.2.2, and some of the compromises this can force us into.

4.4 Interpreting Smoothers: Plots

In a linear regression without interactions, it is fairly easy to interpret the coefficients. The expected response changes by β_i for a one-unit change in the *i*th input variable.

¹⁵Remember that for a matrix v to be "strictly positive", it must be the case that for any vector \vec{a} , $\vec{a} \cdot v\vec{a} > 0$. Covariance matrices are automatically non-negative, so we're just ruling out the case of some weird direction along which the distribution has zero variance.

¹⁶People do sometimes use multivariate Gaussians; we'll glance at this in Chapter 15.



Figure 4.9: An example of a regression surface that would be very hard to learn by piling together interaction terms in a linear regression framework. (Can you guess what the mystery function f is?) — wireframe is from the graphics library lattice.



Figure 4.10: 1000 points sampled from the surface in Figure 4.9, plus independent Gaussian noise (s.d. = 0.05).



Figure 4.11: Gaussian kernel regression of the points in Figure 4.10. Notice that the estimated function will make predictions at arbitrary points, not just the places where there was training data.

The coefficients are also the derivatives of the expected response with respect to the inputs. And it is easy to draw pictures of how the output changes as the inputs are varied, though the pictures are somewhat boring (straight lines or planes).

As soon as we introduce interactions, all this becomes harder, even for parametric regression. If there is an interaction between two components of the input, say x^1 and x^2 , then we can't talk about the change in the expected response for a one-unit change in x^1 without saying what x^2 is. We might *average* over x^2 values, and we'll see next time a reasonable way of doing this, but the flat statement "increasing x^1 by one unit increases the response by β_1 " is just false, no matter what number we fill in for β_1 . Likewise for derivatives; we'll come back to them next time as well.

What about pictures? With only two input variables, we can make wireframe plots like Figure 4.11, or contour or level plots, which will show the predictions for different combinations of the two variables. But what if we want to look at one variable at a time, or there are more than two input variables?

A reasonable way to produce *a* curve for each input variable is to set all the others to some "typical" value, like their means or medians, and to then plot the predicted response as a function of the one remaining variable of interest (Figure 4.12). Of course, when there are interactions, changing the values of the other inputs will change the response to the input of interest, so it's a good idea to produce a couple of curves, possibly super-imposed (Figure 4.12 again).

If there are three or more input variables, we can look at the interactions of any two of them, taken together, by fixing the others and making three-dimensional or contour plots, along the same principles.

The fact that smoothers don't give us a simple story about how each input is associated with the response may seem like a disadvantage compared to using linear regression. Whether it really is a disadvantage depends on whether there really is a simple story to be told, and/or how much big a lie you are prepared to tell in order to keep your story simple.

4.5 Average Predictive Comparisons

Suppose we have a linear regression model

$$Y = \beta_1 X_1 + \beta_2 X_2 + \epsilon \tag{4.24}$$

and we want to know how much Y changes, on average, for a one-unit increase in X_1 . The answer, as you know very well, is just β_1 :

$$[\beta_1(X_1+1) + \beta_2 X_2] - [\beta_1 X_1 + \beta_2 X_2] = \beta_1$$
(4.25)

This is an interpretation of the regression coefficients which you are very used to giving. But it fails as soon as we have interactions:

$$Y = \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \epsilon$$
 (4.26)

Now the effect of increasing X_1 by 1 is

$$[\beta_1(X_1+1) + \beta_2 X_2 + \beta_3(X_1+1)X_2] - [\beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2] = \beta_1 + \beta_3 X_2$$
(4.27)



```
new.frame <- data.frame(x=seq(-3,3,length.out=300),y=median(y.noise))
plot(new.frame$x,predict(noise.np,newdata=new.frame),
    type="l",xlab=expression(x^1),ylab="y",ylim=c(0,1.0))
new.frame$y <- quantile(y.noise,0.25)
lines(new.frame$x,predict(noise.np,newdata=new.frame),lty=2)
new.frame$y <- quantile(y.noise,0.75)
lines(new.frame$x,predict(noise.np,newdata=new.frame),lty=3)</pre>
```

Figure 4.12: Predicted mean response as function of the first input coordinate x^1 for the example data, evaluated with the second coordinate x^2 set to the median (solid), its 25th percentile (dashed) and its 75th percentile (dotted). Note that the changing shape of the partial response curve indicates an interaction between the two inputs. Also, note that the model can make predictions at arbitrary coordinates, whether or not there were any training points there.



curve(exp(7*x)/(1+exp(7*x)),from=-5,to=5,ylab="y")

Figure 4.13: The function of Eq. 4.28, with $\beta = 7$.

The right answer to "how much does the response change when X_1 is increased by one unit?" depends on the value of X_2 ; it's certainly not just " β_1 ".

We also can't give just a single answer if there are nonlinearities. Suppose that the true regression function is this:

$$Y = \frac{e^{\beta X}}{1 + e^{\beta X}} + \epsilon \tag{4.28}$$

which looks like Figure 4.13, setting $\beta = 7$ (for luck). Moving x from -4 to -3 increases the response by 7.57×10^{-10} , but the increase in the response from x = -1 to x = 0 is 0.499. Functions like this are very common in psychology, medicine (dose-response curves for drugs), biology, etc., and yet we cannot sensibly talk about *the* response to a one-unit increase in x. (We will come back to curves which look like this in Chapter 12.)

More generally, let's say we are regressing Y on a vector \vec{X} , and want to assess the impact of one component of the input on Y. To keep the use of subscripts and superscripts to a minimum, we'll write $\vec{X} = (u, \vec{V})$, where u is the coordinate we're really interested in. (It doesn't have to come first, of course.) We would like to know how much the prediction changes as we change u,

$$\mathbf{E}\left[Y|\vec{X} = (u^{(2)}, \vec{v})\right] - \mathbf{E}\left[Y|\vec{X} = (u^{(1)}, \vec{v})\right]$$
(4.29)

and the change in the response per unit change in u,

$$\frac{\mathbf{E}\left[Y|\vec{X}=(u^{(2)},\vec{v})\right] - \mathbf{E}\left[Y|\vec{X}=(u^{(1)},\vec{v})\right]}{u^{(2)} - u^{(1)}}$$
(4.30)

Both of these, but especially the latter, are called the **predictive comparison**. Note that both of them, as written, depend on $u^{(1)}$ (the starting value for the variable of interest), on $u^{(2)}$ (the ending value), and on \vec{v} (the other variables, held fixed during this comparison). We have just seen that in a linear model without interactions, $u^{(1)}$, $u^{(2)}$ and \vec{v} all go away and leave us with the regression coefficient on u. In nonlinear or interacting models, we can't simplify so much.

Once we have estimated a regression model, we can chose our starting point, ending point and context, and just plug in to Eq. 4.29 or Eq. 4.30. (Or problem 9 in problem set 32.) But suppose we do want to boil this down into a single number for each input variable — how might we go about this?

One good answer, which comes from Gelman and Pardoe (2007), is just to average 4.30 over the data¹⁷. More specifically, we have as our **average predictive comparison** for u

$$\frac{\sum_{i=1}^{n} \sum_{j=1}^{n} (\hat{r}(u_{j}, \vec{v_{i}}) - \hat{r}(u_{i}, \vec{v_{i}})) \operatorname{sign}(u_{j} - u_{i})}{\sum_{i=1}^{n} \sum_{j=1}^{n} (u_{j} - u_{i}) \operatorname{sign}(u_{j} - u_{i})}$$
(4.31)

where *i* and *j* run over data points, \hat{r} is our estimated regression function, and the sign function is defined by sign(x) = +1 if x > 0, = 0 if x = 0, and = -1 if x < 0. We use the sign function this way to make sure we are always looking at the consequences of *increasing u*.

The average predictive comparison is a reasonable summary of how rapidly we should expect the response to vary as u changes slightly. But we need to remember that once the model is nonlinear or has interactions, it's just not possible to boil the whole the predictive relationship between u and y into one number. In particular, the value of Eq. 4.31 is going to depend on the distribution of u (and possibly of v), even when the regression function is unchanged. (See Exercise 3.)

4.6 Computational Advice: npreg

The homework will call for you to do nonparametric regression with the np package — which we've already looked at a little. It's a powerful bit of software, but it can take a bit of getting used to. This section is not a substitute for reading Hayfield and Racine (2008), but should get you started.

We'll look at a synthetic-data example with four variables: a quantitative response Y, two quantitative predictors X and Z, and a categorical predictor U, which can be either "A" or "B". The true model is

$$Y = \epsilon + 20X^2 + \begin{cases} Z & \text{if } W = A\\ 10e^Z/(1+e^Z) & \text{if } W = B \end{cases}$$
(4.32)

with $\epsilon \sim \mathcal{N}(0, 0.05)$. Code Example 4 generates some data from this model for us.

The basic function for fitting a kernel regression in np is npreg - conceptually, it's the equivalent of lm. Like lm, it takes a formula argument, which specifies the

¹⁷Actually, they propose something a bit more complicated, which takes into account the uncertainty in our estimate of the regression function, via bootstrapping (Chapter 6).

```
make.demo.df <- function(n) {
  demo.func <- function(x,z,w)
    20*x^2 + ifelse(w=="A",z,10*exp(z)/(1+exp(z))
  }
  x <- runif(n,-1,1)
  z <- rnorm(n,0,10)
  w <- sample(c("A","B"),size=n,replace=TRUE)
  y <- demo.func(x,z,w)+rnorm(n,0,0.05)
  return(data.frame(x=x,y=y,z=z,w=w))
}
demo.df <- make.demo.df(100)</pre>
```

Code Example 4: Generating data from Eq. 4.32.

model, and a data argument, which is a data frame containing the variables included in the formula. The basic idea is to do something like this:

demo.np1 <- npreg(y ~ x + z, data=demo.df)</pre>

The variables on the right-hand side of the formula are the predictors; we use + to separate them. Kernel regression will automatically include interactions between all variables, so there is no special notation for interactions. Similarly, there is no point in either including or excluding intercepts. If we wanted to transform either a predictor variable or the response, as in 1m, we can do so. Run like this, npreg will try to determine the best bandwidths for the predictor variables, based on a sophisticated combination of cross-validation and optimization.

Let's look at the output of npreg:

The main things here are the bandwidths. We also see the root mean squared error on the training data. Note that this is the in-sample root MSE; if we wanted the in-sample MSE, we could do

> demo.np1\$MSE
[1] 14.68851

(You can check that this is the square of the residual standard error above.) If we want the cross-validated MSE used to pick the bandwidths, that's



plot(demo.np1,theta=40,view="fixed")

Figure 4.14: Plot of the kernel regression with just two predictor variables. (See help(npplot) for plotting options.

```
> demo.np1$bws$fval
[1] 25.41536
```

The fitted and residuals functions work on these objects just like they do in lm objects, while the coefficients and confint functions do not. (Why?)

The predict function also works like it does for 1m, expecting a data frame containing columns whose names match those in the formula used to fit the model:

```
> predict(demo.np1, newdata=data.frame(x=-1,z=5))
[1] 21.45737
```

With two predictor variables, there is a nice three-dimensional default plot (Figure 4.14).

Kernel functions can also be defined for categorical and ordered variables. These can be included in the formula by wrapping the variable in factor() or ordered(), respectively:

```
19:55 Thursday 22<sup>nd</sup> January, 2015
```

No. Unordered Categorical Explanatory Vars.: 1

Again, there's no point, or need, to indicate interactions. Including the extra variable, not surprisingly, improves the cross-validated MSE:

> demo.np3\$bws\$fval
[1] 6.535231

With three or more predictor variables, we'd need a four-dimensional plot, which is hard. Instead, the default is to plot what happens as we sweep one variable with the others held fixed (by default, at their medians; see help(npplot) for changing that), as in Figure 4.15. We get something parabola-ish as we sweep X (which is right), and something near a step function as we sweep Z (which is right when W = B), so we're not doing badly for estimating a fairly complicated function of three variables with only 100 samples. We could also try fixing W at one value or another and making a perspective plot — Figure 4.16.

The default optimization of bandwidths is *extremely* aggressive. It keeps adjusting the bandwidths until the changes in the cross-validated MSE are very small, or the changes in the bandwidths themselves are very small. The "tolerances" for what count as "very small" are controlled by arguments to npreg called tol (for the bandwidths) and ftol (for the MSE), which default to about 10^{-8} and 10^{-7} , respectively. With a lot of data, or a lot of variables, this gets *extremely* slow. One can often make npreg run much faster, with no real loss of accuracy, by adjusting these options. A decent rule of thumb is to start with tol and ftol both at 0.01. One can use the bandwidth found by this initial coarse search to start a more refined one, as follows:

```
> bigdemo.df <- make.demo.df(1e3)
> system.time(demo.np4 <- npreg(y~x+z+factor(w),data=bigdemo.df,tol=0.01,
    ftol=0.01))
    user system elapsed
75.694 0.331 77.144</pre>
```

In words, it took 77 seconds altogether for R to run npreg, mostly of which was spent on our job, with just a little on background system tasks. The result of the run is stored in demo.np4:



plot(demo.np3)

Figure 4.15: Predictions of demo.np3 as each variable is swept over its range, with the others held at their medians.

The bandwidths have all shrunk (as they should), and the cross-validated MSE is also much smaller (0.87 versus 6.5 before). Figure 4.16 shows the estimated regression surfaces for both values of the categorical variable.

The package also contains a function, npregbw, which takes a formula and a data frame, and just optimizes the bandwidth. This is called automatically by npreg, and many of the relevant option are documented in its help page. One can also use the output of npregbw as an argument to npreg, in place of a formula.

4.7 Further Reading

Simonoff (1996) is a good practical introduction to kernel smoothing and related methods. Wasserman (2006) provides more theory. Li and Racine (2007) is a detailed treatment of nonparametric methods for econometric problems, overwhelming focused on kernel regression and kernel density estimation (which we'll get to in Chapter 16); Racine (2008) summarizes.

Kernel regression was introduced, independently, by Nadaraya (1964) and Watson (1964); both were inspired by kernel density estimation.

4.8 Exercises

1. Suppose we use a uniform ("boxcar") kernel extending over the region (-h/2, h/2). Show that

$$\mathbf{E}\left[\hat{r}(0)\right] = \mathbf{E}\left[r(X)\left|-\frac{h}{2} < X < \frac{h}{2}\right]$$
(4.33)

$$= r(0) + r'(0) \mathbf{E} \left[X \left| -\frac{h}{2} < X < \frac{h}{2} \right]$$

$$+ \frac{r''(0)}{2} \mathbf{E} \left[X^2 \left| -\frac{h}{2} < X < \frac{h}{2} \right] + o(h^2)$$
(4.34)

Show that $\mathbb{E}\left[X \left| -\frac{h}{2} < X < \frac{h}{2}\right.\right] = O(r'(0)f'(0)h^2)$, and that $\mathbb{E}\left[X^2 \left| -\frac{h}{2} < X < \frac{h}{2}\right.\right] = O(h^2)$. Conclude that the over-all bias is $O(h^2)$.

- 2. Use Eqs. 4.21, 4.17 and 4.16 to show that the excess risk of the kernel smoothing, when the bandwidth is selected by cross-validation, is also $O(n^{-4/5})$.
- 3. Generate 1000 data points where X is uniformly distributed between -4 and 4, and $Y = e^{7x}/(1 + e^{7x}) + \epsilon$, with ϵ Gaussian and with variance 0.01. Use non-parametric regression to estimate $\hat{r}(x)$, and then use Eq. 4.31 to find the average predictive comparison. Now re-run the simulation with X uniform on the interval [0,0.5] and re-calculate the average predictive comparison. What happened?



```
W=B
```



```
x.seq <- seq(from=-1,to=1,length.out=50)
z.seq <- seq(from=-30,to=30,length.out=50)
grid.A <- expand.grid(x=x.seq,z=z.seq,w="A")
grid.B <- expand.grid(x=x.seq,z=z.seq,w="B")
yhat.A <- predict(demo.np4,newdata=grid.A)
yhat.B <- predict(demo.np4,newdata=grid.B)
par(mfrow=c(1,2))
persp(x=x.seq,y=z.seq,z=matrix(yhat.A,nrow=50),theta=40,main="W=A",
xlab="x",ylab="z",zlab="y",ticktype="detailed")
persp(x=x.seq,y=z.seq,z=matrix(yhat.B,nrow=50),theta=40,main="W=B",
xlab="x",ylab="z",zlab="y",ticktype="detailed")
```

Figure 4.16: The regression surfaces learned for the demo function at the two different values of the categorical variable. Note that holding z fixed, we always see a parabolic shape as we move along x (as we should), while whether we see a line or something close to a step function at constant x depends on w, as it should.