
Additive Models

8.1 Additive Models

The **additive model** for regression is that the conditional expectation function is a sum of **partial response** functions, one for each predictor variable. Formally, when the vector \vec{X} of predictor variables has p dimensions, x_1, \dots, x_p , the model says that

$$\mathbb{E}[Y|\vec{X} = \vec{x}] = \alpha + \sum_{j=1}^p f_j(x_j) \quad (8.1)$$

This includes the linear model as a special case, where $f_j(x_j) = \beta_j x_j$, but it's clearly more general, because the f_j s can be arbitrary nonlinear functions. The idea is still that each input feature makes a separate contribution to the response, and these just add up (hence "partial response function"), but these contributions don't have to be strictly proportional to the inputs. We do need to add a restriction to make it identifiable; without loss of generality, say that $\mathbb{E}[Y] = \alpha$ and $\mathbb{E}[f_j(X_j)] = 0$.¹

Additive models keep a lot of the nice properties of linear models, but are more flexible. One of the nice things about linear models is that they are fairly straightforward to interpret: if you want to know how the prediction changes as you change x_j , you just need to know β_j . The partial response function f_j plays the same role in an additive model: of course the change in prediction from changing x_j will generally depend on the level x_j had before perturbation, but since that's also true of reality that's really a feature rather than a bug. It's true that a set of plots for f_j s takes more room than a table of β_j s, but it's also nicer to look at, conveys more information, and imposes fewer systematic distortions on the data.

Of course, none of this would be of any use if we couldn't actually estimate these models, but we can, through a clever computational trick which is worth knowing for its own sake. The use of the trick is also something they share with linear models, so we'll start there.

¹ To see why we need to do this, imagine the simple case where $p = 2$. If we add constants c_1 to f_1 and c_2 to f_2 , but subtract $c_1 + c_2$ from α , then nothing *observable* has changed about the model. This degeneracy or lack of identifiability is a little like the way collinearity keeps us from defining true slopes in linear regression. But it's less harmful than collinearity because we can fix it with this convention.

[[TODO:
 Re-
 organize:
 bring curse
 of dimen-
 sionality
 up, then
 additive
 models
 as com-

8.2 Partial Residuals and Back-fitting

8.2.1 Back-fitting for Linear Models

The general form of a linear regression model is

$$\mathbb{E}[Y|\vec{X} = \vec{x}] = \beta_0 + \vec{\beta} \cdot \vec{x} = \sum_{j=0}^p \beta_j x_j \quad (8.2)$$

where x_0 is always the constant 1. (Adding this fictitious constant variable lets us handle the intercept just like any other regression coefficient.)

Suppose we don't condition on all of \vec{X} but just one component of it, say X_k . What is the conditional expectation of Y ?

$$\mathbb{E}[Y|X_k = x_k] = \mathbb{E}[\mathbb{E}[Y|X_1, X_2, \dots, X_k, \dots, X_p] | X_k = x_k] \quad (8.3)$$

$$= \mathbb{E}\left[\sum_{j=0}^p \beta_j X_j | X_k = x_k\right] \quad (8.4)$$

$$= \beta_k x_k + \mathbb{E}\left[\sum_{j \neq k} \beta_j X_j | X_k = x_k\right] \quad (8.5)$$

where the first line uses the law of total expectation², and the second line uses Eq. 8.2. Turned around,

$$\beta_k x_k = \mathbb{E}[Y|X_k = x_k] - \mathbb{E}\left[\sum_{j \neq k} \beta_j X_j | X_k = x_k\right] \quad (8.6)$$

$$= \mathbb{E}\left[Y - \left(\sum_{j \neq k} \beta_j X_j\right) | X_k = x_k\right] \quad (8.7)$$

The expression in the expectation is the k^{th} **partial residual** — the (total) residual is the difference between Y and its expectation, the partial residual is the difference between Y and what we expect it to be *ignoring* the contribution from X_k . Let's introduce a symbol for this, say $Y^{(k)}$.

$$\beta_k x_k = \mathbb{E}[Y^{(k)}|X_k = x_k] \quad (8.8)$$

In words, if the over-all model is linear, then the partial residuals are linear. And notice that X_k is the only input feature appearing here — if we could somehow get hold of the partial residuals, then we can find β_k by doing a simple regression, rather than a multiple regression. Of course to get the partial residual we need to know all the other β_j s...

This suggests the following estimation scheme for linear models, known as the **Gauss-Seidel algorithm**, or more commonly and transparently as **back-fitting**; the pseudo-code is in Example 17.

“You say This is an iterative approximation algorithm. Initially, we look at how far each
 ‘vicious circle’, I ² As you learned in baby prob., this is the fact that $\mathbb{E}[Y|X] = \mathbb{E}[\mathbb{E}[Y|X, Z]|X]$ — that we can always
 say ‘it- condition more variables, provided we then average over those extra variables when we’re done.
 erative improvement’.”

Given: $n \times (p + 1)$ inputs \mathbf{x} (0th column all 1s)
 $n \times 1$ responses \mathbf{y}
small tolerance $\delta > 0$
center \mathbf{y} and each column of \mathbf{x}

$$\hat{\beta}_j \leftarrow 0 \text{ for } j \in 1 : p$$

until (all $|\hat{\beta}_j - \gamma_j| \leq \delta$) {
for $k \in 1 : p$ {
 $y_i^{(k)} = y_i - \sum_{j \neq k} \hat{\beta}_j x_{ij}$
 $\gamma_k \leftarrow$ regression coefficient of $y^{(k)}$ on $x_{\cdot k}$
 $\hat{\beta}_k \leftarrow \gamma_k$
}
}
 $\hat{\beta}_0 \leftarrow (n^{-1} \sum_{i=1}^n y_i) - \sum_{j=1}^p \hat{\beta}_j n^{-1} \sum_{i=1}^n x_{ij}$
Return: $(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$

CODE EXAMPLE 17: *Pseudocode for back-fitting linear models. Assume we make at least one pass through the **until** loop. Recall from Chapter 1 that centering the data does not change the β_j s; this way the intercept only has to be calculated once, at the end. [[ATTN: Fix horizontal lines]]*

point is from the global mean, and do a simple regression of those deviations on the first input variable. This then gives us a better idea of what the regression surface really is, and we use the deviations from *that* surface in a simple regression on the next variable; this should catch relations between Y and X_2 that weren't already caught by regressing on X_1 . We then go on to the next variable in turn. At each step, each coefficient is adjusted to fit in with what we have already guessed about the other coefficients — that's why it's called "back-fitting". It is not obvious³ that this will ever converge, but it (generally) does, and the fixed point on which it converges is the usual least-squares estimate of β .

Back-fitting is rarely used to fit linear models these days, because with modern computers and numerical linear algebra it's faster to just calculate $(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}$. But the cute thing about back-fitting is that it doesn't actually rely on *linearity*.

8.2.2 Backfitting Additive Models

Defining the partial residuals by analogy with the linear case, as

$$Y^{(k)} = Y - \left(\alpha + \sum_{j \neq k} f_j(x_j) \right) \quad (8.9)$$

a little algebra along the lines of §8.2.1 shows that

$$\mathbb{E}[Y^{(k)} | X_k = x_k] = f_k(x_k) \quad (8.10)$$

³ Unless, I suppose, you're Gauss.

Given: $n \times p$ inputs \mathbf{x} $n \times 1$ responses \mathbf{y} small tolerance $\delta > 0$ one-dimensional smoother \mathcal{S} $\hat{\alpha} \leftarrow n^{-1} \sum_{i=1}^n y_i$ $\hat{f}_j \leftarrow 0$ for $j \in 1 : p$ until (all $ \hat{f}_j - g_j \leq \delta$) { for $k \in 1 : p$ { $y_i^{(k)} = y_i - \sum_{j \neq k} \hat{f}_j(x_{ij})$ $g_k \leftarrow \mathcal{S}(y^{(k)}) \sim x_{\cdot k}$ $g_k \leftarrow g_k - n^{-1} \sum_{i=1}^n g_k(x_{ik})$ $\hat{f}_k \leftarrow g_k$ } } Return: $(\hat{\alpha}, \hat{f}_1, \dots, \hat{f}_p)$	
--	--

CODE EXAMPLE 18: *Pseudo-code for back-fitting additive models. Notice the extra step, as compared to back-fitting linear models, which keeps each partial response function centered.*

If we knew how to estimate arbitrary one-dimensional regressions, we could now use back-fitting to estimate additive models. But we have spent a lot of time learning how to use smoothers to fit one-dimensional regressions! We could use nearest neighbors, or splines, or kernels, or local-linear regression, or anything else we feel like substituting here.

Our new, improved back-fitting algorithm in Example 18. Once again, while it's not obvious that this converges, it does. Also, the back-fitting procedure works well with some complications or refinements of the additive model. If we know the function form of one or another of the f_j , we can fit those parametrically (rather than with the smoother) at the appropriate points in the loop. (This would be a **semiparametric** model.) If we think that there is an interaction between x_j and x_k , rather than their making separate additive contributions for each variable, we can smooth them together; etc.

There are actually *two* packages standard packages for fitting additive models in R: **gam** and **mcmc**. Both have commands called **gam**, which fit **generalized** additive models — the generalization is to use the additive model for things like the probabilities of categorical responses, rather than the response variable itself. If that sounds obscure right now, don't worry — we'll come back to this in Chapters 11–12 after we've looked at generalized linear models. §8.4 below illustrates using one of these packages to fit an additive model.

8.3 The Curse of Dimensionality

Before illustrating how additive models work in practice, let's talk about why we'd want to use them. So far, we have looked at two extremes for regression models; additive models are somewhere in between.

On the one hand, we had linear regression, which is a parametric method (with $p+1$ parameters). Its weakness is that the true regression function μ is hardly ever linear, so even with infinite data linear regression will always make systematic mistakes in its predictions — there's always some approximation bias, bigger or smaller depending on how non-linear μ is. The strength of linear regression is that it converges very quickly as we get more data. Generally speaking,

$$MSE_{\text{linear}} = \sigma^2 + a_{\text{linear}} + O(n^{-1}) \quad (8.11)$$

where the first term is the intrinsic noise around the true regression function, the second term is the (squared) approximation bias, and the last term is the estimation variance. Notice that the rate at which the estimation variance shrinks doesn't depend on p — factors like that are all absorbed into the big O .⁴ Other parametric models generally converge at the same rate.

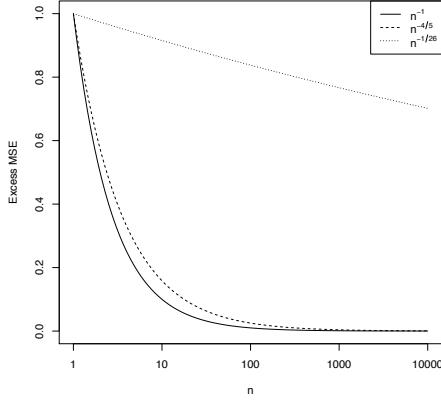
At the other extreme, we've seen a number of completely nonparametric regression methods, such as kernel regression, local polynomials, k -nearest neighbors, etc. Here the limiting approximation bias is actually *zero*, at least for any reasonable regression function μ . The problem is that they converge more slowly, because we need to use the data not just to figure out the coefficients of a parametric model, but the sheer shape of the regression function. We saw in Chapter 4 that the mean-squared error of kernel regression in one dimension is $\sigma^2 + O(n^{-4/5})$. Splines, k -nearest-neighbors (with growing k), etc., all attain the same rate. But in p dimensions, this becomes (Wasserman, 2006, §5.12)

$$MSE_{\text{nonpara}} - \sigma^2 = O(n^{-4/(p+4)}) \quad (8.12)$$

There's no ultimate approximation bias term here. Why does the rate depend on p ? Well, to hand-wave a bit, think of kernel smoothing, where $\hat{\mu}(\vec{x})$ is an average over y_i for \vec{x}_i near \vec{x} . In a p dimensional space, the volume within ϵ of \vec{x} is $O(\epsilon^p)$, so the probability that a training point \vec{x}_i falls in the averaging region around \vec{x} gets exponentially smaller as p grows. Turned around, to get the same number of training points per \vec{x} , we need exponentially larger sample sizes. The appearance of the 4s is a little more mysterious, but can be resolved from an error analysis of the kind we did for kernel regression in Chapter 4⁵. This slow rate isn't just

⁴ See Appendix C if you are not familiar with “big O ” notation.

⁵ Remember that in one dimension, the bias of a kernel smoother with bandwidth h is $O(h^2)$, and the variance is $O(1/nh)$, because only samples falling in an interval about h across contribute to the prediction at any one point, and when h is small, the number of such samples is proportional to nh . Adding bias squared to variance gives an error of $O(h^4) + O((nh)^{-1})$, solving for the best bandwidth gives $h_{\text{opt}} = O(n^{-1/5})$, and the total error is then $O(n^{-4/5})$. Suppose for the moment that in p dimensions we use the same bandwidth along each dimension. (We get the same end result with more work if we let each dimension have its own bandwidth.) The bias is still $O(h^2)$, because the Taylor expansion still goes through. But now only samples falling into a region of volume $O(h^p)$



```

curve(x^(-1),from=1,to=1e4,log="x",xlab="n",ylab="Excess MSE")
curve(x^(-4/5),add=TRUE,lty="dashed")
curve(x^(-1/26),add=TRUE,lty="dotted")
legend("topright",legend=c(expression(n^{-1}),
  expression(n^{-4/5}),expression(n^{-1/26})),
  lty=c("solid","dashed","dotted"))

```

Figure 8.1 Schematic of rates of convergence of MSEs for parametric models ($O(n^{-1})$), one-dimensional nonparametric regressions or additive models ($O(n^{-4/5})$), and a 100-dimensional nonparametric regression ($O(n^{-1/26})$). Note that the horizontal but not the vertical axis is on a logarithmic scale.

a weakness of kernel smoothers, but turns out to be the best any nonparametric estimator can do.

For $p = 1$, the nonparametric rate is $O(n^{-4/5})$, which is of course slower than $O(n^{-1})$, but not all that much, and the improved bias usually more than makes up for it. But as p grows, the nonparametric rate gets slower and slower, and the fully nonparametric estimate more and more imprecise, yielding the infamous **curse of dimensionality**. For $p = 100$, say, we get a rate of $O(n^{-1/26})$, which is not very good at all. (See Figure 8.1.) Said another way, to get the same precision with p inputs that n data points gives us with one input takes $n^{(4+p)/5}$ data points. For $p = 100$, this is $n^{20.8}$, which tells us that matching the error of $n = 100$ one-dimensional observations requires $O(4 \times 10^{41})$ hundred-dimensional observations.

So completely unstructured nonparametric regressions won't work very well in high dimensions, at least not with plausible amounts of data. The trouble is that

around x contribute to the prediction at x , so the variance is $O((nh^p)^{-1})$. The best bandwidth is now $h_{\text{opt}} = O(n^{-1/(p+4)})$, yielding an error of $O(n^{-4/(p+4)})$ as promised.

there are just *too many* possible high-dimensional functions, and seeing only a trillion points from the function doesn't pin down its shape very well at all. [[ATTN:

This is where additive models come in. Not every regression function is additive, More so they have, even asymptotically, some approximation bias. But we can estimate mathematically each f_j by a simple one-dimensional smoothing, which converges at $O(n^{-4/5})$, matical almost as good as the parametric rate. So overall

$$MSE_{\text{additive}} - \sigma^2 = a_{\text{additive}} + O(n^{-4/5}) \quad (8.13)$$

explanation in appendix?

Since linear models are a sub-class of additive models, $a_{\text{additive}} \leq a_{\text{lm}}$. From a purely predictive point of view, the only time to prefer linear models to additive models is when n is so small that $O(n^{-4/5}) - O(n^{-1})$ exceeds this difference in approximation biases; eventually the additive model will be more accurate.⁶

8.4 Example: California House Prices Revisited

As an example, we'll look at data on median house prices across Census tracts from the data-analysis assignment in §A.13. This has both California and Pennsylvania, but it's hard to visually see patterns with both states; I'll do California, and let you replicate this all on Pennsylvania, and even on the combined data.

Start with getting the data:

```
housing <- read.csv("http://www.stat.cmu.edu/~cshalizi/ADAFaEPoV/data/calif_penn_2011.csv")
housing <- na.omit(housing)
calif <- housing[housing$STATEFP == 6, ]
```

(How do I know that the STATEFP code of 6 corresponds to California?)

We'll fit a linear model for the log price, on the thought that it makes some sense for the factors which raise or lower house values to multiply together, rather than just adding.

```
calif.lm <- lm(log(Median_house_value) ~ Median_household_income + Mean_household_income +
  POPULATION + Total_units + Vacant_units + Owners + Median_rooms + Mean_household_size_owners +
  Mean_household_size_renters + LATITUDE + LONGITUDE, data = calif)
```

This is very fast — about a fifth of a second on my laptop.

Here are the summary statistics⁷:

```
print(summary(calif.lm), signif.stars = FALSE, digits = 3)
##
## Call:
## lm(formula = log(Median_house_value) ~ Median_household_income +
##     Mean_household_income + POPULATION + Total_units + Vacant_units +
##     Owners + Median_rooms + Mean_household_size_owners + Mean_household_size_renters +
##     LATITUDE + LONGITUDE, data = calif)
##
```

⁶ Unless the best additive approximation to μ is linear; then the linear model has no more bias and less variance.

⁷ I have suppressed the usual stars on “significant” regression coefficients, because, as discussed in Chapter ??, those aren't really the most important variables, and I have reined in R's tendency to use far too many decimal places.

```
predlims <- function(preds, sigma) {
  prediction.sd <- sqrt(preds$e.fit^2 + sigma^2)
  upper <- preds$fit + 2 * prediction.sd
  lower <- preds$fit - 2 * prediction.sd
  lims <- cbind(lower = lower, upper = upper)
  return(lims)
}
```

CODE EXAMPLE 19: *Calculating quick-and-dirty prediction limits from a prediction object (`preds`) containing fitted values and their standard errors, plus an estimate of the noise level. Because those are two (presumably uncorrelated) sources of noise, we combine the standard deviations by “adding in quadrature”.*

```
## Residuals:
##      Min    1Q Median     3Q    Max
## -3.855 -0.153  0.034  0.189  1.214
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)
## (Intercept)              -5.74e+00  5.28e-01 -10.86 < 2e-16
## Median_household_income 1.34e-06  4.63e-07   2.90  0.0038
## Mean_household_income   1.07e-05  3.88e-07  27.71 < 2e-16
## POPULATION               -4.15e-05 5.03e-06 -8.27 < 2e-16
## Total_units                8.37e-05 1.55e-05   5.41  6.4e-08
## Vacant_units               8.37e-07 2.37e-05   0.04  0.9719
## Owners                   -3.98e-03 3.21e-04 -12.41 < 2e-16
## Median_rooms              -1.62e-02 8.37e-03  -1.94  0.0525
## Mean_household_size_owners 5.60e-02 7.16e-03   7.83  5.8e-15
## Mean_household_size_renters -7.47e-02 6.38e-03 -11.71 < 2e-16
## LATITUDE                  -2.14e-01 5.66e-03 -37.76 < 2e-16
## LONGITUDE                 -2.15e-01 5.94e-03 -36.15 < 2e-16
##
## Residual standard error: 0.317 on 7469 degrees of freedom
## Multiple R-squared:  0.639, Adjusted R-squared:  0.638
## F-statistic: 1.2e+03 on 11 and 7469 DF,  p-value: <2e-16
```

Figure 8.2 plots the predicted prices, ± 2 standard errors, against the actual prices. The predictions are not all that *accurate* — the RMS residual is 0.317 on the log scale (i.e., 37% on the original scale), but they do have pretty reasonable coverage; about 96% of actual prices fall within the prediction limits⁸. On the other hand, the predictions *are* quite precise, with the median of the calculated

⁸ Remember from your linear regression class that there are two kinds of confidence intervals we might want to use for prediction. One is a confidence interval for the *conditional mean* at a given value of x ; the other is a confidence interval for the *realized values of Y* at a given x . Earlier examples have emphasized the former, but since we don’t know the true conditional means here, we need to use the latter sort of intervals, prediction intervals proper, to evaluate coverage. The `predlims` function in Code Example 19 calculates a rough prediction interval by taking the standard error of the conditional mean, combining it with the estimated standard deviation, and multiplying by 2. Strictly speaking, we ought to worry about using a t -distribution rather than a Gaussian here, but with 7469 residual degrees of freedom, this isn’t going to matter much. (Assuming Gaussian noise is likely to be more of a concern, but this is only meant to be a rough cut anyway.)

standard errors being 0.011 on the log scale (i.e., 1.1% in dollars). This linear model *thinks* it knows what's going on.

Next, we'll fit an additive model, using the `gam` function from the `mgcv` package; this automatically sets the bandwidths using a fast approximation to leave-one-out CV called **generalized cross-validation**, or **GCV** (§3.4.3).

```
system.time(calif.gam <- gam(log(Median_house_value) ~ s(Median_household_income) +
  s(Mean_household_income) + s(POPULATION) + s(Total_units) + s(Vacant_units) +
  s(Owners) + s(Median_rooms) + s(Mean_household_size_owners) + s(Mean_household_size_renters) +
  s(LATITUDE) + s(LONGITUDE), data = calif))
##    user    system   elapsed
##  3.452   0.144   3.614
```

(That is, it took about five seconds total to run this.) The `s()` terms in the `gam` formula indicate which terms are to be smoothed — if we wanted particular parametric forms for some variables, we could do that as well. (Unfortunately we can't just write `MedianHouseValue ~ s(.)`, we have to list all the variables on the right-hand side.⁹) The smoothing here is done by splines (hence `s()`), and there are lots of options for controlling the splines, or replacing them by other smoothers, if you know what you're doing.

Figure 8.3 compares the predicted to the actual responses. The RMS error has improved (0.27 on the log scale, or 130%, with 96% of observations falling with ± 2 standard errors of their fitted values), at only a fairly modest cost in the claimed precision (the median standard error of prediction is 0.02, or 2.1%). Figure 8.4 shows the partial response functions.

It makes little sense to have latitude and longitude make separate additive contributions here; presumably they interact. We can just smooth them together¹⁰:

```
calif.gam2 <- gam(log(Median_house_value) ~ s(Median_household_income) + 
  s(Mean_household_income) + s(POPULATION) + s(Total_units) + s(Vacant_units) + s(Owners) + s(Median_rooms) + 
  s(Mean_household_size_owners) + s(Mean_household_size_renters) + s(LONGITUDE, 
  LATITUDE), data = calif)
```

This gives an RMS error of ± 0.25 (log-scale) and 96% coverage, with a median standard error of 0.021, so accuracy is improving (at least in sample), with little loss of precision.

Figures 8.6 and 8.7 show two different views of the joint smoothing of longitude and latitude. In the perspective plot, it's quite clear that price increases specifically towards the coast, and even more specifically towards the great coastal cities. In the contour plot, one sees more clearly an inward bulge of a negative, but not too very negative, contour line (between -122 and -120 longitude) which embraces Napa, Sacramento, and some related areas, which are comparatively more developed and more expensive than the rest of central California, and so

⁹ Alternately, we could use Kevin Gilbert's `formulaTools` functions — see
<https://gist.github.com/kgilbert-cmu>.

¹⁰ If the two variables which interact have very different magnitudes, it's better to smooth them with a `te()` term than an `s()` term, but here they are comparable. See §8.5 for more, and
`help(gam.models)`.

```

graymapper <- function(z, x = calif$LONGITUDE, y = calif$LATITUDE, n.levels = 10,
  breaks = NULL, break.by = "length", legend.loc = "topright", digits = 3,
  ...) {
  my.greys = grey(((n.levels - 1):0)/n.levels)
  if (!is.null(breaks)) {
    stopifnot(length(breaks) == (n.levels + 1))
  }
  else {
    if (identical(break.by, "length")) {
      breaks = seq(from = min(z), to = max(z), length.out = n.levels +
        1)
    }
    else {
      breaks = quantile(z, probs = seq(0, 1, length.out = n.levels + 1))
    }
  }
  z = cut(z, breaks, include.lowest = TRUE)
  colors = my.greys[z]
  plot(x, y, col = colors, bg = colors, ...)
  if (!is.null(legend.loc)) {
    breaks.printable <- signif(breaks[1:n.levels], digits)
    legend(legend.loc, legend = breaks.printable, fill = my.greys)
  }
  invisible(breaks)
}

```

CODE EXAMPLE 20: *Map-making code.* In its basic use, this takes vectors for x and y coordinates, and draws gray points whose color depends on a third vector for z , with darker points indicating higher values of z . Options allow for the control of the number of gray levels, setting the breaks between levels automatically, and using a legend. Returning the break-points makes it easier to use the same scale in multiple maps. See online for commented code.

more expensive than one would expect based on their distance from the coast and San Francisco.

If you worked through problem set A.13, you will recall that one of the big things wrong with the linear model is that its errors (the residuals) are highly structured and very far from random. In essence, it totally missed the existence of cities, and the fact that houses cost more in cities (because land costs more there). It's a good idea, therefore, to make some maps, showing the actual values, and then, by way of contrast, the residuals of the models. Rather than do the plotting by hand over and over, let's write a function (Code Example 20).

Figures 8.8 and 8.9 show that allowing for the interaction of latitude and longitude (the smoothing term plotted in Figures 8.6–8.7) leads to a much more *random* and less systematic clumping of residuals. This is desirable in itself, even if it does little to improve the mean prediction error. Essentially, what that smoothing term is doing is picking out the existence of California's urban regions, and their distinction from the rural background. Examining the plots of the interaction term should suggest to you how inadequate it would be to just put in a $\text{LONGITUDE} \times \text{LATITUDE}$ term in a linear model.

Including an interaction between latitude and longitude in a spatial problem is

pretty obvious. There are other potential interactions which might be important here — for instance, between the two measures of income, or between the total number of housing units available and the number of vacant units. We could, of course, just use a completely unrestricted nonparametric regression — going to the opposite extreme from the linear model. In addition to the possible curse-of-dimensionality issues, however, getting something like `npreg` to run with 7000 data points and 11 predictor variables requires a lot of patience. Other techniques, like nearest neighbor regression (§1.5.1) or regression trees (Ch. 13), may run faster, though cross-validation can be demanding even there.

8.5 Interaction Terms and Expansions

One way to think about additive models, and about (possibly) including interaction terms, is to imagine doing a sort of Taylor series or power series expansion of the true regression function. The zero-th order expansion would be a constant:

$$\mu(x) \approx \alpha \quad (8.14)$$

The best constant to use here would just be $\mathbb{E}[Y]$. (“Best” here is in the mean-square sense, as usual.) A purely additive model would correspond to a first-order expansion:

$$\mu(x) \approx \alpha + \sum_{j=1}^p f_j(x_j) \quad (8.15)$$

Two-way interactions come in when we go to a second-order expansion:

$$\mu(x) \approx \alpha + \sum_{j=1}^p f_j(x_j) + \sum_{j=1}^p \sum_{k=j+1}^p f_{jk}(x_j, x_k) \quad (8.16)$$

(Why do I limit k to run from $j+1$ to p ? rather than from 1 to p ?) We will, of course, insist that $\mathbb{E}[f_{jk}(X_j, X_k)] = 0$ for all j, k . If we want to estimate these terms in R, using `mrgcv`, we use the syntax `s(xj, xk)` or `te(xj, xk)`. The former fits a thin-plate spline over the (x_j, x_k) plane, and is appropriate when those variables are measured on similar scales, so that curvatures along each direction are comparable. The latter uses a tensor product of smoothing splines along each coordinate, and is more appropriate when the measurement scales are very different¹¹.

There is an important ambiguity here: for any j , with additive partial-response function f_j , I could take any of its interactions, set $f'_{jk}(x_j, x_k) = f_{jk}(x_j, x_k) + f_j(x_j)$ and $f'_j(x_j) = 0$, and get exactly the same predictions under all circumstances. This is the parallel to being able to add and subtract constants from the first-order functions, provided we made corresponding changes to the intercept term. We therefore need to similarly fix the two-way interaction functions.

A natural way to do this is to insist that the second-order f_{jk} function should

¹¹ For the distinction between thin-plate and tensor-product splines, see §7.4. If we want to interact a continuous variable x_j with a categorical x_k , `mrgcv`'s syntax is `s(xj, by=xk)` or `te(xj, by=xk)`.

be uncorrelated with (“orthogonal to”) the first-order functions f_j and f_k ; this is the analog to insisting that the first-order functions all have expectation zero. The f_{jk} s then represent purely interactive contributions to the response, which could not be captured by additive terms. If this is what we want to do, the best syntax to use in `mgcv` is `ti`, which specifically separates the first- and higher-order terms, e.g., `ti(xj) + ti(xk) + ti(xj, xk)` will estimate three functions, for the additive contributions and their interaction.

An alternative is to just *pick* a particular f_{jk} , and absorb f_j into it. The model then looks like

$$\mu(x) \approx \alpha + \sum_{j=1}^p \sum_{k=j+1}^p f_{jk}(x_j, x_k) \quad (8.17)$$

We can also mix these two approaches, if we specifically do not want additive or interactive terms for certain predictor variables. This is what I did above, where I estimated a single second-order smoothing term for both latitude and longitude, with no additive components for either.

Of course, there is nothing special about two-way interactions. If you’re curious about what a three-way term would be like, and you’re lucky enough to have data which amenable to fitting it, you could certainly try

$$\mu \approx \alpha + \sum_{j=1}^p f_j(x_j) + \sum_{j=1}^p \sum_{k=j+1}^p f_{jk}(x_j, x_k) + \sum_{j,k,l} f_{jkl}(x_j, x_k, x_l) \quad (8.18)$$

(How should the indices for the last term go?) More ambitious combinations are certainly possible, though they tend to become a confused mass of algebra and indices.

Geometric interpretation

It’s often convenient to think of the regression function as living in a big (infinite-dimensional) vector space of functions. Within this space, the *constant* functions form a linear sub-space¹², and we can ask for the projection of the true regression function on to that sub-space; this would be the best approximation¹³ to μ as a constant. This is, of course, the expectation value. The *additive* functions of all p variables also form a linear sub-space¹⁴, so the right-hand side of Eq. 8.15 is just the projection of μ on to that space, and so forth and so on. When we insist on having the higher-order interaction functions be uncorrelated with the additive functions, we’re taking the projection of μ on to the space of all functions *orthogonal* to the additive functions.

¹² Because if f and g are two constant functions, $af + bg$ is also a constant, for any real numbers a and b .

¹³ Remember that projecting a vector on to a linear sub-space finds the point in the sub-space closest to the original vector. This is equivalent to minimizing the (squared) bias.

¹⁴ By parallel reasoning to the previous footnote.

Selecting interactions

There are two issues with interaction terms. First, the curse of dimensionality returns: an order- q interaction term will converge at the rate $O(n^{-4/(4+q)})$, so they can dominate the over-all uncertainty. Second, there are lots of possible interactions ($\binom{p}{q}$, in fact), which can make it very demanding in time and data to fit them all, and hard to interpret. Just as with linear models, therefore, it can make a lot of sense to selectively examine interactions based on subject-matter knowledge, or residuals of additive models.

Varying-coefficient models

In some contexts, people like to use models of the form

$$\mu(x) = \alpha + \sum_{j=1}^p x_j f_j(x_{-j}) \quad (8.19)$$

where f_j is a function of the non- j predictor variables, or some subset of them. These **varying-coefficient** functions are obviously a subset of the usual class of additive models, but there are occasions where they have some scientific justification¹⁵. These are conveniently estimated in `mgcv` through the `by` option, e.g., `s(xk, by=xj)` will estimate a term of the form $x_j f(x_k)$.¹⁶

8.6 Closing Modeling Advice

With modern computing power, there are very few situations in which it is actually better to do linear regression than to fit an additive model. In fact, there seem to be only two good reasons to prefer linear models.

1. Our data analysis is guided by a credible scientific theory which asserts linear relationships *among the variables we measure* (not others, for which our observables serve as imperfect proxies).
2. Our data set is so massive that either the extra processing time, or the extra computer memory, needed to fit and store an additive rather than a linear model is prohibitive.

Even when the first reason applies, and we have good reasons to believe a linear theory, the truly scientific thing to do would be to *check* linearity, by fitting a flexible non-linear model and seeing if it looks close to linear. (We will see formal tests based on this idea in Chapter 9.) Even when the second reason applies, we would like to know how much bias we're introducing by using linear predictors, which we could do by randomly selecting a subset of the data which is small enough for us to manage, and fitting an additive model.

In the vast majority of cases when users of statistical software fit linear models, neither of these justifications applies: theory doesn't tell us to expect linearity,

¹⁵ They can also serve as a “transitional object” when giving up the use of purely linear models.

¹⁶ As we saw above, `by` does something slightly different when given a categorical variable. How are these two uses related?

and our machines don't compel us to use it. Linear regression is then employed for no better reason than that users know how to type `lm` but not `gam`. You now know better, and can spread the word.

8.7 Further Reading

Simon Wood, who wrote the `mgcv` package, has a nice book about additive models and their generalizations, Wood (2006); at this level it's your best source for further information. Buja *et al.* (1989) is a thorough theoretical treatment.

The expansions of §8.5 are sometimes called “functional analysis of variance” or “functional ANOVA”. Making those ideas precise requires exploring some of the geometry of infinite-dimensional spaces of functions (“Hilbert space”). See Wahba (1990) for a treatment of the statistical topic, and Halmos (1957) for a classic introduction to Hilbert spaces.

Historical notes

Ezekiel (1924) seems to be the first publication advocating the use of additive models as a general method, which he called “curvilinear multiple correlation”. His paper was complete with worked examples on simulated data (with known answers) and real data (from economics)¹⁷. He was explicit that any reasonable smoothing or regression technique could be used to find what we'd call the partial response functions. He also gave a successive-approximation algorithm for estimate the over-all model: start with an initial guess about all the partial responses; plot all the partial residuals; refine the partial responses simultaneously; repeat. This differs from back-fitting in that the partial response functions are updating in parallel within each cycle, not one after the other. This is a subtle difference, and Ezekiel's method will often work, but can run into trouble with correlated predictor variables, when back-fitting will not.

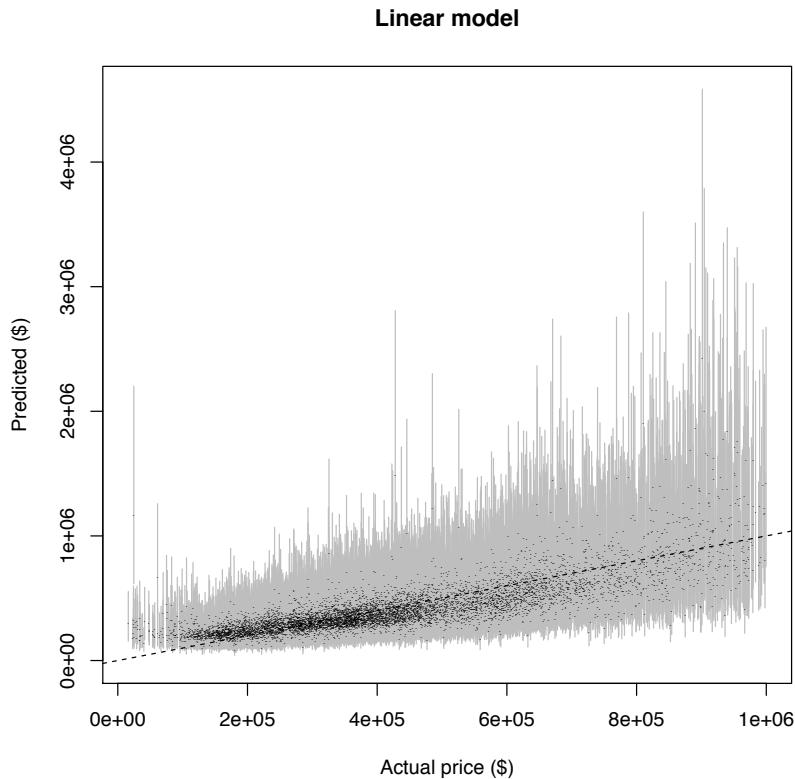
The Gauss-Seidel or backfitting algorithm was invented by Gauss in the early 1800s during his work on least squares estimation in linear models; he mentioned it in letters to students, described it as something one could do “while half asleep”, but never published it. Seidel gave the first published version in 1874. (For all this history, see Benzi 2009.) I am not sure when the connection was made between additive statistical models and back-fitting.

Exercises

- 8.1 Repeat the analyses of California housing prices with Pennsylvania housing prices. Which partial response functions might one reasonably hope would stay the same? Do they? (How can you tell?)

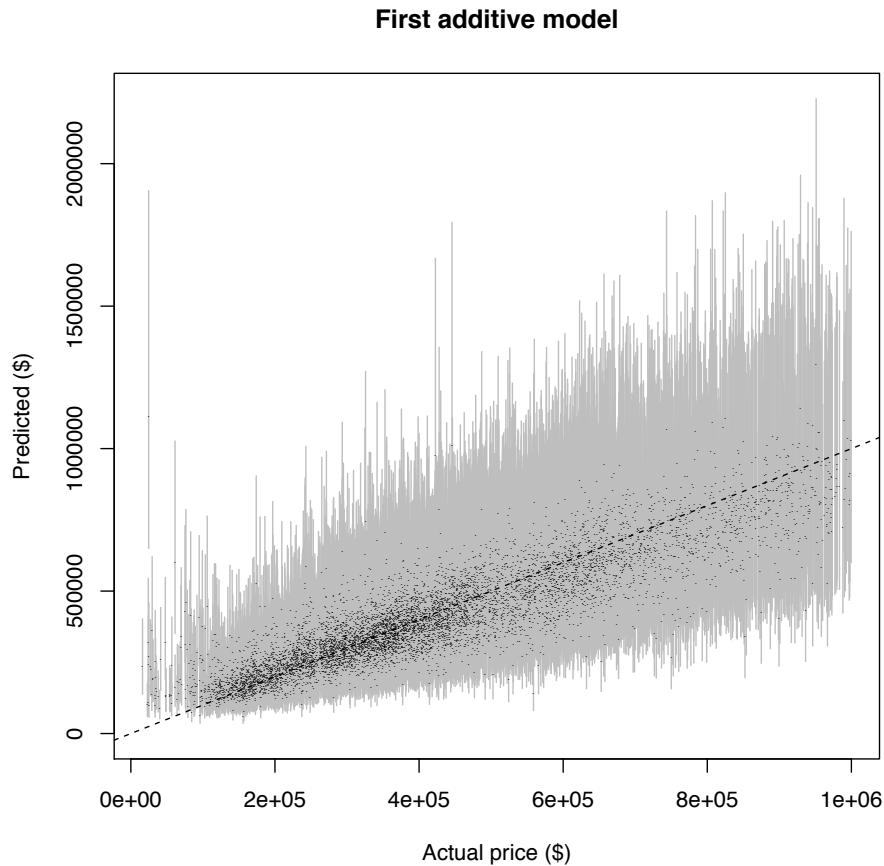
¹⁷ “Each of these curves illustrates and substantiates conclusions reached by theoretical economic analysis. Equally important, they provide definite quantitative statements of the relationships. The method of . . . curvilinear multiple correlation enable[s] us to use the favorite tool of the economist, *caeteris paribus*, in the analysis of actual happenings equally as well as in the intricacies of theoretical reasoning” (p. 453).

- 8.2 *Additive?* For general p , let $\|\vec{x}\|$ be the (ordinary, Euclidean) length of the vector \vec{x} . Is this an additive function of the (ordinary, Cartesian) coordinates? Is $\|\vec{x}\|^2$ an additive function? $\|\vec{x} - \vec{x}_0\|$ for a fixed \vec{x}_0 ? $\|\vec{x} - \vec{x}_0\|^2$?
- 8.3 *Additivity vs. parallelism*
1. Take any additive function f of p arguments x_1, x_2, \dots, x_p . Fix a coordinate index i and a real number c . Prove that $f(x_1, x_2, \dots, x_i, \dots, x_p) - f(x_1, x_2, \dots, x_i + c, \dots, x_p)$ depends only on x_i and c , and not on the other coordinates.
 2. Suppose $p = 2$, and continue to assume f is additive. Consider the curve formed by plotting $f(x_1, x_2)$ against x_1 for a fixed value of x_2 , and the curve formed by plotting $f(x_1, x_2)$ against x_1 with x_2 fixed at a different value, say x'_2 . Prove that the curves are parallel, i.e., that the vertical distance between them is constant.
 3. For general p and additive f , consider the surfaces formed by the f by varying all but one of the coordinates. Prove that these surfaces are always parallel to each other.
 4. Is the converse true? That is, do parallel regression surfaces imply an additive model?



```
plot(calif$Median_house_value, exp(preds.lm$fit), type = "n", xlab = "Actual price ($)",  
     ylab = "Predicted ($)", main = "Linear model", ylim = c(0, exp(max(predlims.lm))))  
segments(calif$Median_house_value, exp(predlims.lm[, "lower"]), calif$Median_house_value,  
        exp(predlims.lm[, "upper"])), col = "grey")  
abline(a = 0, b = 1, lty = "dashed")  
points(calif$Median_house_value, exp(preds.lm$fit), pch = 16, cex = 0.1)
```

Figure 8.2 Actual median house values (horizontal axis) versus those predicted by the linear model (black dots), plus or minus two *predictive* standard errors (grey bars). The dashed line shows where actual and predicted prices are equal. Here `predict` gives both a fitted value for each point, and a standard error for that prediction. (Without a `newdata` argument, `predict` defaults to the data used to estimate `calif.lm`, which here is what we want.) Predictions are exponentiated so they're comparable to the original values (and because it's easier to grasp dollars than log-dollars).



```

plot(calif$Median_house_value, exp(preds.gam$fit), type = "n", xlab = "Actual price ($)",
     ylab = "Predicted ($)", main = "First additive model", ylim = c(0, exp(max(predlims.gam))))
segments(calif$Median_house_value, exp(predlims.gam[, "lower"]), calif$Median_house_value,
        exp(predlims.gam[, "upper"]), col = "grey")
abline(a = 0, b = 1, lty = "dashed")
points(calif$Median_house_value, exp(preds.gam$fit), pch = 16, cex = 0.1)
  
```

Figure 8.3 Actual versus predicted prices for the additive model, as in Figure 8.2. Note that the `sig2` attribute of a model returned by `gam()` is the estimate of the noise variance around the regression surface (σ^2).

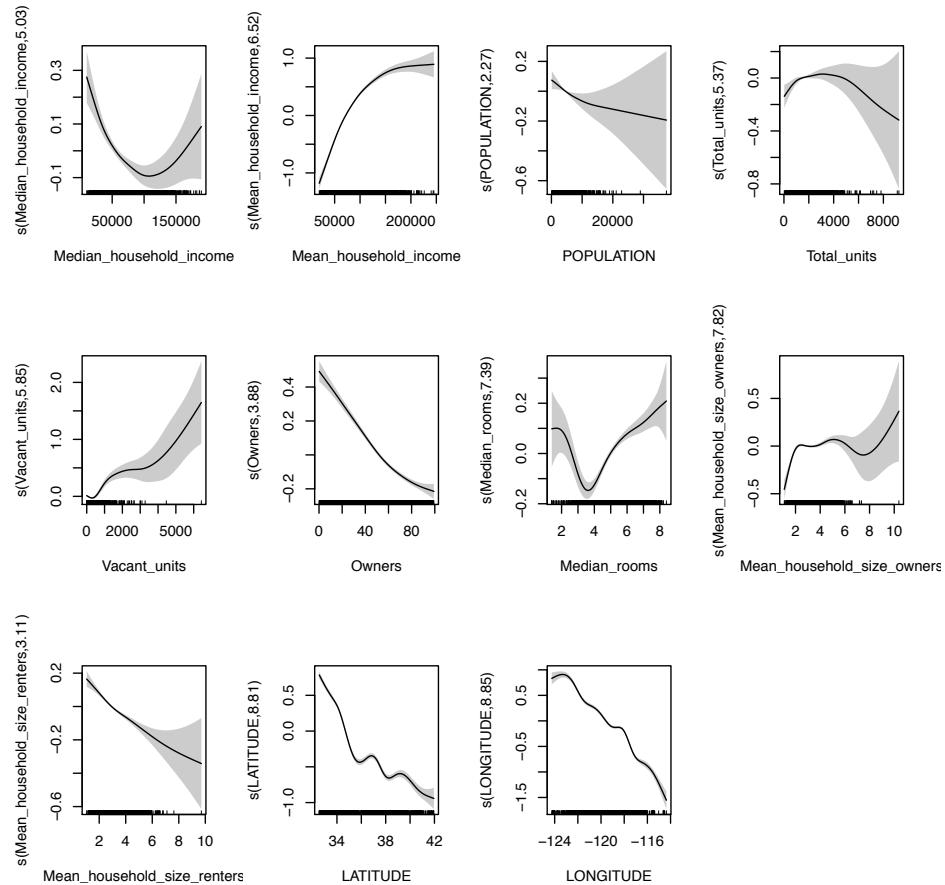


Figure 8.4 The estimated partial response functions for the additive model, with a shaded region showing ± 2 standard errors. The tick marks along the horizontal axis show the observed values of the input variables (a **rug plot**); note that the error bars are wider where there are fewer observations. Setting `pages=0` (the default) would produce eight separate plots, with the user prompted to cycle through them. Setting `scale=0` gives each plot its own vertical scale; the default is to force them to share the same one. Finally, note that here the vertical scales are logarithmic.

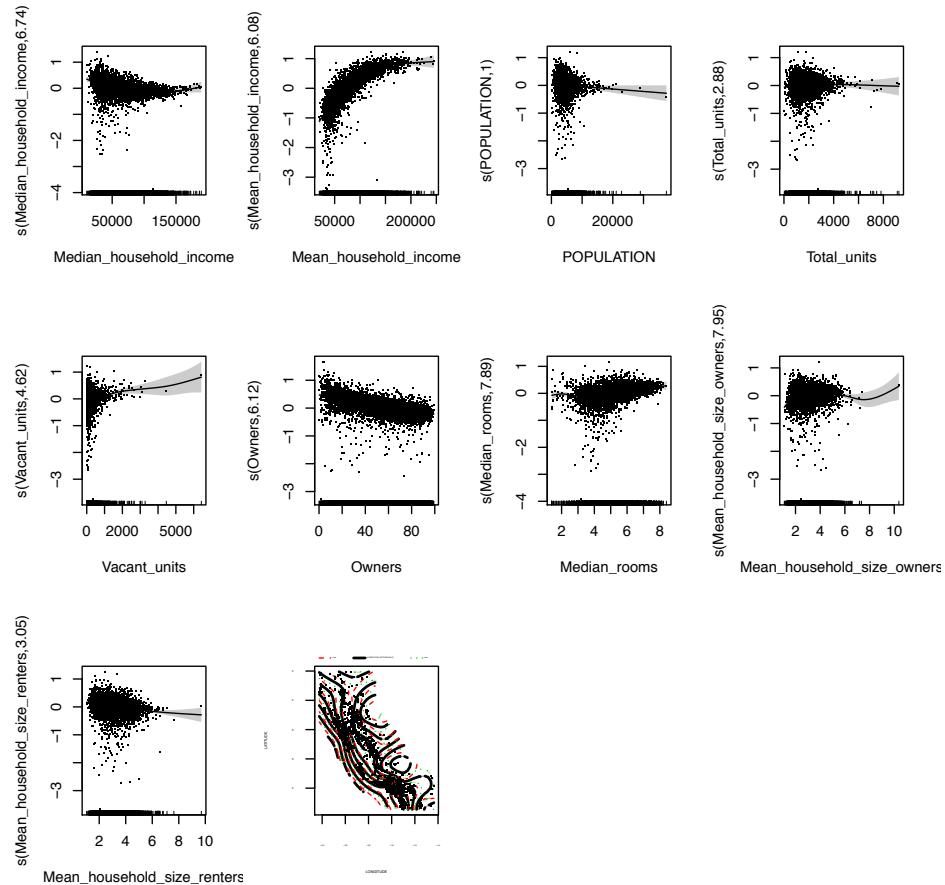
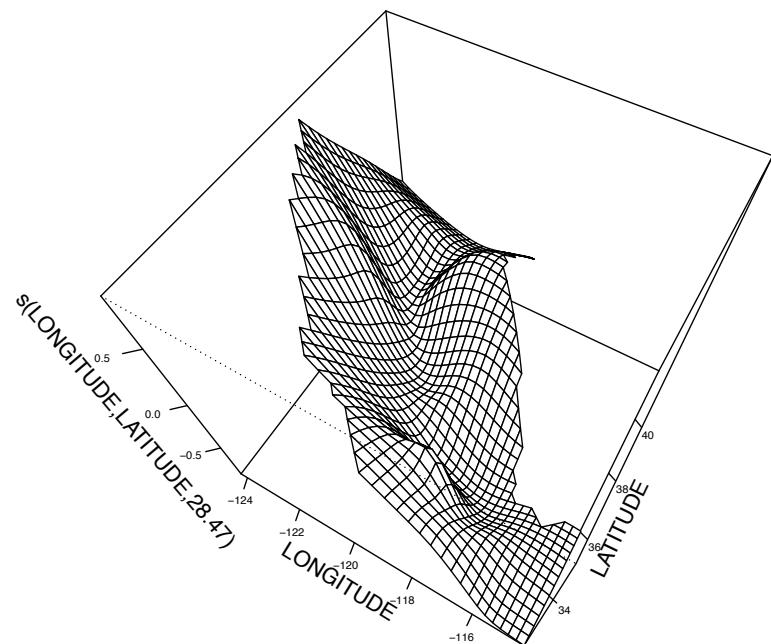


Figure 8.5 Partial response functions and partial residuals for `addfit2`, as in Figure 8.4. See subsequent figures for the joint smoothing of longitude and latitude, which here is an illegible mess. See `help(plot.gam)` for the plotting options used here.



```
plot(calif.gam2, select = 10, phi = 60, pers = TRUE, ticktype = "detailed",
      cex.axis = 0.5)
```

Figure 8.6 The result of the joint smoothing of longitude and latitude.

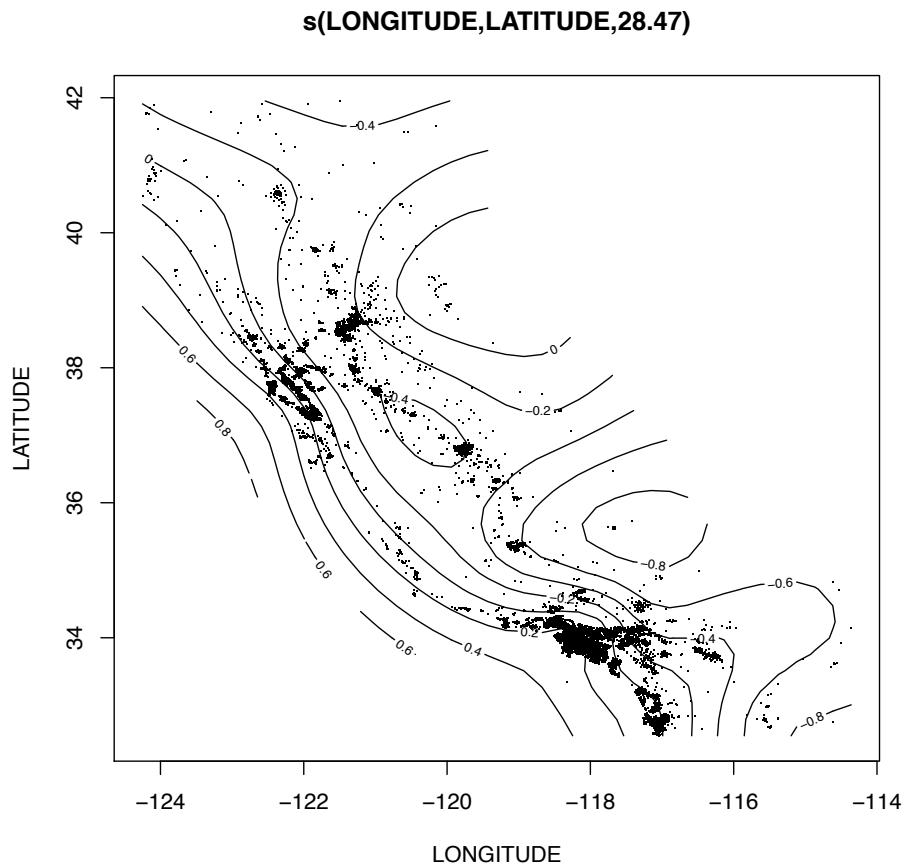
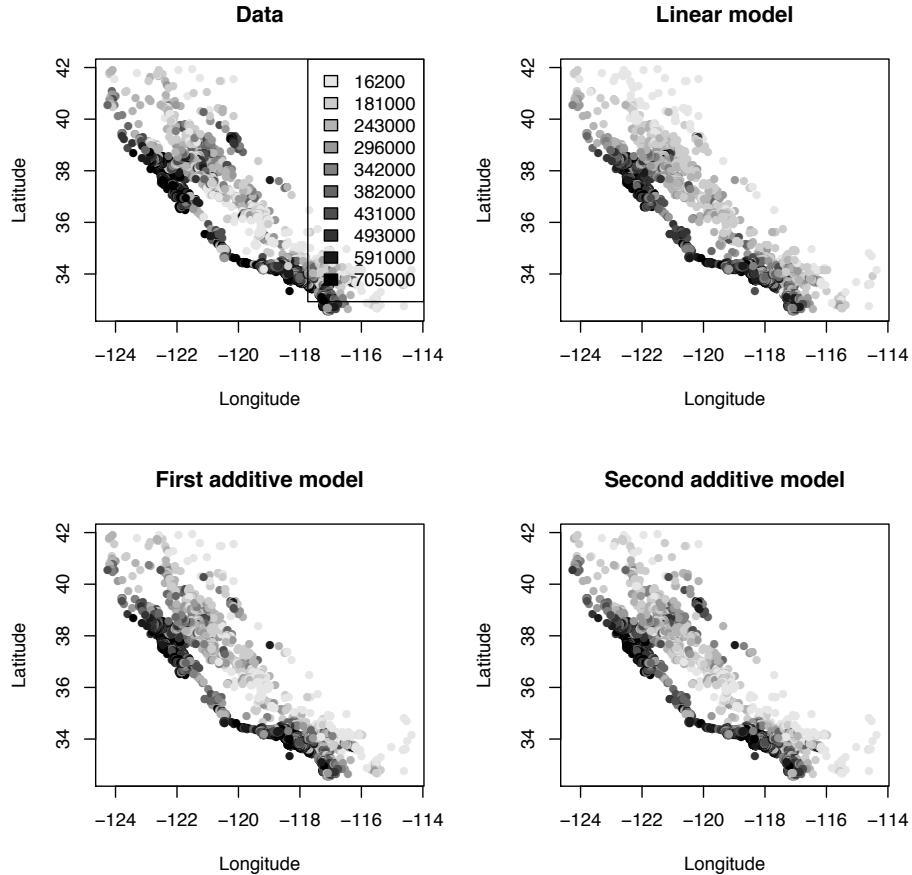


Figure 8.7 The result of the joint smoothing of longitude and latitude. Setting `se=TRUE`, the default, adds standard errors for the contour lines in multiple colors. Again, note that these are log units.



```

par(mfrow = c(2, 2))
calif.breaks <- graymapper(calif$Median_house_value, pch = 16, xlab = "Longitude",
  ylab = "Latitude", main = "Data", break.by = "quantiles")
graymapper(exp(preds.lm$fit), breaks = calif.breaks, pch = 16, xlab = "Longitude",
  ylab = "Latitude", legend.loc = NULL, main = "Linear model")
graymapper(exp(preds.gam$fit), breaks = calif.breaks, legend.loc = NULL, pch = 16,
  xlab = "Longitude", ylab = "Latitude", main = "First additive model")
graymapper(exp(preds.gam2$fit), breaks = calif.breaks, legend.loc = NULL, pch = 16,
  xlab = "Longitude", ylab = "Latitude", main = "Second additive model")
par(mfrow = c(1, 1))
  
```

Figure 8.8 Maps of real prices (top left), and those predicted by the linear model (top right), the purely additive model (bottom left), and the additive model with interaction between latitude and longitude (bottom right). Categories are deciles of the actual prices.

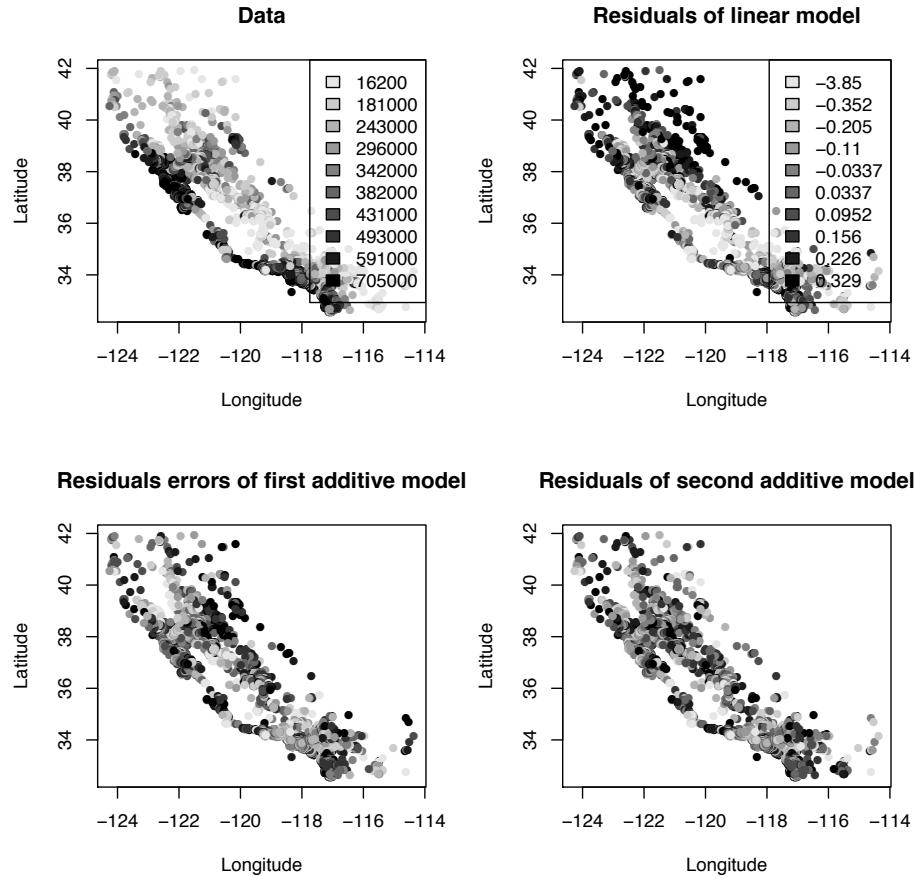


Figure 8.9 Actual housing values (top left), and the residuals of the three models. (The residuals are all plotted with the same color codes.) Notice that both the linear model and the additive model without spatial interaction systematically mis-price urban areas. The model with spatial interaction does much better at having randomly-scattered errors, though hardly perfect. — How would you make a map of the magnitude of regression errors?