

3.11. 1 → Basis functions

is equivalent to find the model equation for $y(x)$:

$$y = p_0 + p_1 x + \dots + p_d x^d + \epsilon$$

→ In linear regression to one variable x , consider the case where the predictor may be some function of x denoted by $\phi(x)$. This means we will transform the variable x using different functions & using the resulting variables as predictors. The functions $\phi_1, \dots, \phi_k(x)$ are called basis functions.

$$\boxed{y = p_0 + p_1 \phi_1(x) + \dots + p_k \phi_k(x) + \epsilon}$$

→ Let's consider 3 types of regressions & their particular choice of basis functions:

- Polynomial regression
- Step function approach
- Piecewise polynomial regression

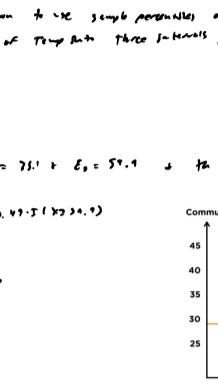
3.11. 2 → Polynomial regression

→ To model the response using \approx w/ d^{th} order/degree polynomial, the model equation is

$$\boxed{y = p_0 + p_1 x + p_2 x^2 + \dots + p_d x^d + \epsilon}$$

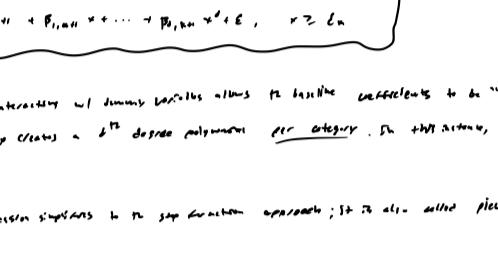
→ Then the basis functions are $b_j(x) = x^j, j=1, \dots, d$

→ Example: $y = 38.9 - 0.76x + 0.004x^2$



3.11. 3 → Step function approach

→ Start by dividing the range where data into intervals or bins. Hence, there are k cutpoints or knots denoted as E_1, E_2, \dots, E_k throughout the range of x . To illustrate, the following diagram shows four knots (blue) in intervals



→ The basis functions here are dummy variables that indicate the interval where an x value is found. Specifically,

$$\boxed{b_j(x) = \begin{cases} 1 & \text{if } x \in [E_j, E_{j+1}) \\ 0 & \text{if } x \in [E_k, \infty) \end{cases}, j=1, \dots, k}$$

→ where $b_j(x)$ is the indicator function. In effect, we have represented the variable x as a vector w/ k components. These dummy variables are called step functions. Notice that the interval $x \in E_i$ does not have a dummy variable, as it is the middle category by definition. As a result,

$$\boxed{\begin{aligned} &\rightarrow p_0 = \text{mean response for } x \in (-\infty, E_1] \\ &\rightarrow p_1 = \text{mean response for } x \in [E_1, E_2], \text{ when } j=1, \dots, k-1 \\ &\rightarrow p_k = \text{mean response for } x \in [E_k, \infty) \end{aligned}}$$

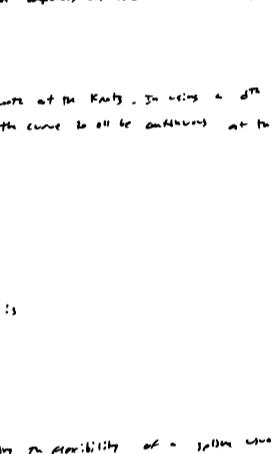
→ While we may naturally choose knot values, it is common to use sample quantiles of x spread uniformly, denoted as E_1, E_2, \dots, E_k where $E_1 = 25^{\text{th}} \text{ sample percentile of temp}$ & $E_k = 59.88^{\text{th}} \text{ sample percentile of temp}$.

$$E_1 = 25^{\text{th}} \text{ sample percentile of temp} \approx$$

$$E_2 = 60^{\text{th}} \text{ sample percentile of temp}$$

→ Using the empirical percentile approach, we set $E_1 = 25.1$ & $E_2 = 59.9$ & the following fit:

$$\begin{aligned} y &= 38.9 - 0.76x + 0.004x^2, 25.1 \leq x \leq 59.9 \\ &= \begin{cases} 38.9 & x < 25.1 \\ 38.9 - 0.76(25.1) & 25.1 \leq x < 59.9 \\ 35.53 & x \geq 59.9 \end{cases} \end{aligned}$$



3.11. 4 → Piecewise polynomial regression

→ Essentially we use step functions as dummy variables & their interactional polynomial functions w/ degree d .

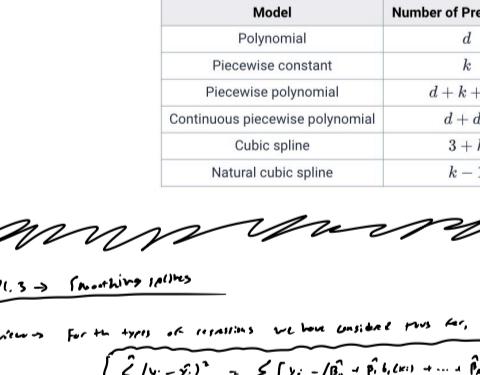
The result is a piecewise polynomial regression, & it is more intuitive to write the model equations:

$$\boxed{y = \begin{cases} p_0 + p_1 x + \dots + p_d x^d + \epsilon, & x < E_1 \\ p_0 + p_1 x + \dots + p_d x^d + \epsilon, & E_1 \leq x < E_2 \\ \vdots \\ p_0 + p_1 x + \dots + p_d x^d + \epsilon, & x \geq E_k \end{cases}}$$

→ Intuitively, the interacting w/ dummy variables allows the baseline coefficients to be "updated" for the associated categories. The step creates a d^{th} degree polynomial per category. In this instance, the categories are the different intervals of x .

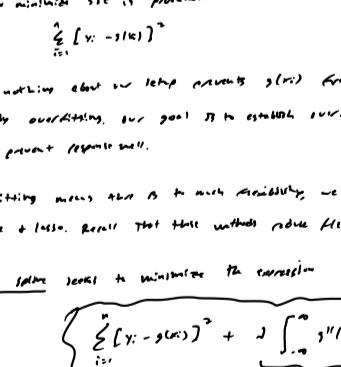
→ When done, the regression appears to be step function approach; it is also called piecewise constant regression for that reason.

3.11. 5 → Examples



3.11. 6 → Natural splines

→ In addition to the three quadratics fit by binning at the knots, we obtain the following fits:



→ This is called a continuous piecewise quadratic regression. In fact, we impose two constraints (one per knot) so that the quadratics are continuous at the knots. Imposing two constraints implies that we drop from $d+2-k$ predictors to $d+k-2$ predictors.

→ Even so, the natural curve goes back and forth near the knots. To smooth a fitted curve, we can impose more knots constraints: the first derivative of the piecewise quadratic must also be continuous at each knot. If we did this, the regression would only consist of $d-k+2$ predictors, since neither the quadratics would be imposed.

→ There is an example to show how we use the $d+k$ predictors.

$$\rightarrow \text{Model equation: } y = p_0 + p_1 x + p_2 x^2 + \epsilon$$

→ Now let's impose the constraint that $p_1 = p_2$. Then the model equation simplifies to

$$\rightarrow y = p_0 + p_1 x + p_2 x^2 + \epsilon$$

→ The result is a fit where the predictor is x & x^2 ; the model features change from $2 \rightarrow 1$ due to the constraint.

→ In general, a constraint is an equation that allows one p to be completely expressed in terms of the other p 's.

3.11. 7 → Splines

→ A regression spline is a continuous piecewise polynomial regression that is smooth at the knots. In using a d^{th} degree polynomial, smoothness is attained by requiring the first $d-1$ derivatives of the curve to all be continuous at the knots.

This means, for example, a cubic spline has

- continuity at the knots
- first derivative continuity at the knots, &
- second derivative continuity at the knots

→ Given that a cubic spline has k constraints per knot, its predictor count is

$$\begin{aligned} p &= 3+k - \underbrace{\sum_{i=1}^k \text{constraints}}_{= 3k} \\ &\rightarrow p = 3k \end{aligned}$$

→ Since a cubic function is sufficiently flexible in a given interval, increasing the flexibility of a spline usually comes from increasing k rather than raising the polynomial order. In other words, the k knots is a flexibility measure whose ideal value can be determined through cross-validation.

→ There are many ways to express the basis functions for a cubic spline; for this class we use the following set of functions: x, x^2, x^3 & the remaining $k-3$ are

$$\boxed{b_j(x) = (x-E_j)_+^3, b_1(x) = (x-E_1)_+^3, \dots, b_{k-3}(x) = (x-E_{k-3})_+^3}$$

→ Known as truncated power basis functions. \Rightarrow If $x > E$, then the basis function $= (x-E)^3$, otherwise = 0

3.11. 8 → Examples

3.11. 9 → Natural splines

→ One disadvantage of fitting polynomials is that (x_i) is typically large near the lower & upper boundaries of the variable. A method to mitigate this for a cubic spline (or spline in general) is to assume that the curve becomes flatter in the intervals $[x_1, E_1]$ & $[x_k, E_k]$. It is true that if E_1 & E_k are boundary knots, then the curve becomes flatter near the lower & upper boundaries of x , so remaining knots are considered as interior knots.

→ This results in a natural cubic spline. Relative to a cubic spline, we want to work w/ the fact that the knot values become flatter near the boundaries of x , i.e. more constraints.

The further apart k are, the less constraints we "need" to impose at the boundaries. Therefore, a natural cubic spline operates with $d+k-4 = k-1$ predictors, although there are fewer basis functions now, their become quite complex & are not shown here.

→ There is an example to show how we use the $d+k$ predictors.

$$\rightarrow \text{Model equation: } y = p_0 + p_1 x + p_2 x^2 + \epsilon$$

→ Now let's impose the constraint that $p_1 = p_2$. Then the model equation simplifies to

$$\rightarrow y = p_0 + p_1 x + p_2 x^2 + \epsilon$$

→ The result is a fit where the predictor is x & x^2 ; the model features change from $2 \rightarrow 1$ due to the constraint.

→ In general, a constraint is an equation that allows one p to be completely expressed in terms of the other p 's.

3.11. 10 → Smoothing splines

→ One disadvantage of fitting polynomials is that (x_i) is typically large near the lower & upper boundaries

of the variable. A method to mitigate this for a cubic spline (or spline in general) is to assume that the curve becomes flatter in the intervals $[x_1, E_1]$ & $[x_k, E_k]$. It is true that if E_1 & E_k are boundary knots, then the curve becomes flatter near the lower & upper boundaries of x , so remaining knots are considered as interior knots.

→ This results in a natural cubic spline. Relative to a cubic spline, we want to work w/ the fact that the knot values become flatter near the boundaries of x , i.e. more constraints.

The further apart k are, the less constraints we "need" to impose at the boundaries. Therefore, a natural cubic spline operates with $d+k-4 = k-1$ predictors, although there are fewer basis functions now, their become quite complex & are not shown here.

→ There is an example to show how we use the $d+k$ predictors.

$$\rightarrow \text{Model equation: } y = p_0 + p_1 x + p_2 x^2 + \epsilon$$

→ Now let's impose the constraint that $p_1 = p_2$. Then the model equation simplifies to

$$\rightarrow y = p_0 + p_1 x + p_2 x^2 + \epsilon$$

→ The result is a fit where the predictor is x & x^2 ; the model features change from $2 \rightarrow 1$ due to the constraint.

→ In general, a constraint is an equation that allows one p to be completely expressed in terms of the other p 's.

3.11. 11 → Effective degrees of freedom

→ In short, $d+k$ is the effective degrees of freedom of the model.

$$\rightarrow d+k$$

→ In general, $d+k$ is the effective degrees of freedom of the model.

$$\rightarrow d+k = \text{df}$$

→ It already mentioned, p_0, \dots, p_d are the free degrees of freedom of the model.

→ The shrinkage parameters λ & γ are the penalty degrees of freedom of the model.

→ The effective degrees of freedom are the sum of the free & penalty degrees of freedom.

$$\boxed{df = \text{df}_0 + \lambda \text{df}_{\text{penalty}}}$$

→ Define the i^{th} diagonal entry of S^{-1} as s_{ii} . As a result, the effective df is

$$\boxed{(df)_i = \sum_{j=1}^n s_{ij}^2}$$

→ It increases from 0 to ∞ , as the shrinkage df from 0 to ∞ , allowing it

→ to work w/ the effective df to estimate the i^{th} parameter p_i .

$$\rightarrow df = \sum_{i=1}^n (df)_i$$

→ Furthermore, recall that p_i doesn't require λ to calculate the i^{th} error when using OLS estimation.

→ In other words, the i^{th} error doesn't need the single fit w/ all the coefficients. That's exactly what we do for smoothing splines when λ is large enough.

$$\boxed{\text{Leverage error} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

→ In summary with a d degree polynomial & k knots, the table below gives the i^{th} predictors for the models we have discussed:

Model	Number of Predictors, p

<tbl_r cells="