Given: TS $(x_1 y_1) \ldots (x_n y_n)$ iid obs of $(X, Y)$

Goal: Generate prediction rule $f(x)$ that predicts $Y$ for query predictor $x$.

Optimal rule $f(x) = E(Y | x)$

Assume single predictor

Approach 1: Kernel smoothing

Approach 2: Expansion based rules

- Choose dictionary of "basis functions" $B_1(x) \ldots B_K(x)$ for which (hopefully)

$$f(x) \approx \sum_{j=1}^{K} a_j B_j(x) \leftarrow$$

How to find $a_1 \ldots a_K$?

Obvious approach: Find $\hat{\underline{a}}$ to minimize resubstitution error:

$$\hat{\underline{a}} = \underset{\underline{a}}{\text{argmin}} \| Y - X \underline{a} \|^2$$

where $X_{ij} = B_j(x_i)$.

Can we do better?

Look at extreme case: $K = n$, $B_1 \ldots B_K$ linearly independent over the data

$(\text{rank}(X) = n)$

$$\hat{y} = Hy = X(X^TX)^{-1}X^Ty$$

$$\hat{a} = (X^TX)^{-1}X^Ty$$
$$= X^{-1}X^TX^{+T}y$$
$$= X^{-1}y$$

$$= \underline{X X^{-1} X^{-T} X^{+}} y$$
$$\quad\quad\quad I$$

$$= y$$

what's the $EEE = E(\frac{1}{n} \| f - Hy \|^2)$

$y_i = f(x_i) + \varepsilon_i \quad E(\varepsilon_i) = 0 \quad V(\varepsilon_i) = \sigma^2 \quad \varepsilon_i$ inde

$$EEE_{est} = \frac{1}{n}\left( \underbrace{\| (I-H)f \|^2}_{\text{squared bias}} + \sigma^2 \underbrace{\boxed{\text{trace } H}}_{} \right) \quad \leftarrow$$

$$\underset{\| \atop 0}{} \qquad\qquad \underset{\text{in HW01 } \text{trace}(W^TW)}{n}$$

$$= EEE_{est} = 0 + \sigma^{2^n} = \sigma^2$$

If we had $K$ basis functions, then

$$EEE_{est} = \frac{1}{n}\left( \| (I-H)f \|^2 + \sigma^2 K \right)$$

Suppose $f(x) = \sum_{i=1}^{n} a_i B_i(x)$ but some of the

$a_i$ are $0$. Then it would be better to remove those basis functions from the dictionary. because this would reduce the variance component of $EEE$ and would not increase the bias component.
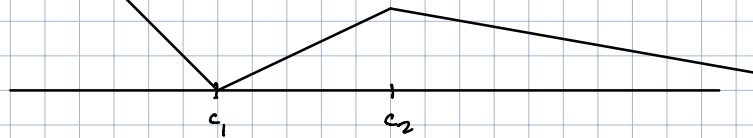
*Turbo*[(x)] Procedure to fit linear splines

Linear spline with knots $c_1 < c_2 \cdots < c_k$ is

- Piecewise linear in $(-\infty, c_1], (c_1, c_2] \cdots (c_k, \infty)$
- Continuous

\* Friedman & Silverman 1989



Any linear spline with knots $c_1 \ldots c_k$ is a linear combination of basis functions

$$B_i(x) = (x - c_i)_+ \quad i = 1 \ldots K$$
$$B_{k+1}(x) = 1$$
$$B_{k+2}(x) = x$$

Given training sample $(x_1, y_1) \ldots (x_n, y_n)$ Turbo finds

- \# of knots
- Knot positions
- slopes and intercepts of linear pieces

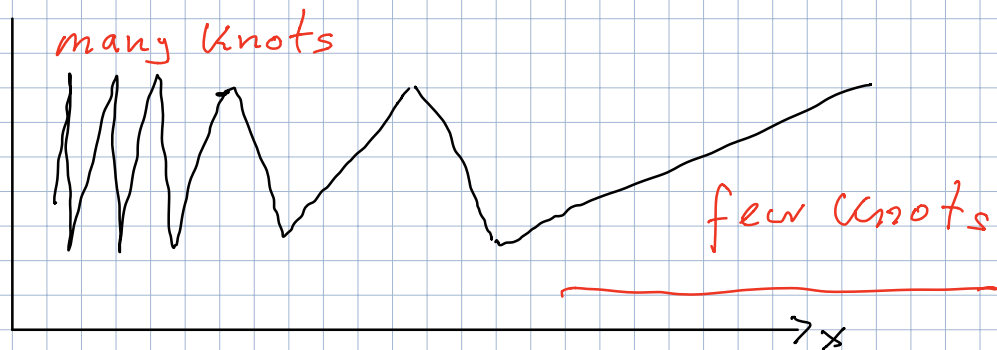to minimize an estimate of the expected squared prediction error.

Hard problem.

Suppose $K = 1$

$$\underset{\underline{a}, \, c_1}{\text{minize}} \quad \sum_{i=1}^{n} \left( y_i - a_1 B_1(x_i) - a_2 - a_3 x_i \right)^2$$

$$B_1(x) = (x - c_1)_+$$

Finding $c_1$ is a nonlinear optimization problem. No explicit solution.



<span style="color:red">many knots</span>

<span style="color:red">few knots</span>

$\rightarrow x$

<span style="color:red">Design choices made by FRS:</span>

For given knot positions choose $\underline{a}$ by least squares

Assume $\quad x_1 < x_2 \qquad < x_n$

Dictionary $\quad B_i = (x - x_i)_+ \quad i = 1 \dots n$

$$B_{n+1} = 1$$

Note: $B_1$ is linear over the range of data

$\Rightarrow$ Don't need an extra linear basis function.

<span style="color:red">Fit model by stepwise forward selection</span>

- Find best model with one basis function using exhaustive search.

$$\hat{1} = \underset{}{\text{argmin}} \quad \underset{}{\text{argmin}} \sum (y_i - a B_i(x_i))^2$$

- Find the best model with two basis functions, given we already have chosen the first one.

$$j_2 = \underset{j}{\arg\min}\ \underset{\underline{a}}{\arg\min} \sum \left( y_i - a_1 B_{j_1}(x_i) - a_2 B_j(x_i) \right)^2$$

⋮ etc