

$$\text{MSE} = \text{bias}^2 + \text{Variance}$$

$$\downarrow$$

$$\text{bias} = \mathbb{E}(\hat{r}(x)) - r(x)$$

$$y = r(x) + \epsilon$$

Non-Parametric Regression In machine learning this is known as "learning a function".

$$Y_i = r(x_i) + \epsilon_i \quad \leftarrow \text{model}$$

We will construct  $\hat{r} = \hat{r}(x)$   $\leftarrow$  estimator of  $r$ , or smoother

Assume that  $\text{Var}(\epsilon_i) = \sigma^2$ , independent of  $x_i$ .

Linear Smoother (not linear regression exactly).

Def:  $\hat{r}$  is a linear smoother if for each  $x$ , there exists a vector  $\vec{l}(x) = \begin{pmatrix} l_1(x) \\ l_2(x) \\ \vdots \\ l_n(x) \end{pmatrix}$  such that

$$\hat{r}(x) = \sum_{i=1}^n l_i(x) Y_i$$

$$\Rightarrow \vec{\hat{r}} = \begin{pmatrix} \hat{r}(x_1) \\ \vdots \\ \hat{r}(x_n) \end{pmatrix} \quad \underline{\text{fitted values}}$$

$$\vec{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}$$

$$\Rightarrow \hat{r} = L Y, \quad \text{where } L_{ij} = l_j(x_i)$$

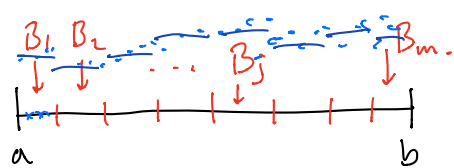
- The matrix  $L$  is called the smoothing matrix, or the "hat" matrix.  $i^{\text{th}}$  row is the "effective kernel" for the estimator of  $r(x_i)$ .
- The effective degrees of freedom =  $\nu = \text{tr}(L)$ .

Exercise Reinterpret linear least squares in this matrix vector formulation. (I.e., find  $L, Y$ ).

It will turn out that most linear smoothers have the property that for all  $x$ ,  $\sum l_i(x) = 1$ . If this is true, then if  $Y_i = c$  for all  $i$ , then  $\hat{r}(x) = \sum l_i(x) Y_i = c$ . (The smoother preserves constants.)

### Example Regressogram

Let  $x_i \in (a, b)$ , and compute  $m$  bins on this interval:



Decide on  $m$ ,  
call  $h = \frac{b-a}{m}$

For  $x \in B_j$ , define  $\hat{r}(x) = \frac{1}{k_j} \sum_{i: x_i \in B_j} Y_i$ ,  $k_j = \# x_i \in B_j$

computing the average of  $Y_i$  over bin  $B_j$ .

$$\Rightarrow \text{For } x \in B_j, \text{ set } l_i(x) = \begin{cases} \frac{1}{k_j} & \text{if } x_i \in B_j \\ 0 & \text{otherwise} \end{cases}$$

$$\Rightarrow \hat{r}(x) = \sum l_i(x) Y_i \quad \text{for } x \in B_j$$

$$\text{and } l(x) = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ y_{k_j} \\ y_{k_j} \\ \vdots \\ y_{k_j} \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \left. \vphantom{\begin{pmatrix} 0 \\ \vdots \\ 0 \\ y_{k_j} \\ y_{k_j} \\ \vdots \\ y_{k_j} \\ 0 \\ \vdots \\ 0 \end{pmatrix}} \right\} i \text{ such that } x_i \in B_j.$$

If  $n=9$ ,  $m=3$ , and  $k_1 = k_2 = k_3 = 3$ .

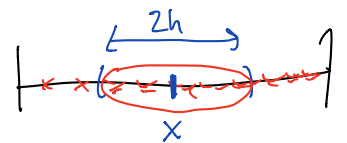
$$\Rightarrow L = \begin{pmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ 0 & 0 & 0 & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ 0 & 0 & 0 & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

$$\Rightarrow v = \text{tr}(L) = 3$$

Ex: Local Smoothing / Averaging

$$\text{For } h > 0, \text{ set } B_x = \{ i : |x_i - x| \leq h \}$$

$$n_x = |B_x|$$



$$\Rightarrow \hat{r}(x) = \frac{1}{n_x} \sum_{i \in B_x} Y_i$$

average our points within a distance of  $h$  from  $x$ .

The key question in both of these examples is how to choose  $h$ ?

- if  $h$  is big, we get a very smooth  $\hat{r}$
- if  $h$  is small, then  $\hat{r}$  looks a lot like  $Y_i$

$h$  is generally referred to as the BANDWIDTH

### Choosing the smoothing parameter

To recall the risk:  $R(h) = \mathbb{E} \left( \frac{1}{n} \sum_i (\hat{r}(x_i) - r(x_i))^2 \right)$

this depends on the unknown  $r$ .

Instead, choose  $h$  to minimize an estimate of  $R$ ,  $\hat{R}(h)$ .

Idea: Use the "training error":  $\frac{1}{n} \sum (Y_i - \hat{r}(x_i))^2$   
residual sum of squares  
RSS.

- Poor estimate of  $R(h)$
- Tends to lead to overfitting since outliers are given equal weight as other data points.

### A better option

Def: "Leave-one-out cross validation"

$$CV = \hat{R}(h) = \frac{1}{n} \sum (Y_i - \hat{r}_{(-i)}(x_i))^2$$

cross validation

$\hat{r}_{(-i)}$  obtained by ignoring the  $i^{\text{th}}$  data point.

For linear smoothers,  $\hat{r}(x) = \sum Y_i l_i(x)$

$\hat{r}_{(-i)}(x) =$  leave one out estimator

$$= \sum_j Y_j l_{j,(-i)}(x)$$

$$\text{where } l_{j,(-i)}(x) = \begin{cases} 0 & \text{if } j=i \\ \frac{l_j(x)}{\sum_{k \neq i} l_k(x)} & \text{if } j \neq i \end{cases} \left. \vphantom{\begin{cases} 0 \\ \frac{l_j(x)}{\sum_{k \neq i} l_k(x)} \end{cases}} \right\} \begin{array}{l} \text{re-weighting} \\ \text{of } l_j \end{array}$$

General idea compute  $\hat{r}$  with only part of the data and then check.

If we use this as  $\hat{R}(h)$ , then what can we say?

$$\begin{aligned} \mathbb{E} \left( Y_i - \hat{r}_{(-i)}(x_i) \right)^2 &= \mathbb{E} \left( \left( Y_i - r(x_i) \right) + \left( r(x_i) - \hat{r}_{(-i)}(x_i) \right) \right)^2 \\ &= \mathbb{E} \left( \left( Y_i - r(x_i) \right)^2 + 2 \left( Y_i - r(x_i) \right) \left( r(x_i) - \hat{r}_{(-i)}(x_i) \right) \right. \\ &\quad \left. + \left( r(x_i) - \hat{r}_{(-i)}(x_i) \right)^2 \right) \end{aligned}$$

model:  $Y_i = r(x_i) + \epsilon_i$

$$= \sigma^2 + \mathbb{E} \left( \left( r(x_i) - \hat{r}_{(-i)}(x_i) \right)^2 \right)$$

$$\text{for large } n \approx \sigma^2 + \mathbb{E} \left( \left( r(x_i) - \hat{r}(x_i) \right)^2 \right)$$

$$\Rightarrow \mathbb{E}(\hat{R}) \approx \underbrace{\sigma^2 + R}_{= \text{predictive risk}} \rightarrow \text{"nearly unbiased" since presumably } R \gg \sigma^2.$$

Is this expensive to do for all  $i$ ?

Then For a linear smoother, the CV <sup>estimator</sup> can be written as

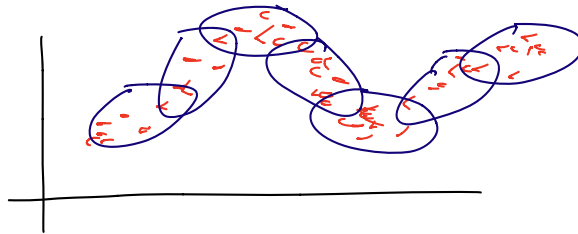
$$\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n \left( \frac{Y_i - \hat{r}(x_i)}{1 - L_{ii}} \right)^2$$

$\hat{L}_{ii} = l_i(x_i)$

Then  $h$  can be chosen by minimizing  $\hat{R}(h)$ .

## Local Regression

Idea: Give more weight to  $x_i, Y_i$  that are near to where you want to evaluate  $\hat{r}$



Definition A kernel is a function  $K = K(x)$

such that  $\int K(x) dx = 1$

$$\int x K(x) dx = 0$$

$$\int x^2 K(x) dx = \sigma_K^2 > 0$$

$< \infty$

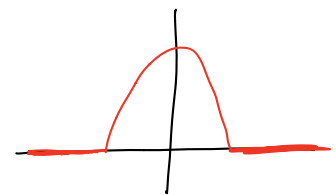
Ex:



$$\frac{1}{2} \mathbb{1}(x) = \begin{cases} 1 & \text{if } |x| \leq 1 \\ 0 & \text{else} \end{cases}$$



$$\frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$



$$\frac{3}{4} (1-x^2) \mathbb{1}(x)$$

Note that: For  $h > 0$ , and some  $x'$ , we have that

$K_n(x-x') = \frac{1}{h} K\left(\frac{x-x'}{h}\right)$  is also a kernel centered at  $x'$ .

Idea Use these "kernels" to smooth out the noisy data  $(x_i, Y_i)$ .

Definition For  $h > 0$ , the Nadaraya-Watson kernel estimator is

$$\hat{r}(x) = \sum_{i=1}^n l_i(x) Y_i$$

where  $l_i(x) = \frac{\frac{1}{h} K\left(\frac{x-x_i}{h}\right)}{\sum_{j=1}^n \frac{1}{h} K\left(\frac{x-x_j}{h}\right)}$

This form is similar to barycentric interpolation formula, but it is not an interpolant.

Goal: ① Choose  $h$  via cross-validation.

② The choice of  $K$  is less important, especially when there is a lot of data.

In order to choose  $h$ , we need to be able to estimate the risk. For our purposes, assume that the  $x_j$ 's are randomly drawn from some density  $f$ .

The risk can then be written as:

$$\begin{aligned} \text{Then } R(\hat{r}, r) &= \frac{h^4}{4} \underbrace{\left( \int x^2 K(x) dx \right)^2}_{\sigma_{K^2}^2} \underbrace{\left( \int (r'' + 2r' \frac{f'}{f})^2 dx \right)}_{\text{bias}^2} \\ &+ \underbrace{\frac{\sigma^2}{nh} \left( \int K^2(x) dx \right) \left( \int \frac{1}{f} dx \right)}_{\text{variance}} + o\left(\frac{1}{n}\right) + o(h^4) \end{aligned}$$

where we have assumed that  $h \rightarrow 0$  and  $nh \rightarrow \infty$ .

$\text{bias}^2 \sim h^4$   
 $\text{variance} \sim \frac{1}{nh}$

Comments (1) "Design bias":  $2r' \frac{f'}{f}$

Since this depends on  $f$ , the bias is sensitive to the distribution of the  $x_i$ 's

(2) To minimize  $R$ , set  $\frac{dR}{dh} = 0$  and solve.

$$\Rightarrow h_* \sim O\left(\frac{1}{n^{4/5}}\right)$$

$$\Rightarrow R \sim O\left(\frac{1}{n^{4/5}}\right) \cdot \left( \begin{array}{l} \text{vs MLE for least} \\ \text{squares, } R \sim \frac{1}{n} \end{array} \right).$$

Note Homoscedasticity

$$y = r(x) + \epsilon$$

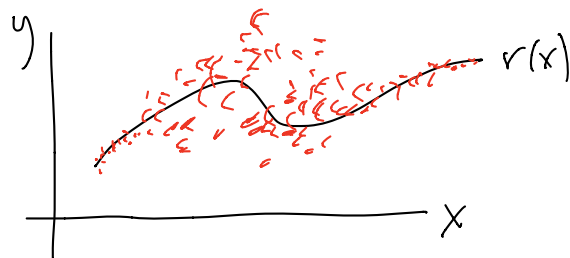
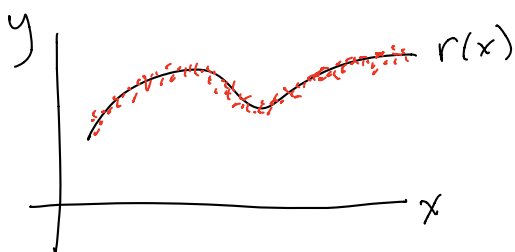
$$\text{Var}(\epsilon) = \sigma^2,$$

and  $\epsilon$  is independent of  $x$  ( $\epsilon|x = \epsilon$ ).

vs. Heteroscedasticity

$$y = r(x) + \sigma(x)\epsilon \quad \begin{array}{l} \swarrow \text{Var}(\epsilon) = 1 \\ \uparrow \text{function of } x \end{array}$$

$$\Rightarrow \text{Var}(\sigma(x)\epsilon) = \sigma^2(x) \text{Var}(\epsilon) = \sigma^2(x).$$



Reformulate the model to separate  $\sigma$  from  $\epsilon$ .

$$y - r = \sigma \epsilon$$

$$(y - r)^2 = \sigma^2 \epsilon^2$$

$$\log(y - r)^2 = \log \sigma^2 + \log \epsilon^2$$

$$\text{Set } z_i = \log(\hat{y}_i - r(x_i))^2$$

$$\delta_i = \log \epsilon^2$$

$$\Rightarrow z_i = \log \sigma^2(x_i) + \delta_i \quad \begin{array}{l} \uparrow \text{independent} \\ \text{of } x_i \end{array}$$

Then our goal is to estimate  $\log \sigma^2$ .



There is a two-step procedure for doing this:

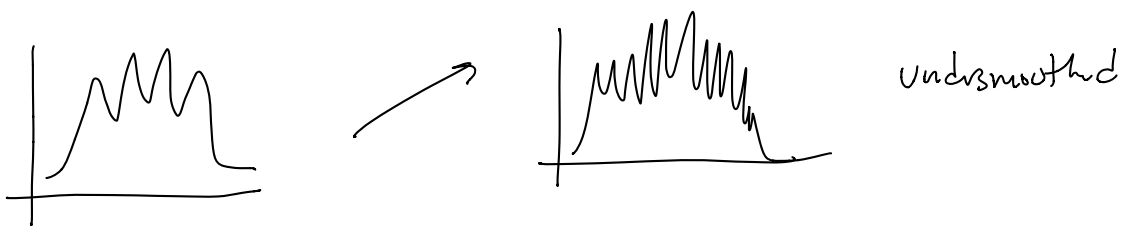
- ① Estimate  $r$  in  $Y_i = r(x_i) + \tilde{\epsilon}_i$  to get  $\hat{r}$ .
- ② Compute  $Z_i = \log(Y_i - \hat{r}(x_i))^2$
- ③ Regress  $Z_i$  on  $x_i$  again to get an estimator  $\hat{q} \approx \log(\sigma^2(x))$ . Set  $\hat{\sigma}^2(x) = e^{\hat{q}(x)} > 0$ .

## Density Estimation

Setup: Observe some data  $X_1, \dots, X_n \sim F$ , and therefore the density is  $f = F'$ .

Goal: Estimate  $f$  using as few assumptions as possible.

Still a smoothing problem =



Let  $\hat{f}$  be our estimate of  $f$ .



One way to measure the error:

$$L = \int (\hat{f}(x) -$$