## Exercise 7.4

$$Err_{in} = \frac{1}{N} \sum_{i=1}^{N} E_{Y_0} \left[ (Y_i^0 - \hat{f}(x_i))^2 \right]$$

$$\overline{err} = \frac{1}{N} \sum_{i=1}^{N} (y_i - \hat{f}(x_i))^2 \qquad y = f(x) + \varepsilon.$$

$$E_Y[op] = \frac{2}{N} \sum_{i=1}^{N} Cov(\hat{y}_i, y)$$

$$E_Y[op] = E_Y \left[ \frac{1}{N} \sum_{i=1}^{N} \left\{ E_{Y_0} \left[ (Y_i^0 - \hat{f}(x_i))^2 \right] - (y_i - \hat{f}(x_i))^2 \right\} \right]$$

$$= E_Y \left[ \frac{1}{N} \sum_{i=1}^{N} \left\{ E_{Y_0}[Y_i^{0^2}] - 2 E_{Y_0}[Y_i^0] \hat{f}(x_i) + \hat{f}(x_i)^2 + \begin{array}{c} +E[\hat{f}(x)] \\ -E[\hat{f}(x)] \end{array} \right. \right.$$
$$\left. \left. - y_i^2 + 2 E[y_i \hat{f}(x_i)] - \hat{f}(x_i)^2 \right\} \right]$$

$$= \frac{1}{N} \sum_{i=1}^{N} E_Y[E_{Y_0}[Y_i^{0^2}]] - 2 E_{Y_0}[Y_i^0] E_Y[\hat{f}(x_i)] - E_Y[y_i^2] +$$
$$+ 2 E_Y[y_i \hat{f}(x_i)]$$

- $E_Y \left[ E_{Y_0}[Y_i^{0^2}] - f(x)^2 + f(x)^2 \right] = \sigma_\varepsilon^2 + f(x)^2$

- $E_Y \left[ y^2 - f(x)^2 + f(x)^2 \right] = \sigma_\varepsilon^2 + f(x)^2$

- $E_{Y_0}[Y_i^0] = f(x) = E_Y[y] \qquad\qquad \hat{f}(x) = \hat{y}$

$$E_Y[op] = \frac{1}{N} \sum_{i=1}^{N} \left\{ -2 E_Y[y] E_Y[\hat{y}] + 2 E_Y[y_i \hat{y}_i] \right.$$
$$= \frac{2}{N} \sum_{i=1}^{N} Cov(y_i, \hat{y}_i)$$

1

## Ex 7.5

For $\hat{y} = Sy$ show that $\sum_{i=1}^{N} Cov(y_i, \hat{y}_i) = trace(S)\sigma^2$

$$\sum_{i=1}^{N} Cov(y_i, \hat{y}_i) = trace\left(Cov(y, \hat{y})\right)$$
$$= trace\left(Cov(y, Sy)\right)$$
$$= trace\left(S \, Cov(y, y)\right)$$
$$= trace\left(S \, Var(y)\right)$$
$$= trace\left(S \, \sigma^2\right)$$
$$= trace(S)\sigma^2$$

Cov matrix $\begin{pmatrix} Cov(y_1, y_1) & Cov(y_1, y_2) & \cdots & Cov(y_1, y_n) \\ Cov(y_2, y_1) & Cov(y_2, y_2) & & \cdot \\ \vdots & & \ddots & \\ Cov(y_n, y_1) & \cdots & & Cov(y_n, y_n) \end{pmatrix}$

$\sum_{i=1}^{rank(M)} m_{ii} = trace \, M$

## Exercise 7.7

Use the approximation: $\dfrac{1}{(1-x)^2} \approx 1 + 2x$

to show similarities between $C_p/AIC$ and GCV

$$C_p = \overline{err} + 2\frac{d}{N}\hat{\sigma}_\varepsilon^2 \qquad\qquad AIC = -\frac{2}{N}\,loglik + 2\frac{d}{N}$$

$$GCV = \frac{1}{N}\sum_{i=1}^{N}\left(\frac{y_i - \hat{f}(x_i)}{1 - trace(S)/N}\right)^2$$

$$C_p = \frac{1}{N}\sum_{i=1}^{N}(y_i - \hat{f}(x_i))^2 + 2\frac{d}{N}\hat{\sigma}_\varepsilon^2$$

$$trace(S) = \frac{\sum Cov(y, \hat{y}_i)}{\sigma^2}$$
$$\frac{1}{\sigma^2}\sum(y_i - \hat{f}(x_i)^2$$

$$GCV = \frac{1}{N}\sum_{i=1}^{N}(y_i - \hat{f}(x_i))^2 \cdot \left(\frac{1}{1 - \frac{trace(S)}{N}}\right)^2$$

$$= \frac{1}{N}\sum_{i=1}^{N}(y_i - \hat{f}(x_i))^2\left(1 + 2\cdot\frac{trace(S)}{N}\right)$$

$$= \frac{1}{N}\sum_{i=1}^{N}(y_i - \hat{f}(x_i))^2 + \frac{2}{N}\,trace(S)\frac{\sum(y_i - \hat{f}(x_i))^2}{N} = \overline{err} + \frac{2}{N}trace(S)\hat{\sigma}_\varepsilon^2$$

Gaussian regression

$$AIC \propto -\frac{2}{N}\,log\left(exp\left\{-\frac{1}{2\sigma_\varepsilon^2}\sum_{i=1}^{N}(y_i - \hat{f}(x_i))^2\right\}\right) + \frac{2d}{N}$$

$$= +\frac{1}{N\sigma^2}\sum_{i=1}^{N}(y_i - \hat{f}(x_i))^2 + \frac{2d}{N}$$

trace(S)

$$\propto \frac{1}{N}\sum_{i=1}^{N}(y_i - \hat{f}(x_i))^2 + \frac{2d}{N}\sigma_\varepsilon^2$$

# Bootstrap methods

→ what is bootstrap
→ how to use bootstrap for error estimation    $\widehat{E[Err_\tau]}$

IDEA: generate bootstrap sample from the empirical distribution
computed on original sample
→ by resampling with replacement from the original sample

- suppose $\tau = \{ (x_1, y_1), \dots, (x_n, y_n) \}$

- by resampling, $\tau_1^* = \{ (x_1^*, y_1^*), \dots, (x_n^*, y_n^*) \}$

- repeat for $B$ large, $\tau_1^*, \tau_2^*, \dots, \tau_B$

Based on the generated bootstrap sample (which mimic new experiments)
we can estimate any aspect of the distribution of a map

## Example

original sample $\tau = \{ z_1, z_2, z_3, z_4 \} = \{ 1, 3, 4, 6 \}$

generate $B$ bootstrap sample      $\tau_1^* = \{ z_1^*, z_2^*, z_3^*, z_4^* \}$
by resampling with replacement                   $| $
from $\tau$                                              $= \{ 1, 4, 1, 6 \}$

$\tau_2^* = \{ 4, 4, 3, 3 \}$

$\vdots$

$\tau_B^* = \{ 1, 1, 1, 1 \}$

$(x, y)$      $Cov(x, y)$

$\tau$  $(x_1; y_1), (x_2; y_2), \dots, (x_n; y_n)$
  $(4; 3), (1; 3), \dots, (4; 2)$

$\tau_1^* = \{ (1; 3), (4; 2), \dots, (4; 2) \}$

## Bootstrap approach for prediction error estimation

WRONG APPROACH
- estimate our $\hat{f}(x)$ from each bootstrap sample
- evaluate how well $\hat{f}_b^*(x)$ estimate $y$

$$\widehat{Err}_b^{\text{WRONG}} = \frac{1}{B} \sum_{b=1}^{B} \left( \frac{1}{N} \sum_{i=1}^{N} L\left(y, \hat{f}_b^*(x_i)\right) \right)$$

⚠ : training and test set are <u><u>not</u></u> independent

$$E[\widehat{Err}_b^{\text{WRONG}}] = 0.184 \quad , \quad 1NN \quad , \quad Y \perp X$$

$$\begin{array}{l} \bullet \; y_i \in \tau_b^* \longrightarrow \text{error} = 0 \\ \bullet \; y_i \notin \tau_b^* \longrightarrow \text{error} = 0.5 \end{array}$$

$$= 0.5 \times \underline{Pr[y_i \notin \tau_b^*]} + 0 \times \underline{Pr[y_i \in \tau_b^*]}$$
$$\qquad\qquad \underset{0.368}{}$$

$$Pr[\text{observation } i \notin \underline{\text{bootstrap sample } b}]$$

$$Pr[\tau_{b[i]}^* \neq y_i] = \frac{N-1}{N} \qquad \Rightarrow Pr[y_i \notin \tau_b^*] = \left(\frac{N-1}{N}\right)^N$$
$$\underset{\text{same for all positions}}{} \qquad\qquad\qquad\qquad |$$
$$\qquad\qquad\qquad\qquad\qquad = e^{-1} \approx 0.368$$

$$E[\widehat{Err}_b^{\text{WRONG}}]\Big|_{\substack{X \perp Y \\ 1NN}} \approx 0.5 \times e^{-1} = 0.184$$

An important fact
$$Pr[\text{observation } i \text{ belongs to a bootstrap sample } b] = 1 - e^{-1}$$
$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad |$$
$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \approx \boxed{0.632}$$

$$\text{approximation } 1 - \left(\frac{N-1}{N}\right)^N$$
$$\text{of}$$

CORRECT APPROACH

$\tau = \{z_1, \ldots, z_N\}$

$\tau_s^* = \{z_1^*, \ldots, z_N^*\}$ resampling with replacement

→ there are original observations which are included more than once
⟹ there are original observations which are not included at all
also
↳ these can be used as a test set as they are not used in the training process.

$$\widehat{Err}^{(1)} = \frac{1}{N} \sum_{i=1}^{N} \frac{1}{|C^{-i}|} \sum_{b \in C^{-i}} L\left(y_i, \hat{f}_b^*(x_i)\right)$$

where $|C^{-i}|$ is the number of bootstrap samples that do NOT contain $i$

## Issues

→ the average number of unique observations in the training set is 0.632N → not so far from 0.5N, that is the value related to 2-fold CV

→ similar issues of training-set-size bias than 2-fold CV
→ result in a small overestimation of the error

To solve the issue, the .632 estimator has been developed

$$\widehat{Err}^{(.632)} = 0.368 \, \overline{err} + 0.632 \, \widehat{Err}^{(1)}$$

In general, it works well, but in some case it fails, like in our 1NN
$x \perp y$

$$\overline{err} = 0 \quad \rightarrow \quad \widehat{Err}^{(.632)} = 0.632 \, \widehat{Err}^{(1)}$$

Further improvements ".632+ estimator"

- based on the quantity $\gamma$, the no-information-error rate
  error that we obtain when inputs and class label are independent
  $\hat{\gamma}$ is computed by permuting $x$ and $y$ separately, we compute
  the prediction error for each combination of $y_i$ and $x_j$

$$\hat{\gamma} = \frac{1}{N} \sum_{i=1}^{N} \frac{1}{N} \sum_{j=1}^{N} L\left(y_i, \hat{f}(x_j)\right)$$

- $\hat{\gamma}$ is used to compute the overfitting rate

$$\hat{R} = \frac{\widehat{Err}^{(1)} - \overline{err}}{\hat{\gamma} - \overline{err}} \qquad 0 \le R \le 1$$
$$\qquad\qquad\qquad\qquad\qquad \hookleftarrow \text{no overfitting}$$

- finally

$$\widehat{Err}^{(632+)} = (1 - \hat{w}) \, \overline{err} + \hat{w} \, \widehat{Err}^{(1)}$$

$$\text{where } w = \frac{0.632}{1 - 0.368 \hat{R}}$$

# Generalized Additive Models
- extensions of the (generalized) linear model

## Linear model
 - powerfull tool
 - can be used in several cases (regression, classification, ...)

## Main limitation
 - it suppose linear effects, often not true in reality
   ($\hat{\beta}$ is the increments in $y$ when the corresponding $x$ increase of $1$)

In the context of regression, the (generalized) additive model has the form

$$E[Y | x_1, ..., x_p] = \alpha + f_1(x_1) + ... + f_p(x_p)$$

where
  $Y$ is the outcome
  $x_j$ are the predictors
  $f_j$ is a function which describe the effect of $x_j$
- we already saw that we can use $f_j(x_j) = x_j^2$, $f_j(x_j) = \log x_j$
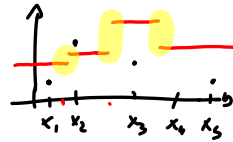- we can be more general, and use a nonparametric function
  (splines, kernel, ... )

  splines $\to$ §5.2


  kernel $\to$ §6.1 , §6.2


- cubic splines $\to$ bottom right of fig 5.2
- natural splines $\to$ since the estimation outside the observations range
                        can be dangerous, the line is forced to be
                        linear outside the range

6

Kernel method
- extension of k-nn



2·nn (red)
- use the $y_i$ of the 2 nearest neighbours
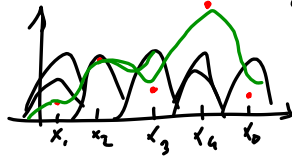- costant
- ugly and unnecessary discontinuities

kernel method
- weights through a kernel
$$D\left(\frac{|x-x_0|}{h_\lambda(x_0)}\right)$$

different kinds of kernels          $h, \lambda$ are parameters
                                    $h$ controls the shape
                                    $\lambda$ the width of the kernel



---

**GAM**

for functional effect, we can

- $E[Y|x] = \alpha + \beta_1 f_1(x_1) + \beta_2 f_2(x_2) + \dots + \beta_p f_p(x_p)$
  
  $\leftarrow \log(x_i)$
  
  → the least square estimator approach is usable

- $\overline{E[Y|x]} = \alpha + f_1(x_1) + \dots + f_p(x_p)$
  
  → backfitting algorithm

__Generalized__ Additive Model
  └→ extending the GLM   (STk 3100)

GLM   $g(\mu(x)) = \alpha + \beta^T X$         extending the linear model to
                                              all exponential families sampling
  └→ link function                            models

  e.g.  logistic model

      $g() = $ logit              $\mu(x) = P[Y = 1 | X = x]$

      $\log\left(\frac{\mu(x)}{1-\mu(x)}\right) = \alpha + \beta_1 x_1 + \dots + \beta_p x_p$

GAM   $g(\mu(x)) = \alpha + \sum_{j=1}^{p} f_j(x_j)$

additive logistic regression :  $\log\left(\frac{\mu(x)}{1-\mu(x)}\right) = \alpha + \sum_{j=1}^{p} f_j(x_j)$

Advantages of GAM:
- flexibility, due to $f$ (we can capture non-linear effects)
- interpretability, due to the additivity (not so different from
   the usual interpretation of GLM)
Note: not all effect need to be non-linear/linear

- semiparametric model   $g(\mu(x)) = \underbrace{X^T\beta}_{\text{parametric}} + \underbrace{f(z)}_{\text{non-parametric}}$

  e.g. semiparametric model : Cox model
       $\lambda(t) = \underbrace{\lambda_0(t)}_{\text{non-parametric}} \underbrace{\exp(X^T\beta)}_{\text{parametric part}}$

7