

## Exercise 10.2

$$f^*(x) = \arg \min_{f(x)} E_{Y|x} [e^{-Y f(x)}]$$

$$\frac{\partial}{\partial f(x)} E_{Y|x} [e^{-Y f(x)}] = E_{Y|x} [-Y e^{-Y f(x)}]$$

$$E_{Y|x} [-Y e^{-Y f(x)}] = 0 \quad Y = \begin{cases} -1 & \Pr[Y = -1|x] \\ 1 & \Pr[Y = 1|x] \end{cases}$$

$$(-1) e^{-(+1) f(x)} \Pr[Y = -1|x] - (1) e^{-(-1) f(x)} \Pr[Y = 1|x] = 0$$

$$e^{f(x)} \Pr[Y = -1|x] = e^{-f(x)} \Pr[Y = 1|x] \quad \text{multiply both sides by } e^{f(x)}$$

$$e^{2f(x)} \Pr[Y = -1|x] = \Pr[Y = 1|x]$$

$$e^{2f(x)} = \frac{\Pr[Y = 1|x]}{\Pr[Y = -1|x]}$$

$$f(x) = \frac{1}{2} \log \left( \frac{\Pr[Y = 1|x]}{\Pr[Y = -1|x]} \right)$$

## Exercise 10.4

$$X = (x_1, \dots, x_{10}) \quad x_j \sim \mathcal{N}(0, 1) \quad x_j \text{ iid} \quad N = 2000$$

$$Y_i = \begin{cases} 1 & \text{if } \sum_{j=1}^{10} x_j^2 > \chi_{10}^2(0.5) \\ -1 & \text{if } \sum_{j=1}^{10} x_j^2 \leq \chi_{10}^2(0.5) \end{cases}$$

$$\begin{matrix} \rightarrow & \begin{matrix} 1 \\ 2 \\ \vdots \\ N \end{matrix} \\ \text{obs} & \left\{ \begin{matrix} x_{11} & x_{12} & \dots & x_{110} \\ x_{21} & x_{22} & \dots & x_{210} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{N10} \end{matrix} \right\} \end{matrix} \quad \text{apply } (x^2, \_, \text{sum})$$

Statistical boosting  $\begin{cases} \text{gradient boosting} \\ \text{likelihood-based boosting} \end{cases}$

From the previous lecture:

AdaBoost: classifier

- weak estimator
- loss function
- iteratively apply a weak estimator to modifications of the data in order to minimize a loss function

AdaBoost: - stump  
- tree  
- ...

↳ AdaBoost: weight more misclassified observations

for classification  $\rightarrow$  AdaBoost  $e^{-Yf(u)}$

from classification to regression

- loss function:  $RSS$   $\xrightarrow{\text{LS estimator}}$

- weak estimator:  $\lambda (X^T X)^{-1} X^T y$   $\lambda \rightarrow 0$   
 $0 < \lambda < 1$   $\xrightarrow{\text{weak estimator}}$

penalty  
parameter

$\rightarrow$  makes our LSE "weak" default = 0.1

- modification of the data:  $y \rightarrow u$  residuals

focusing on the not explain  
part of the variation  
outcome

## $L_2$ Boost algorithm for linear regression

① Initialization: initialize the regression coefficient estimate  $\hat{\beta}^{[0]} = (0, \dots, 0)$

(first modification of the data:  $v = y - X\beta = y - 0 = y$ )

② for  $m$  from 1 to  $m\_stop$

a) fit the weak estimator to the modification of the data

$$\hat{b}^{[m]} = \frac{1}{n} (X^T X)^{-1} X^T v \quad \lambda = 0.1$$

b) update the estimate  $\hat{\beta}^{[m]} = \hat{\beta}^{[m-1]} + \hat{b}^{[m]}$

c) modify the data:  $v = y - X^T \hat{\beta}^{[m]}$

③ final estimate

$$\hat{\beta}_{L_2\text{Boost}} = \sum_{m=1}^{m\_stop} \hat{b}^{[m]} = \hat{\beta}^{[m\_stop]}$$

Note:

$$m \rightarrow \infty, \quad \hat{\beta}_{L_2\text{Boost}} \rightarrow \hat{\beta}_{OLS}$$

need of an early stop (find the "right"  $m\_stop$ ) in order to not overfit (to find the best balance between bias and variance for the prediction error)

-  $m\_stop$  is the crucial tuning parameter

if it is too small: too much bias (our model does not explain the outcome variation)

if it is too big: too much variance (we overfit the data)

→ complexity parameter

boosting has a second tuning parameter,  $\lambda$  (is not so important, because smaller values  $\Rightarrow$  more steps (iterations)  
larger values  $\Rightarrow$  less steps

$X$  must be centred  $E[X_i] = 0$   
(it is ok to standardize)

$L_2$  Boost algorithm in general

- the goal is to minimize the loss function. At each step we want to identify the direction of the greatest decrease of the loss function

→ negative gradient :  $-\frac{\partial L(y, f(x))}{\partial f(x)}$

e.g.  $\frac{\partial \sum_{i=1}^n (y_i - x_i^T \beta)^2}{\partial x_i^T \beta} = \frac{\partial \sum (y - x^T \beta)}{\partial \text{residuals}}$

$\frac{1}{n} \exp\left\{-\frac{1}{2}(y - x^T \beta)\right\}$

In general:

① Initialization :  $\hat{f}(x) \equiv 0$  or  $\hat{f}(x) = \bar{y}$

② for  $m$  from 1 to  $m_{\text{stop}}$

(a) derive  $v = -\frac{\partial L(y, \hat{f}(x))}{\partial \hat{f}(x)}$

(b) fit our weak estimator :  $\hat{\beta} = g(v, x, v)$

(c) update the estimate  $\hat{f}^{(m)} = \hat{f}^{(m-1)} + \hat{\beta}$

③ Finalization  $\hat{f}_{\text{boost}} = \hat{f}^{(m_{\text{stop}})}$

GAM :  $\hat{g} = \sum_{j=1}^p f_j(x_j)$

$L_2$  Boost for High Dimensional Data

- one of the advantages of boosting is that we can handle HAD
- componentwise version of boosting

## Componentwise boosting

• Linear <sup>Gaussian</sup> regression model

① Initialization  $\hat{\beta}_j^{(0)} \equiv 0 \quad j = 1, \dots, p$   
 $(u = y - X\hat{\beta}_j^{(0)} = y)$

② for  $m$  from 1 to  $m_{\text{stop}}$

a) compute possible updates for each dimension of the regression coefficient vector separately  
 (fit a weak estimator on each dimension of  $X$ )

$$b_j = \frac{\sum_{i=1}^n x_i^{(j)} u_i}{\sum_{i=1}^n x_i^{(j)2}} \quad j = 1, \dots, p$$

b) select the best update among the  $p$  possibilities

$$j^* : \arg \min_j \sum_{i=1}^n (u_i - x_i^{(j)} b_j)^2$$

c) update the  $j^*$ -th regression coefficient

$$\hat{\beta}^{(m)} = \hat{\beta}^{(m-1)} + (0, \dots, 0, b_{j^*}, 0, \dots, 0)$$

d) modify the data  $u = y - X\hat{\beta}^{(m)}$

fit a GAM in R with boosting

$$y = f_1(x_1) + f_2(x_2) + \dots + f_p(x_p)$$

splines  
tree

→ function `mboost` of the package `mboost`

$$\mu(y) = X^T \beta \rightarrow \text{glmboost}$$

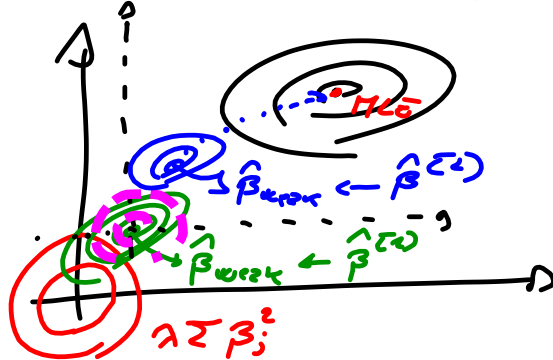
Gaussian regression  
logistic regression  
⋮

## Likelihood-based boosting

- fully statistical approach  $\rightarrow$  likelihood-based
- weak estimator:

standard estimator  $\hat{\beta}_{MLE} = \arg\max_{\beta} \ell(\beta)$

weak estimator  $\hat{\beta}_{weak} = \arg\max_{\beta} \ell_{pen}(\beta)$



ridge penalty  
 $\hat{\beta}_{weak} = (X^T X + \lambda I)^{-1} X^T y$

same meaning of the parameter  $\lambda$

## Boosting Ridge (Gaussian regression)

- ① initialize  $\hat{\beta}^{(0)} = (0, \dots, 0) \rightarrow v = y - X\hat{\beta}^{(0)}$
- ② fit the weak estimator  $\hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$

$$\frac{\partial \ell_{pen}(\beta)}{\partial \beta} = 0 \quad \ell_{pen} = \frac{1}{2} (y - X\beta)^T (y - X\beta) + \frac{\lambda}{2} \beta^T \beta$$

kernel of 2 Gaussian log-likelihood

$$\exp \left\{ -\frac{1}{2} (y - X\beta)^T (y - X\beta) \right\}$$

$$\frac{\partial \ell_{pen}}{\partial \beta} = -X^T (y - X\beta) + \lambda \beta = 0$$

$$-X^T y + X^T X \beta + \lambda \beta = 0$$

$$\beta (X^T X + \lambda I) = X^T y$$

$$\hat{\beta} = \hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$$

$$(b) \hat{\beta}^{(n)} = \hat{\beta}^{(n-1)} + \hat{\beta}$$

(c) modification of the data  
 (add an offset in the log-likelihood)

$$\frac{1}{2} (y - X\hat{\beta}^{(n)} - X\beta) (y - X\hat{\beta}^{(n)} - X\beta)^T$$

$$3 \hat{\beta}_{L_{Boost}} = \sum_{n=1}^{n_{step}} \hat{\beta}^{(n)}$$

weak estimator is the ridge estimator  
 Gaussian regression  $\rightarrow$  likelihood-based boosting  $\rightarrow$  some model  
 gradient boosting  $\rightarrow$  when we choose  $\lambda$  and  $\lambda$  in the right way

### Exercise

$$\hat{y} = X \hat{\beta}_{\text{boost}} = \sum_{m=1}^{m_{\text{step}}} S(I-S)^{m-1} y = (I - (I-S)^{m_{\text{step}}}) y$$

where  $S = X(X^T X + \lambda I)^{-1} X^T$

using  $\hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$  Boosting with ridge estimator

Show through SVD that Boosting ridge and ridge regression provide different shrinkage effect

ridge regression :  $\frac{d_j^2}{d_j^2 + \lambda}$

Boosting Ridge :  $(1 - (1 - \frac{d_j^2}{d_j^2 + \lambda})^{m_{\text{step}}})$