

Exercise 3.6

$$P(\beta|y) \propto p(y|\beta)P(\beta)$$

where

$$p(y|\beta) \leftrightarrow \mathcal{N}(X\beta, \sigma^2 I)$$

$$P(\beta) \leftrightarrow \mathcal{N}(0, \tau I)$$

$$\frac{1}{\sqrt{2\pi} \sigma} \text{const}$$

$$P(\beta|y) \propto \text{const.} \exp\left\{-\frac{1}{2\sigma^2} (y - X\beta)^T (y - X\beta)\right\} \exp\left\{-\frac{1}{2\tau} \beta^T \beta\right\}$$

$$\propto \exp\left\{(y - X\beta)^T (y - X\beta) + \frac{\sigma^2}{\tau} \beta^T \beta\right\}$$

Exercise 3.10

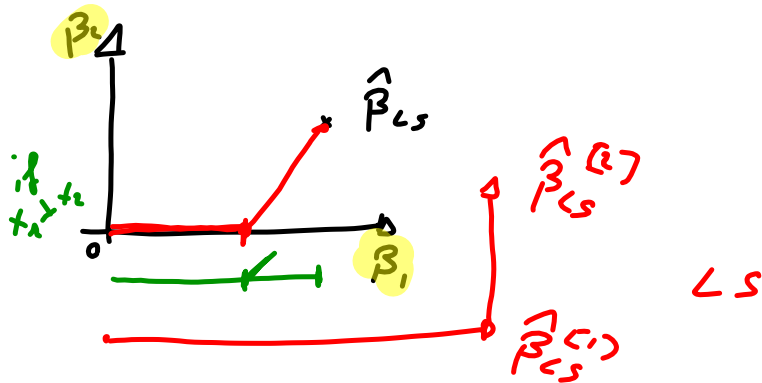
$$F = \frac{(RSS_0 - RSS_1) / (\overset{=1}{p_1 - p_0})}{RSS_1 / (N - p_1 - 1)}$$

→ find the  $\beta_j$  s.t.  $\beta_j = 0$  lead to the smallest  $RSS_0 - RSS_1$

we know (ex 3.1) ,  $F_{1, N-p-1} \stackrel{d}{=} \varepsilon_j^2$

⇒ the  $\beta_j$  which, when set equal to 0, increases the least the RSS is that with smallest  $\varepsilon_j^2$

## Least angle regression



we update  $\hat{\beta}_1$  until

$$\text{cor}(x_1, r) < \text{cor}(x_2, r)$$

then

we update both  $(\hat{\beta}_1, \hat{\beta}_2)$

$$r = y - \bar{y} - \sum_{j=1}^p x_j \hat{\beta}_j$$

- start with  $r = y - \bar{y}$
- check the largest correlation between  $x_j$  and  $r$
- update the regression coefficient of  $x_j$  ( $\hat{\beta}_j$ ) until  $\langle x_j, r \rangle$

- first we update  $\hat{\beta}_1$

Suppose that  $\beta_j$  are ordered by importance, i.e., effect of  $x_i$  is larger than effect of  $x_j$ , i.e.,

- we reach a point where we add also  $\hat{\beta}_2$  to our solution  
we are updating  $\{\hat{\beta}_1, \hat{\beta}_2\}$

- we reach a point where the update is related to  $\{\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3\}$

⋮

once a coefficient enters in the set of active regression coefficients, it stays

FORWARD REGRESSION

LAR

1<sup>st</sup> step  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1^{LS} x_1$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1^{LS} x_1$$

- less can be seen as a special case of a modification of LAR, for which once the estimate of a regression coeff. reach 0, it is excluded by the set of active regression coefficient.

## METHODS USING DERIVED INPUT DIRECTIONS

- Principal component regression
- Partial least squares

### Principal component regression

IDEA: inputs have different variabilities in different directions

→ directions with largest variability provide more information

↓ principal components  
linear combinations of  $X$  based on directions of largest variability

$$z_m = X V_m \quad \text{eigen vectors}$$

$z_1$  = direction with the largest variability

$z_2$  = " " the 2<sup>nd</sup> largest variability st.  $z_1 \perp z_2$

$\vdots$

$z_p$  = " " the  $p$ <sup>th</sup> largest variability  $z_p \perp z_1, \dots, z_{p-1}$

### Model

$$y = \theta_0 + \sum_{m=1}^M \theta_m z_m + \varepsilon$$

$$\hat{\theta}_0 = \bar{y}$$

$$\hat{\theta}_m = \frac{\langle z_m, y \rangle}{\langle z_m, z_m \rangle}$$

### NOTE

principal component analysis is scale dependent

↓

IMPORTANT to first standardize  $X$

for each  $X_j$

$$x_{ij}^{st} = \frac{(x_{ij} - \bar{x}_j)}{sd(X_j)}$$

$$\begin{aligned}
 \hat{y} &= \hat{\theta}_0 + \sum_{m=1}^M \hat{\theta}_m z_m \\
 &= \hat{\theta}_0 + \sum_{m=1}^M \hat{\theta}_m X v_m \\
 &= \hat{\theta}_0 + X \sum_{m=1}^M \hat{\theta}_m v_m \\
 &= \hat{\theta}_0 + X \hat{\beta}_{PCR}
 \end{aligned}$$

$$z_m = X v_m$$

if  $M=p$ , principal component regression provides the least square estimates, because all  $p$  principal components span the space of  $X$

- idea: use  $M < p$  principal components, to exclude directions with less information

we obtain results similar to ridge regression

$$\hat{\beta}_{\text{ridge}} = \sum_{j=1}^p u_j \frac{d_j^2}{d_j^2 + \lambda} u_j^T y$$

$$\hat{\beta}_{LS} = \sum_{j=1}^{\textcircled{p}} u_j u_j^T y$$

$$\hat{\beta}_{PCR} = \sum_{j=1}^p u_j \underbrace{\mathbb{1}[j \leq M]} u_j^T y = \sum_{j=1}^{\textcircled{M}} u_j u_j^T y$$

## Partial least squares

- in the construction of principal components, we do not take into account  $y$

↓  
- in the construction of derived input directions, we also consider  $y$

as for principal component regression it is important to first standardize  $X$

1<sup>st</sup> step:  $\hat{\phi}_{1j} = \frac{\langle x_j, y \rangle}{\langle x_j, x_j \rangle}$

1<sup>st</sup> pls direction  
 $\underline{e}_1 = \sum_{j=1}^p \hat{\phi}_{1j} x_j$

$\hat{\theta}_1 = \frac{\langle e_1, y \rangle}{\langle e_1, e_1 \rangle}$

$x_j^{(2)} = x_j - \frac{\langle e_1, x_j \rangle}{\langle e_1, e_1 \rangle} e_1$

$j = 1, \dots, p$

$\hat{\phi}_{2j} = \frac{\langle x_j^{(2)}, y \rangle}{\langle x_j^{(2)}, x_j^{(2)} \rangle}$

$\underline{e}_2 = \sum_{j=1}^p \hat{\phi}_{2j} x_j^{(2)}$

$\hat{\theta}_2 = \frac{\langle e_2, y \rangle}{\langle e_2, e_2 \rangle}$

Differences:

pcr: derived input directions are the principal components of  $X$ , constructed by looking at the variability within  $X$

pls: directions take into consideration both the variability and the correlation with  $y$

Mathematically

PCR

$\max_{\alpha} \text{Var}(X\alpha) \quad \text{s.t.} \quad \|\alpha\| = 1$

$\alpha^T S_V e = 0 \quad e = 1, \dots, m-1$

sample covariance matrix

PLS

$\max_{\alpha} \text{Cor}^2(y, X\alpha) \text{Var}(X\alpha)$

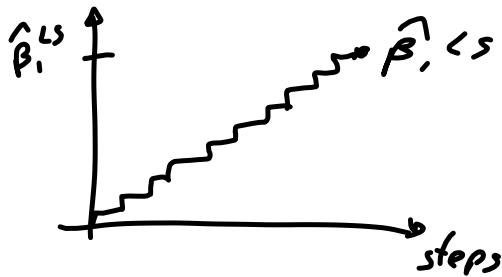
uncorrelated directions

In practice, the variance component tends to dominate the correlation one

s.t.  $\|\alpha\| = 1$

$\alpha^T S_{\hat{y}} e = 0 \quad e = 1, \dots, m-1$

→ similar results!

Incremental forward **stagewise** regression

similar to LAR, but each update involves only one parameter each time

$$\hat{\beta}_j^{(m)} = \hat{\beta}_j^{(m-1)} + \underbrace{\delta}_{\text{step}}$$

→ we will go back when we will talk about boosting

## → grouped LASSO

different construction of penalty to take into account structures in the data

- dummy variables related to the same categorical variable
- working with genetic data, group genes belonging to the same pathway

$$\hat{\beta}_{gL} = \argmin_{\beta} \left\{ \sum_{i=1}^n (y_i - \hat{\beta}_0 - \sum_{\ell=1}^L x_{i\ell} \beta_{\ell}) + \lambda \sum_{\ell=1}^L \sqrt{p_{\ell}} \|\beta_{\ell}\|_2 \right\}$$

$\uparrow$   
 group size

$\|\cdot\|_2$  = Euclidean distance  
 it is 0  $\Leftrightarrow$  all components are 0

- adaptive LASSO } → 2-step procedure  
 - relaxed LASSO }

- SCAD  $L(\beta) + J(\beta)$

$$\frac{dJ(\beta)}{d\beta} = \lambda \operatorname{sign}(\beta) \left[ \mathbb{1}(|\beta| \leq 1) + \frac{(a^2 - |\beta|)}{(a-1)\lambda} \mathbb{1}(|\beta| > 1) \right]$$