# Methods of regularization



— bias
— variance

add small bias, reduce variance

ordinary
- least squares estimator
  BLUE (Best Linear Unbiased
                    Estimator)
                         $\hat{\beta}$

  minimize the variance

accepting a small increase in the
bias in order to reduce the variance
to get the smallest possible prediction error

- model selection : reduce the variance by
  removing those variables that are not
  contributing much to the bias-reduction
- same idea for PCR

Penalized (regularized) regression → we add to our usual loss (deviance) a term
                                      (penalty) which force the estimator to be a little
                                      bit biased, but obtain smaller variance

Notes : - we now assume to have centred response $y_i - \bar{y}$, because we do not
          want to penalize the intercept
        - for reasons which will be clear soon, we will work with standardized X
          $$x_i^\# = \frac{x_i - \bar{x}}{se(x)}$$

# Ridge regression

Consider a linear model $y = X\beta + \varepsilon$
Ridge regression finds the regression coefficient estimates by minimizing

$$D_{ridge}(\beta, \lambda) = \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2 = \| y - X\beta \|_2^2 + \lambda \beta^T \beta$$

tuning parameter
(penalty parameter)

$L_2$ penalty

The minimizer of this loss function is

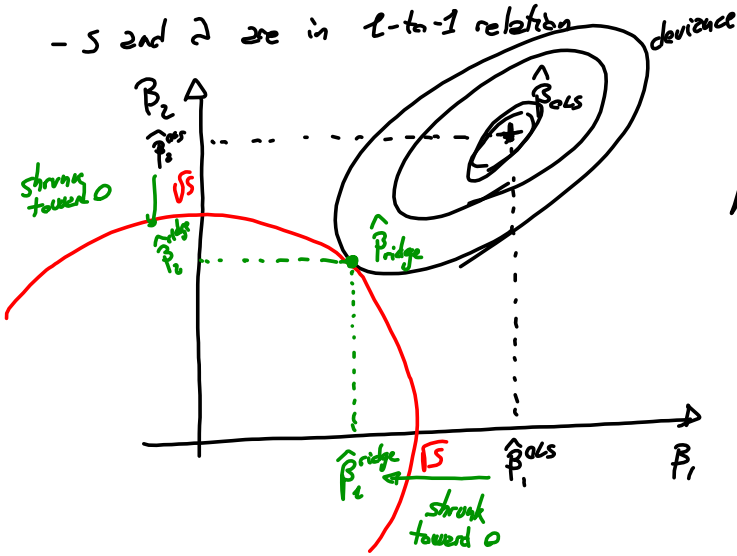$$\hat{\beta}_{ridge}(\lambda) = \left( X^T X + \lambda I \right)^{-1} X^T y$$

$\boxed{\lambda > 0}$

Note that $\lambda$ is a very important parameter that controls the amount of penalization
- $\lambda = 0$ → no penalty, so $\hat{\beta}_{ridge} = \hat{\beta}_{ols}$
- $\lambda = +\infty$ → each small deviation of $\hat{\beta}$ from 0 are strongly penalized, so $y - \bar{y} = 0$
        $\hat{\beta}_{ridge}^{(j)} = 0 \ \forall j$
- normally, it is computed by cross-validation
- even for small $\lambda$, we reduce problems of collinearity
- allow regression of the case $p > n$

Alternative formulation. ~~Ridge~~

Minimize $\sum_{i=1}^{n} (y_i - x_i^T \beta)^2$   subject to   $\sum_{j=1}^{p} \beta_j^2 \le s$

– $s$ and $\lambda$ are in 1-to-1 relation



simple case, two variables → two regression coefficients

Note: due to correlation, it is not necessarily true that

$$\lambda_a > \lambda_b \not\Rightarrow |\hat{\beta}_j(\lambda_a)| < |\hat{\beta}_j(\lambda_b)|$$
$$s_a < s_b$$

Bias

$E[\hat{\beta}_{ols}] = \beta$   ordinary least squares estimator is unbiased   BLUE

$E[\hat{\beta}_{ridge}] = E[(X^TX + \lambda I)^{-1} X^T y]$

$\phantom{E[\hat{\beta}_{ridge}]} = E[(X^TX + \lambda I)^{-1} (X^TX)(X^TX)^{-1} X^T y]$

$\phantom{E[\hat{\beta}_{ridge}]} = E[(I_p + \lambda(X^TX)^{-1})^{-1} (X^TX)^{-1} X^T y]$

$\underbrace{\phantom{= E[(I_p + \lambda(X^TX)^{-1})^{-1}}}_{W(\lambda)} \underbrace{\phantom{(X^TX)^{-1} X^T y]}}_{\hat{\beta}_{ols}}$

$\phantom{E[\hat{\beta}_{ridge}]} = \underline{W(\lambda) E[\hat{\beta}_{ols}] = W(\lambda)\beta}$

$(X^TX + \lambda I)^{-1}(X^TX)$
$([X^TX + \lambda I](X^TX)^{-1})^{-1}$

the ridge estimator is unbiased only if $W(\lambda) = I$, i.e., $\lambda = 0$

$Var(\hat{\beta}_{ridge}) = Var(W(\lambda) \hat{\beta}_{ols})$

$\phantom{Var(\hat{\beta}_{ridge})} = W(\lambda) Var(\hat{\beta}_{ols}) W(\lambda)^T$     $Var(\hat{\beta}_{ols}) = \sigma^2 (X^TX)^{-1}$

$\phantom{Var(\hat{\beta}_{ridge})} = \sigma^2 W(\lambda) (X^TX)^{-1} W(\lambda)^T$

$Var(\hat{\beta}_{ols}) - Var(\hat{\beta}_{ridge}) = \sigma^2 [(X^TX)^{-1} - W(\lambda)(X^TX)^{-1}W(\lambda)^T]$

$\phantom{Var} = \sigma^2 W(\lambda)[W(\lambda)^{-1}(X^TX)^{-1}(W(\lambda)^T)^{-1} - (X^TX)^{-1}] W(\lambda)^T$

$\phantom{Var} = \sigma^2 W(\lambda)[(I_p + \lambda(X^TX)^{-1})(X^TX)^{-1}(I_p + \lambda(X^TX)^{-1}) - (X^TX)^{-1}] W(\lambda)^T$
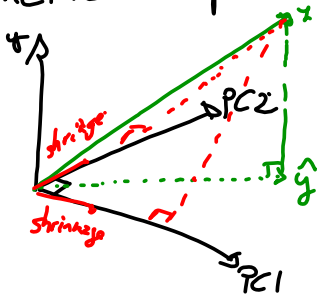
$\phantom{Var} = \sigma^2 W(\lambda)[\cancel{(X^TX)^{-1}} + \lambda(X^TX)^{-2} + \lambda(X^TX)^{-2} + \lambda^2(X^TX)^{-3} - \cancel{(X^TX)^{-1}}] W(\lambda)^T$

$\phantom{Var} = \sigma^2 W(\lambda)[2\lambda(X^TX)^{-2} + \lambda^2(X^TX)^{-3}] W(\lambda)^T \quad > 0$

all are quadratic form                    $Var(\hat{\beta}_{ridge}^{(\lambda)}) \le Var(\hat{\beta}_{ols})$
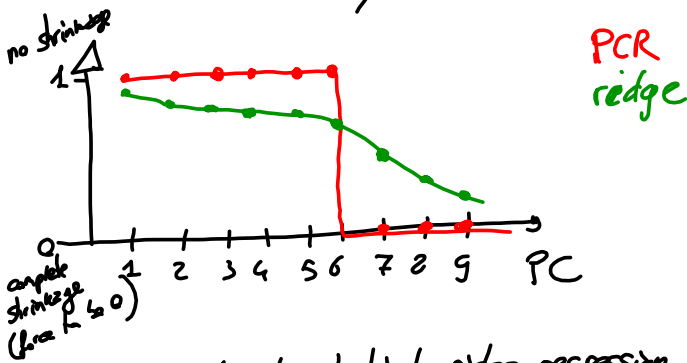
Geometric interpretation related to PCA



$$\hat{\beta}_1^{ols} = \frac{\|y\|}{\|pc1\|}$$

$$\hat{\beta}_1^{ridge} = \frac{\|y\|}{\|pc1\| + \lambda} \hat{\beta}_1^{ols}$$

$$\hat{\beta}_2^{ols} = \frac{\|y\|}{\|pc2\|}$$

$$\hat{\beta}_2^{ridge} = \frac{\|y\|}{\|pc2\| + \lambda} \hat{\beta}_2^{ols}$$

- ridge regression projects the response on the principal components
- shrink the <u>low-variance components</u> more than the <u>high-variance components</u>
   likely to be noise                    likely to be signal



PCR
ridge

The general idea behind ridge regression, is that we minimize a loss function of the form

$$\sum_{i=1}^{n} (y_i - x_i^T \beta)^2 + \lambda \underline{P(\beta)}$$

→ specifically for ridge regression, $P(\beta) = \beta^T \beta = \sum_{j=1}^{P} \beta_j^2$

$P(\beta)$ can but must not be the $L_2$ norm → we can use the $L_1$ norm
from the sum of squares to the sum of absolute values → $\sum_{j=1}^{P} |\beta_j|$

Minimize $\sum_{i=1}^{n} (y_i - x_i^T \beta)^2 + \lambda \sum_{j=1}^{P} |\beta_j|$
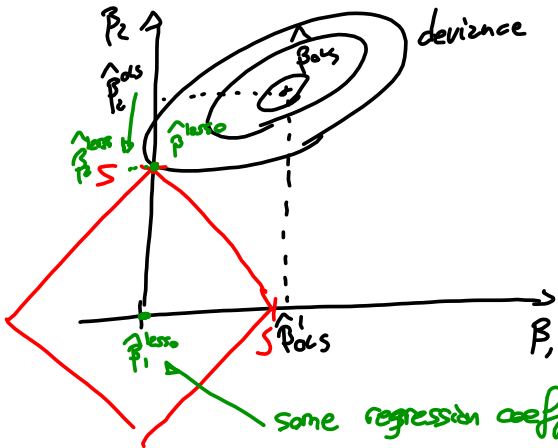
LASSO (Least Angle <u>Shrinkage</u> and <u>Selection Operator</u>)

The most interesting feature of LASSO is that it forces some estimates to be exactly
equal to O → intrinsic <u>variable selection</u>                → for a suitable $\lambda$
The same considerations done for $\lambda$ of ridge regression are valid for LASSO

Using the alternative form for LASSO as well

$$\hat{\beta}_{lasso} = \arg\min_{\beta} \sum_{i=1}^{n} (y_i - x_i^T \beta)^2 \quad \text{subject to} \quad \sum_{j=1}^{p} |\beta_j| \le s$$

deviance
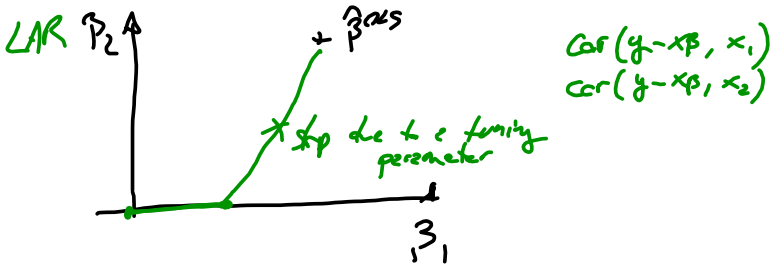
some regression coefficients estimates are forced to be 0

Advantages
- shrinkage
- automatic variable selection

Disadvantages
- no close form for $\hat{\beta}_{lasso}$ due to the non-differentiability of the $L_1$ loss
- we need to rely on numerical computations
  $\hookrightarrow$ in practice, very good algorithms based on LAR, of which lasso is a special case

LAR

$cor(y - x\beta, x_1)$
$cor(y - x\beta, x_2)$

stop due to a tuning parameter

$$\sum_{i=1}^{n} \left( y_i - \beta_0 + \sum_{j=1}^{p} x_j \beta_j \right) + \sum_{j=1}^{p} |\beta|$$

intercept is not penalized (no gain in variance by penalizing this term)

- variable selection
- ridge
- lasso

# Prediction of quantitative variables

Back to the initial problem, predict $y$ using $x$

$$y = f(x) \quad \longrightarrow \text{until now : parametric approach}$$
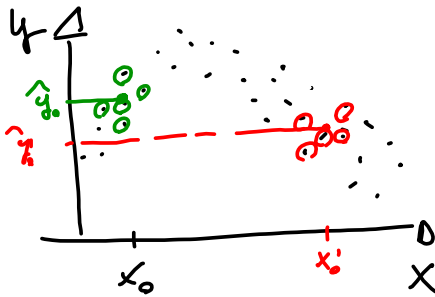
$$y = f(x ; \beta)$$

$$\uparrow$$

family of functions, indicized by $\boxed{\beta}$

$$\hat{y} = f(x ; \beta)\Big|_{\beta = \hat{\beta}}$$

- in the class of parametric functions, find the best by selecting the best $\hat{\beta}$
- simple, easy to compute

Alternative: do not restrict to a parametric form, base our estimation on the data only
(plus some regularity conditions) $\longrightarrow$ NON-PARAMETRIC APPROACH



$$\hat{y}_0 = \hat{f}(x_0)$$

$$K=5 = \frac{1}{5} \sum_{x_i \in N_5(x_0)} y_i$$

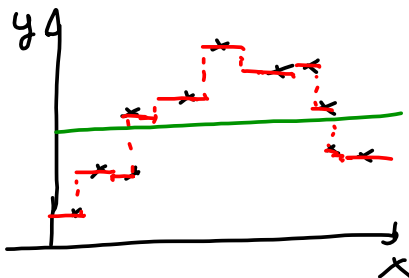$$\hat{y}_0 = \frac{1}{K} \sum_{x_i \in N_k(x_0)} y_i$$

where $N_k(x_0)$ denotes the set of the $K$ closest points to $x_0$

We assume that the values of the closest points are similar to the one of interest, and we base our estimate on those response $\rightarrow$ simple mean

Basically, we are assuming that $f(x)$ is constant close to $x_0$     $f(x) = \beta_0$

$K$ is the tuning parameter, that tells how many neighbours to include in the estimate
$\rightarrow$ smaller value, more complex model, until the extreme $K=1$, in which each $x_0$ is
estimated though its closest neighbour
$\rightarrow$ larger values, simpler model, until the extreme $K=n$, in which $x_0$ is estimated with
all the observations : $f(x) = \bar{y}$



$k=1$
$k=n$

- complexity is inverse of $K$
- how to choose $K$
  - CROSS-VALIDATION