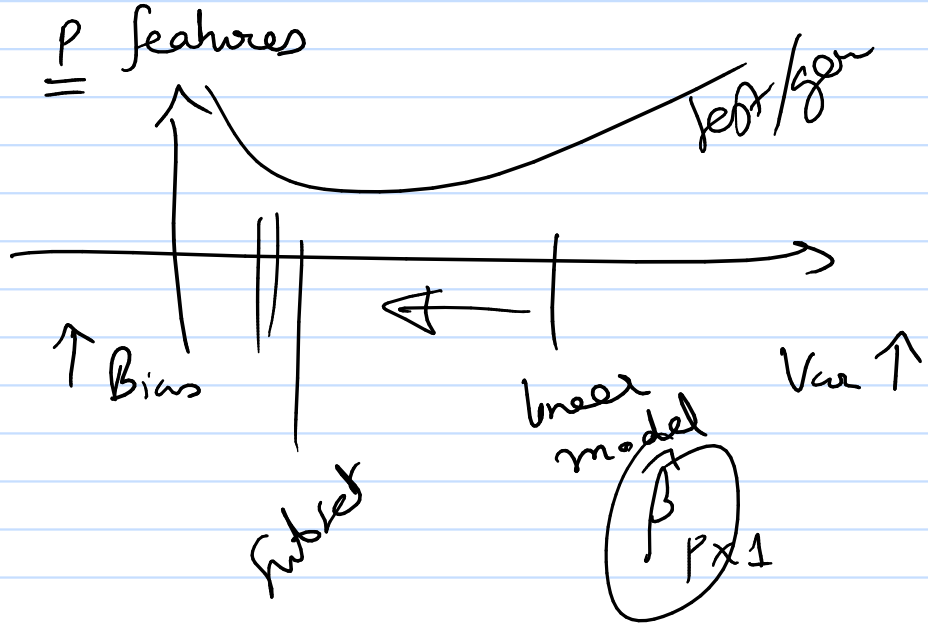


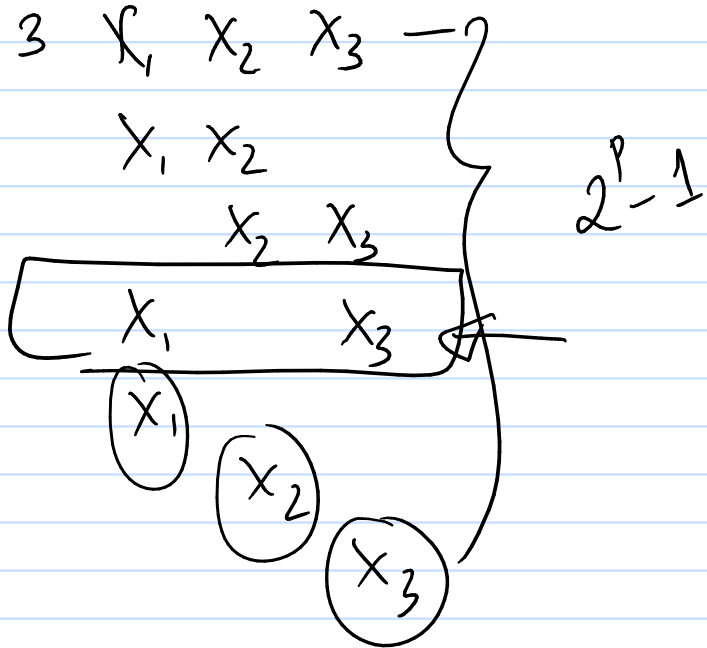
# ① Subset selection : why?



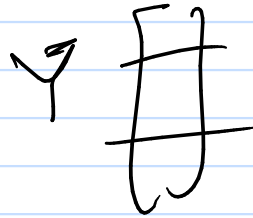
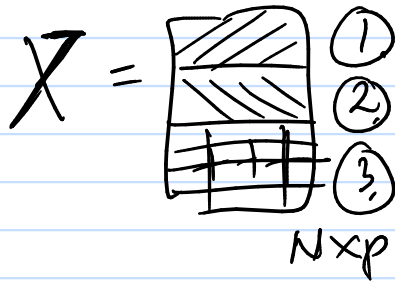
$p$  features

1, 2, 3, ...,  $p$   
↑  
 $p$   $\binom{p}{2}$  . . . . . 1

How to find the best subset



Recall



for each choice:

for each fold  $i$ :

get a model using remaining folds

error (fold  $i$ )

estimated  
EPE = Sum (errors) across folds  $\times \frac{1}{3}$

① LASSO

Criteria: RSS( $\beta$ ) +  $\lambda \cdot \sum_{j=1}^p |\beta_j|$

$10^{-2} = 0.01$

$\|\beta\|_1$

minim subset selection

$$\beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4$$

$\uparrow$   $\uparrow$   
0 0

② Ridge regression

$$RSS(\beta) + \lambda \sum_{j=1}^p \beta_j^2$$

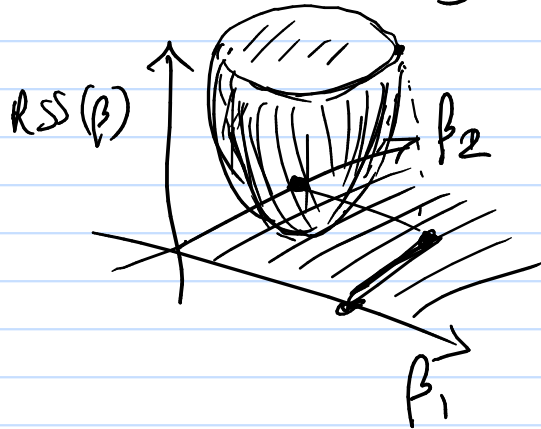
$\|\beta\|_2^2$

LASSO

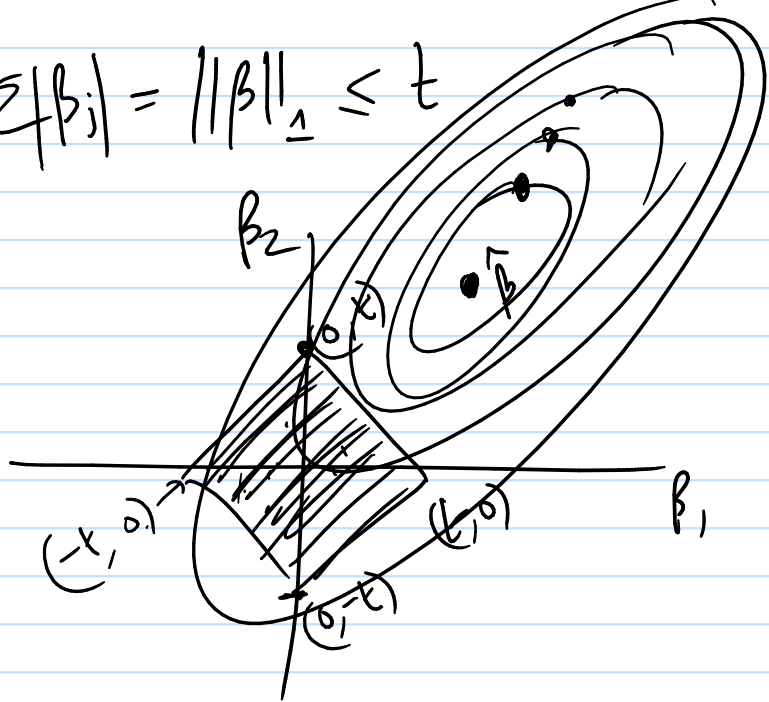
$$\min_{\beta} \text{RSS}(\beta) \\ \text{st. } \|\beta\|_1 \leq t$$

RR  $\min_{\beta} \text{RSS}(\beta)$   
st.  $\|\beta\|_2^2 \leq t$

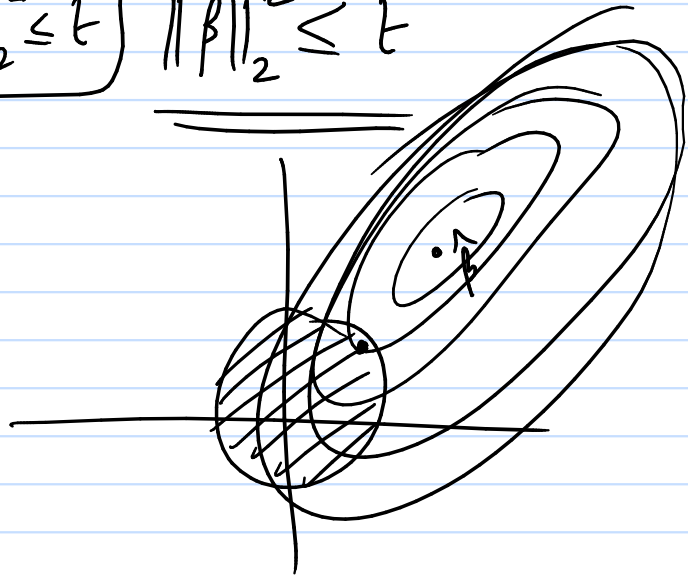
$\beta_1$   $\beta_2$   $\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$   $p=2$

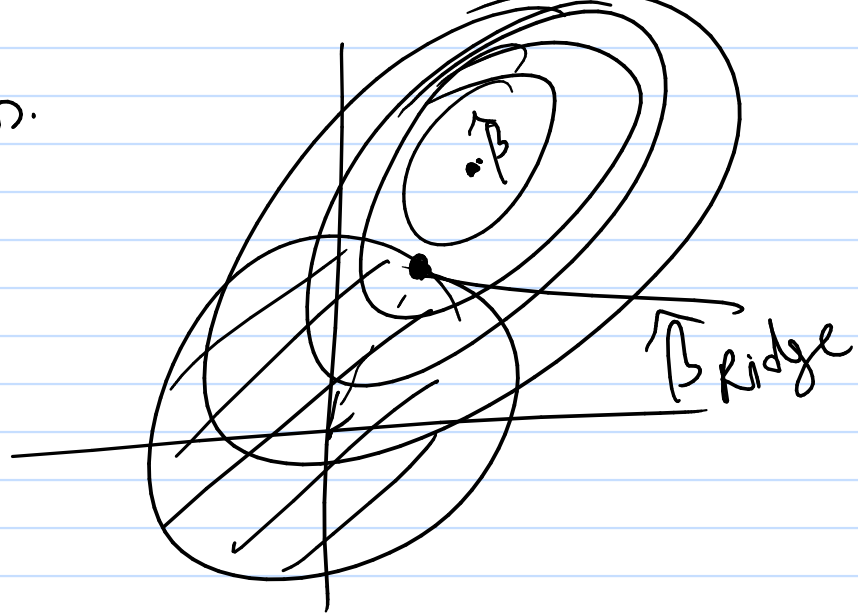
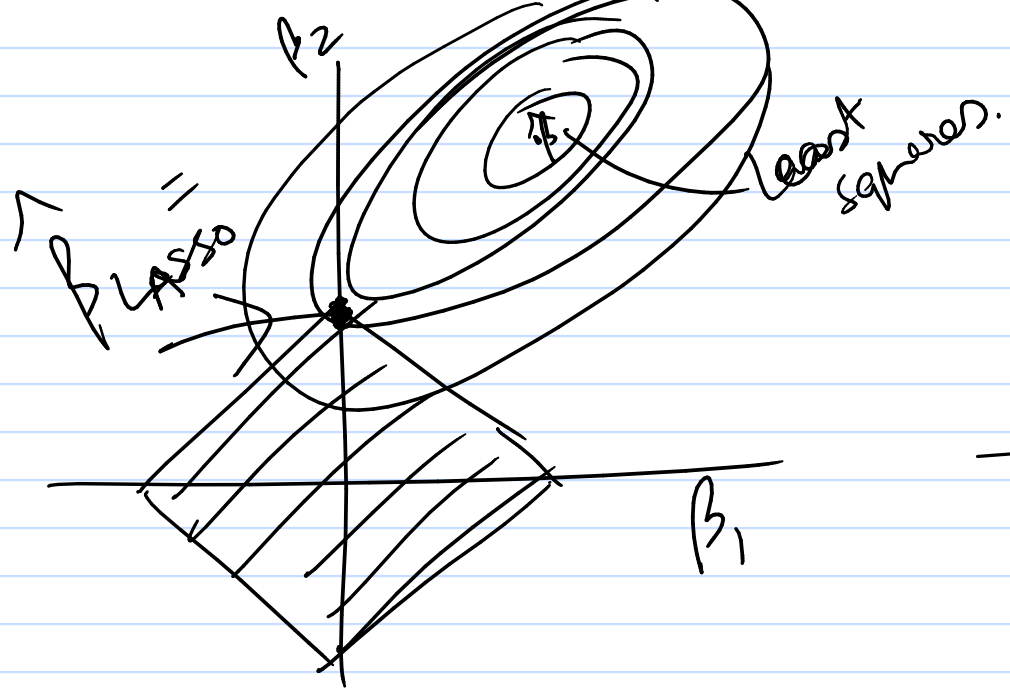


$$\sum |\beta_j| = \|\beta\|_1 \leq t$$



$$\boxed{\beta_1^2 + \beta_2^2 \leq t} \quad \underline{\underline{\|\beta\|_2^2 \leq t}}$$





$$P(B|A) = \frac{P(B \cap A)}{P(A)}$$

$$P(Z|W) = \frac{P(Z, W)}{P(W)}$$

$$P(Z=z | W=w) = \frac{P(W=w, Z=z)}{P(W=w)}$$

Bayes rule:

$$P(Z|W) \cdot P(W) = P(W|Z) \cdot P(Z) = P(Z, W)$$

$$P(Z|W) = \frac{P(W|Z) \cdot P(Z)}{P(W)}$$



$\beta \sim \text{RV}$  Prior:  $\begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} \sim \mathcal{N}\left(0, \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}\right)$   $X$  fixed

Likelihood:  $Y | \beta; X \sim \mathcal{N}(\beta^T X, \sigma^2)$

Then

$P(\underbrace{Y_1=y_1, Y_2=y_2, \dots, Y_N=y_N}_{\text{data}} | X) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - \beta^T x_i)^2}{2\sigma^2}\right)$

$P(\beta | \underbrace{Y_1=y_1, Y_2=y_2, \dots, Y_N=y_N}_{\text{data}}) \leftarrow \underline{\text{Normal}}$

Ridge:  $\min_{\beta} \sum_{i=1}^N (y_i - \beta^T x_i)^2 + \lambda \cdot \sum_{j=1}^p \beta_j^2$

(MAP)

$P(\beta | \text{data}) \propto P(\text{data} | \beta) \cdot P(\beta)$  : Maximum a-posteriori estimate

$\rightarrow \max_{\beta} P(\text{data} | \beta) \cdot P(\beta)$

$\min_{\beta} \log P(\text{data} | \beta) + \log P(\beta)$

$\log \frac{1}{(\sqrt{2\pi})^2} \exp\left(-\frac{1}{2\sigma^2} (\beta_1^2 + \beta_2^2)\right)$

$= \lambda \sum_{j=1}^p \beta_j^2 + \dots$

$$\sum_{i=1}^N (y_i - \beta^T x_i)^2 \propto -\log P(\text{data} | \beta)$$

RSS( $\beta$ )  $\longleftrightarrow$  maximizing likelihood

min  $-\log$  likelihood.

$$-\log P(\text{data} | \beta) = -\log \left[ P(Y_1 = y_1 | \beta) \times P(Y_2 = y_2 | \beta) \cdot \dots \times P(Y_N = y_N | \beta) \right]$$

$$\log(a \times b) = \log a + \log b$$

$$= -\log P(Y_1=y_1|\beta) - \log P(Y_2=y_2|\beta) - \dots$$

$$= \sum_{i=1}^N -\log P(Y_i=y_i|\beta) \quad \frac{1}{\sqrt{2\pi\tau}} \exp\left(-\frac{(y_i-\beta^T x_i)^2}{2\tau^2}\right)$$

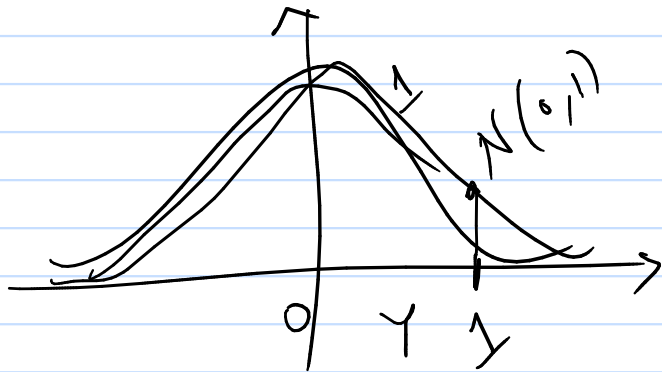
The term  $P(Y_i=y_i|\beta)$  is circled. The Gaussian function is decomposed into two parts:  $\frac{1}{\sqrt{2\pi\tau}}$  (labeled 'a') and  $\exp\left(-\frac{(y_i-\beta^T x_i)^2}{2\tau^2}\right)$  (labeled 'b').

$$\propto \sum_{i=1}^N \frac{(y_i - \beta^T x_i)^2}{2\tau^2}$$

The entire expression is circled.

$$\log \frac{1}{\sqrt{2\pi\tau}} - \frac{(y_i - \beta^T x_i)^2}{2\tau^2}$$

The first part of the decomposition is circled.



$$\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right) = P(Y=y)$$

$$P(Y=1) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1^2}{2}\right)$$