

POPULATION

SAMPLE

$$Y = \beta' X + \epsilon$$

$$y_i = \tilde{\beta}' \tilde{X}_i + \tilde{\epsilon}_i$$

$$E(Y^2) = E(\beta' X)^2 + E(\epsilon^2)$$



$$E_n y_i^2 = E_n (\beta' \tilde{X}_i)^2 +$$

$$E_n \tilde{\epsilon}_i^2$$

$$\underline{MSE}_{\text{pop}} = E \epsilon^2$$



$$\underline{MSE}_{\text{sample}} = \frac{1}{n} \sum \tilde{\epsilon}_i^2$$

$$E(Y - \beta' X)^2$$



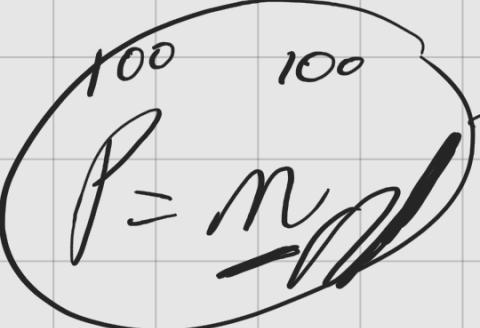
$$\underline{R^2}_{\text{pop}} = 1 - \frac{E \epsilon^2}{E Y^2}$$

$$\underline{R^2}_{\text{sample}} = \frac{1 - \frac{1}{n} \sum \tilde{\epsilon}_i^2}{\frac{1}{n} \sum y_i^2}$$

$$\in [0, 1]$$

$$\in [0, 1]$$

$$X \sim N(0, I_p) \quad Y \sim N(0, 1)$$

If $P = n$  \rightarrow OVERFITTING $R^2_{\text{sample}} = 1 \checkmark$

$$P = \frac{n}{2}$$

$$R^2_{\text{sample}} \approx 0,5$$

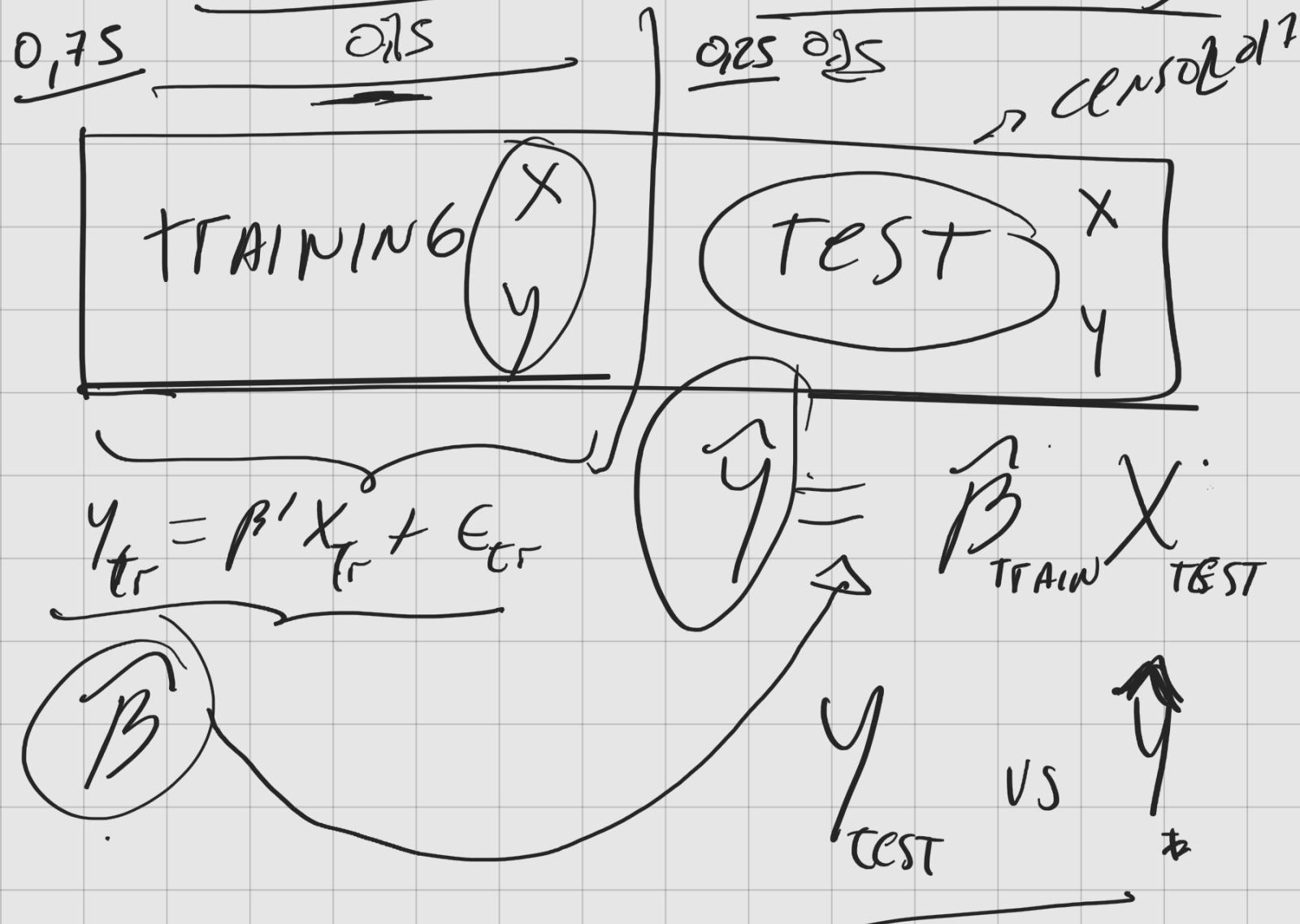
$$P = \frac{n}{20}$$

$$R^2_{\text{sample}} \approx 0,05$$

$$\text{MSE}_{\text{Adjusted}} = \frac{1}{n-p} \sum_{i=1}^n \hat{e}_i^2$$

$$R^2_{\text{Adjusted}} = 1 - \frac{\sum_{i=1}^n \hat{e}_i^2}{\sum_{i=1}^n y_i^2}$$

* SAMPLE SPLITTING



- ① ~~Sample~~ SPLITTING
- ② Training data $\rightarrow f(x)$

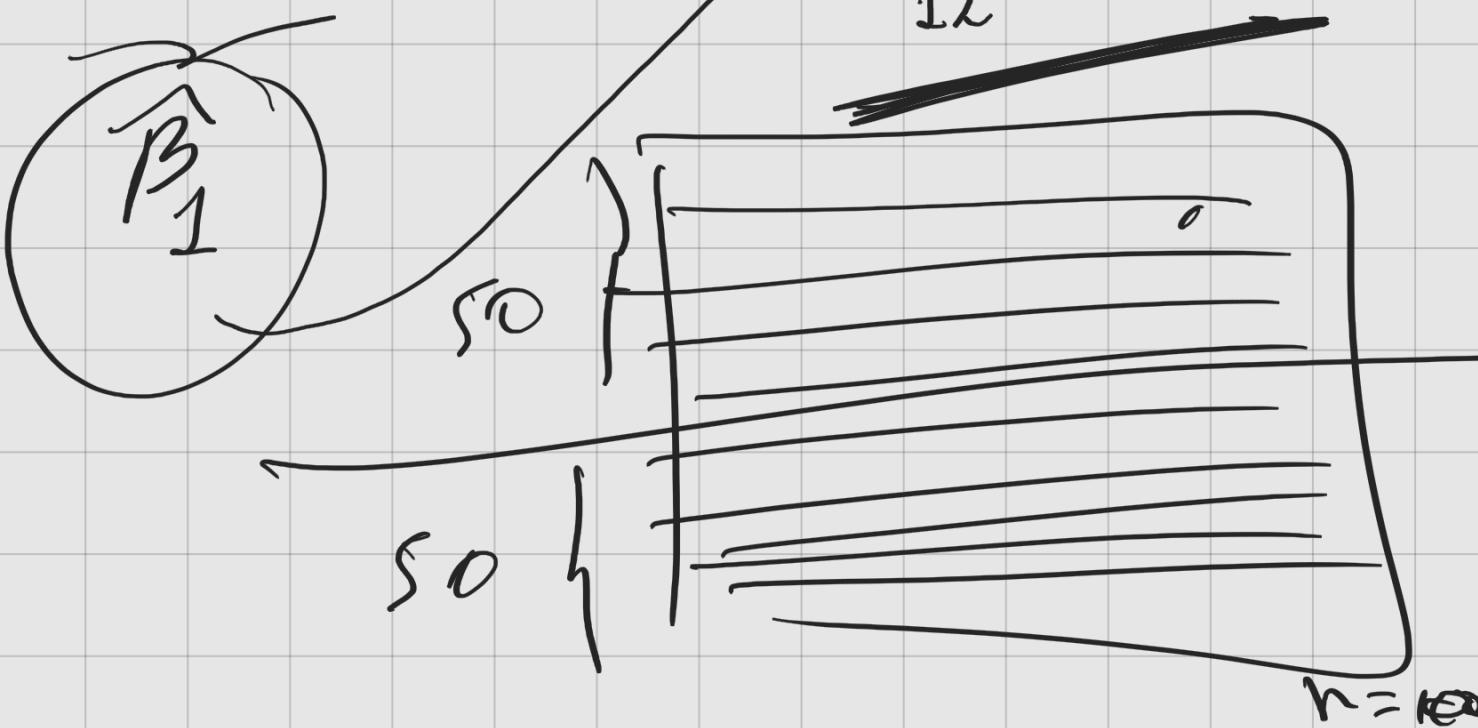
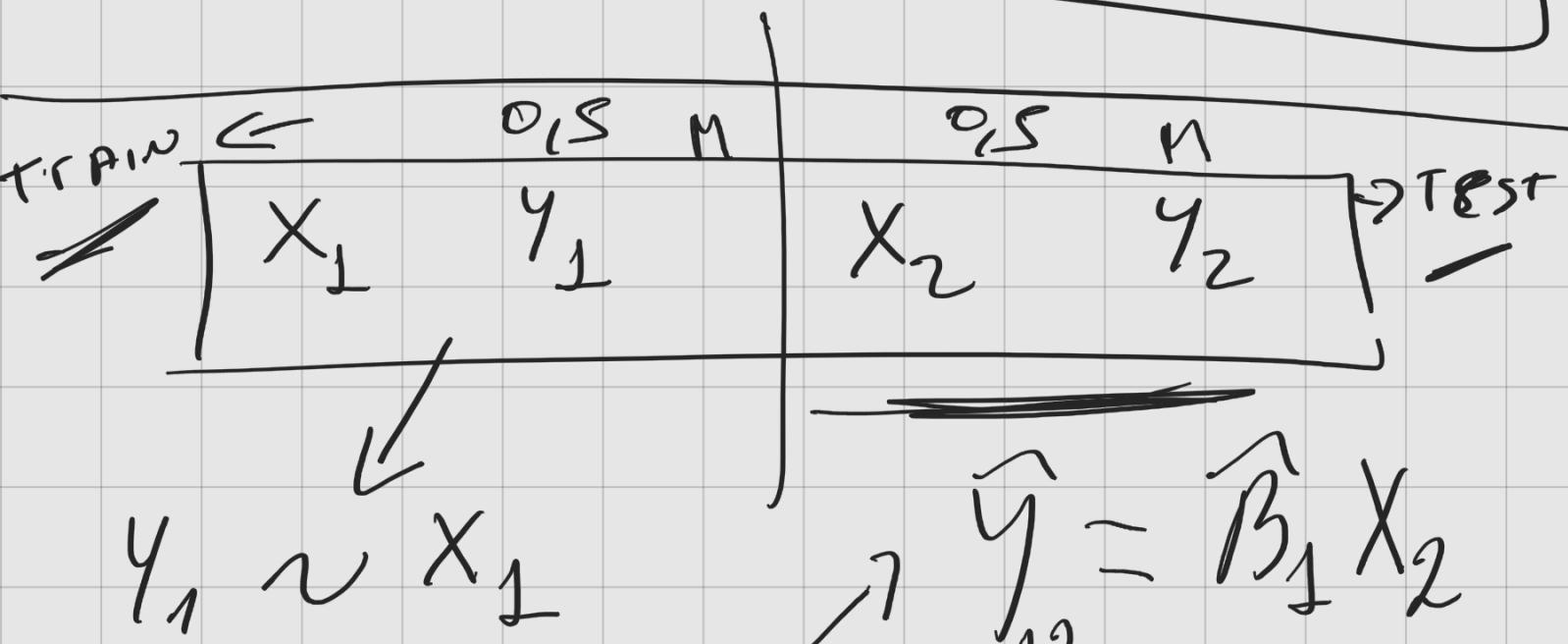
$$f(x) = \beta' x$$

$$\text{MSE}_{\text{test}} = \frac{1}{n} \sum_{k \in V} (y_k - \hat{f}(x_k))^2$$

\Rightarrow ~~test~~ ^{test} _{sample}

$MSE_{test} \Rightarrow$ OUT OF SAMPLE MSE

$$R^2_{test} = 1 - \frac{MSE_{test}}{\frac{1}{m} \sum_{k=1}^m y_k^2}$$



$$MSE_{OS} = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_{i2})^2$$

$$R^2_{OS} = 1 - \frac{MSE_{OS}}{\frac{1}{m} \sum (y_i)^2}$$

HDS / LASSO

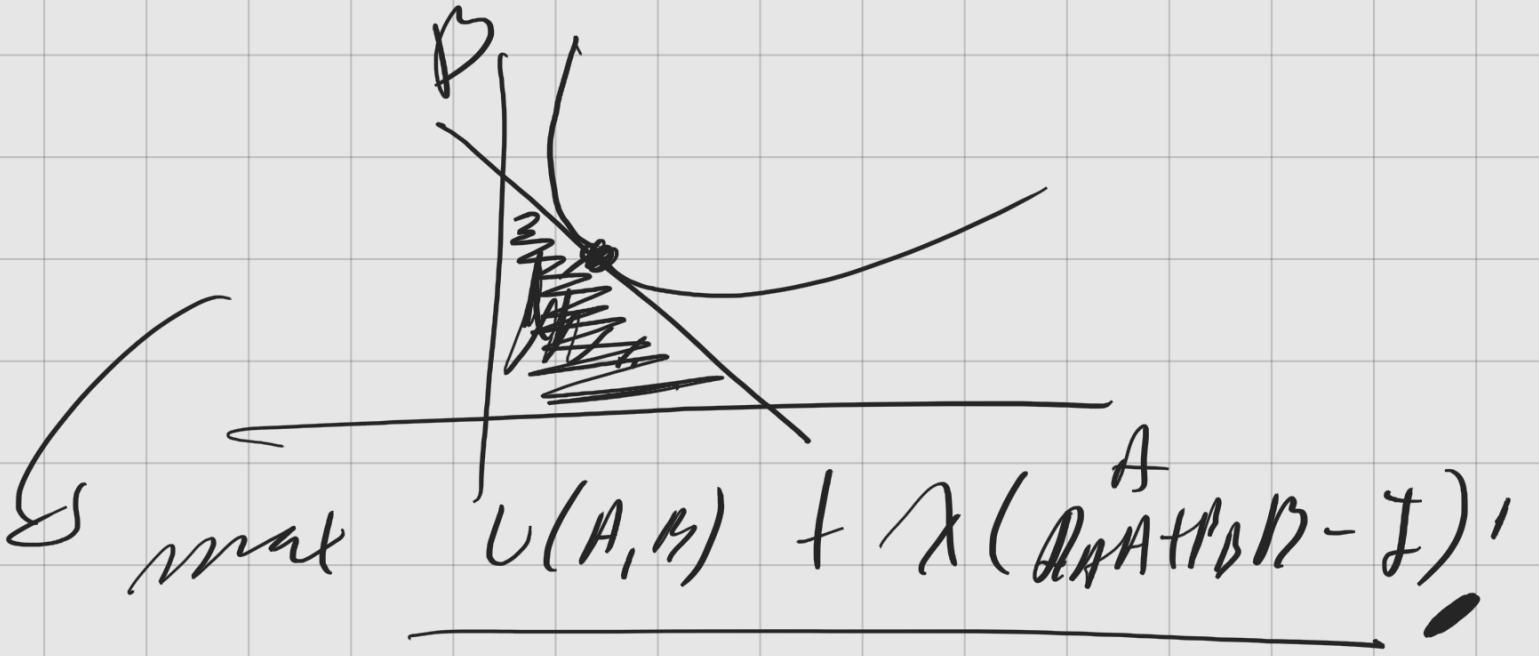
$$* \min_{b \in \mathbb{R}^P} \frac{\sum_i (y_i - b' x_i)^2}{1} + \lambda \frac{\sum_{j=1}^P |b_j|}{2}$$

① MSE

$$\textcircled{2} \quad \lambda \sum_{j=1}^P |b_j|$$

$$* \min \delta(K, l) \text{ S.t } y = f(K, l)$$

$$* \max U(A, B) \text{ s.t. } J = P_A A + P_B B$$



$$\min \sum_i (y_i - b' x_i)^2 \quad \text{S.T.} \quad \textcircled{b}$$

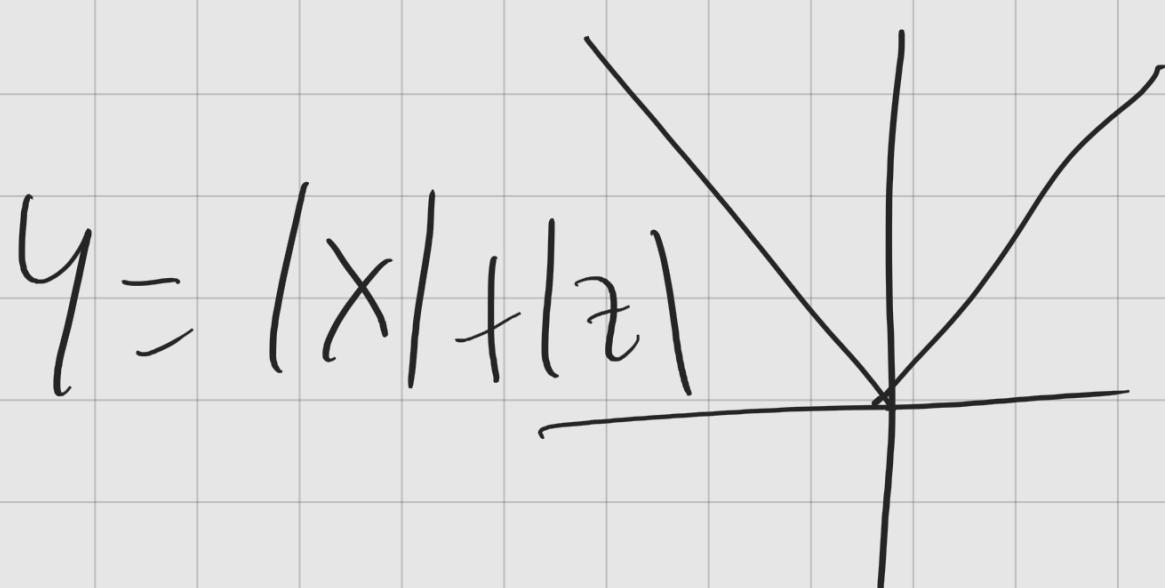
$$\sum_{j=1}^P |b_j| \leq S$$

$$\min_b \sum_i (y_i - b' x_i)^2 + \lambda \left[\sum_{j=1}^P |b_j| - S \right]$$

$$\textcircled{D} \quad \min \sum_i (y_i - b' x_i)^2 + \lambda \sum_{j=1}^P |b_j| \quad \textcircled{1}$$

$$x_1, x_2 = y$$

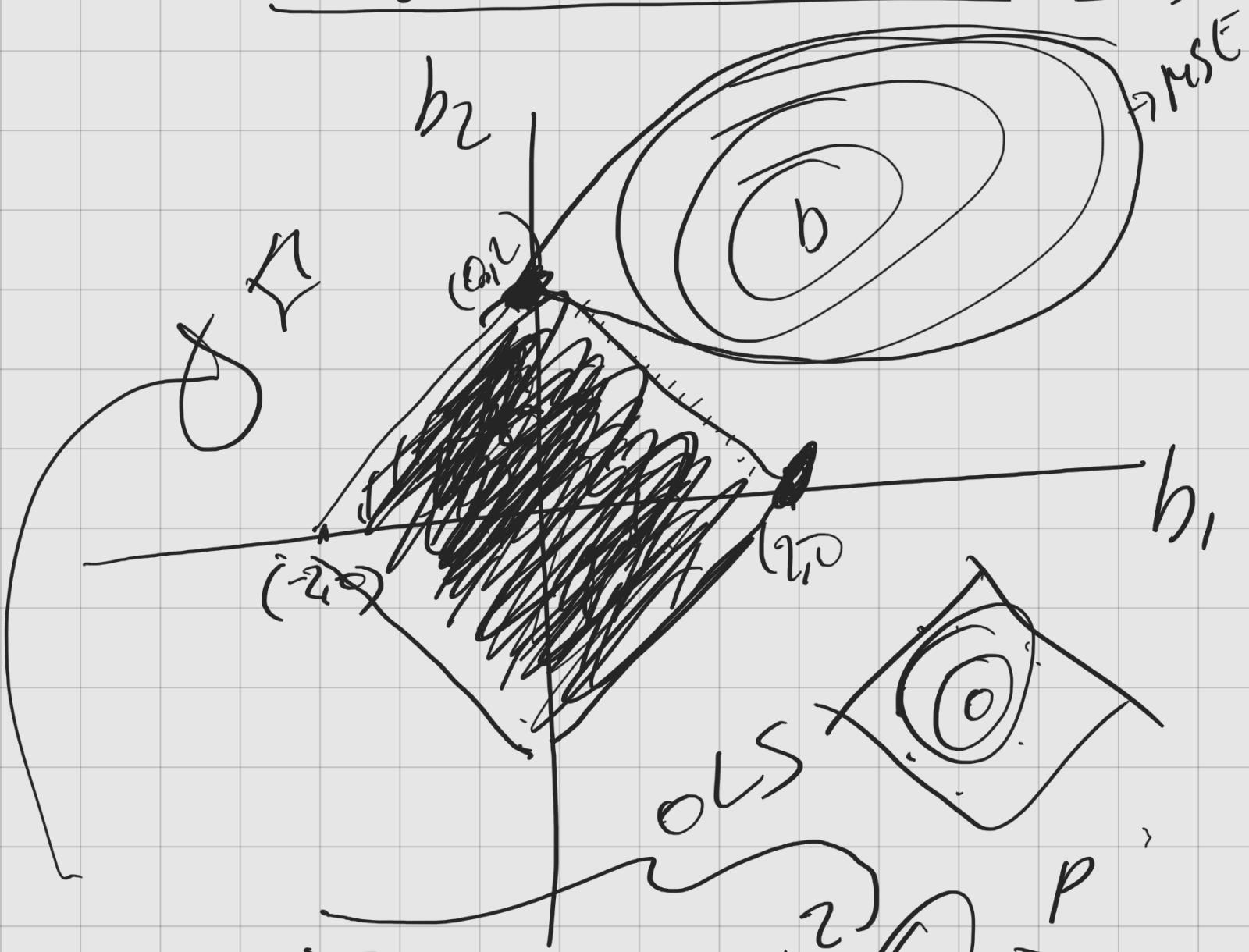
②



$$\sum_{j=1}^p |b_j|$$

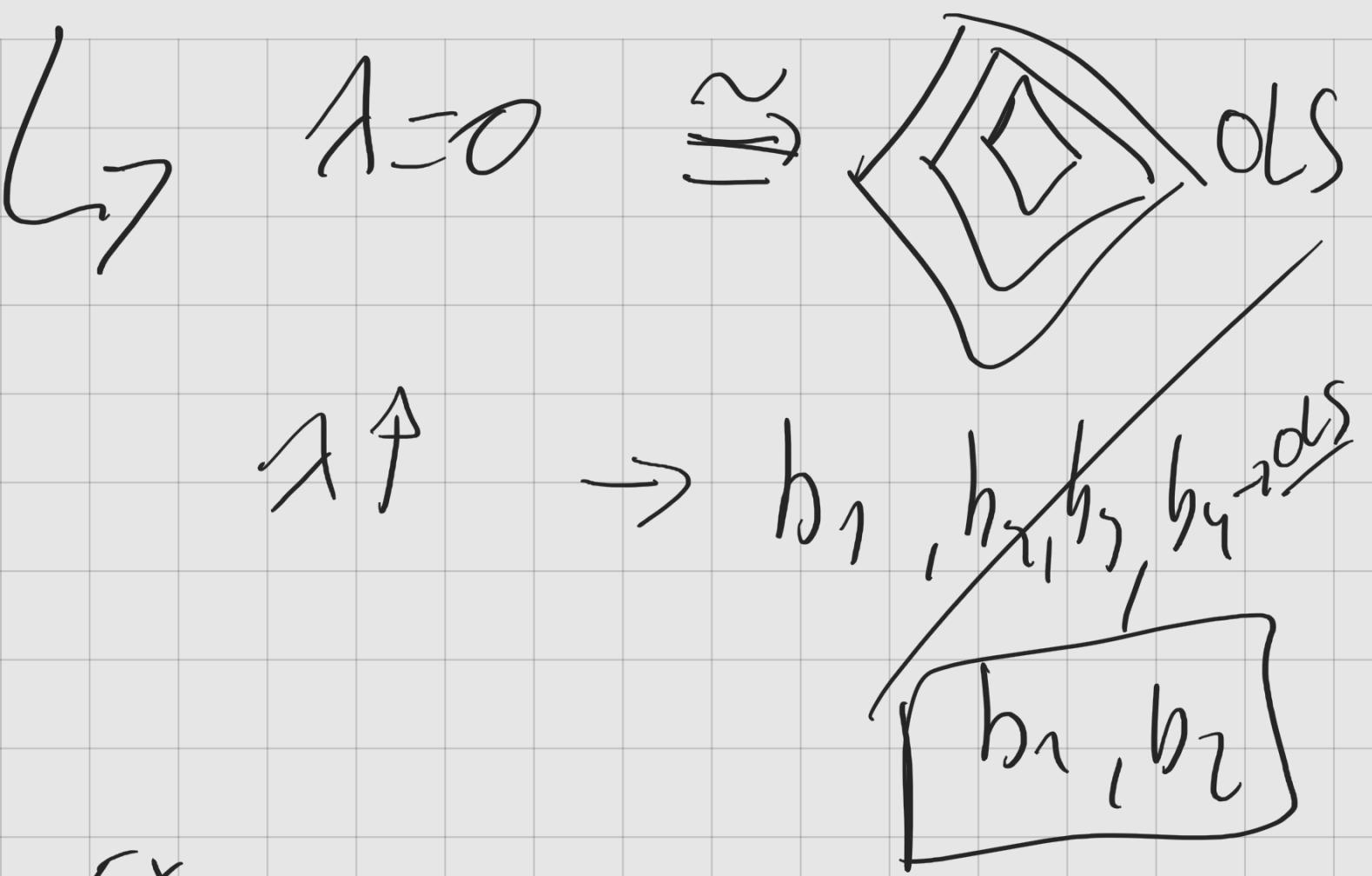
$$P=2$$

$$\Rightarrow |b_1| + |b_2| \leq 2$$



$$\min_{b \in \mathbb{R}^p} \sum_i (y_i - b' x_i)^2 + (\lambda) \sum_{j=1}^p |b_j|$$

CASSO



Ex-
 $\hat{Y} = \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \dots + \hat{\beta}_n X_n$

$\hat{Y} = 1,5X_1 + 2,3X_2 + \dots + 0,07X_n$

↓
 $0,02X_{n-1}$

↑
 $0,07$

$$Ab \rightarrow O \rightarrow \begin{cases} 0,007 \\ 0,002 \end{cases} \quad \boxed{\text{OLS}}$$

$$\begin{cases} n = 100 \\ p = 100 \end{cases} \xrightarrow{\text{OLS}} \text{OVERF.} \quad \text{h} = 0,003$$

$$\begin{cases} \text{QUASSO} \\ \hookrightarrow P = 25 \end{cases}$$

* B, CH \rightarrow 2013 *

$$H = 2c \bar{\sigma} \sqrt{2n \log(2P/\delta)}$$

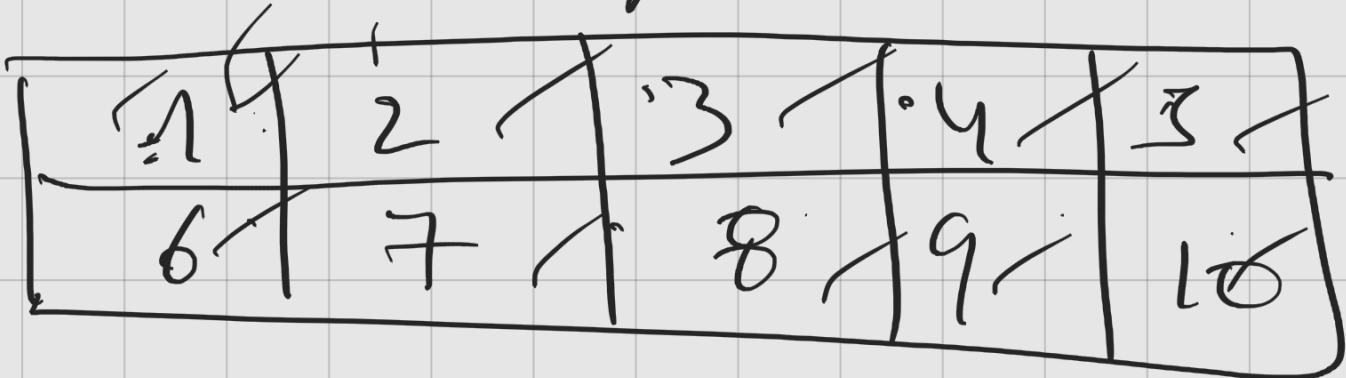
$$\bar{\sigma} = \sqrt{E\epsilon^2}$$

$$c = 1,1$$

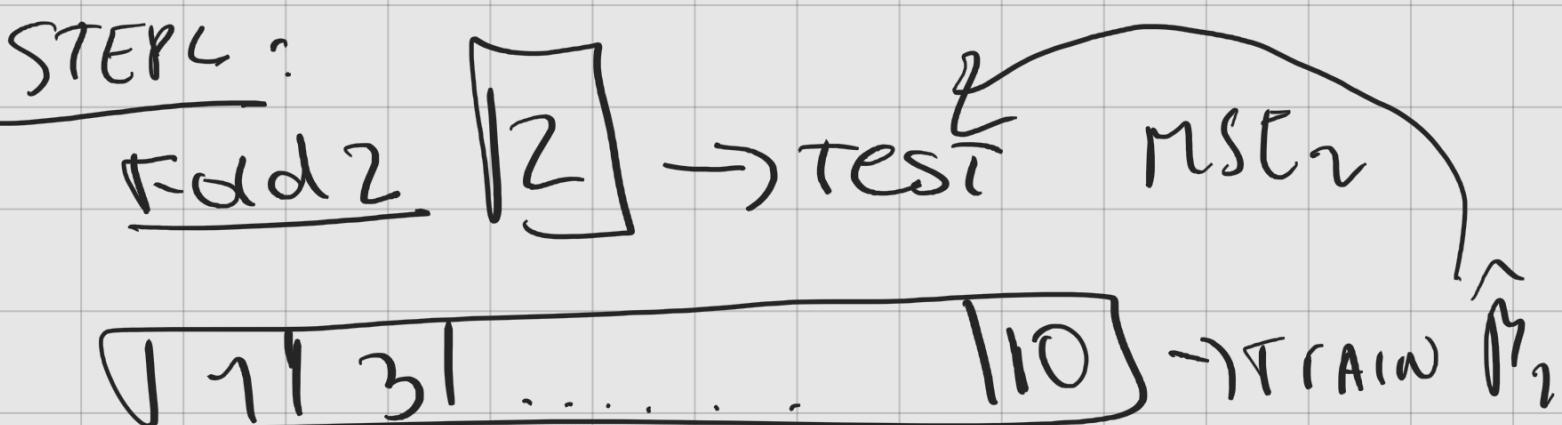
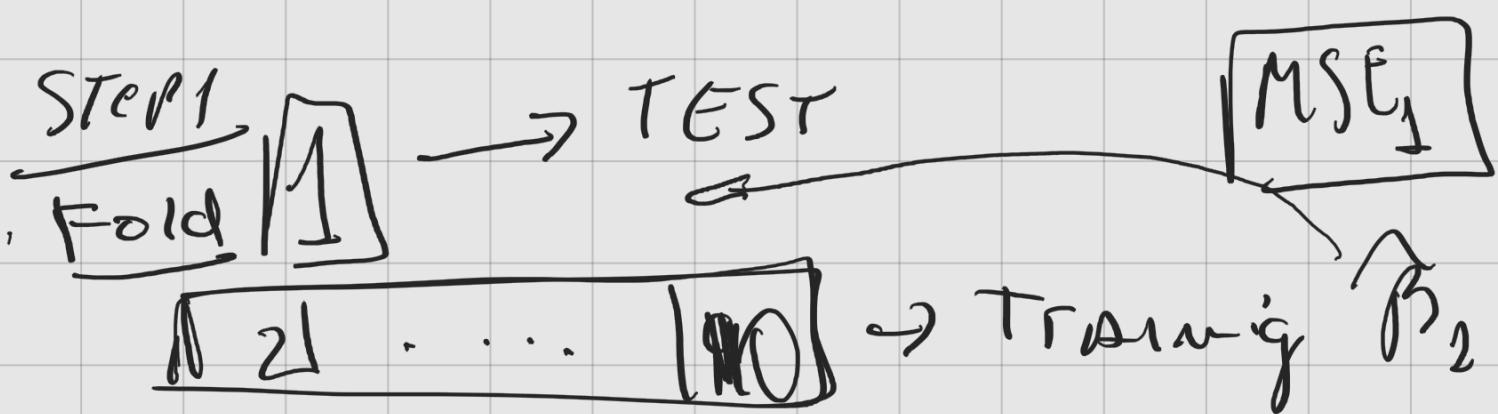
$$\delta = 0,05 \rightarrow \underline{\text{C.J}}$$

K-Fold Cross Validation

$K = \# \text{ Samples}$



$$k = \{0,01; 0,02; 0,03; \dots; 0,09\}$$



$$\boxed{\text{STEP 10}} \rightarrow \sqrt{\text{MSE}_{10}}$$

Ridge

$$\sum (y_i - \beta' x)^2 + \lambda_1 |b_i|^2$$

LASSO

$$\sum (y_i - \beta' x)^2 + \lambda_2 |b_i|$$

$$\sum (y_i - \beta' x)^2 + \alpha \lambda_1 |b_i|^2 +$$

$$(\alpha) \lambda_2 |b_i|$$

$$|b_1| + |b_2| + \dots + |b_n| \leq s$$

$$\underbrace{|b_1|^2 + |b_2|^2 + \dots + |b_n|^2}_{\leq s^2}$$

