

# Normal Equations & Geo Interpretations

## 1. Simple linear regression model

Data  $(x_i, Y_i)$ ,  $i=1, \dots, n$

random  $\mathbb{E}[e_i] = 0$

model:  $Y_i = \beta_0 + \beta_1 x_i + e_i$

↑  
dependent var./observation

↑  
intercept

↑  
indep. var./predictor

↑  
noise: measurement noise, biological variation

$\theta = [\beta_0, \beta_1]^T$  is the param vector/weights

$\mathbb{E}[Y_i] = \beta_0 + \beta_1 x_i = \text{lin. comb. of unknowns } \beta\text{'s,}$   
w/ known coeff  $(1, x_i)$ .

$$\underset{\sim}{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}_{n \times 1} \quad X = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}_{n \times 2}, \quad \underset{\sim}{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}_{2 \times 1}, \quad \underset{\sim}{e} = \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix}_{n \times 1}$$

$\underbrace{\quad}_{\underset{\sim}{1}} \quad \underbrace{\quad}_{\underset{\sim}{x}}$

$$\underset{\sim}{Y} = X \underset{\sim}{\beta} + \underset{\sim}{e} \quad \text{"Matrix-vector form"}$$

$\uparrow$   
 data matrix

$$E[\underset{\sim}{Y}] = X \underset{\sim}{\beta} = \begin{bmatrix} \underset{\sim}{1} & \underset{\sim}{x} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \beta_0 \underset{\sim}{1} + \beta_1 \underset{\sim}{x}$$

## 2. (Multiple) Linear regression model

$$y_i = \sum_{j=1}^p x_{ij} \beta_j + e_i, \quad i=1, \dots, n.$$

$$\underset{\sim}{Y}_{n \times 1} = X_{n \times p} \underset{\sim}{\beta}_{p \times 1} + \underset{\sim}{e}_{n \times 1} \quad \text{vector of rand. elements.}$$

Common Assumptions:

$$\mathbb{E}[\underset{\sim}{e}] = \underset{\sim}{0}$$

$$\text{Var}(e_i) = \sigma^2, \quad i=1, \dots, n.$$

$$\mathbb{E}[e_i^2] = \text{Var}(e_i) = \sigma^2$$

$$\text{Cov}(e_i, e_{i'}) = \mathbb{E}[e_i e_{i'}] = 0, \quad \forall i \neq i'$$

$$\text{VarCov}(\underset{\sim}{e}) = \sigma^2 \mathbb{I} = \mathbb{E}[\underset{\sim}{e} \underset{\sim}{e}^T]$$

$$\text{VarCov}(\underset{\sim}{Y}) = \text{VarCov}(\underset{\sim}{e}) = \sigma^2 \mathbb{I}$$

$$\text{VarCov}(u) = \mathbb{E}[(u - \mathbb{E}u)(u - \mathbb{E}u)^T]$$

$\int_{n \times 1} \xrightarrow{\quad} 1 \times n$

Ex: (Regression in narrow sense)  $Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + e_i, i=1, \dots, 50$ .

$Y_i$ : grade

$x_{i1}$ : time spent on hw

$x_{i2}$ : time spent on review

$$\begin{bmatrix} Y_1 \\ \vdots \\ Y_{50} \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} \\ \vdots & \vdots & \vdots \\ 1 & x_{50,1} & x_{50,2} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} + \begin{bmatrix} e_1 \\ \vdots \\ e_{50} \end{bmatrix}$$

Ex: [Analysis of Variance (ANOVA)]  $Y_{ij} = \mu + \alpha_i + e_{ij}, i=1, 2, 3; j=1, 2$ .

$$\underset{\sim}{Y} = \begin{bmatrix} y_{11} \\ y_{12} \\ \hline y_{21} \\ y_{22} \\ \hline y_{31} \\ y_{32} \end{bmatrix} \quad X = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ \hline 1 & 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 \\ \hline 1 & 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 1 & 1 \end{bmatrix} \quad \underset{\sim}{\beta} = \begin{bmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{bmatrix}$$

$\mu$  = "overall effect"  
 $\alpha_i$  = "effect of  $i^{\text{th}}$  treatment"

$$\hat{\theta}_1 = (\alpha_1 - \alpha_2) ?$$

$$\hat{\theta}_2 = (\alpha_1 - (\alpha_2 + \alpha_3)/2)$$

rank(X) = 3

"Linear (independent)", "vector space", "basis", "column/row rank".

# Linear Algebra Review

Hayes 2.3.2 ; Scheffe App. I

Given  $\{\underline{v}_1, \dots, \underline{v}_n\}$ .  $\alpha_1 \underline{v}_1 + \dots + \alpha_n \underline{v}_n = \underline{0} \Rightarrow$

- $\alpha_i = 0, \forall i$ : "linearly independent"
- not all  $\alpha_i$  are 0: "linearly dependent".

For "lin dependent" case, may write

$$\underline{v}_1 = \beta_2 \underline{v}_2 + \dots + \beta_n \underline{v}_n$$

Ex:  $\underline{v}_1 = [1 \ 2 \ 1]^T$   $\underline{v}_2 = [1 \ 0 \ 1]^T$

$$\alpha_1 \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} + \alpha_2 \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} = \underline{0} \quad \begin{cases} \alpha_1 + \alpha_2 = 0 \\ 2\alpha_1 + 0 = 0 \\ \alpha_1 + \alpha_2 = 0 \end{cases} \Rightarrow \begin{cases} \alpha_1 = 0 \\ \alpha_2 = 0 \end{cases} \Rightarrow \text{"lin. indep."}$$

Ex:  $\underline{v}_1 = [1 \ 2 \ 1]^T$   $\underline{v}_4 = [-2 \ -4 \ -2]^T$

$$\underline{v}_4 = -2 \cdot \underline{v}_1 \Rightarrow \text{"lin dependent"}$$

Ex:  $v_1, v_2, v_3 = [0 \ 1 \ 0]^T$

$$v_1 = v_2 + 2v_3 \Rightarrow \text{"lin dep."}$$

Def: Vector space: A set of all vectors that are lin. comb. of

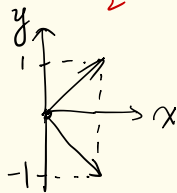
$$\{\underline{v}_i\}_{i=1}^n, \text{ i.e., } V = \left\{ \underline{v} = \sum_{i=1}^n \alpha_i \underline{v}_i, \alpha_i \in \mathbb{R} \right\}.$$

$\underline{v}_i$ 's are said to span the vector space  $V$ , i.e.,  $V = \text{span}\{\underline{v}_1, \dots, \underline{v}_n\}$ .

Ex:  $V^{(1)} = \left\{ \alpha_1 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \alpha_2 \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \alpha_i \in \mathbb{R} \right\}$



$V^{(2)} = \left\{ r_1 \begin{bmatrix} 1 \\ 1 \end{bmatrix} + r_2 \begin{bmatrix} 1 \\ -1 \end{bmatrix}, r_i \in \mathbb{R} \right\}$



Def: A basis for  $V$  is a set of lin indep vectors that span  $V$ .

Ex:  $\left\{ \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ -1 \end{bmatrix} \right\} \checkmark$

$\left\{ \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\} \checkmark$   ~~$\rightarrow$~~

$\left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\} \checkmark$

$\left\{ \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 2 \end{bmatrix} \right\} \times$

Def: The dimension of vector space  $V$  is the # of vectors in any basis of  $V$ . Column/row rank: dim of column/row vector space, respectively.

Ex: What's the rank of  $V = \left\{ \underline{v} = \alpha_1 \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} + \alpha_2 \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} + \alpha_3 \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \alpha_i \in \mathbb{R} \right\}$ ?

Ans1: Basis =  $\left\{ \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \right\}$  or  $\left\{ \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \right\}$  or  $\left\{ \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \right\}$ .

$$\text{rank}(V) = 2.$$

Ans2:  $V = \left\{ \underline{v} = \alpha_1 (\underline{v}_2 + 2\underline{v}_3) + \alpha_2 \underline{v}_2 + \alpha_3 \underline{v}_3 \right.$   
 $\left. = (\alpha_1 + \alpha_2) \underline{v}_2 + (\alpha_1 + \alpha_3) \underline{v}_3 \right\}, \underline{v}_2 \perp \underline{v}_3 \Rightarrow \text{rank}(V) = 2.$

### 3. Geometric interpretation of LS:

Problem Setup:  $\underline{Y} = X\underline{\beta} + \underline{e}$ , where  $X \triangleq [\underline{x}_1, \dots, \underline{x}_p]$

Estimate  $\underline{\beta}$  such that  $J(\underline{\beta}) = \|\underline{Y} - X\underline{\beta}\|^2$  is minimized.

$$J(\underline{\beta}) = \sum_{i=1}^n \left( Y_i - \sum_{j=1}^p x_{ij} \beta_j \right)^2$$

This is called "least-squares".

Claims: 1.  $\hat{\underline{\beta}}$  always exists, but  
2. not always unique.



Goal Estimate  $\beta$ ,  $\sigma^2$ ; "error (variance) of estimators,  
 e.g.,  $\text{Var}(\hat{\beta}_0)$ ,  $\text{Var}(\hat{\beta}_1)$ , ... ;  $\text{VarCov}(\hat{\beta})$ ,  $\text{VarCov}(\hat{y})$ .

Method 1: 
$$\frac{\partial J}{\partial \beta_k} = \sum_{i=1}^n 2(Y_i - \sum_{j=1}^p X_{ij} \beta_j) \frac{\partial}{\partial \beta_k} \underbrace{\left( -(\dots + X_{ik} \beta_k + \dots) \right)}_{-X_{ik}}$$

$$J(\beta) = \sum_{i=1}^n \left( Y_i - \sum_{j=1}^p X_{ij} \beta_j \right)^2$$

$$= \begin{cases} 0 \\ \beta_j = \hat{\beta}_j \end{cases}, \quad k=1, \dots, p$$

$$\Leftrightarrow \sum_i Y_i X_{ik} = \sum_i \sum_j X_{ij} \hat{\beta}_j X_{ik} \Leftrightarrow \boxed{X^T Y = X^T X \hat{\beta}} \quad \begin{array}{l} \text{Normal} \\ \text{Equation} \\ \text{(NE)} \end{array}$$

where  $X^T Y = \left[ \sum_{i=1}^n X_{ik} Y_i \right]_{p \times 1}$ ,  $X^T X = \left[ \sum_{i=1}^n X_{ij} X_{ik} \right]_{p \times p}$ ,  $\hat{\beta} = (X^T X)^{-1} X^T Y$

$$X^T X \hat{\beta} = \left[ \sum_{j=1}^p \left( \sum_{i=1}^n X_{ij} X_{ik} \right) \hat{\beta}_j \right]_{p \times 1}$$

$$\text{Method 2: } \nabla_{\beta} J(\beta) = 0$$

$\beta = \hat{\beta}$

$$J(\beta) = \|\tilde{y} - X\beta\|^2$$

$$\nabla_{\beta} J(\beta) = 2[-X^T(\tilde{y} - X\beta)] = 0$$

$\beta = \hat{\beta}$

$$X^T \tilde{y} = X^T X \hat{\beta}$$

L.S. procedure: Find a vector in  $\mathcal{C}(X)$  which is as close as possible to  $\underline{y}$ .

Claim: If  $(\underline{y} - X\hat{\beta}) \perp \mathcal{C}(X)$ , then  $\hat{\beta}$  solves N.E.

Proof:  $(\underline{y} - \hat{\underline{y}}) \perp \mathcal{C}(X)$

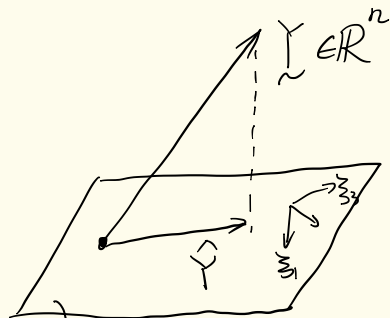
$$\Leftrightarrow (\underline{y} - \hat{\underline{y}}) \perp X\underline{b}, \forall \underline{b} \in \mathbb{R}^p$$

$$\Leftrightarrow \sum_j^T (\underline{y} - \hat{\underline{y}}) = 0, j=1, \dots, p$$

$$\Leftrightarrow \left[ \sum_1^T, \dots, \sum_p^T \right] (\underline{y} - X\hat{\beta}) = \underline{0}$$

$$\Leftrightarrow X^T \underline{y} = X^T X \hat{\beta}$$

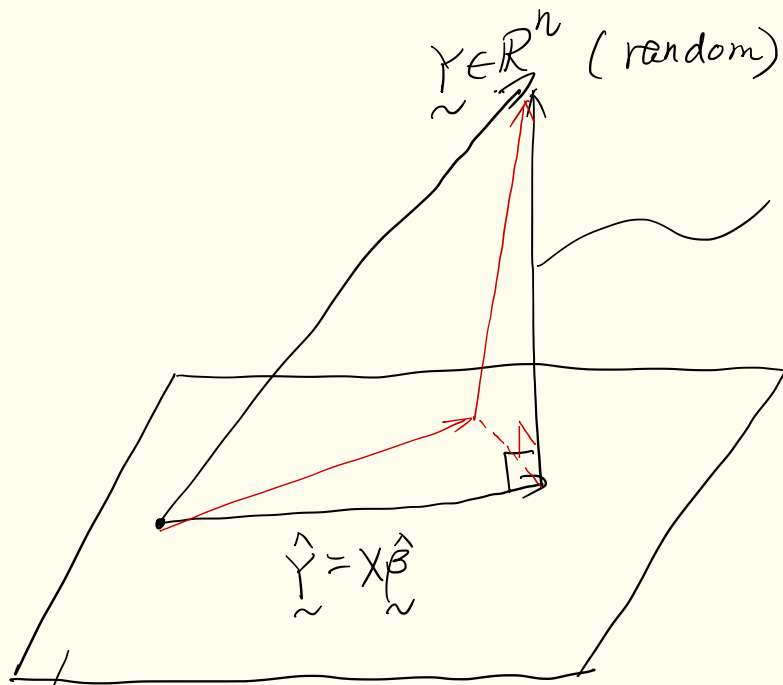
Regardless the rank of  $X^T X$ , there's a soln to the problem.



$$\mathcal{C}(X) = \{X\underline{b}, \underline{b} \in \mathbb{R}^p\}$$

$$X \stackrel{\text{def}}{=} \begin{bmatrix} \sum_1^T & \sum_2^T & \dots & \sum_p^T \end{bmatrix}$$

$$\dim(\mathcal{C}(X)) = r \leq p$$



$$\begin{aligned}
 \text{length}^2 &= \|Y - \hat{Y}\|^2 \\
 &= \|Y - X\hat{\beta}\|^2 \\
 &= \sum_{i=1}^n \left( Y_i - \sum_{j=1}^p x_{ij} \hat{\beta}_j \right)^2 \\
 &\quad (\text{random})
 \end{aligned}$$

Column space of  $X$  (fixed), namely,

$$C(X) \subset \mathbb{R}^n$$

def ||

$$\{Xb, b \in \mathbb{R}^p\}$$

If rank(X) = p = p ①  $\hat{\beta} = (X^T X)^{-1} X^T Y$  is unique soln.

$$E[\hat{\beta}] = E[(X^T X)^{-1} X^T Y] = (X^T X)^{-1} X^T X \beta = \beta \quad (\text{unbiased})$$

$$\begin{aligned} \text{VarCov}(\hat{\beta}) &= (X^T X)^{-1} X^T \underbrace{\text{VarCov}(Y)}_{\sigma^2 I} X (X^T X)^{-1} \\ &= \sigma^2 (X^T X)^{-1} \end{aligned}$$

Note:  $\text{Var}(aX) = a^2 \text{Var}(X)$

$$\text{VarCov}(a^T W) = a^T \text{VarCov}(W) a$$

$$\textcircled{2} \hat{Y} = X \hat{\beta} = X (X^T X)^{-1} X^T Y = H Y$$

H: "hat" matrix, "orthogonal projector".  $H^n = H$ . Why?

$$\text{VarCov}(\hat{Y}) = H \text{Cov}(Y) H^T = \sigma^2 H H^T = \sigma^2 H^2 = \sigma^2 H = \sigma^2 \cdot X (X^T X)^{-1} X^T$$

Ex:  $X = \begin{bmatrix} 2 & 1 \\ 1 & 2 \\ 1 & 1 \end{bmatrix}$   $Y = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$   $Y = X\beta + e$

$$\hat{\beta} = (X^T X)^{-1} X^T Y = \left( \begin{bmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ 1 & 2 \\ 1 & 1 \end{bmatrix} \right)^{-1} \left( \begin{bmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right)$$

$$= \begin{bmatrix} 6 & 5 \\ 5 & 6 \end{bmatrix}^{-1} \begin{bmatrix} 4 \\ 4 \end{bmatrix} = \begin{bmatrix} 6 & -5 \\ -5 & 6 \end{bmatrix} \cdot \frac{1}{11} \cdot \begin{bmatrix} 4 \\ 4 \end{bmatrix}$$

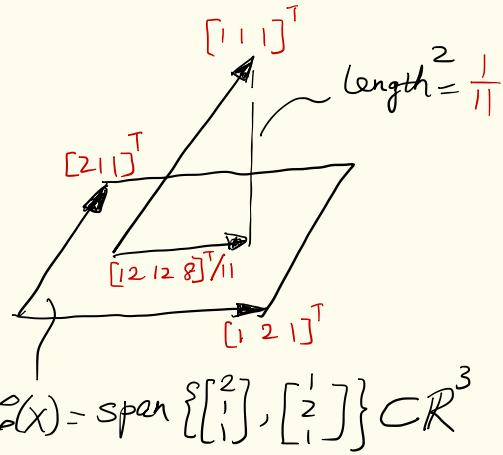
$$= \frac{4}{11} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\text{VarCov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1} = \frac{\sigma^2}{11} \begin{bmatrix} 6 & -5 \\ -5 & 6 \end{bmatrix}$$

$$H = X (X^T X)^{-1} X^T = \begin{bmatrix} 2 & 1 \\ 1 & 2 \\ 1 & 1 \end{bmatrix} \frac{1}{11} \begin{bmatrix} 6 & -5 \\ -5 & 6 \end{bmatrix} \begin{bmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \end{bmatrix} = \frac{1}{11} \begin{bmatrix} 10 & -1 & 3 \\ -1 & 10 & 3 \\ 3 & 3 & 2 \end{bmatrix}$$

$$\hat{Y} = H Y = \frac{1}{11} \begin{bmatrix} 12 \\ 12 \\ 8 \end{bmatrix} \neq Y$$

$$\text{LS err: } \| \tilde{Y} - \hat{Y} \|^2 = \frac{1}{11^2} \left\| \begin{bmatrix} 12-11 \\ 12-11 \\ 8-11 \end{bmatrix} \right\|^2 = \frac{1}{11^2} (1+1+9) = \frac{1}{11}$$



If  $\text{rank}(X) \triangleq r < p$  (Not full column rank), what does  $\hat{\beta}$  look like?

Ans:  $r < p \Leftrightarrow X$  has a non-trivial null space.

$\Leftrightarrow \underline{c} \in \mathbb{R}^p$  s.t.  $X \underline{c} = \underline{0} = \sum_{i=1}^p c_i \underline{\Sigma}_i$ . Then

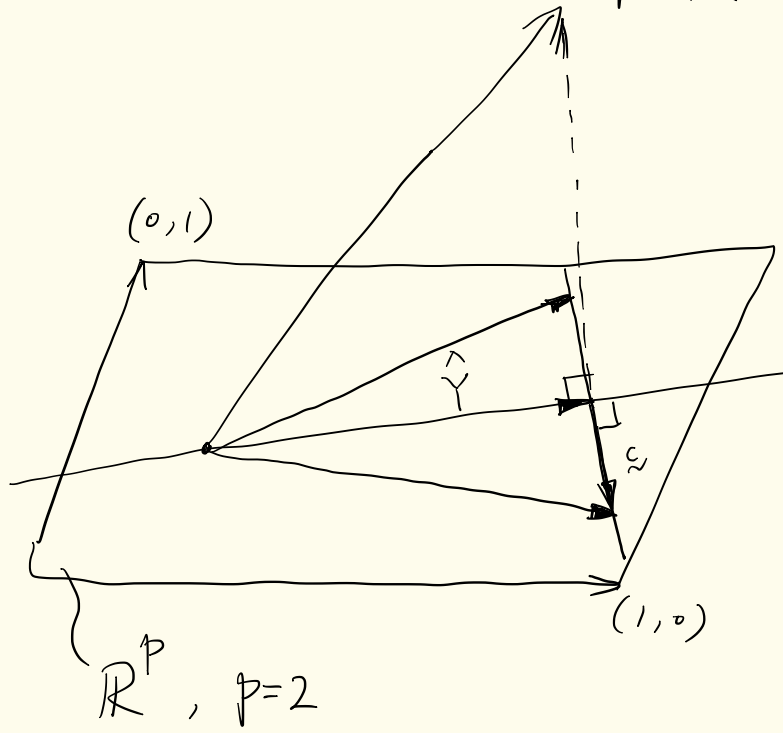
$$\exists \underline{c} \neq \underline{0}$$

$$X \hat{\beta} = X \hat{\beta} + X \underline{c} = X \underbrace{(\hat{\beta} + \underline{c})}_{\tilde{\beta}}$$

$\tilde{\beta} = \hat{\beta} + k \underline{c}$ ,  $k \in \mathbb{R}$  are all LS estimates.

However,  $\hat{y}$  is unique:  $\hat{y} = X(\hat{\beta} + k \underline{c}) = X \hat{\beta} + \underbrace{k X \underline{c}}_{\underline{0}} = X \hat{\beta}$ .

$Y \in \mathbb{R}^n, n=100.$



$\mathcal{C}(X) \subset \mathbb{R}^r,$   
 $r=1$