

# Lecture 6: Least Squares + Subspaces

①

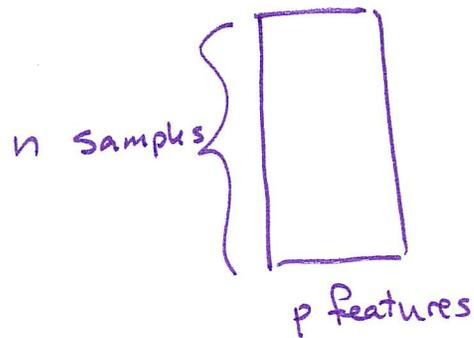
have features and labels for  $i=1, \dots, n$   
 $\underline{x}_i \in \mathbb{R}^p$        $y_i \in \mathbb{R}$

want to predict label  $\hat{y}_i = \langle \underline{x}_i, \underline{w} \rangle = \underline{x}_i^T \underline{w}$

want  $w$  to minimize  $\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \|y - \hat{y}\|_2^2$

$$\hat{y} = X \underline{w}, \text{ where } X = \begin{bmatrix} | & | & & | \\ \underline{x}_1 & \underline{x}_2 & \dots & \underline{x}_n \\ | & | & & | \end{bmatrix}^T$$

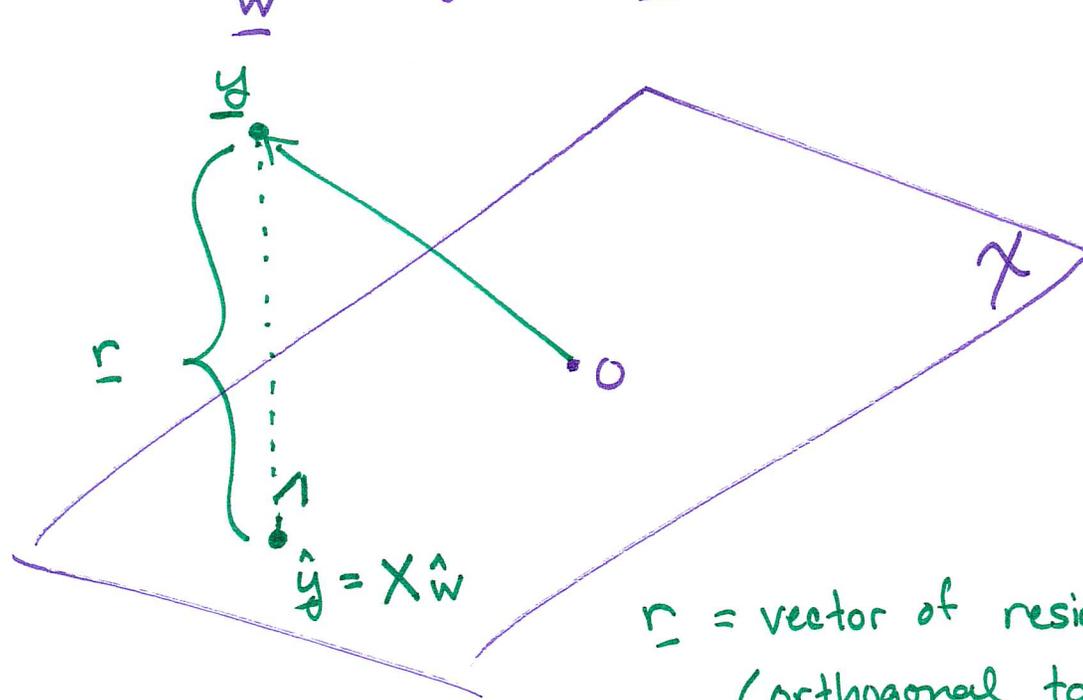
$\hat{y}$  = weighted sum  
of the  
columns of  $X$



$$X \in \mathbb{R}^{n \times p}$$

Last time, assume  $n \geq p$  and  $p$  columns of  $X$  are linearly independent.

$$\hat{\underline{w}} = \arg \min_{\underline{w}} \| \underline{y} - X \underline{w} \|_2^2 \implies \hat{\underline{w}} = (X^T X)^{-1} X^T \underline{y}$$



$\mathcal{X} = \text{span}(\text{columns of } X)$   
 ~~$(x_1, x_2, \dots, x_p)$~~   
 = all vectors that can be written as weighted sum of columns of  $X$

$\underline{r} =$  vector of residuals  $\in \mathbb{R}^n$   
 (orthogonal to  $\mathcal{X}$ )  
 $\underline{r}^T X = 0$ , or  $X^T \underline{r} = 0$

Review of vector calculus :

$y = f(\underline{x})$  maps  $\underline{x} \in \mathbb{R}^n$  to  $y \in \mathbb{R}$

$$\nabla_x f = \begin{bmatrix} df/dx_1 \\ df/dx_2 \\ \vdots \\ df/dx_n \end{bmatrix}$$

a)  $f(\underline{x}) = \langle \underline{c}, \underline{x} \rangle = \underline{x}^T \underline{c} = c_1 x_1 + c_2 x_2 + \dots + c_n x_n$

$$\nabla_x f = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{bmatrix} = \underline{c}$$

b)  $f(\underline{x}) = \|\underline{x}\|^2 = \underline{x}^T \underline{x}$   
 $= x_1^2 + x_2^2 + \dots + x_n^2$

$$\nabla_x f = \begin{bmatrix} 2x_1 \\ 2x_2 \\ \vdots \\ 2x_n \end{bmatrix} = 2\underline{x}$$

(  $f(\underline{x}) = \underline{x}^T Q \underline{x}$  where  $Q = I$  )

$$\begin{bmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{bmatrix}$$

$$c) f(\underline{x}) = \underline{x}^T Q \underline{x}$$

④

assume  $Q$  is symmetric:  $Q = Q^T = \frac{1}{2}(Q + Q^T)$

$$= \sum_{i=1}^n \sum_{j=1}^n x_i Q_{ij} x_j$$

$$[\nabla_x f]_k = df/dx_k \Rightarrow \frac{d}{dx_k} \cancel{x_i} Q_{ij} x_j = \begin{cases} 2Q_{ij} x_i & i=j=k \\ Q_{ij} x_j & i=k, i \neq j \\ Q_{ij} x_i & j=k, i \neq j \end{cases}$$

$$\nabla_x f = (Q + Q^T) x = 2Q x$$

d)  $\hat{\underline{w}} = \arg \min_{\underline{w}} \underbrace{\|y - X\underline{w}\|_2^2}$

$$f(\underline{w}) = \underline{y}^T \underline{y} - 2\underline{w}^T X^T \underline{y} + \underbrace{\underline{w}^T X^T X \underline{w}}_Q$$

$$\nabla_{\underline{w}} f(\underline{w}) = -2X^T \underline{y} + 2X^T X \underline{w} = 0$$

$$\Rightarrow \hat{\underline{w}} = (X^T X)^{-1} X^T \underline{y}$$

key  $X^T X$  is positive-definite :  ~~$\underline{w}^T (X^T X) \underline{w} > 0$~~   
 for all  $\underline{w} \neq 0$

A matrix  $Q$  is positive-definite (p.d.)  $\iff Q \succ 0$  ⑥

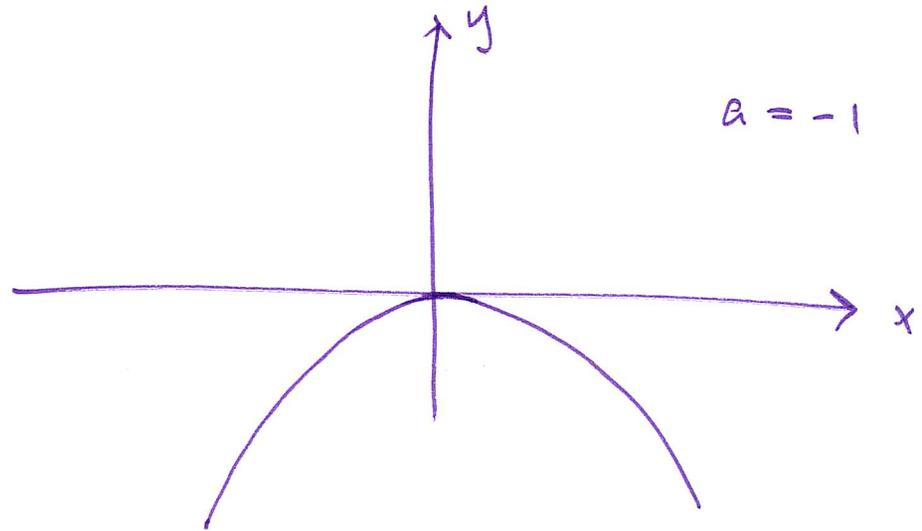
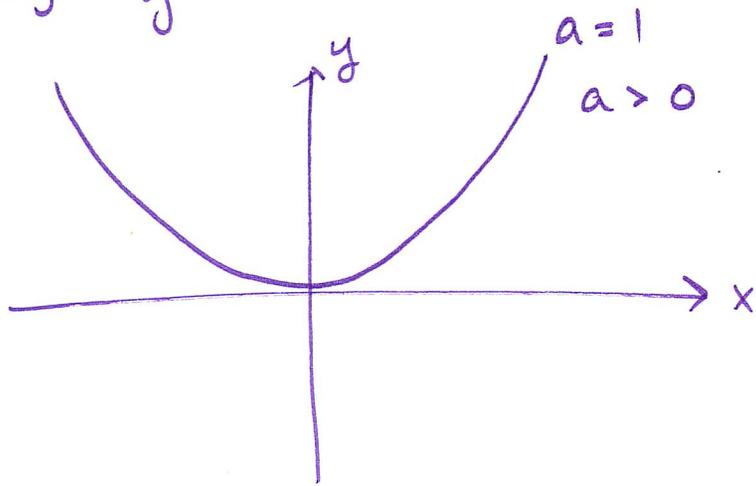
if  $\underline{x}^T Q \underline{x} > 0$  for all  $\underline{x} \neq 0$

A matrix  $Q$  is positive semi-definite (psd)  $\iff Q \succeq 0$

if  $\underline{x}^T Q \underline{x} \geq 0$  for all  $\underline{x} \neq 0$

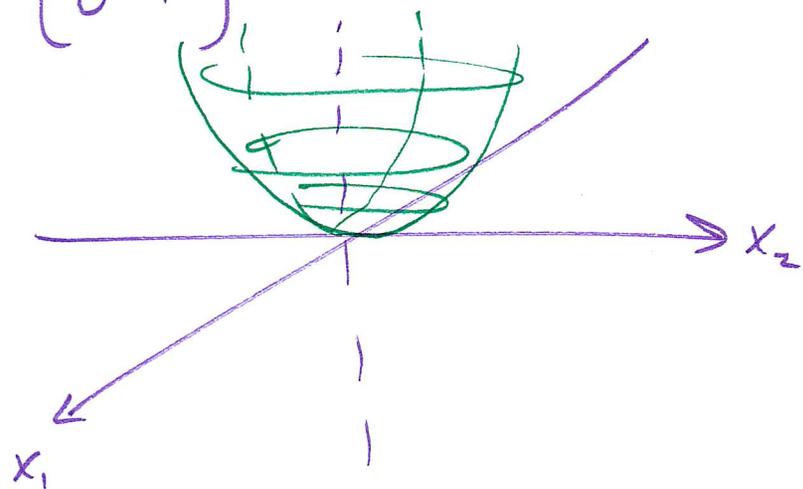
do not mean all  
elements of  $Q \succeq 0$

e.g.  $y = ax^2$

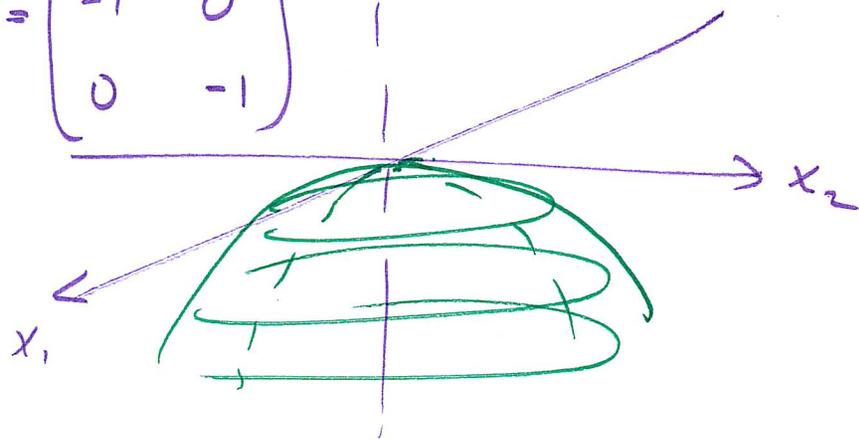


$$Q = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$y = x^T Q x$$



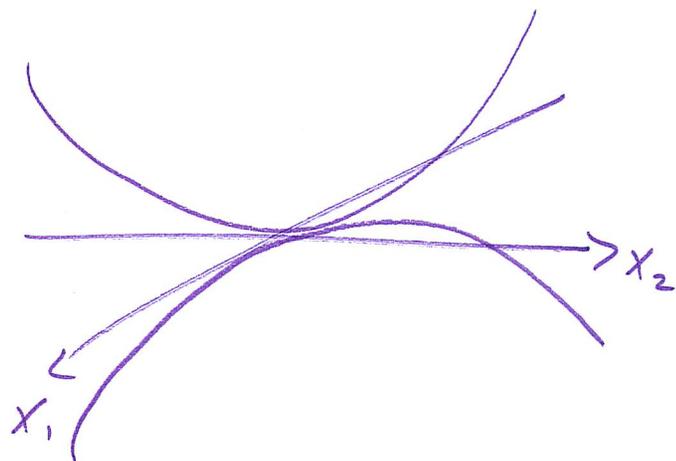
$$Q = \begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix}$$



$Q$  is pos. def.

$$x^T Q x = x_1^2 + x_2^2$$

$$Q = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}$$



(7)

## Properties of Positive Definite Matrices:

1) if  $P \succ 0$  and  $Q \succ 0$ , then  $P+Q \succ 0$

2) if  $Q \succ 0$  and  $\alpha > 0$ , then  $\alpha Q \succ 0$

$$\underline{x}^T Q \underline{x} > 0, \text{ then } \underline{x}^T (\alpha Q) \underline{x} = \alpha (\underline{x}^T Q \underline{x}) > 0$$

3) for any  $A$ ,  $A^T A \succeq 0$  and  $AA^T \succeq 0$

if columns of  $A$  are linearly independent, then  $A^T A \succ 0$

$$\underline{x}^T A^T A \underline{x} \succeq 0$$

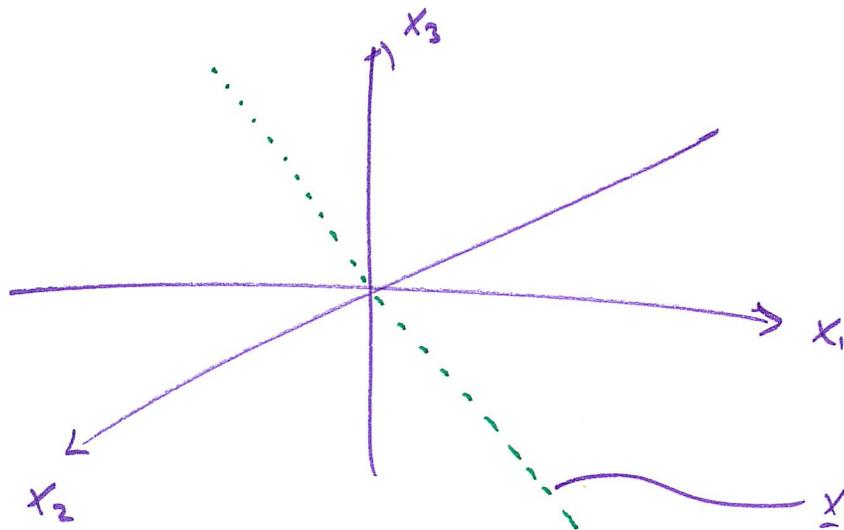
4) if  $A \succ 0$  then  $A^{-1}$  exists

5)  $A \succ B$  ~~is~~ means  $A - B \succ 0$

Subspaces — very common in machine learning

(9)

e.g.  $\underline{x}_i \in \mathbb{R}^3$



$\underline{x}_i$ 's are points on a 1-d subspace of  $\mathbb{R}^3$

$$S = \{ \underline{x} \in \mathbb{R}^3 : x_1 = x_2 = x_3 \}$$

$$\underline{x} = \alpha \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \Rightarrow S = \text{span} \left\{ \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right\}$$

$$\text{Basis for } S = \left\{ \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} / \sqrt{3} \right\}$$

$$\dim(S) = 1$$

A set of points  $S \subseteq \mathbb{R}^n$  (every point in  $S$  is also in  $\mathbb{R}^n$ ) (10)

is a Subspace if

(i)  $0 \in S$  ( $S$  contains origin)

(ii) if  $\underline{x}, \underline{y} \in S$ , then  $\underline{x} + \underline{y} \in S$

(iii) if  $\underline{x} \in S$ ,  $\alpha \in \mathbb{R}$ , then  $\alpha \underline{x} \in S$

ex 1  $n = 3$ . Subspace  $S = \{ \underline{x} \in \mathbb{R}^3 : x_1 = x_2 = x_3 \}$

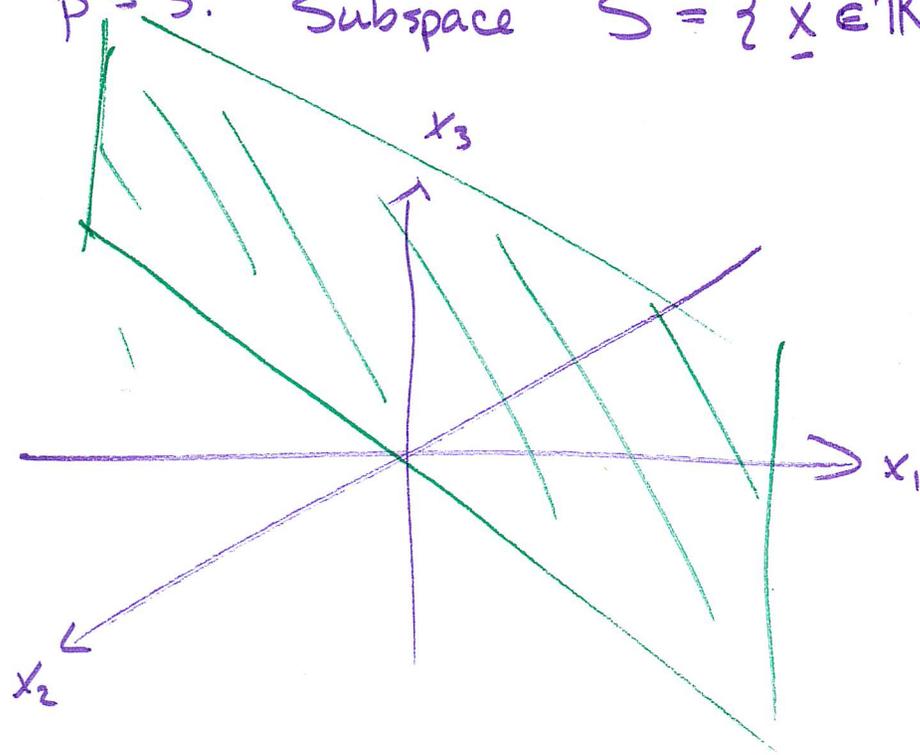
$$\underline{x} + \underline{y} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \\ x_3 + y_3 \end{bmatrix} \in S$$

$$\alpha \underline{x} = \begin{bmatrix} \alpha x_1 \\ \alpha x_2 \\ \alpha x_3 \end{bmatrix} \in S$$

ex 2

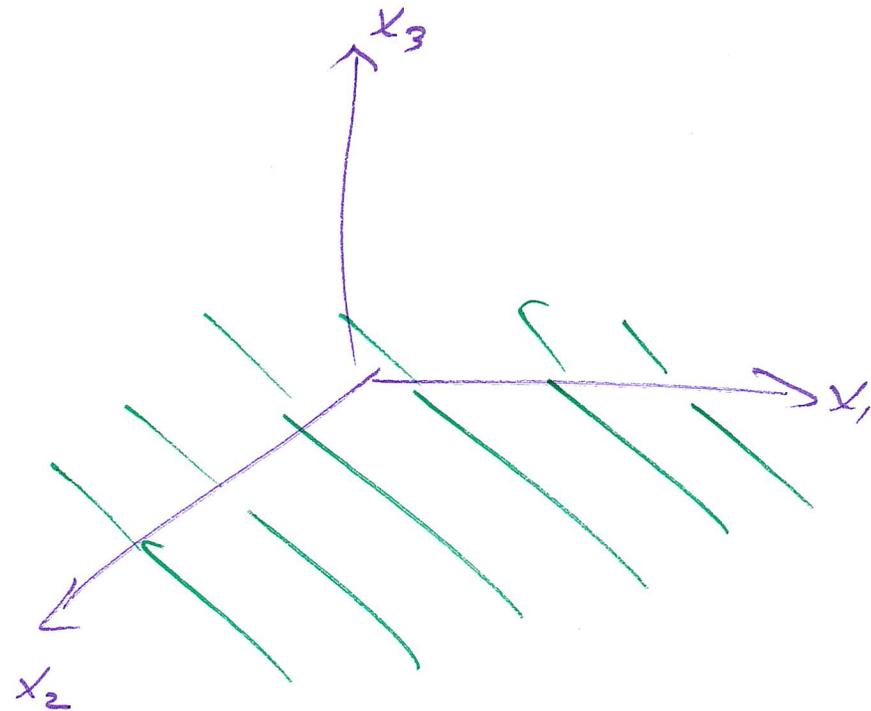
$p = 3$ . Subspace  $S = \{ \underline{x} \in \mathbb{R}^3 : x_1 = x_2 \}$

(11)



ex 3

$S = \{ \underline{x} \in \mathbb{R}^3 : x_3 = 0 \}$



ex 41

(12)

we are given points  $\underline{x}_1, \underline{x}_2, \underline{x}_3 \in \mathbb{R}^n$

$$\text{Span}\{\underline{x}_1, \underline{x}_2, \underline{x}_3\} = \left\{ \underline{y} \in \mathbb{R}^n : \underline{y} = \alpha_1 \underline{x}_1 + \alpha_2 \underline{x}_2 + \alpha_3 \underline{x}_3 \right. \\ \left. \text{for some } \alpha_1, \alpha_2, \alpha_3 \in \mathbb{R} \right\}$$

↓  
 = set of  $\underline{y}$  vecs = linear combo.  
 of  $\underline{x}_1, \underline{x}_2, \underline{x}_3$   
 is a subspace

if  $X = [\underline{x}_1 \quad \underline{x}_2 \quad \underline{x}_3]$ , then range(X) = span(cols(X))

eg. a

$$X = \left[ \begin{array}{c} \text{user} \\ \text{movie} \end{array} \right] = \left[ \begin{array}{c} \text{T} \\ \text{W} \end{array} \right]$$

model: each user's taste profile lies in a subspace spanned by cols of T (representative taste profiles)

= r representative taste profiles

W = r weights for each user

How to represent a subspace?

(13)

- as ~~the~~ span of a set of vectors  
⇒ can be hard to interpret, hard to compute with, redundant
- as the span of a set of linearly indep. vectors
- as the span of a set of orthonormal vectors (BASIS of subspace)

$$S = \text{span} \{ \underline{u}_1, \underline{u}_2, \dots, \underline{u}_r \}$$

where  $\underline{u}_i$ 's are orthonormal

- orthogonal:  $\underline{u}_i^T \underline{u}_j = 0$  if  $i \neq j$

- normalized  $\underline{u}_i^T \underline{u}_i = \|\underline{u}_i\| = 1$  for all  $i$

⇒ dimension ("size") of  $S = r$

e.g.  $S = \{ \underline{x} \in \mathbb{R}^3 : x_3 = 0 \}$

all  $\underline{x} \in S$  have form

$$\underline{x} = \begin{bmatrix} \alpha \\ \beta \\ 0 \end{bmatrix} = \alpha \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + \beta \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$\text{Basis} = \left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \right\}$$

$$\dim(S) = 2$$

$$U = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$$

= basis matrix

$$\text{rank}(U) = 2$$

