

Principal Component Analysis (PCA)

- Math
- Intuition
- Code
- Limitations

@end, delivery case-study
(Q&A)



Math

max
 u

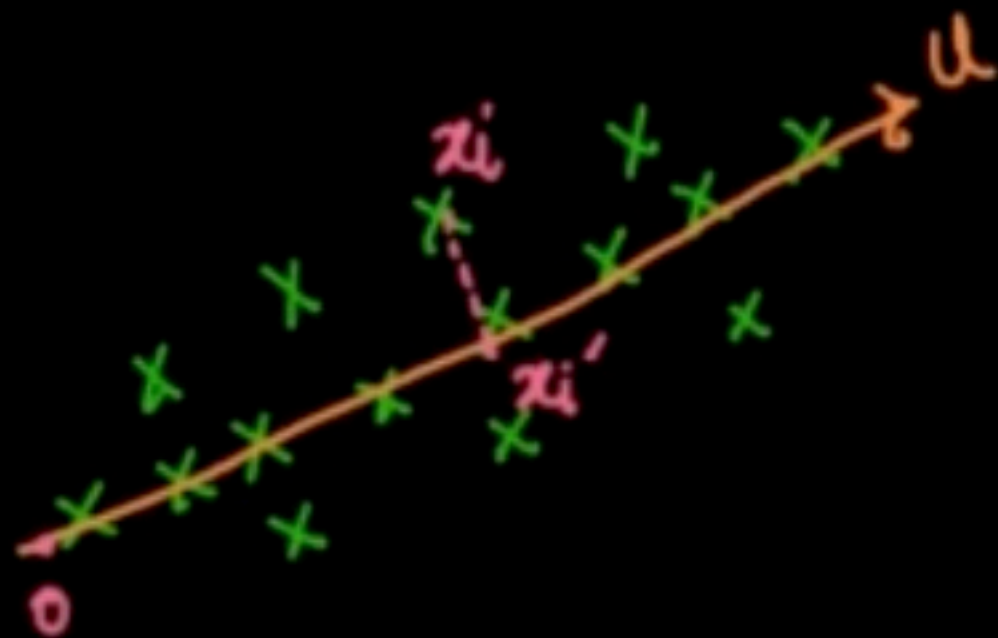
$$\frac{1}{n} \sum (u^T x_i - u^T \bar{x})^2$$

} → let us
simplify this

s.t $\|u\|=1$

Gradient

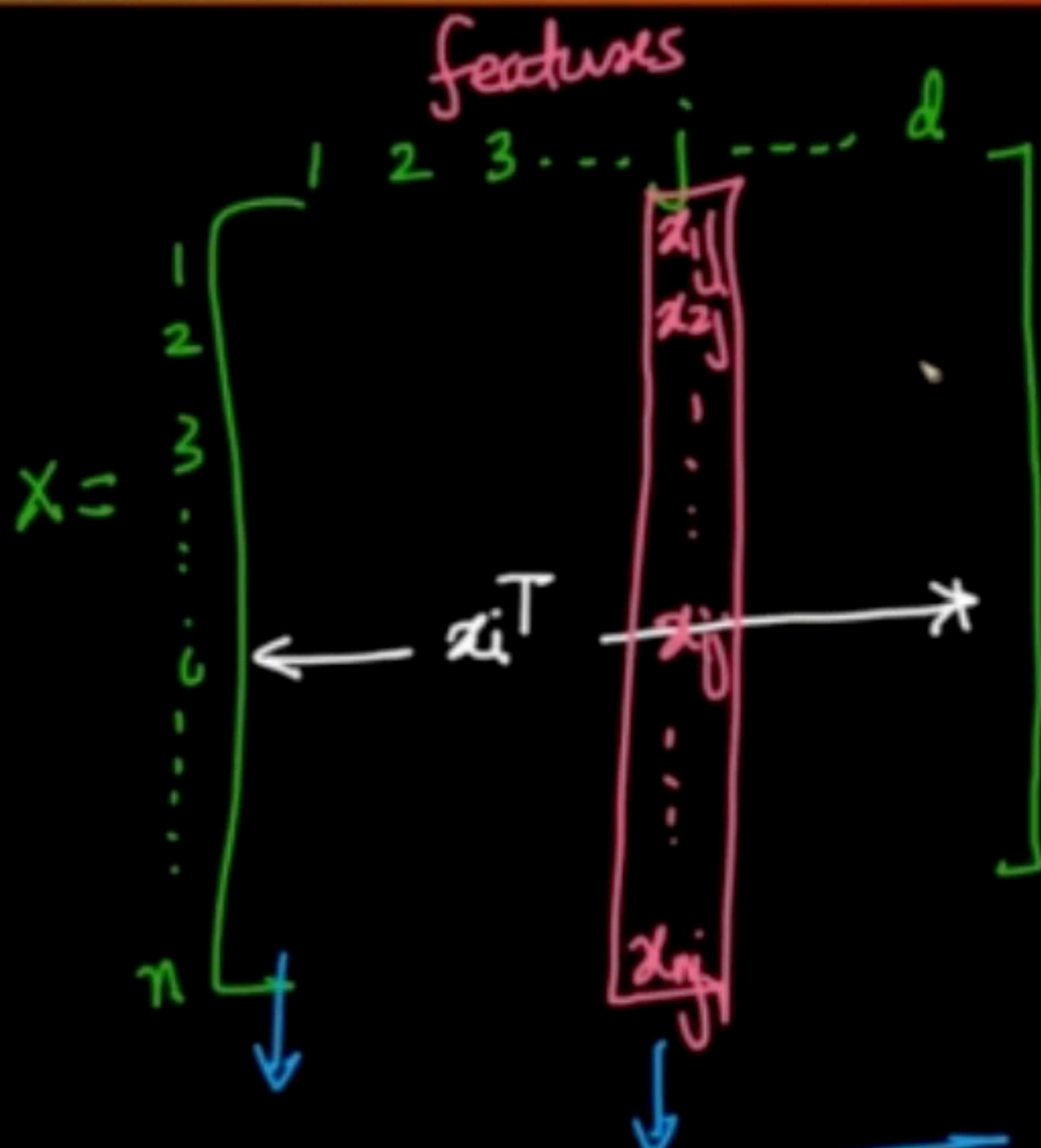
$$\bar{x} = \text{mean of } x_i\text{'s}$$



$$\text{mean of } x_i\text{'s} = u^T \bar{x}$$

Detail

Data matrix

 $\mathbb{R}^d \ni \bar{x}$

Column Standardization

① Mean Centering

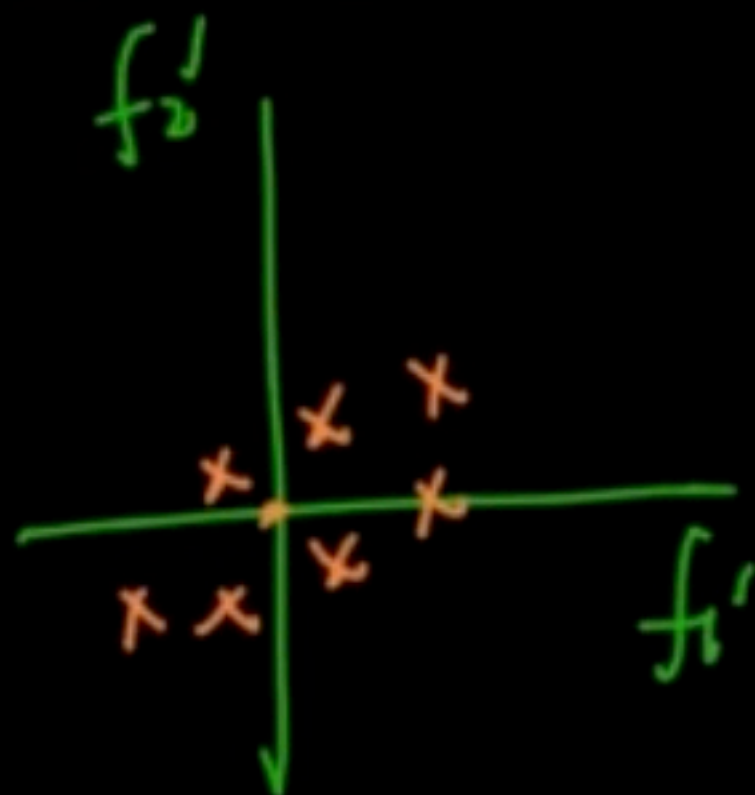
$x_i - \bar{x}$

$\mu_i = 1/\sigma_i$

② Variance Scaling



mean
 →
 Centring



e.g.: $A = N(\mu, \sigma^2)$

$\frac{A - \mu}{\sigma} \sim N(0, 1)$

variance scaling:

$$X = \begin{matrix} & \begin{matrix} 1 & 2 & \dots & j & \dots & d \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ i \\ \vdots \\ n \end{matrix} & \left[\begin{array}{cccccc} & & & x_{1j} & & \\ & & & x_{2j} & & \\ & & & \vdots & & \\ & & & x_{ij} & & \\ & & & \vdots & & \\ & & & x_{nj} & & \end{array} \right] \end{matrix}$$

\downarrow

$$\sigma_1 \quad \sigma_2 \quad \dots \quad \sigma_j \quad \dots \quad \sigma_d$$

$\rightarrow i^{th} \text{ pt}; j^{th} \text{ feature}$

$$x_{ij}^{new} = \frac{x_{ij} - \bar{x}_j}{\sigma_j}$$

Mean-centering
& Variance scaling \rightarrow Standardization

optimization faster?

every

Column's
Columns

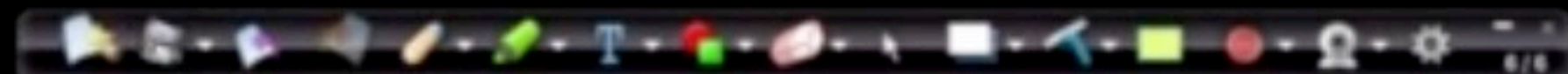
mean = 0

var = 1 = std dev

Simplify Math

copy

SCALER



X^{raw}

 pre-processing

 mean centering

 var-scaling

$X \rightarrow$ optimization problem

$$\max_u \frac{1}{n} \sum_{i=1}^n \left(u^T x_i - \cancel{u^T \bar{x}} \right)^2$$

$\rightarrow [0, 0, 0, \dots, 0]^T$

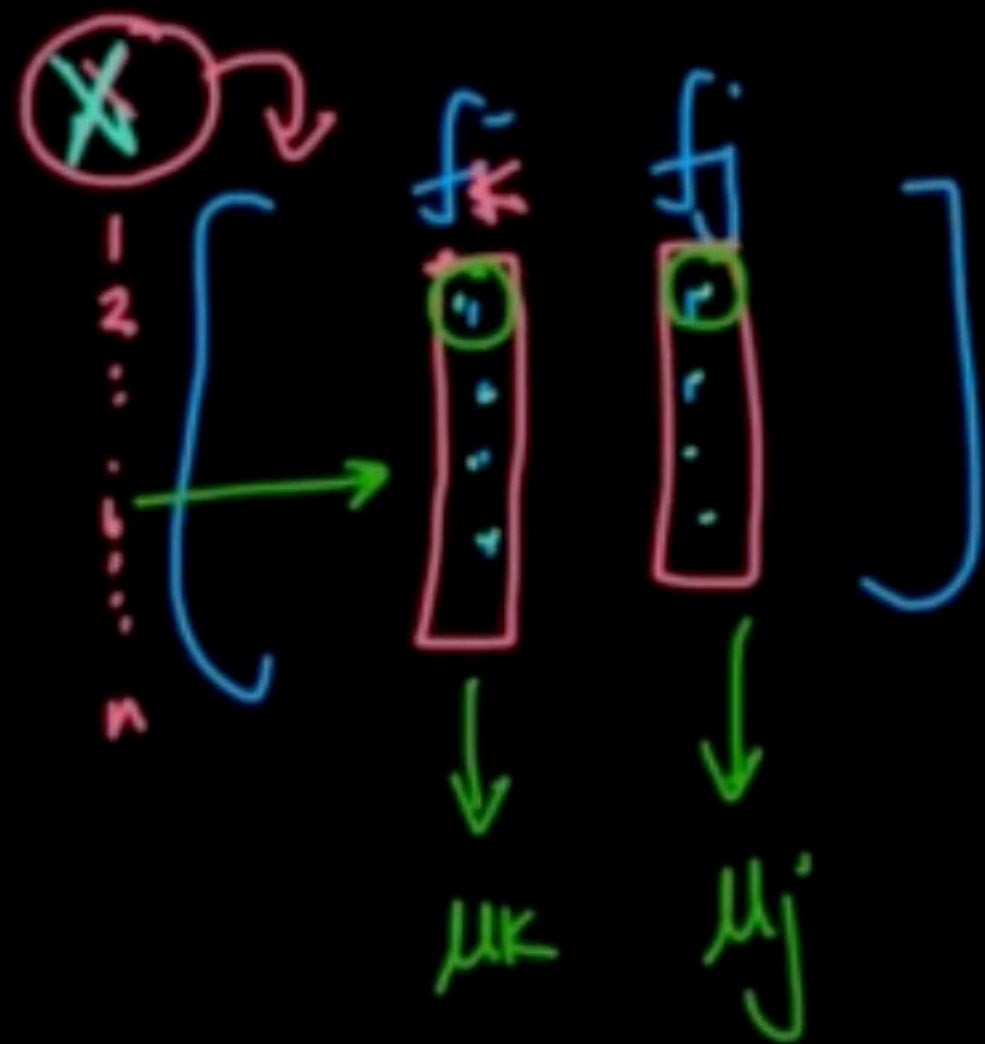
$$\text{s.t. } \|u\| = 1$$

let

$$S_{ij} = \text{Covariance}(f_k, f_j)$$

$$= \frac{1}{n} \sum_{i=1}^n (x_{ik} - \bar{x}_k) (x_{ij} - \bar{x}_j)$$

$$S_{ij} = \frac{1}{n} \sum_{i=1}^n (x_{ik} \cdot x_{ij})$$



(8)

x_1	x_2	x_3	-	-	-	x_{10}
10	20	30	40	-	-	62

$$\bar{x} = \frac{1}{n} \sum x_i$$

$$x_i' = x_i - \bar{x}$$

$$\text{mean}(x_i') = 0$$

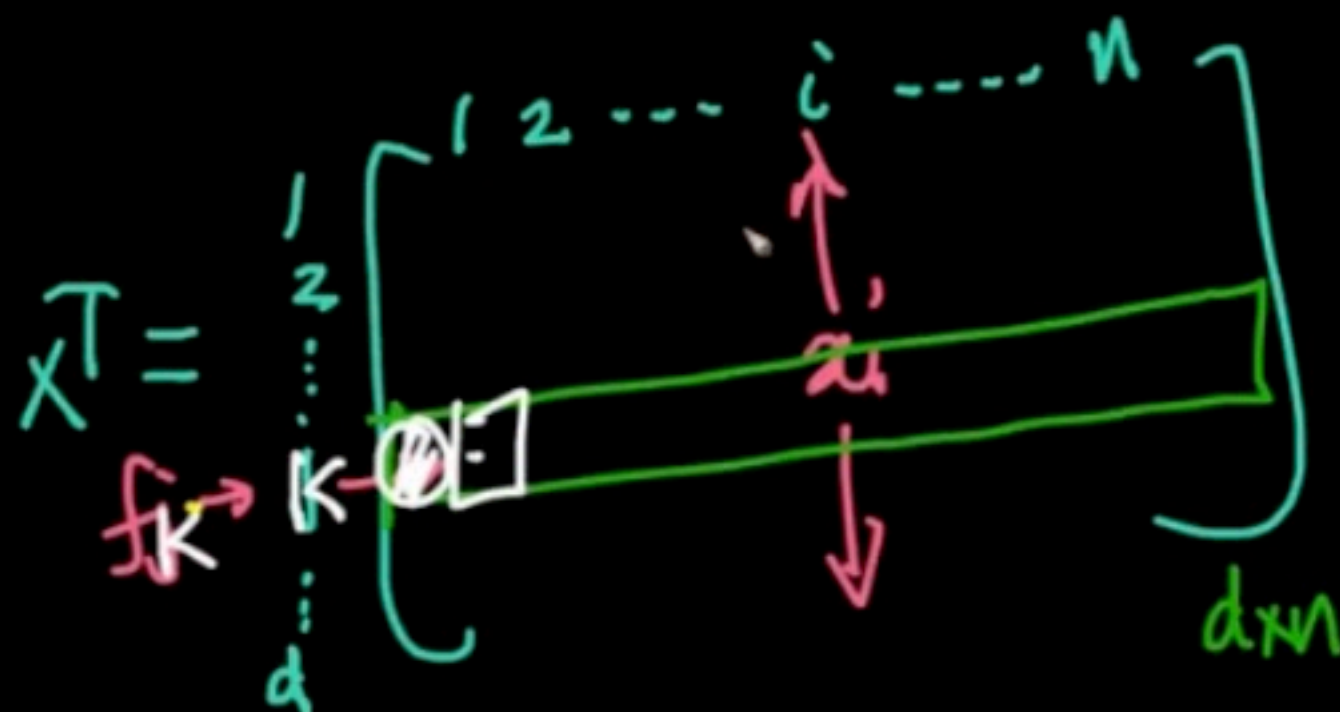
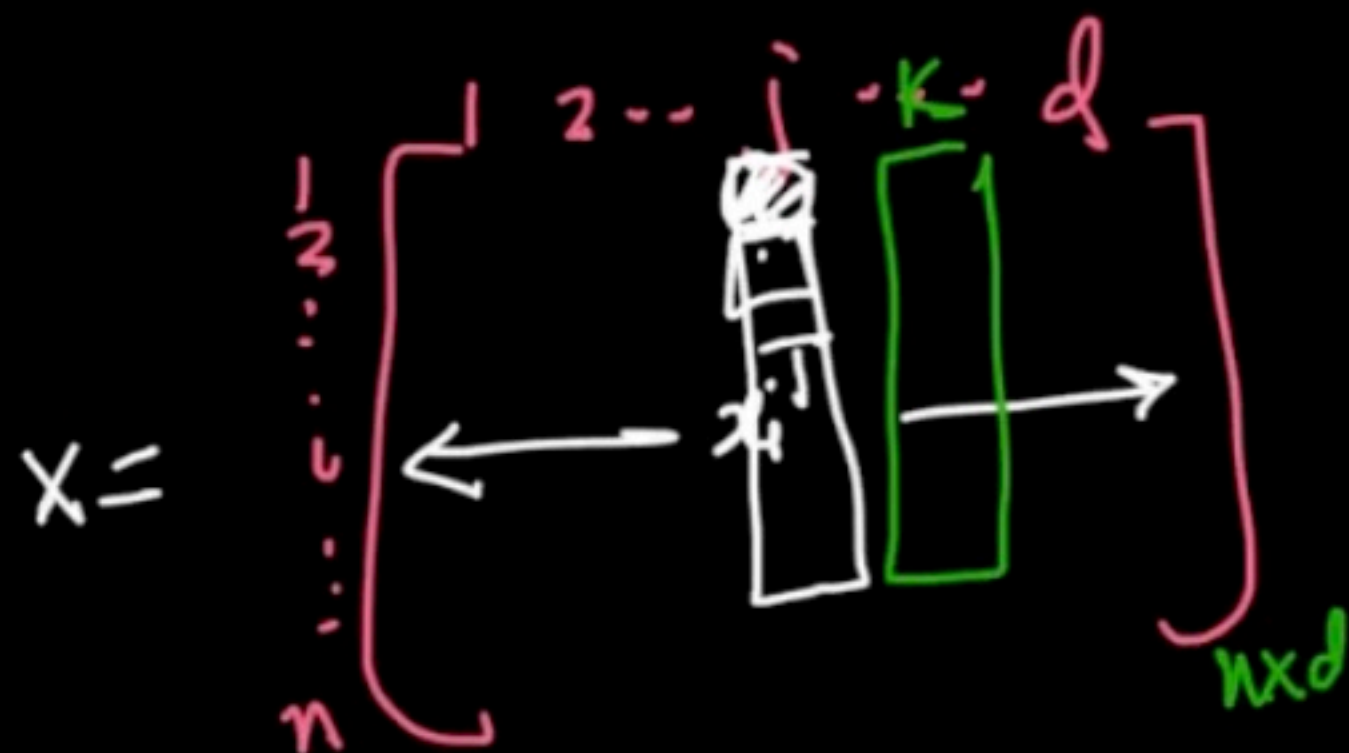


$$S_{ij} = \frac{1}{n} \sum_{i=1}^n \underline{x_{ik}} \cdot \underline{x_{ij}}$$

kth row
jth col

$$S_{d \times d} = \frac{1}{n} \cdot X_{d \times n}^T \cdot X_{n \times d}$$

$$S_{kj} = \frac{1}{n} \cdot \sum_{i=1}^n x_{ik} \cdot x_{ij}$$



Optimizing:

max
 u

$$\frac{1}{n} \sum_{i=1}^n (u^T x_i)^2$$

st $\|u\| = 1$

$\underbrace{\qquad\qquad\qquad}_{\qquad\qquad\qquad} \rightarrow u^T u = 1$



$$\frac{1}{n} \sum_{i=1}^n (u^T x_i)^2$$

↓
scalar

$$\underbrace{u^T}_{1 \times d} \cdot \underbrace{X^T}_{d \times n} = A_{1 \times n}$$

$$u \in \mathbb{R}^d$$

$$\begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_d \end{bmatrix}_{d \times 1}$$

$$A_i = u^T x_i$$

$$A = \begin{bmatrix} \text{---} \\ 1 & 2 & \dots & i & \dots & n \\ \text{---} \end{bmatrix}$$

↓
 $u^T x_i$

$$\left[\leftarrow u^T \rightarrow \right] \begin{bmatrix} 1 & 2 & \dots & i & \dots & n \\ \uparrow \\ x_i \\ \downarrow \end{bmatrix}$$

$$B_{n \times 1} = X_{n \times d} u_{d \times 1} = \frac{1}{2} \begin{bmatrix} \vdots \\ i \\ \vdots \end{bmatrix} \leftarrow x_i^T \rightarrow \begin{bmatrix} \uparrow \\ u \\ \downarrow \end{bmatrix}$$

$$\frac{1}{2} \begin{bmatrix} \vdots \\ i \\ \vdots \end{bmatrix} \rightarrow x_i^T u$$

$n \times 1$

SCALER

$$Res_{|x|} = (u^T x^T) \cdot (xu)$$

$A_{1 \times n}$ $B_{n \times 1}$

$$u^T x_i = x_i^T u$$

$$\begin{bmatrix} 1 & 2 & \dots & i & \dots & n \end{bmatrix}$$

$(u^T x_i)$

$$\begin{bmatrix} \vdots \\ x_i^T u \\ \vdots \end{bmatrix}$$

$(x_i^T u)$

$$= \sum_{i=1}^n (u^T x_i)$$

$$= \sum_{i=1}^n (x_i^T u)$$

$$= \sum_{i=1}^n (u^T x_i)^2$$

Optimzn:

$$\max_u : \underbrace{(u^T x^T)(x u)}_{\text{blue bracket}}$$

$$\text{s.t. } u^T u = 1$$

$$= \frac{1}{n} u^T \underbrace{X^T X}_{S} u$$

Diagram illustrating the derivation of the covariance matrix S from the data matrix X . The expression $\frac{1}{n} u^T X^T X u$ is shown, with a red bracket over $X^T X$ and a blue bracket under it. Arrows point from these brackets to the term $u^T S u$ below, where S is the covariance matrix.

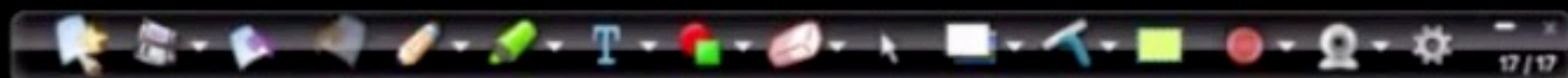
opt:

$$\begin{cases} \text{Max}_u & u^T S u \\ \text{s.t.} & u^T u = 1 \end{cases}$$

The optimization problem is defined as maximizing $u^T S u$ subject to the constraint $u^T u = 1$. A blue arrow points from the text "COV of data." to the matrix S in the objective function.

X : data matrix

\downarrow
 S



$$\begin{array}{c} \min \\ \cancel{\max} \end{array} - \underline{u}^T \overset{\downarrow}{S} \underline{u} + \lambda (\underline{u}^T \underline{u} - 1)$$

①

$$\cancel{s.t. \underline{u}^T \underline{u} = 1}$$

②

$$\min_{\underline{u}} \underbrace{-\underline{u}^T S \underline{u}}_{\rightarrow} + \lambda (\underline{u}^T \underline{u} - 1) \rightarrow \mathcal{L}(\underline{u}, \lambda)$$

$$\frac{\partial \mathcal{L}}{\partial \underline{u}} = 0 \Rightarrow$$

$$\min_{\text{max}} - \underline{u}^T \overset{\downarrow}{S} \underline{u} + \lambda (\underline{u}^T \underline{u} - 1)$$

①

$$\text{s.t. } \underline{u}^T \underline{u} = 1$$

$$\min_u - \underline{u}^T \underline{S} \underline{u} + \lambda (\underline{u}^T \underline{u} - 1) \rightarrow \mathcal{L}(u, \lambda)$$

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial u} = 0 \end{cases} \Rightarrow -2(\underline{u}^T \underline{S})^T + 2\lambda \underline{u} = 0$$

\downarrow
 $\underline{S}^T \underline{u}$

Gradient
descent

$$\frac{d}{du} \left(-\overbrace{u^T}^A S \underbrace{u}_A \right) = -S^T u - S u$$

$$= -2Su$$

$$\frac{d Au}{du} = A^T$$

$$\frac{d(u^T A)}{du} = A$$

$$-2Su + 2\lambda u = 0$$

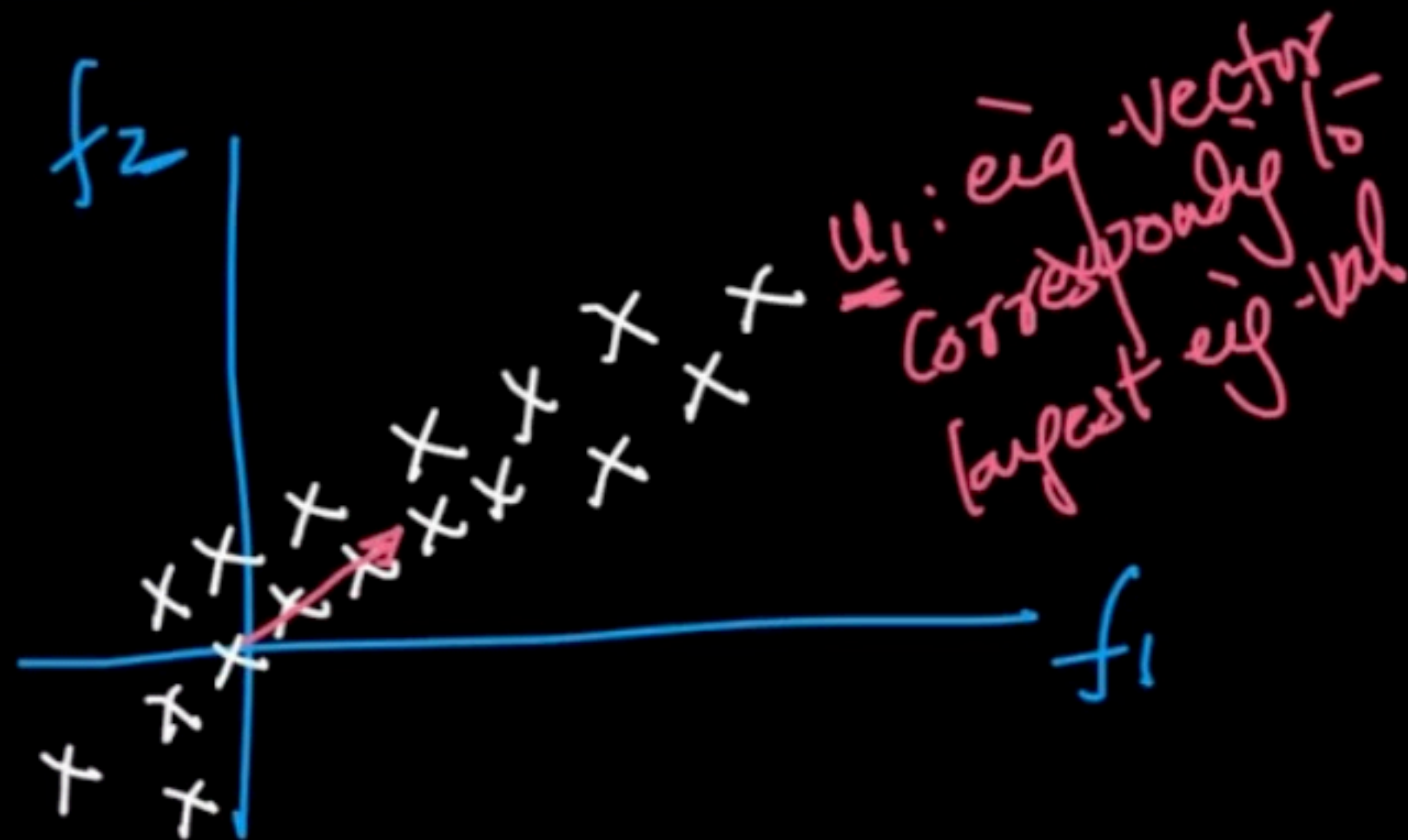
$$\Rightarrow \cancel{2} \lambda u = \cancel{2} Su$$

$\left\{ \begin{array}{l} u\text{'s are eig-vec of } S \\ \lambda\text{'s are eig-val of } S \end{array} \right.$

① Data matrix
pre-processing: mean centering & var-scaling
 $X_{n \times d}$

② Compute $S_{d \times d} = \text{Cov. matrix of } X$

③ eig-vec & eig val of S
 $\swarrow \quad \searrow$
 $u_1, u_2, u_3, \dots, u_d \quad \lambda_1, \lambda_2, \lambda_3, \dots, \lambda_d$



$X_{n \times d}$ ✓



$S_{d \times d}$ ✓



✓

$\lambda_1, \lambda_2, \dots, \lambda_d$

u_1, u_2, \dots, u_d

largest

Viz

100-dim \longrightarrow 2-dim

X

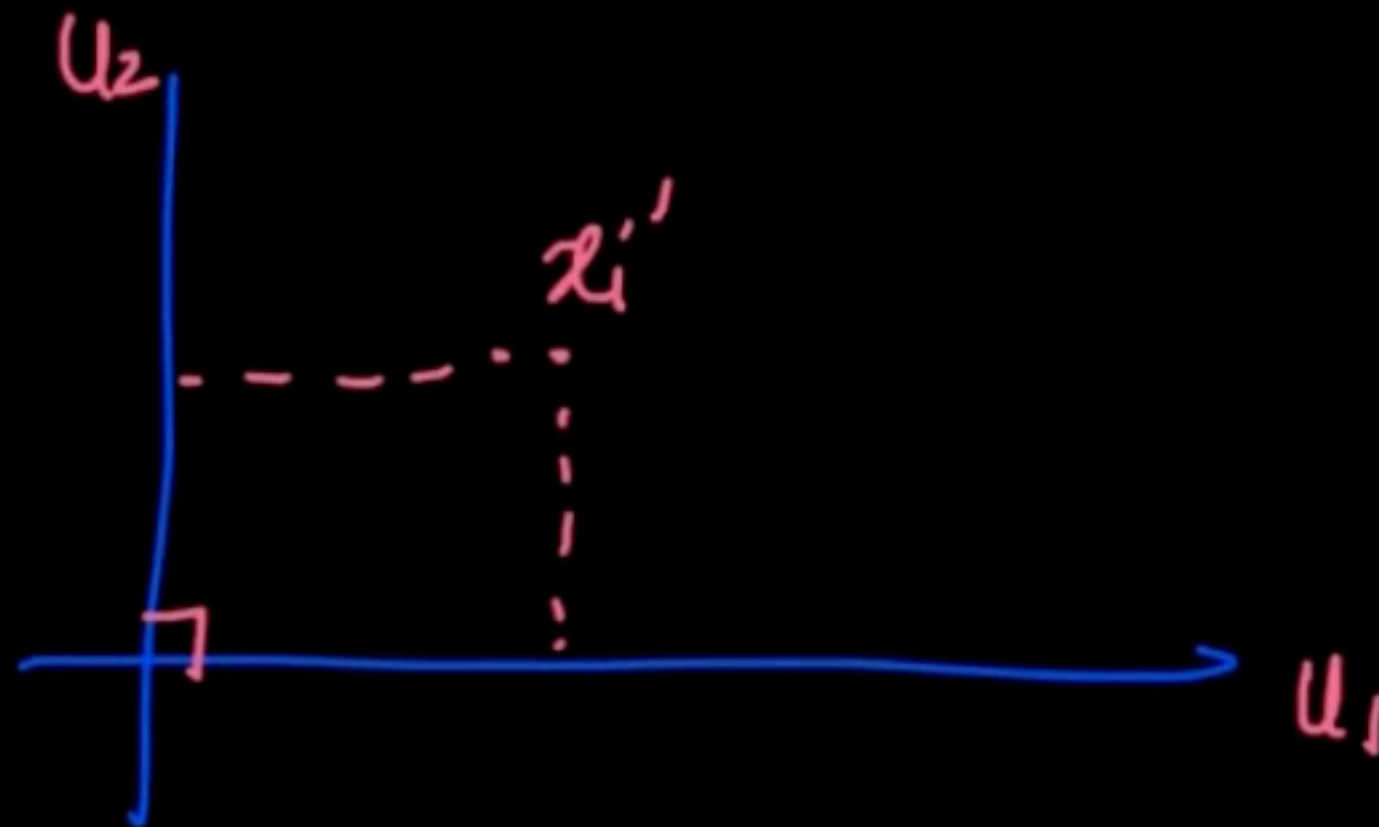
\hookrightarrow S

\longrightarrow

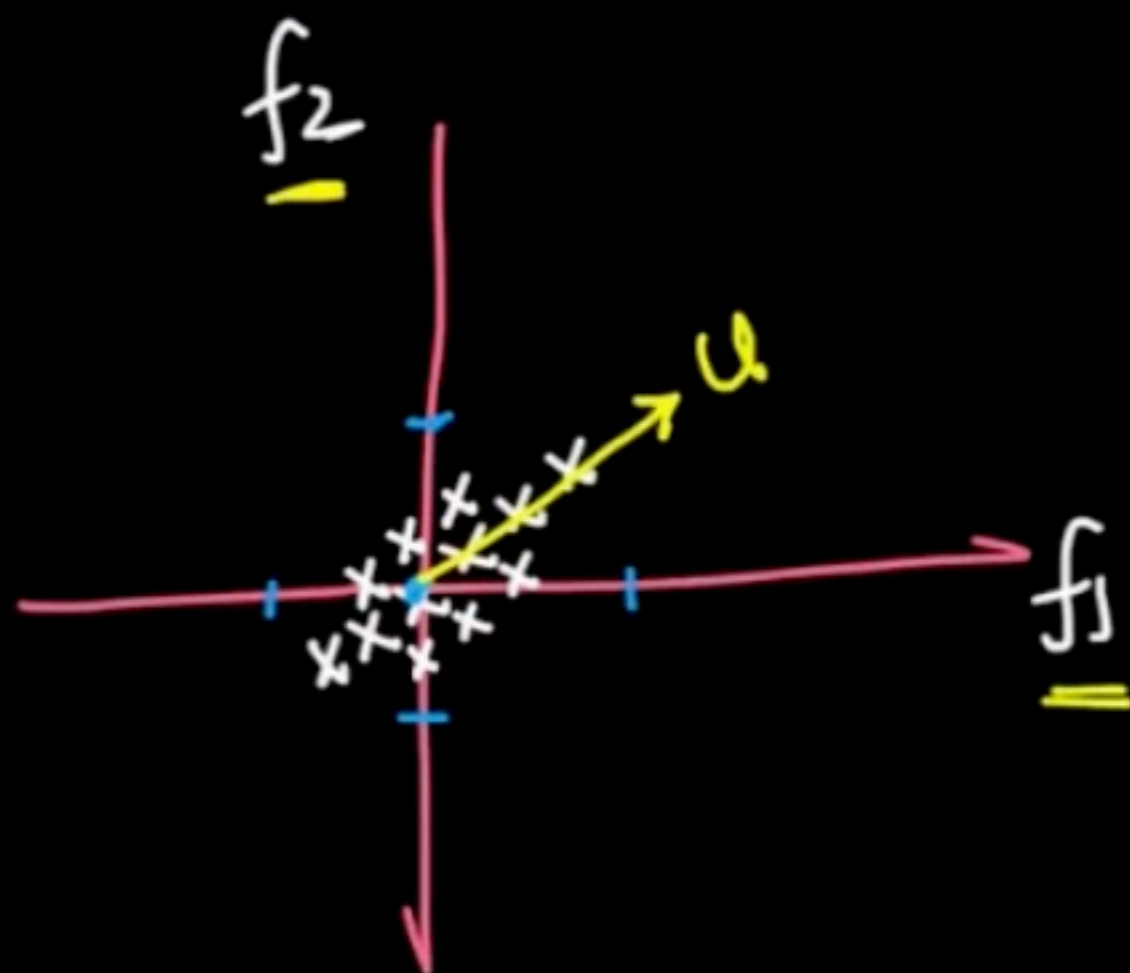
$\lambda_1 > \lambda_2 > \dots > \lambda_n$
 u_1 u_2 \dots u_n

u_2
 u_1

$$x_i \cdot u_1 = \begin{pmatrix} x_{i1}' \\ x_{i2}' \end{pmatrix}$$
$$x_i \cdot u_2 = \begin{pmatrix} x_{i1}' \\ x_{i2}' \end{pmatrix}$$



(Q)



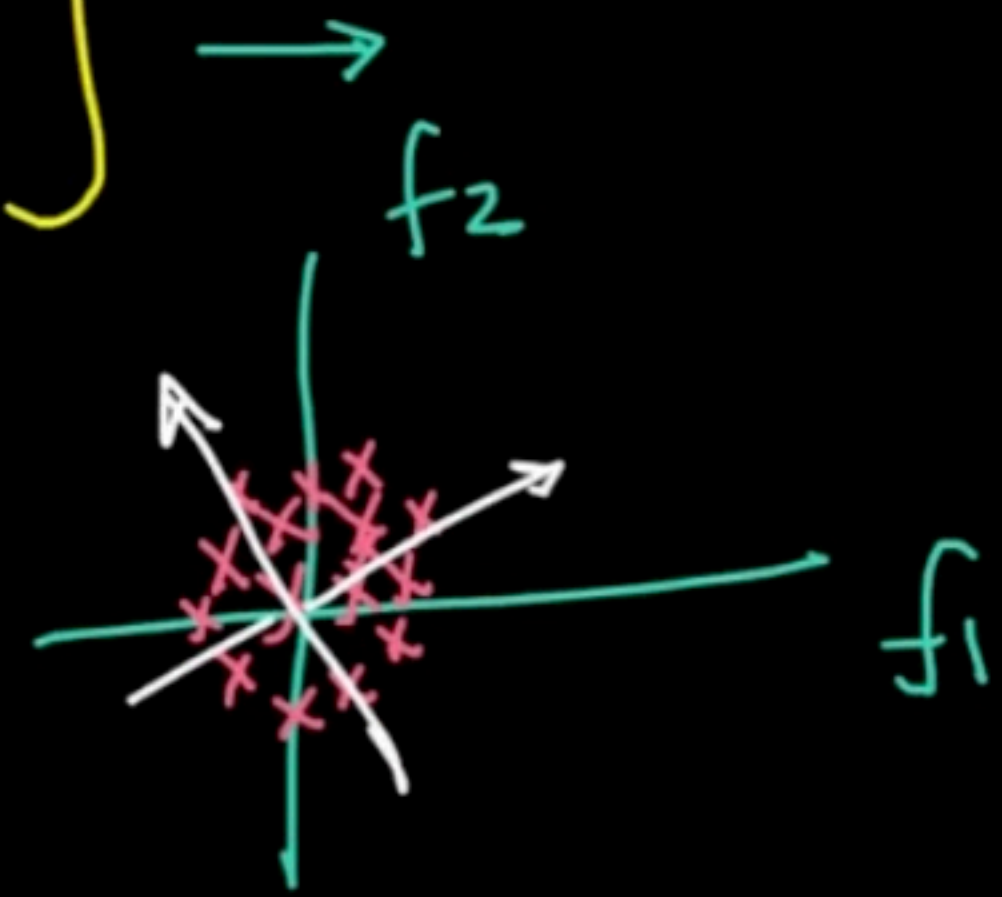
✓ mean centered
✓ var scaling

(Q)

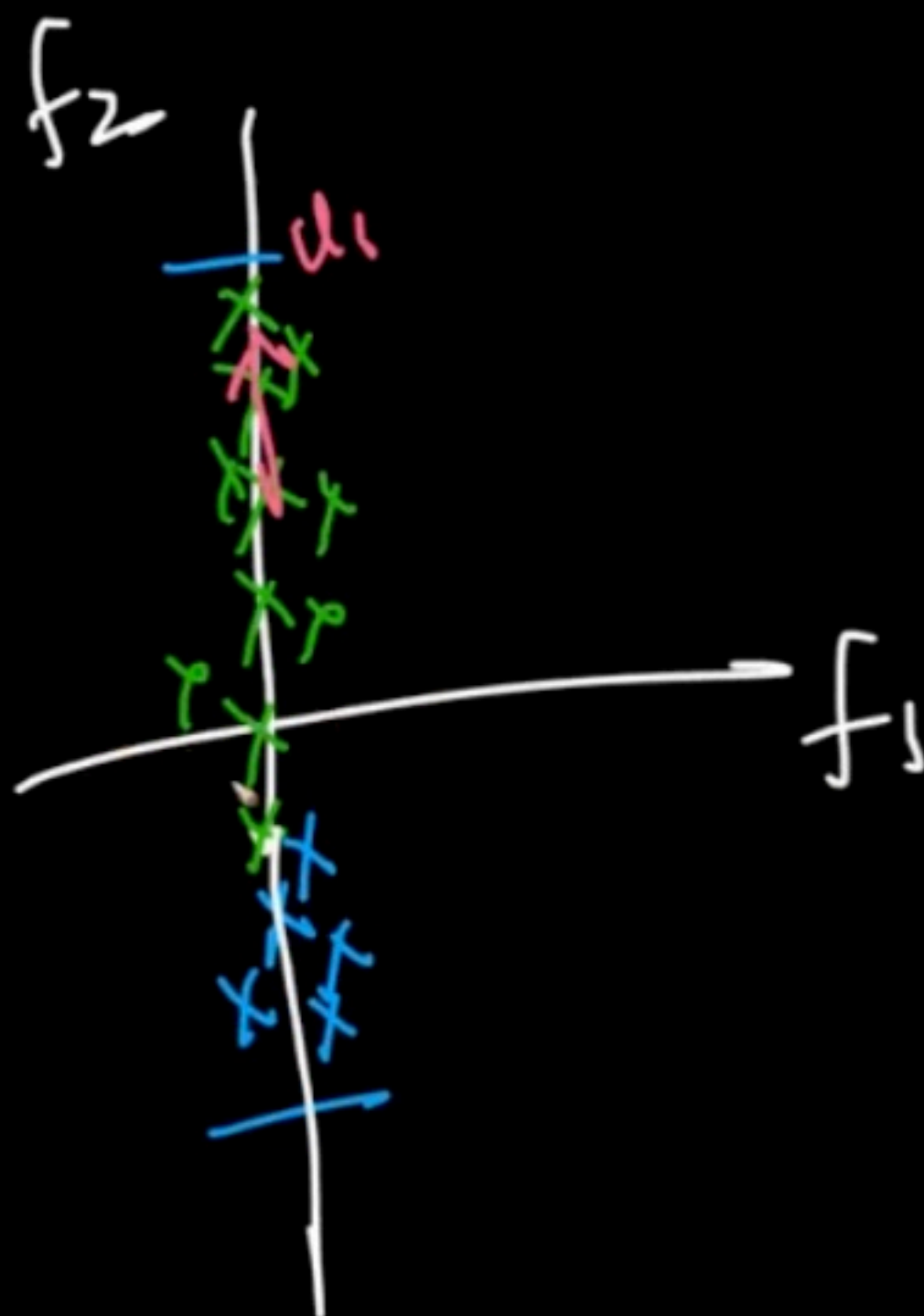
Identity matrix

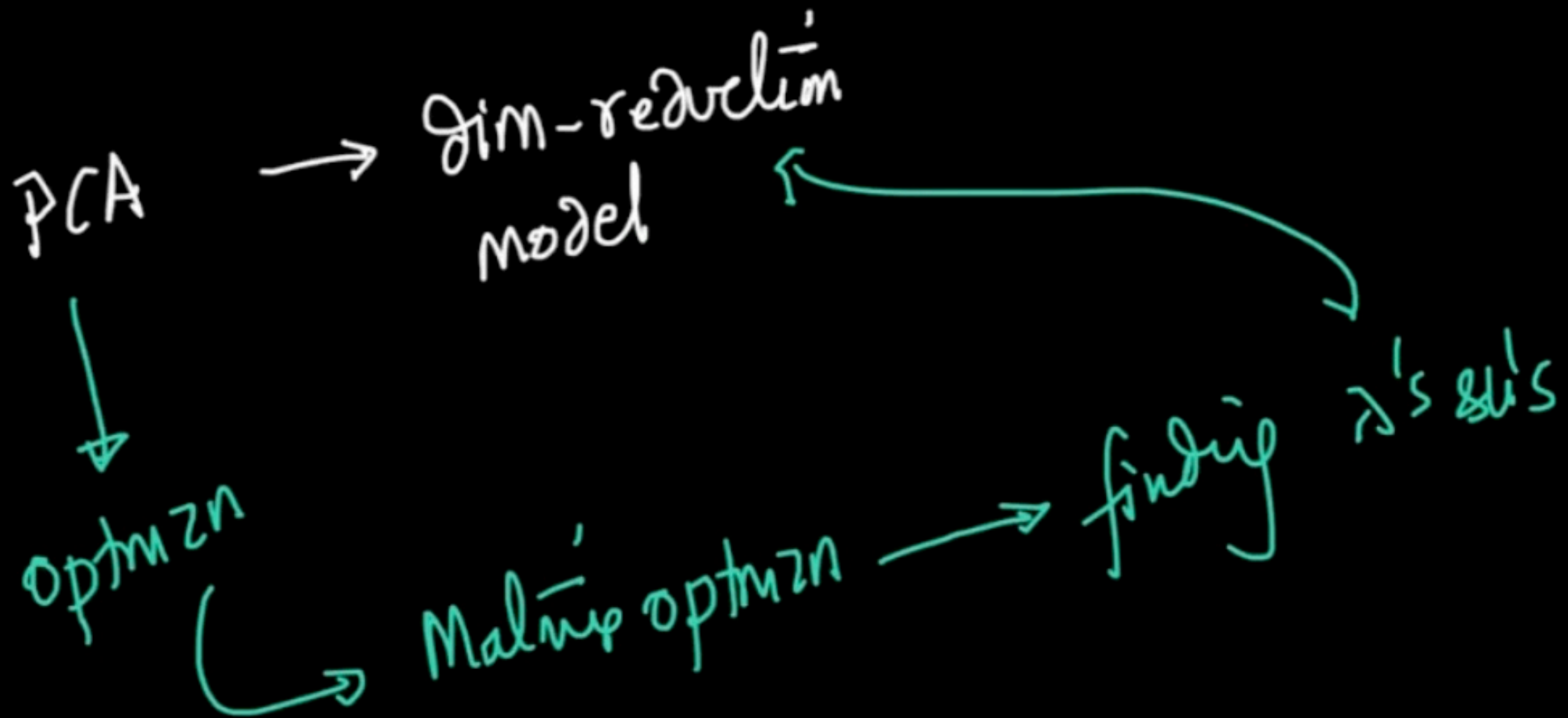
$$S = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

u_1, u_2, \dots
principal components



(Q)





(Q)

100-d

→ Classification or regression

Linear Reg
or
Logistic Reg
or
KNN

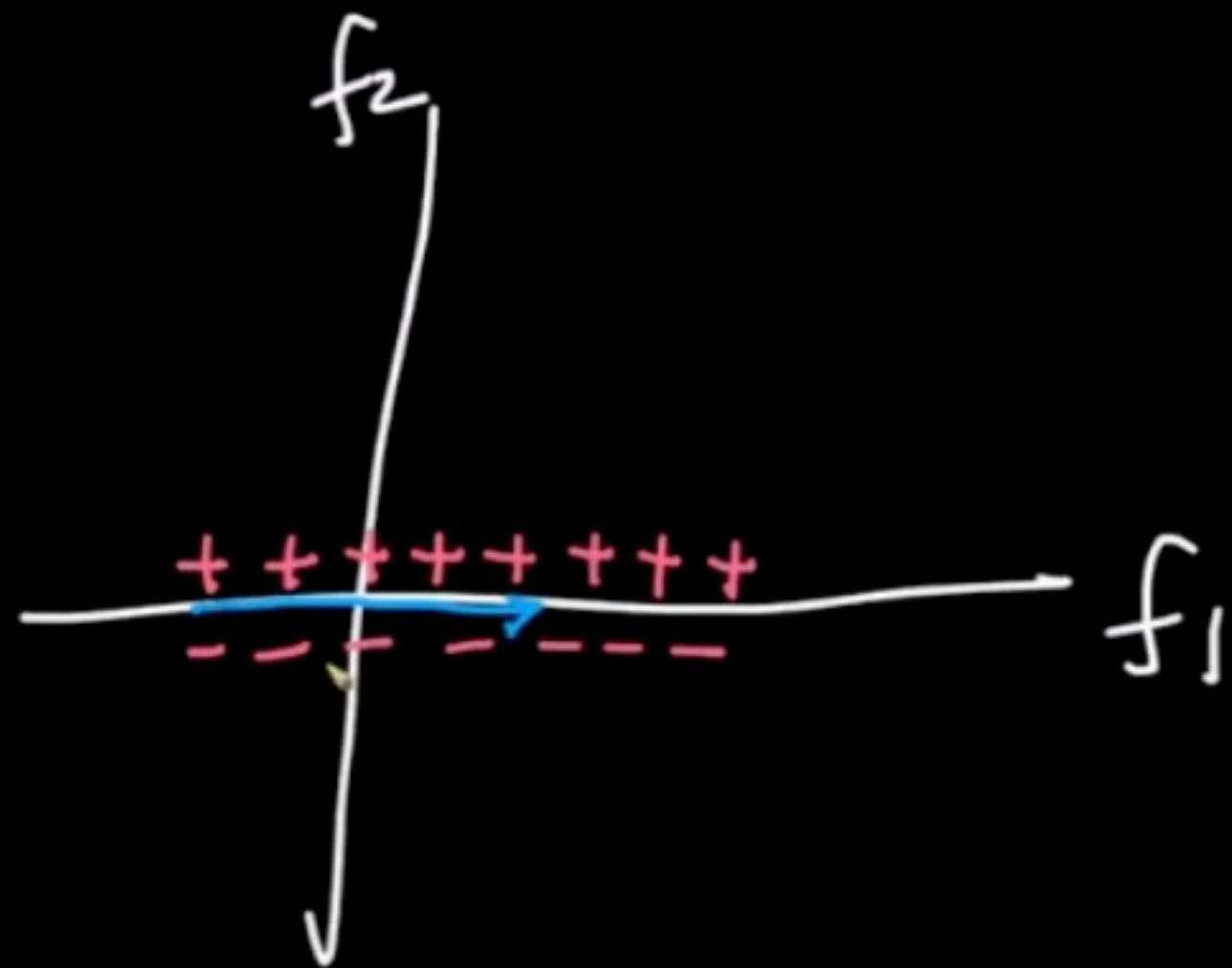
5-dim

→ ML

($u_1, u_2 \dots u_5$)

↓
loss of information

↖
(you can)
(why?)

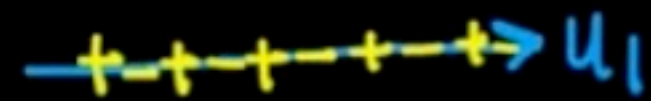


2D

SCALER

does not use class labels

PCA



↓
Classfn?

(Q) Feature map \rightarrow Classification or regression

Linear Reg
or
Logistic Reg
or
KNN

100-d

5-dim
($u_1, u_2 \dots u_5$)

\rightarrow ML

(you can)
(why?)

loss of information