

# Principal Component Analysis

[HOR-1]

- Reachability Distance
- Comparison of outlier detectors
- PCA

## Reachability Distance :

$$RD(A, B) = \max(KD_B, \text{dist}(A, B))$$

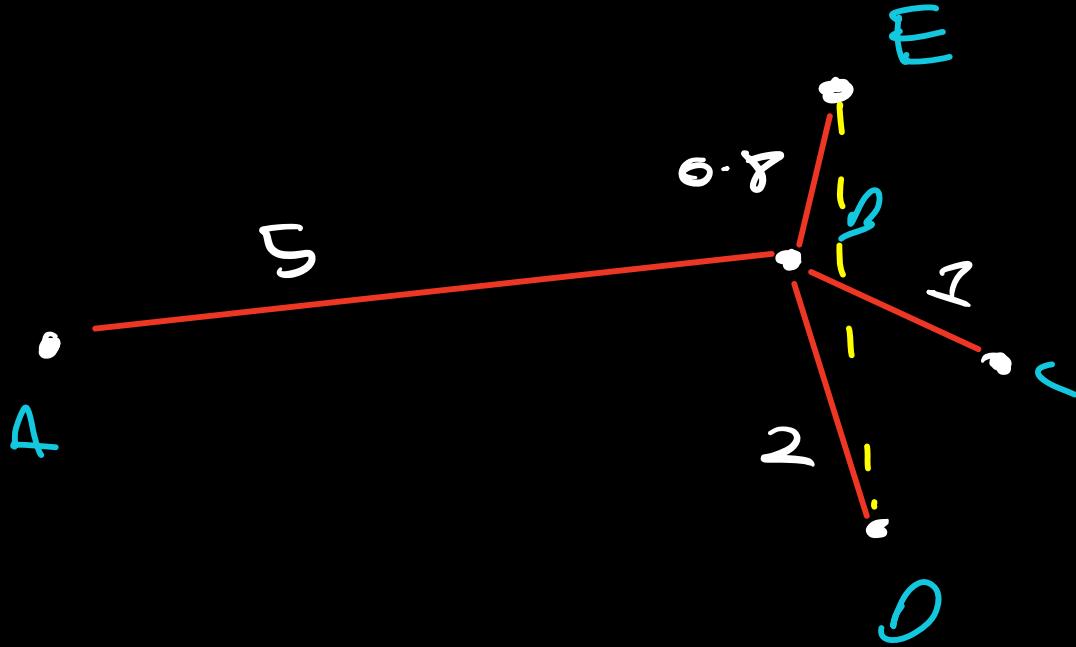
$$\text{Density}(A) = \frac{1}{K} \sum RD(A, \text{Neighbours})$$

This means that we will use distance between A, B but it has to be min equal to  $K\text{-dist}(B)$   
in other words!

RDC me , you) :

if you are in my top ' $k$ ' but I  
am not in your top ' $k$ ' then  
let's use actual distance but if  
we both are in each other's top ' $k$ '  
then let's use your  $k$ -dist

Note: RD is not symmetric.



$$RD(A, B) = 5 \quad \text{not } 1$$

$$RD(C, B) = 2$$

$$RD(E, B) = 2$$

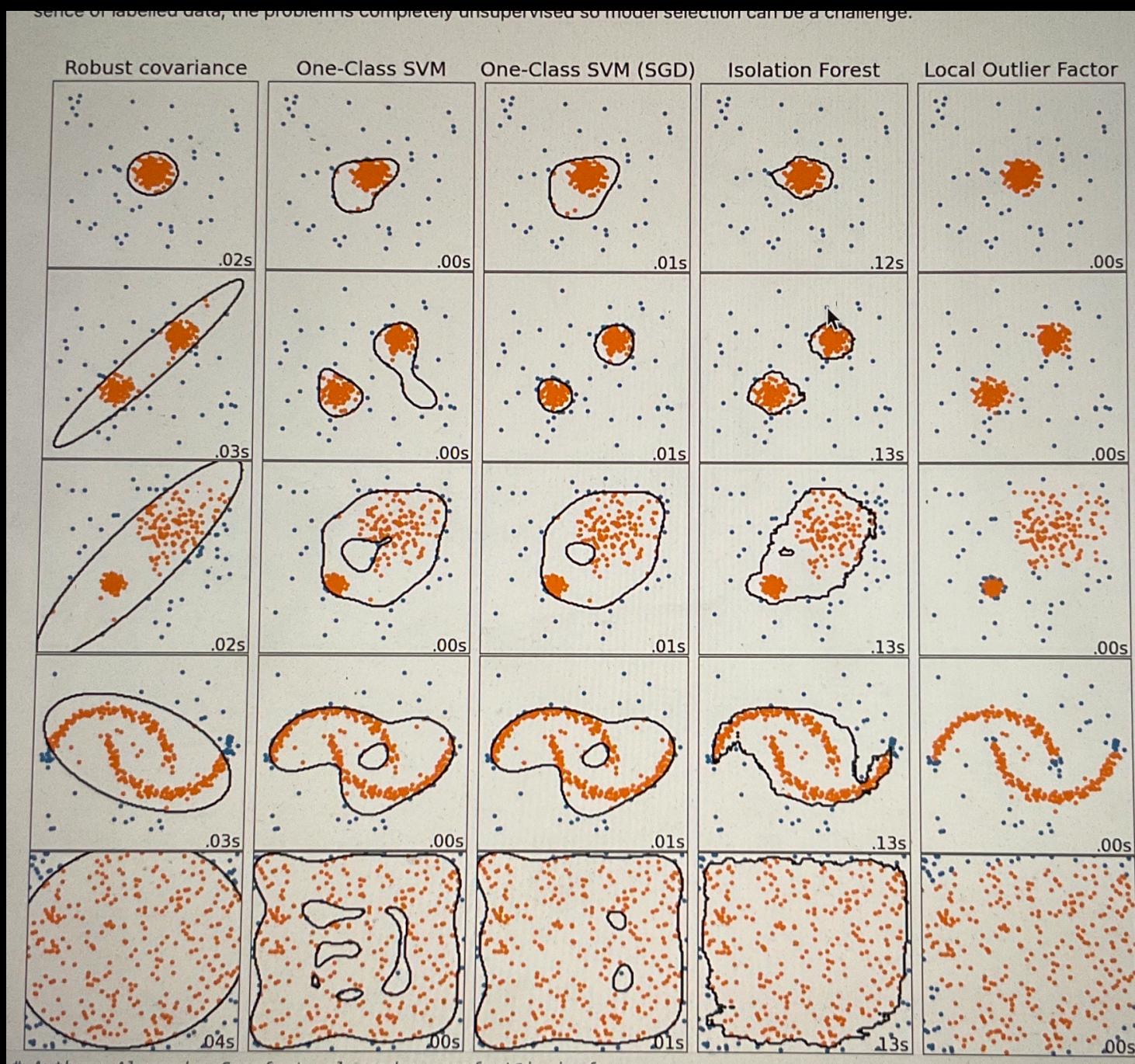
$$RD(D, B) = 2$$

$$RD(B, E) = \underline{\underline{2.5}}$$

This is less intuitive but has some benefits when  $1^{\text{st}}$  neighbour and  $k^{\text{th}}$  neighbour distance may be very different.

# Comparison

In the absence of labelled data, the problem is completely unsupervised so model selection can be a challenge.



# Principal Component Analysis

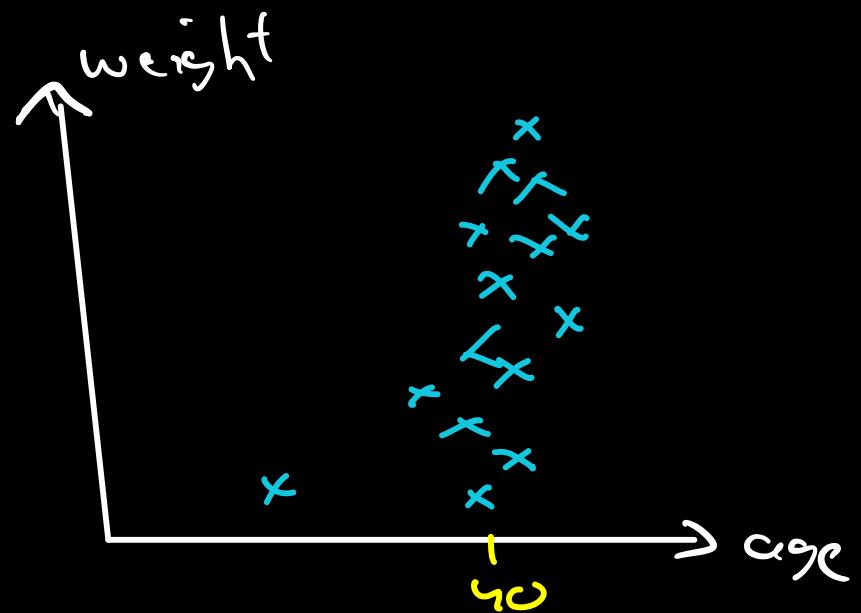
Let's take a few examples:

Diabetes prediction:

W	A	O
:	:	Y
:	:	~Y
:	:	Y
!	!	Y
!	!	Y
!	!	Y

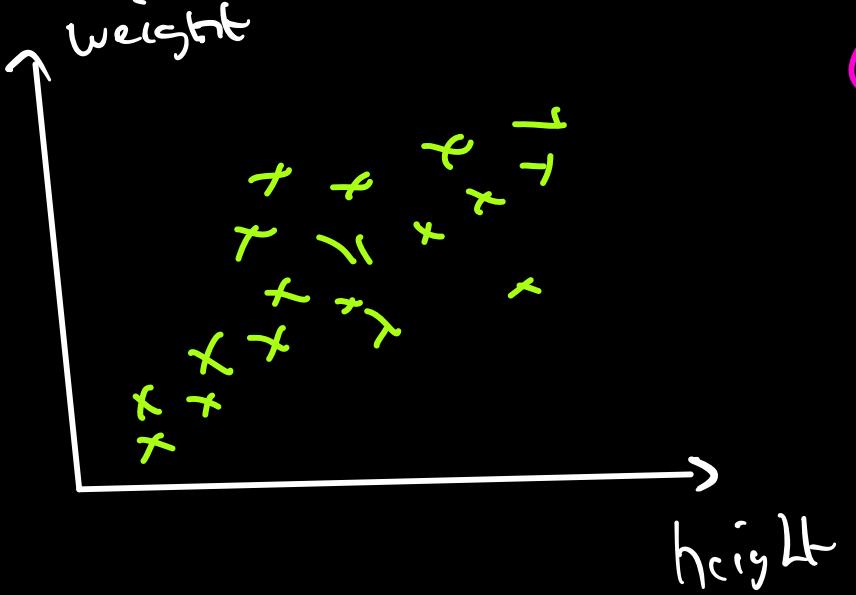
$$f_1 = \text{weight}$$

$$f_2 = \text{age}$$

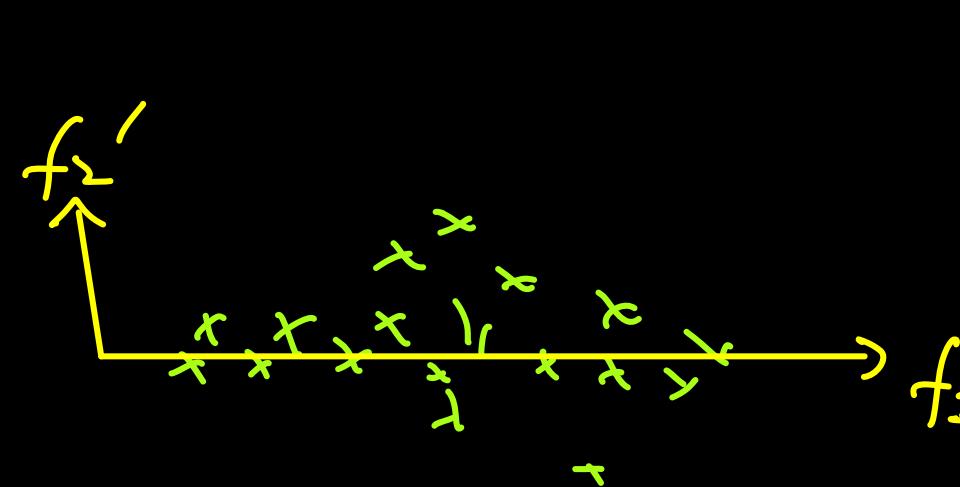
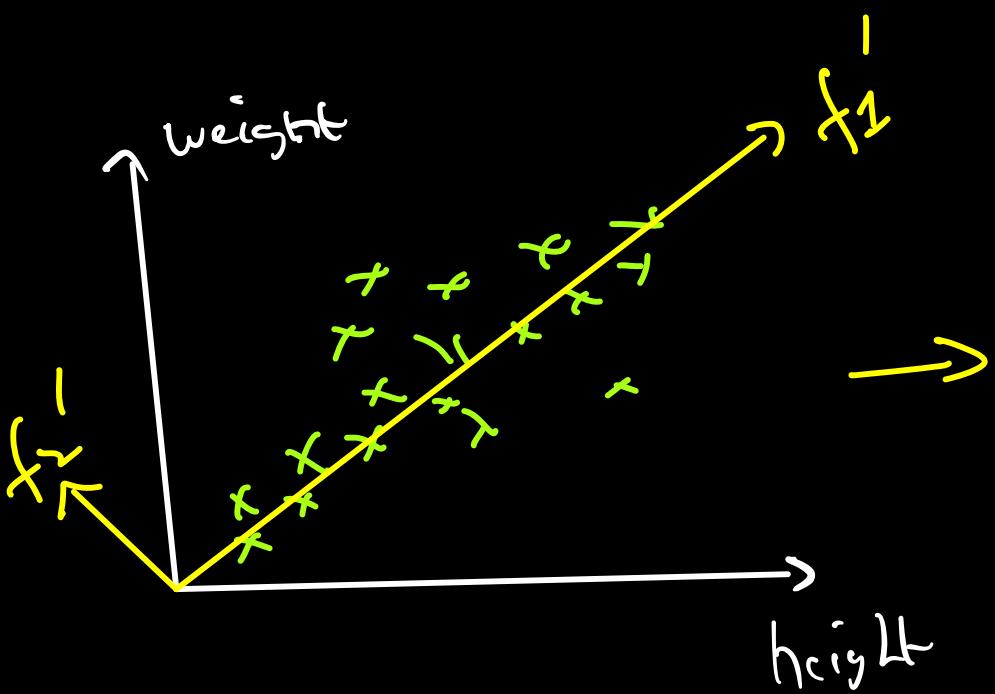


Q: If you could only use 1 feature  
which would you use?

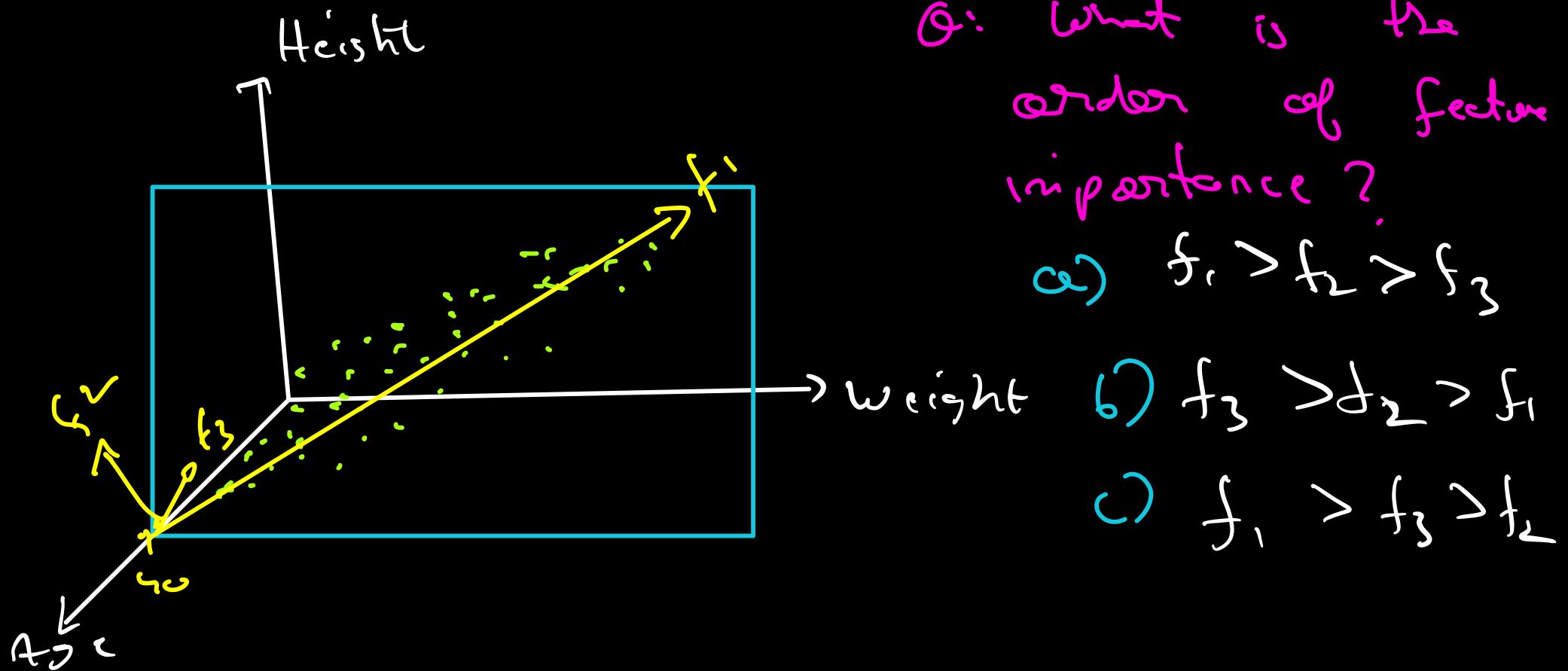
- a) Age
- b) weight
- c) A + w



- Q: Which feature?
- weight
  - height
  - $w + h$



- Q: Which features?
- $f_1$
  - $f_2$
  - $f_1 + f_2$
  - $f_1 - f_2$



Q: What is the order of feature importance?

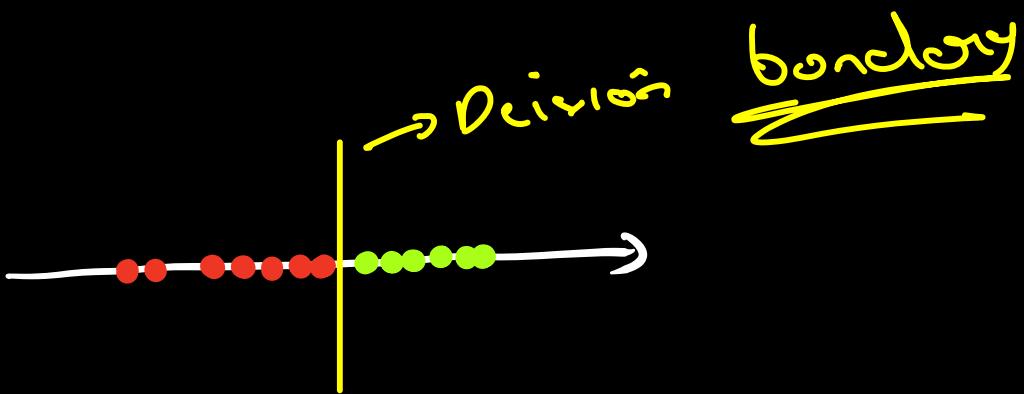
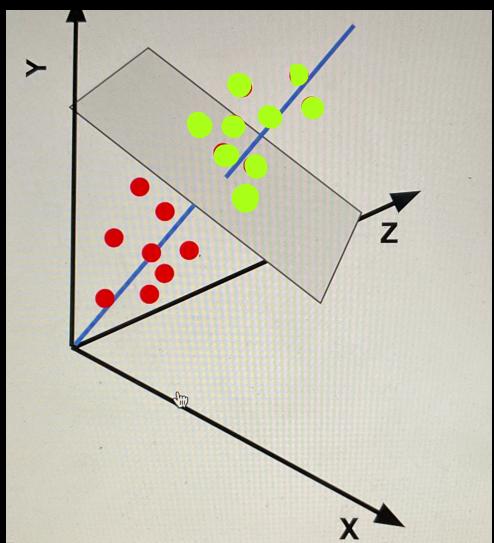
- a)  $f_1 > f_2 > f_3$
- b)  $f_3 > f_2 > f_1$
- c)  $f_1 > f_3 > f_2$

### Goal:

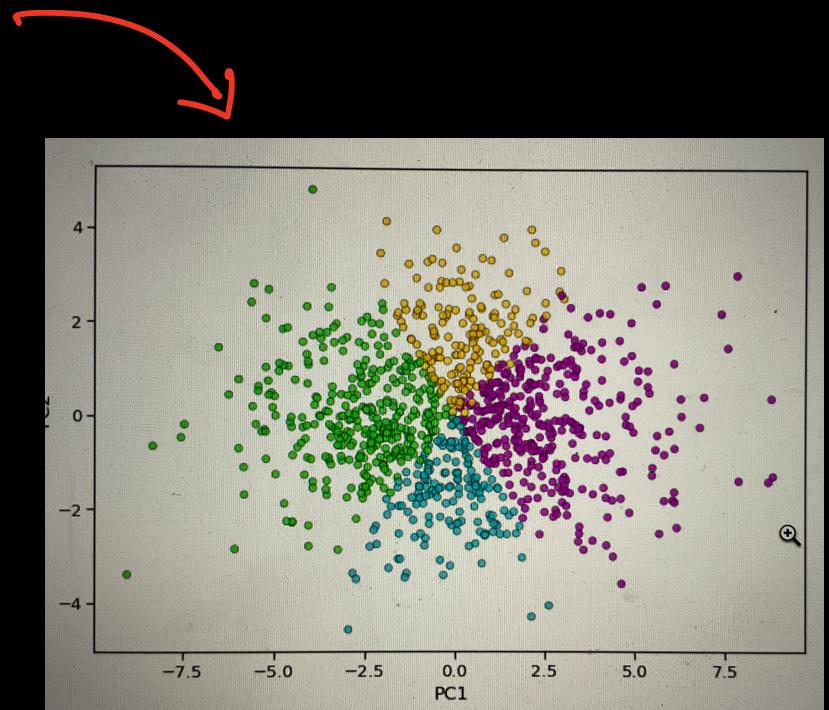
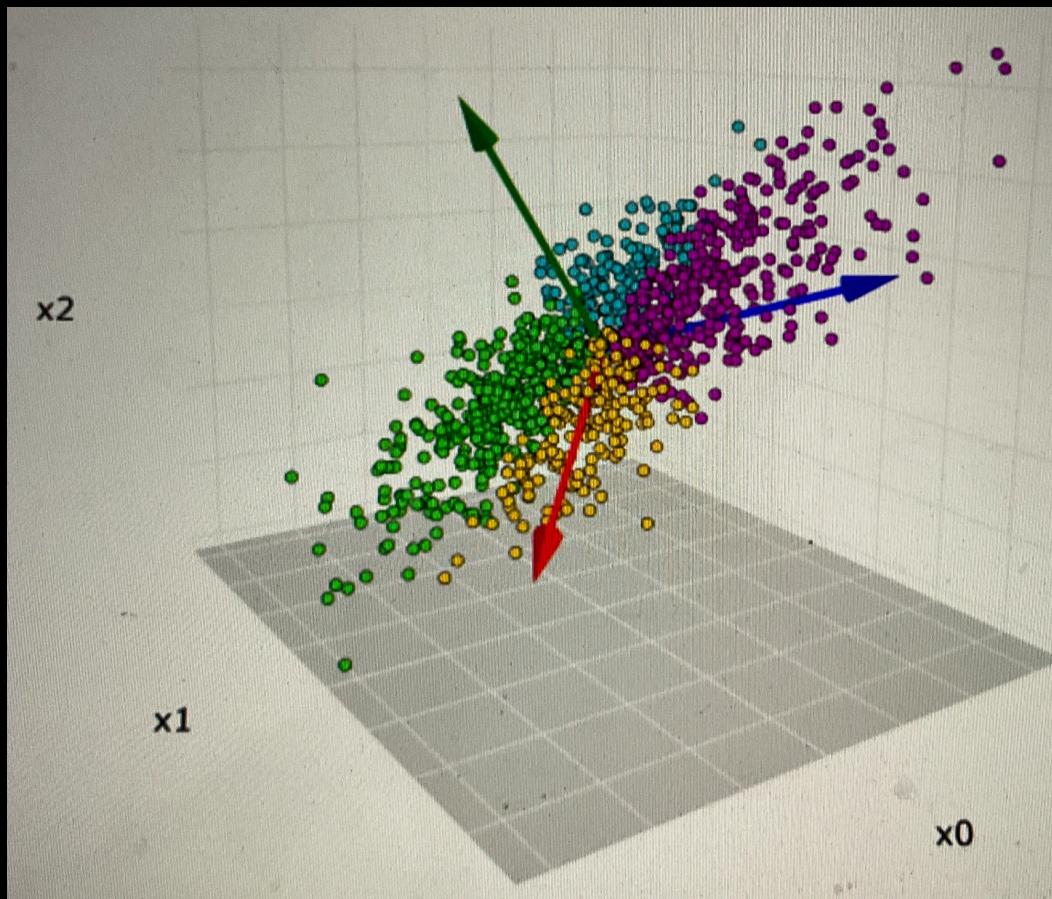
- We see that all the features don't carry equal amounts of info.
- Sometimes in ML, we need to reduce the number of features.

- Visualisation
- compression (space saving)
- Faster ML training

3D → 1D

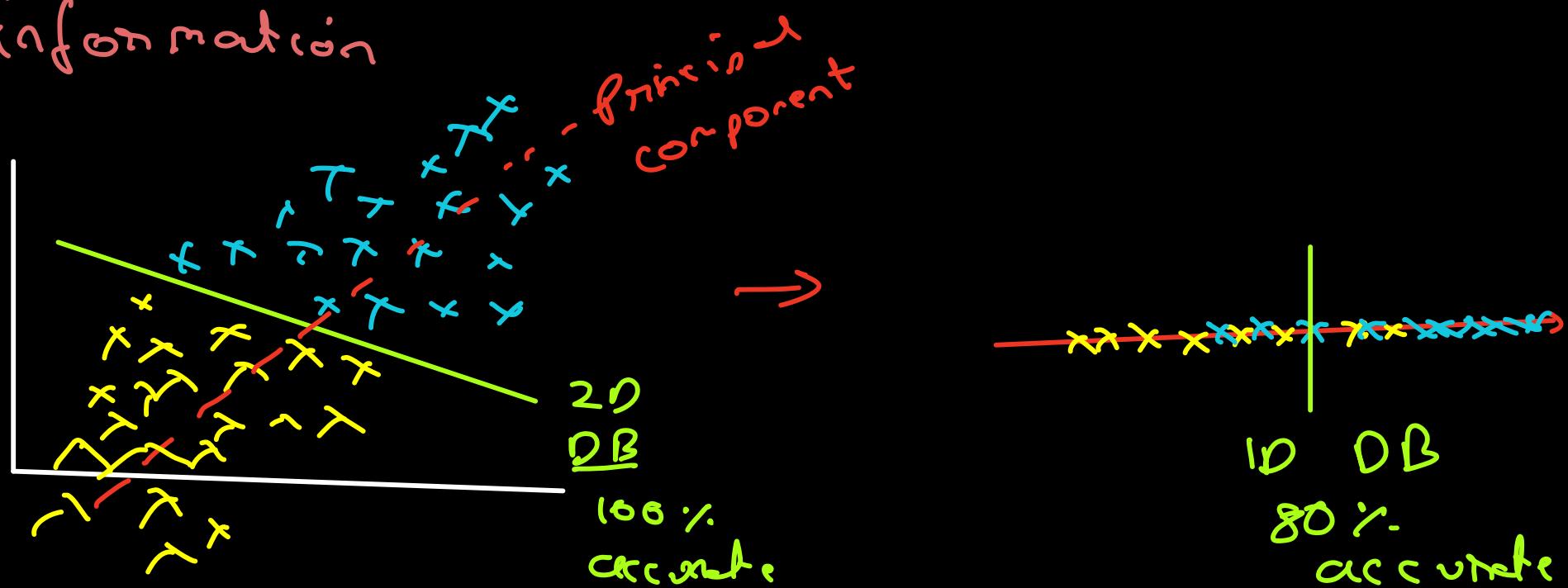


3D - 2D



Principal component analysis refers to the act of finding new axis to represent the data so that a few principal components may contain most information.

However, this may lead to some loss of information

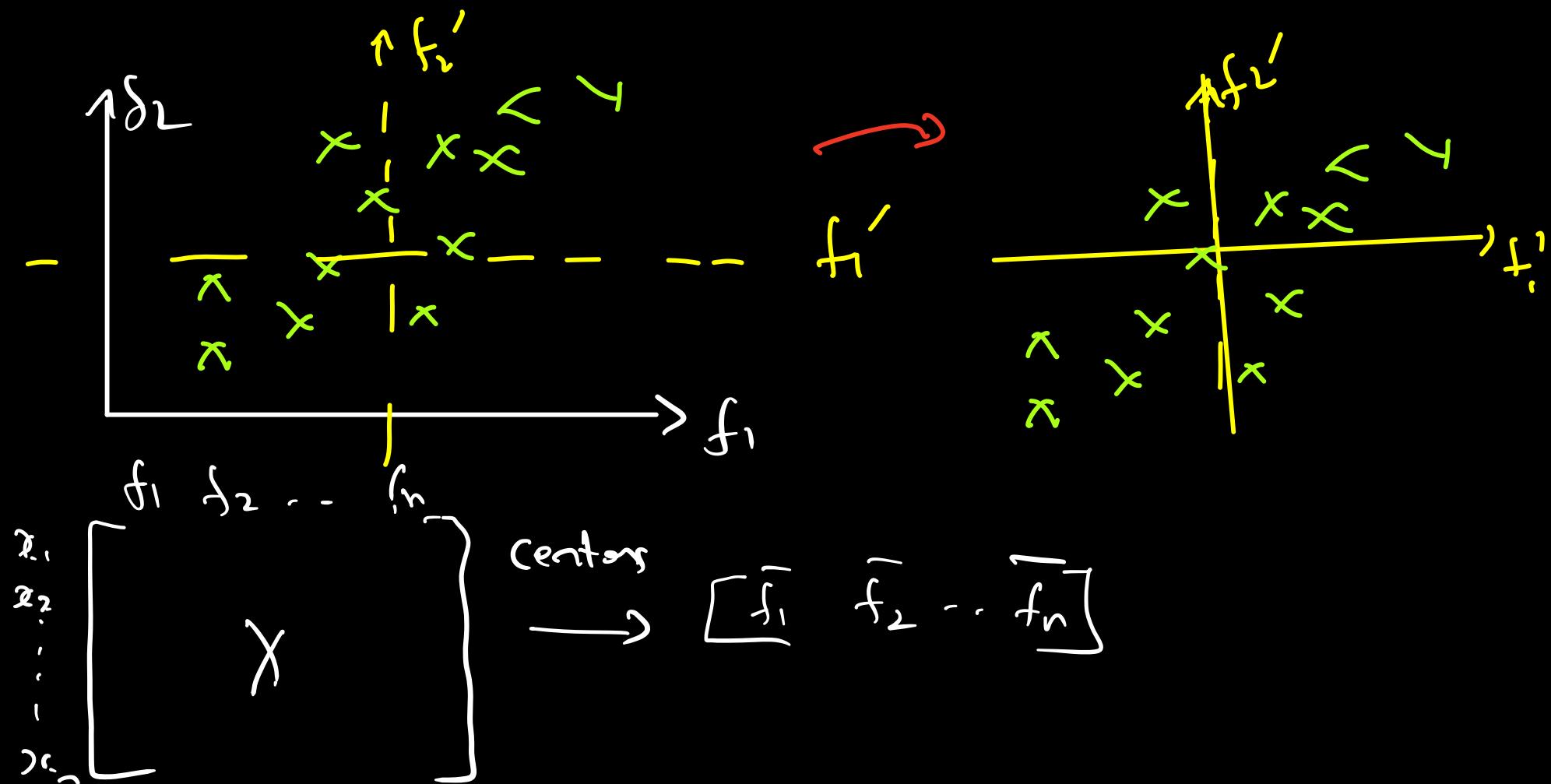


# PCA Mathematics

→ Change of axis

→ Retain most information in few axis.

Step-1 : Centering



$$X - X_{\text{mean}} = X_c$$

## Step 2: Standardize / Scaling

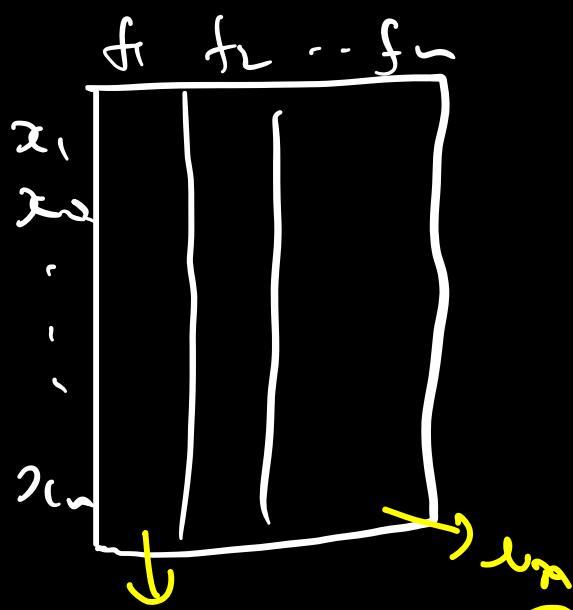
We need to make different features such as height (4-6) ft and weight (40 - 180) kg comparable.

Best is to use Standard Scaling because it is less affected by outliers

Combining step 1 and 2

$$X' = \frac{X - X_{\text{mean}}}{X_{\text{std}}}$$

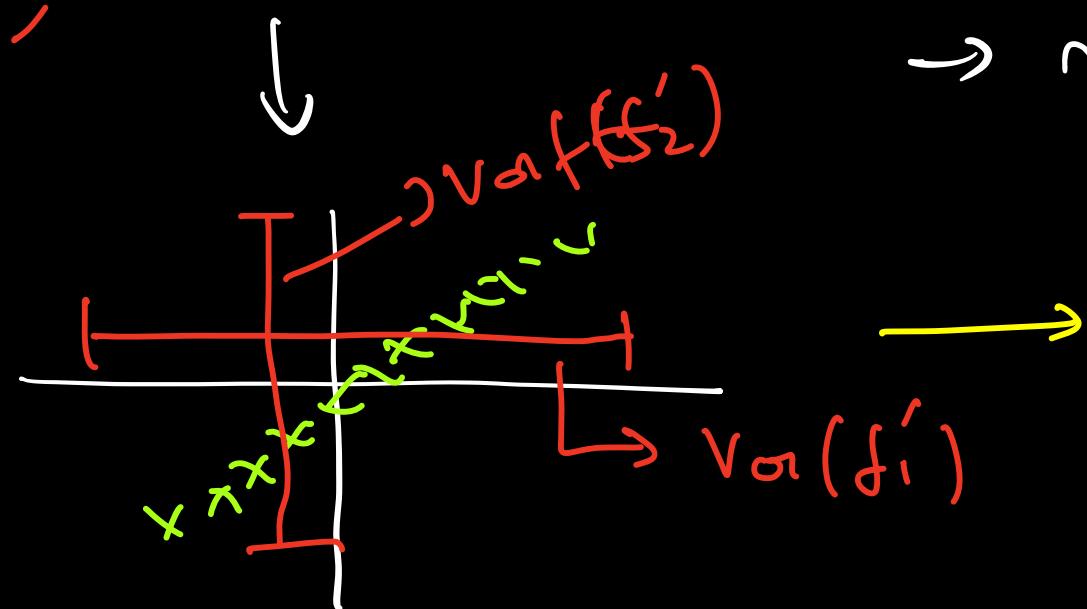
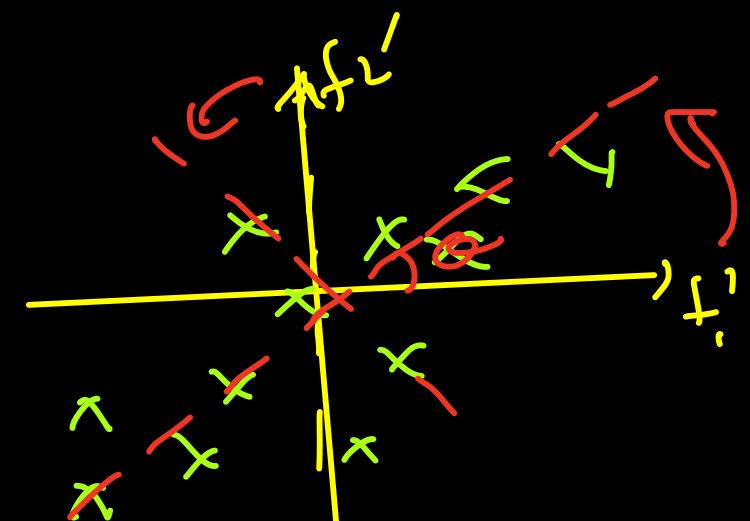
→ use sklearn



$u_1, v_1$

$u_2, v_2$

### Step - 3: Rotate

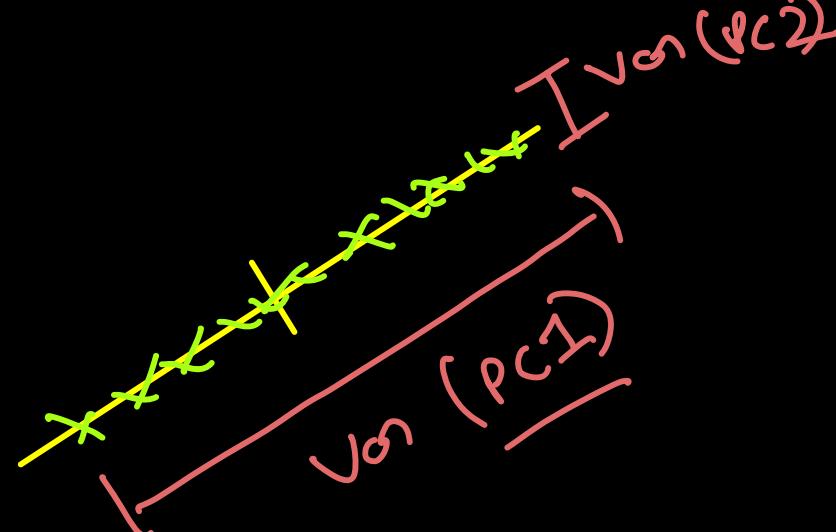


How to rotate?

Q: Any suggestions on how to determine the best axis?

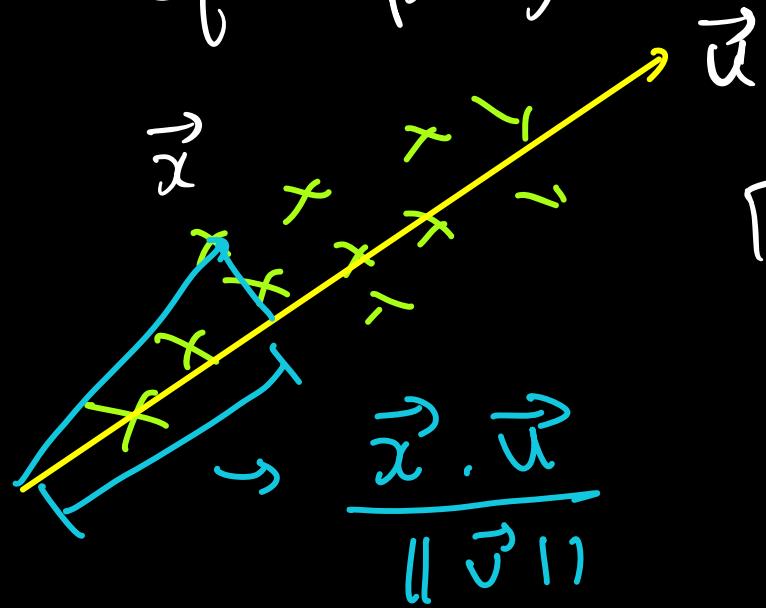
→ maximise variance

→ max length of projection



$$v_{\text{on}}(p(1)) > v_{\text{on}}(g_1') > v_{\text{on}}(f_2') > v_{\text{on}}(p(2))$$

We can also think in terms of length of projection:



Best  $\vec{v}$  will be where

$$\max_{\vec{v}} \sum_{i=1}^n \frac{\vec{x}_i \cdot \vec{v}}{\|\vec{v}\|} \cdot \frac{1}{n}$$

However, projections can also be -ve

$$\max_{\vec{v}} \sum_{i=1}^n \frac{|\vec{x}_i \cdot \vec{v}|}{\|\vec{v}\|} \cdot \frac{1}{n}$$

But this is not differentially

$$\max_{\vec{v}} \frac{1}{n} \sum_{i=1}^n (\vec{x}_i \cdot \vec{v})^2 \rightarrow \text{still hard to diff}$$

↓

constraint:  $\|\vec{u}\|=1$

$$\boxed{\max_{\vec{v}, \lambda} \frac{1}{n} \sum_{i=1}^n (\vec{x}_i \cdot \vec{v})^2 + \lambda (\|\vec{v}\| - 1)}$$

Lagrangian loss function for P.C.A

This can be solved using gradient descent.

However there is also a direct linear algebra way.

Math Ahead !!  $\rightarrow$  Not required to remember/  
reproduce

- We will use some identities w/o proof
- We will learn a new concept

(dimension)  
scaled square  
 $\downarrow$  (homogeneous)

Variance  $\downarrow$

$$\max_{\vec{u}, \lambda} \frac{1}{n} \sum_{i=1}^n (\vec{x}_i \cdot \vec{u})^2 + \lambda (\|u\|^2 - 1)$$

$$\max_{\vec{u}, \lambda} \frac{\left( \sum_i \vec{x}_i \cdot \vec{u} \right)^2}{n} + \lambda (\|u\|^2 - 1)$$

$$\text{identity : } A^T = A^T A = \|A\|^2$$

if  $A$  is invertible square matrix

$$\rightarrow \max_{\vec{u}, \lambda} \frac{(\vec{x}\vec{u})^T (\vec{x}\vec{u})}{n} + \lambda(\vec{u}^T \vec{u} - 1)$$

$$\text{identity } (AB)^T = B^T A^T$$

$$\rightarrow \max_{\vec{u}, \lambda} \vec{u}^T \frac{\vec{x}^T \vec{x}}{n} u + \lambda(\vec{u}^T \vec{u} - 1)$$

$$\rightarrow \max_{\vec{u} - \lambda} \vec{u}^T V \vec{u} + \lambda (\vec{u}^T \vec{u} - 1)$$

V matrix:

$$\frac{\vec{x}^T \vec{x}}{n} = \frac{1}{n} \left[ \begin{matrix} x \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{matrix} \right]^T \left[ \begin{matrix} x \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{matrix} \right]$$

$d \times n \quad \cdot \quad n \times d$

Covariance matrix

$$= \begin{bmatrix} \text{var}(f_1) & \text{cov}(f_1, f_2) & \dots & \text{cov}(f_1, f_d) \\ \text{cov}(f_2) & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \text{var}(f_d) \end{bmatrix}$$

↳ This matrix has pairwise covering of all features [similar to pd. cov()]

So now we have

$$L = u^T V u + \lambda(u^T u) - \lambda$$

$$\frac{\partial L}{\partial u} = 2Vu + 2\lambda u$$

$$\frac{\partial L}{\partial \lambda} = u^T u - 1$$

Now we set both to 0

$$\therefore u^T u = 1$$

identities

$$\frac{2A^T B A}{A} = 2BA$$

$$\frac{2A^T A}{2A} = 2A$$

this is like

$$\frac{2x \cdot x}{2x} = 2x$$

but in matrix

$$\nabla u = \underbrace{-\lambda}_{\downarrow} u$$

$$\nabla u = \underbrace{x' u}_{\substack{\downarrow \\ \text{matrix}}} \rightarrow \text{vector}$$

matrix      vector      scalar

Interesting

→ we have  $\nabla$ , which is just  $\text{ppcov}(x)$

→  $\vec{v}$  must be such that when I multiply  $\nabla$  to it, it is equivalent to multiplying some scalar to it.

# Eigen Values and Vectors

For any matrix  $A$ , there exists\* one or more vectors, such that

$$A \vec{x} = \lambda \vec{x} \quad \begin{matrix} \text{eigen vector} \\ \downarrow \\ \text{eigen value} \end{matrix}$$

Geometrically

$$\text{Let's say: } \vec{x} = \begin{bmatrix} 2 \\ 5 \end{bmatrix}, \quad A = \begin{bmatrix} -6 & 3 \\ 4 & 5 \end{bmatrix}$$

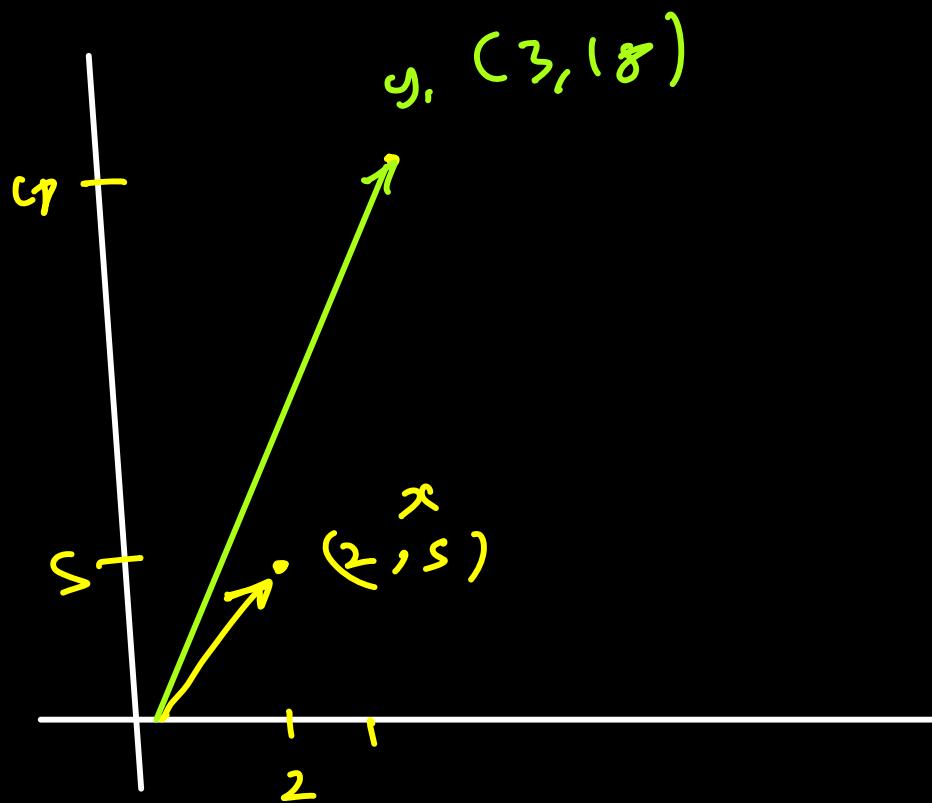
when we multiply  $A \vec{x} = \vec{y}$ , then

$\vec{y}$  has a magnitude and direction, but if the direction of  $\vec{x}$  and  $\vec{y}$  are same, then  $\vec{x}$  is said to be an eigen vector of A with eigen value

$$\lambda = \frac{\|\vec{y}\|}{\|\vec{x}\|}$$

There can be multiple eigen vectors, which are always orthogonal (+) to each other.

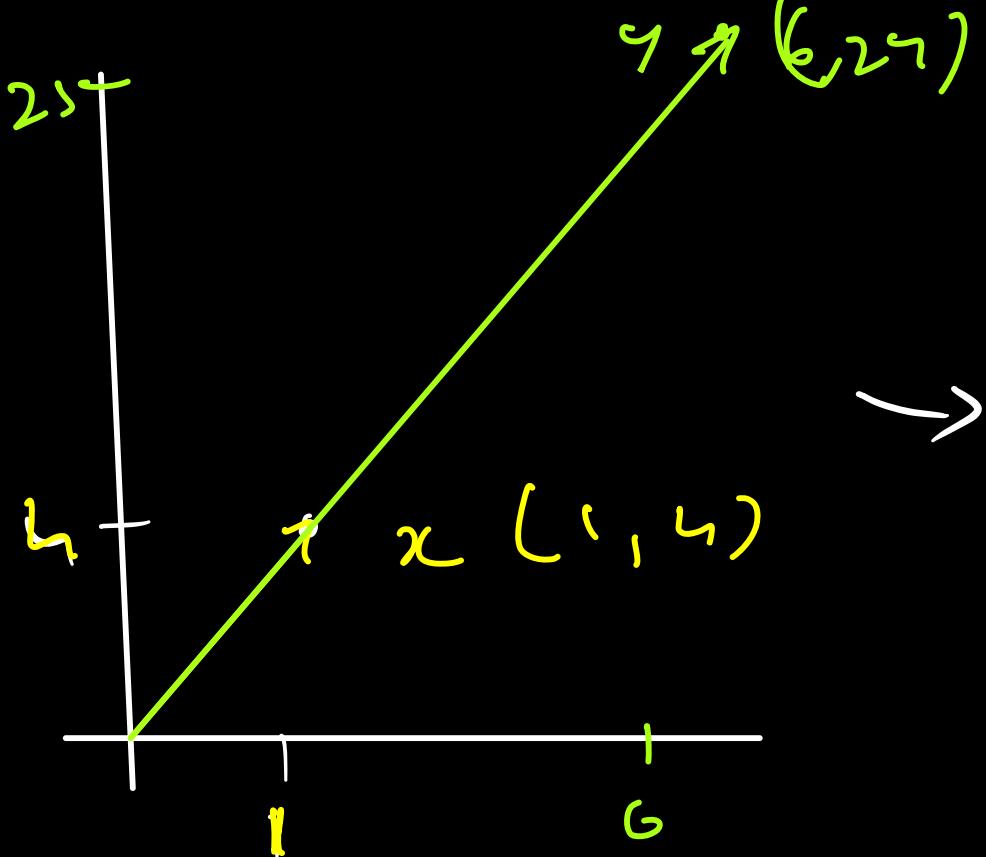
Eg:  $\begin{bmatrix} -6 & 3 \\ 4 & 5 \end{bmatrix} \begin{bmatrix} 2 \\ 5 \end{bmatrix} = \begin{bmatrix} -12 + 15 \\ 8 + 10 \end{bmatrix} = \begin{bmatrix} 3 \\ 18 \end{bmatrix}$



But, if  $x = 1, 4$

$$\begin{bmatrix} -6 & 3 \\ 4 & 5 \end{bmatrix} \begin{bmatrix} 1 \\ 4 \end{bmatrix} = \begin{bmatrix} 6 \\ 24 \end{bmatrix} = 6 \begin{bmatrix} 1 \\ 4 \end{bmatrix}$$

↓  
 eigen  
 value      ↓  
 eigen  
 vector



$$y = \underline{\underline{6 \cdot x}}$$

$\rightarrow x$  is such a vector that will not change direction when multiplied by  $\underline{\underline{\lambda}}$

We can find eigen values and vectors very easily in numpy.

$\rightarrow \underline{\underline{\text{Code}}}$