

Principal Component Analysis

[Dimensionality Reduction]

- Doubt
- Variants of GRD
- PCA

Doubt in last class:

→ gradient does not change:

$$\frac{\partial L}{\partial w_i} = - \sum y_i x_{1i}$$

This is not the complete gradient !!

$$\frac{\partial L}{\partial w_i} = - \sum y_i x_{1i} + \frac{\lambda}{2} (\sqrt{w_i^2 + \dots})$$

$$\frac{\partial L}{\partial w_i} = - \sum y_i x_{1i} + \frac{1}{2 \|w\|} \cdot 2 w_i = \underline{F(w)}$$

Hence gradient does change.

Variants of Gradient Descent

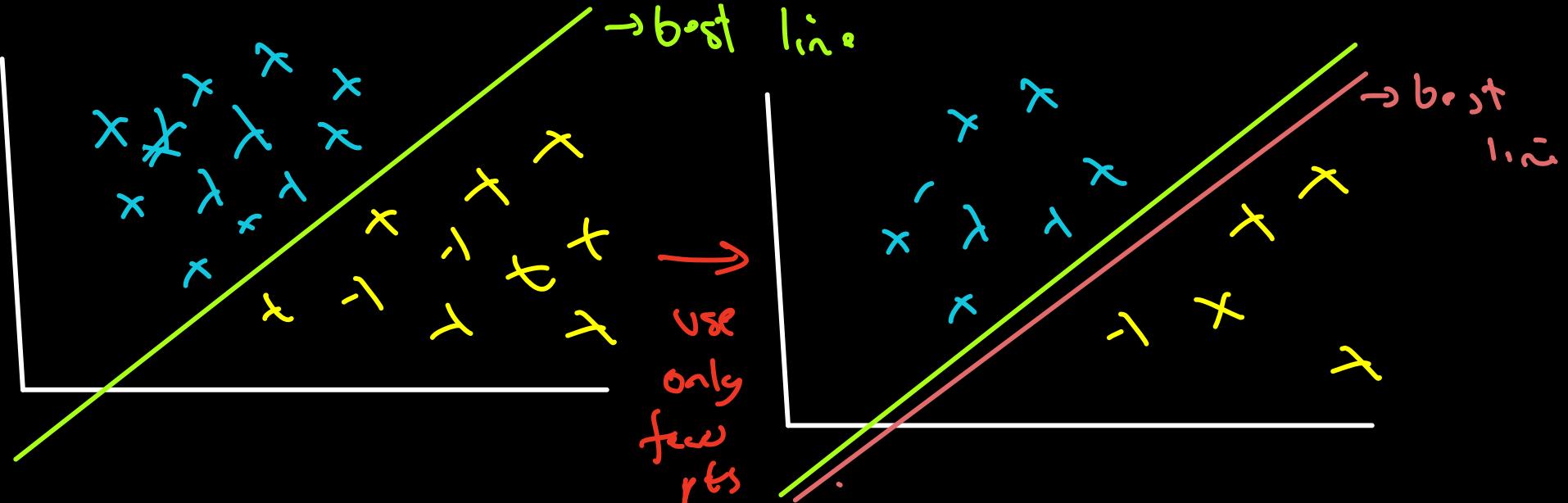
$$\text{GD: } \theta^{t+1} = \theta^t - n \frac{\partial f}{\partial \theta}$$

$$\rightarrow \theta^{t+1} = \theta^t - n \sum_{i=1}^n \frac{\partial f(x_i)}{\partial \theta}$$

Use all data points
for 1 update.

→ high compute time

Q: Let's say I have 1M data points,
can I just update θ with 0.5 M points?



→ But both best lines are similar !!

Idea:

For a single update use only a few data points (batch) randomly and do many updates.

- ↳ This is called
- Very similar accuracy
- Huge speed improvement

Batch Gradient
Descent

Batch GD:

$$\theta^{t+1} = \theta^t - n \sum_{i \in B} \frac{\partial f(x_i)}{\partial \theta}; B \text{ is a random sample of } N$$

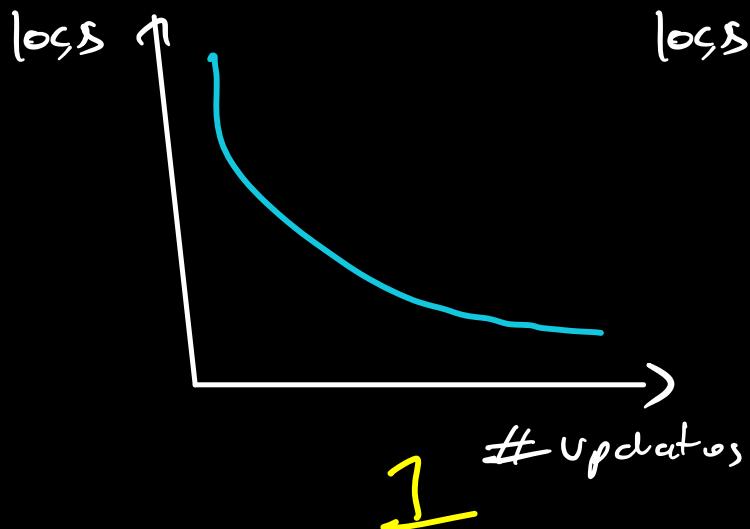
Q: Could I just use 1 point?

Stochastic GD:

$$\theta^{t+1} = \theta^t - n \frac{\partial f(x_k)}{\partial \theta}; k \text{ is random number from } 1-n$$

→ Faster but too many updates are needed
→ less preferred.

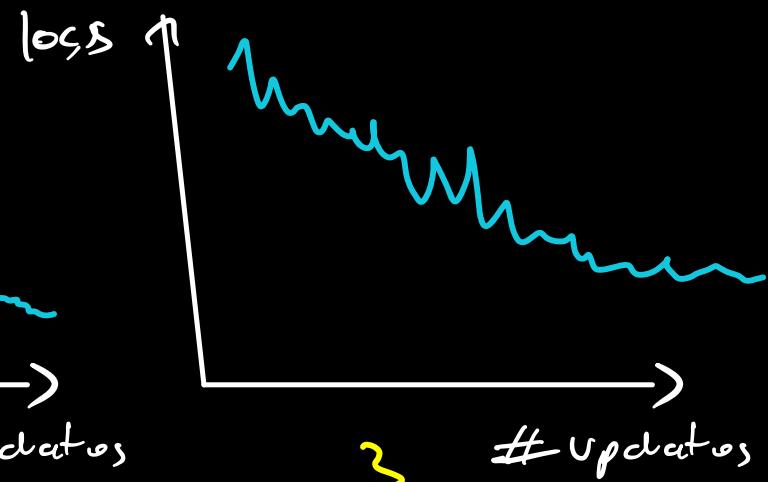
Quiz: Match loss trajectory of respective variant



1 #updates



2 #updates



3 #updates

- a) 1 - sgd b) 1 → gd c) 1 → gd
2 → bsgd 2 → sgd 2 → bsgd
3 → gd 3 → bsgd 3 → sgd

Principal Component Analysis

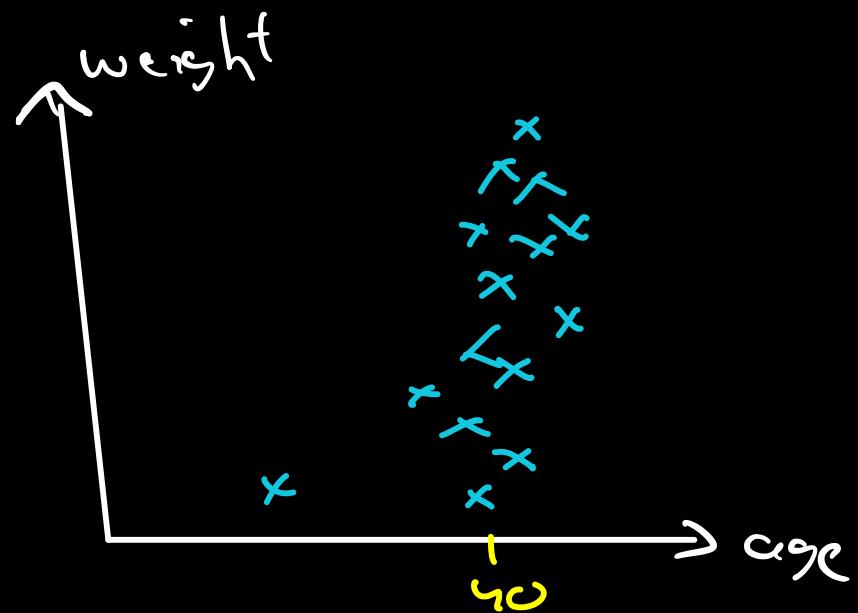
Let's take a few examples:

Diabetes prediction:

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

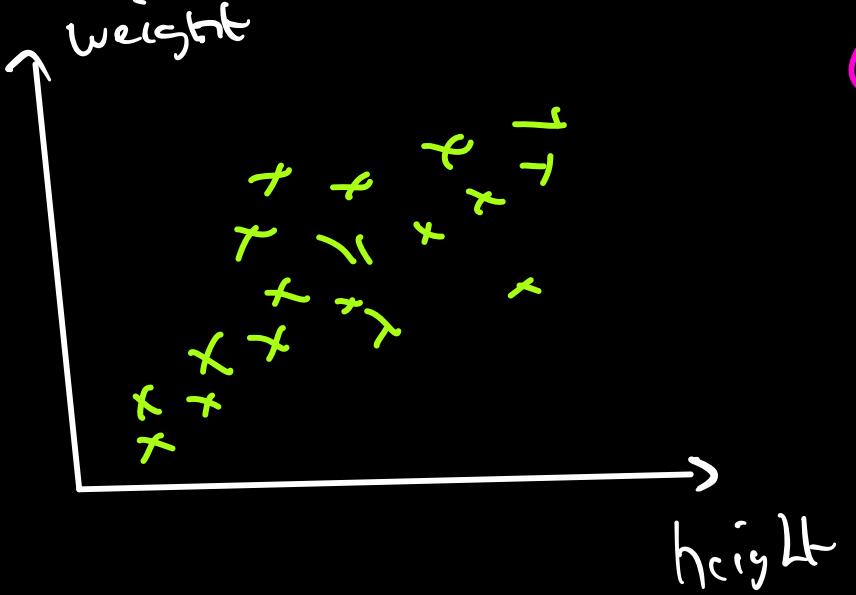
$$f_1 = \text{weight}$$

$$f_2 = \text{age}$$

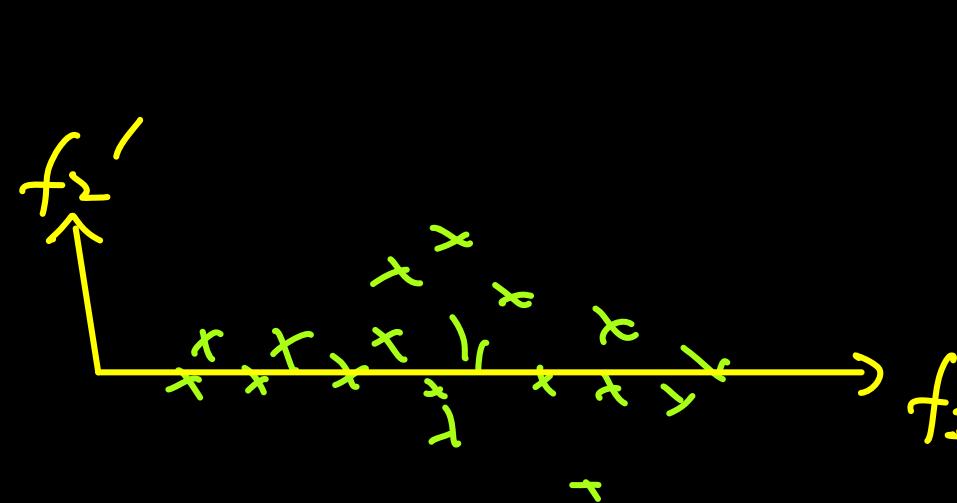
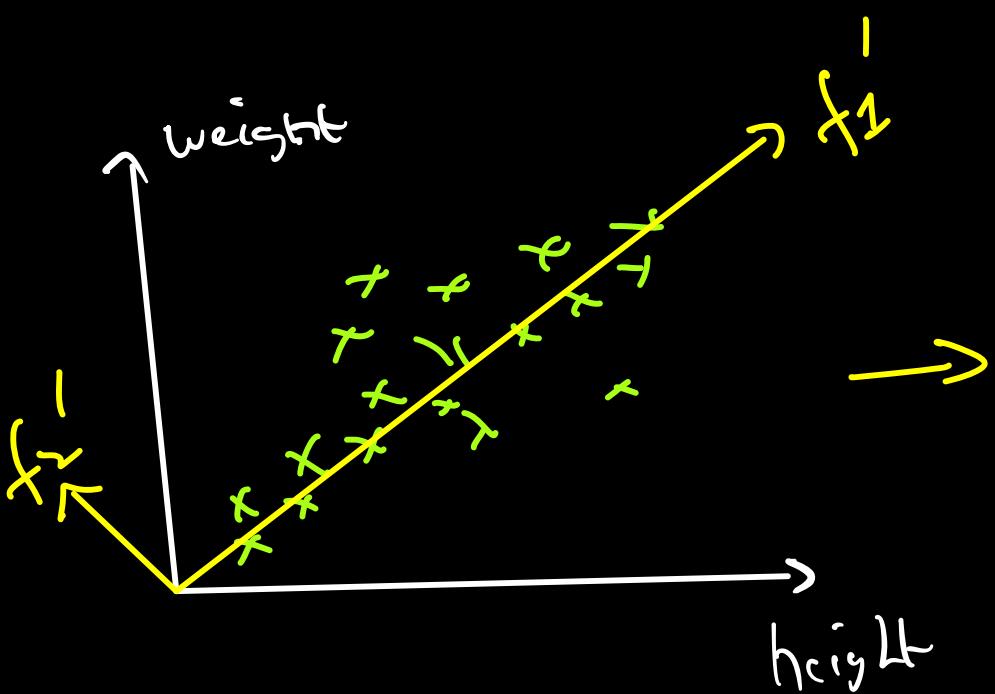


Q: If you could only use 1 feature
which would you use?

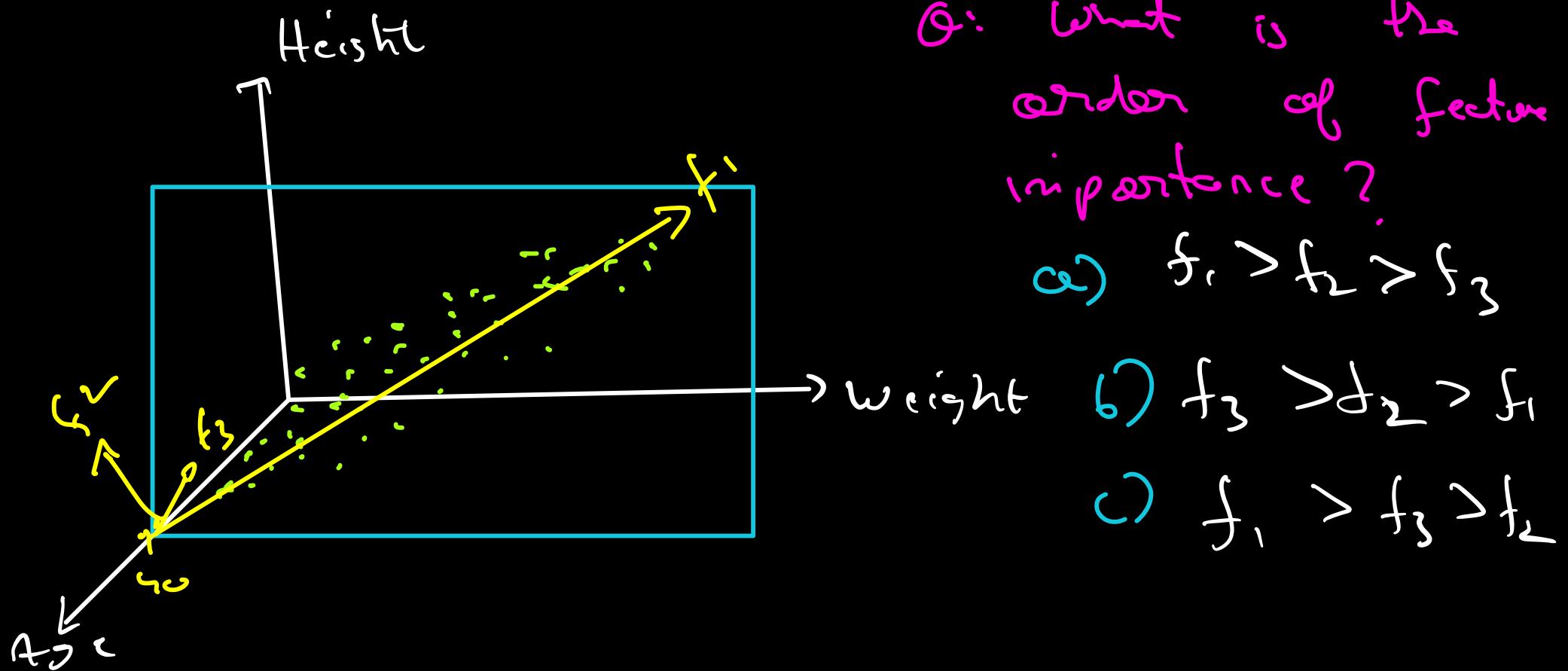
- a) Age
- b) weight
- c) A + w



- 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?

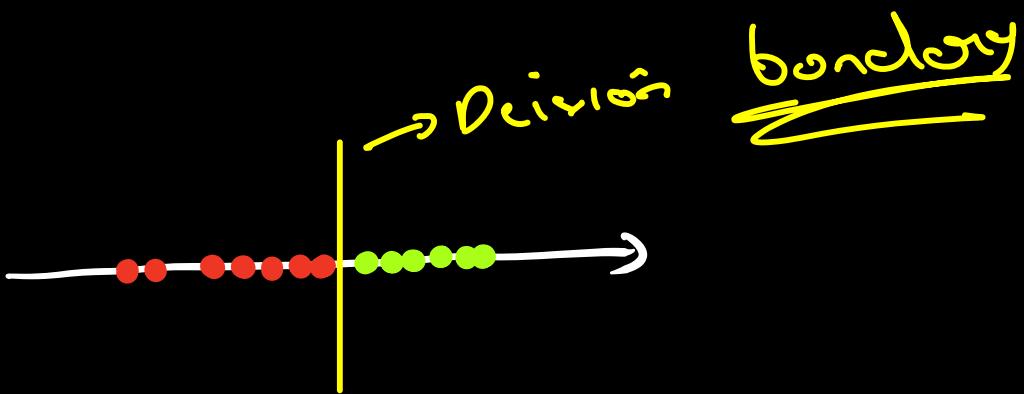
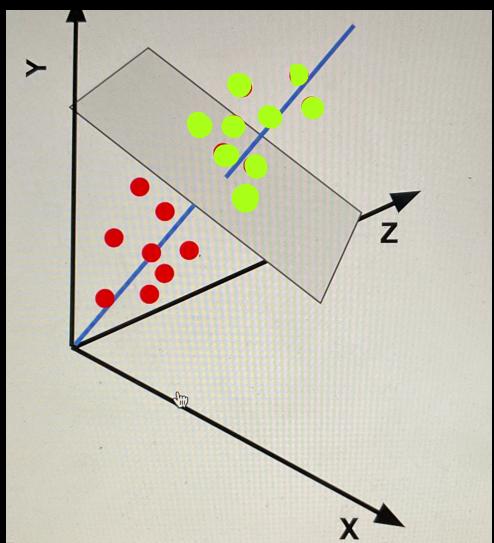
- a) $f_1 > f_2 > f_3$
- b) $f_3 > f_2 > f_1$
- c) $f_1 > f_3 > f_2$

Goal:

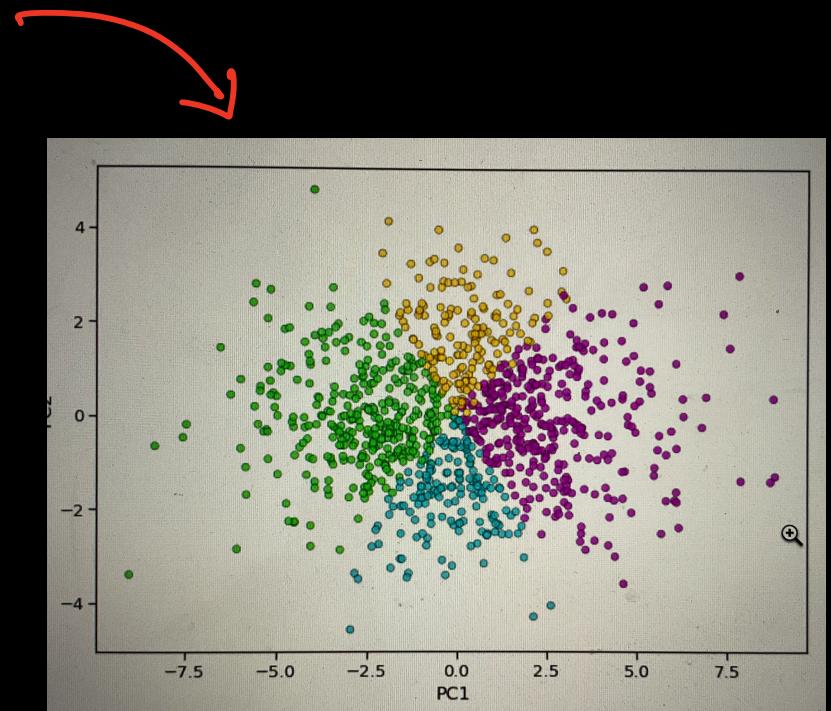
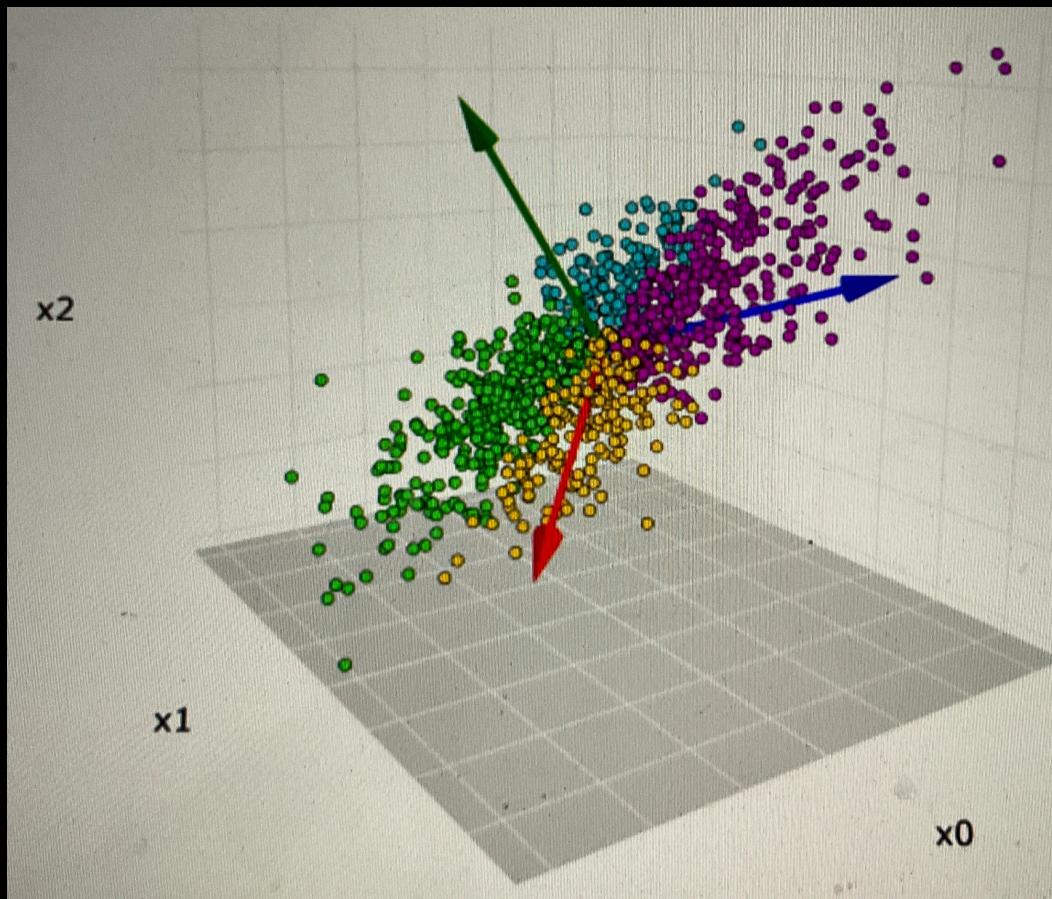
- We see that all the features don't carry equal amounts of info.
- Sometimes in ML, we need to reduce the number of features.

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

3D → 1D

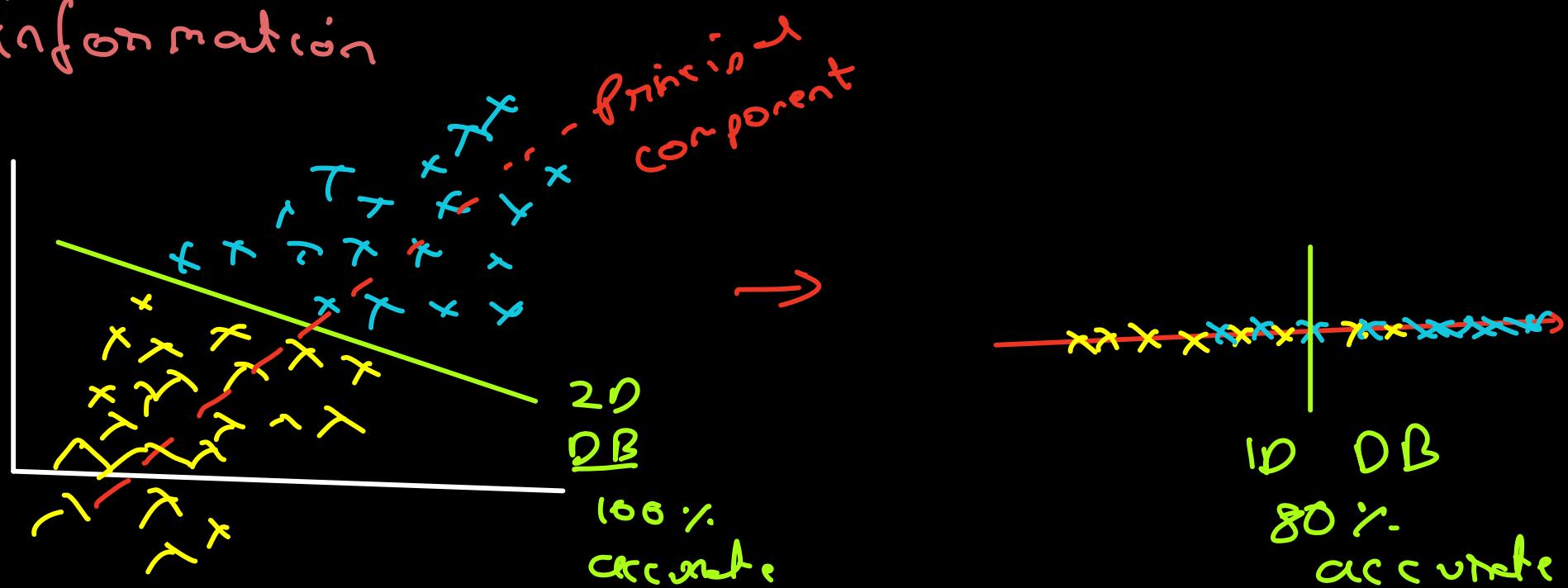


3D - 2D



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

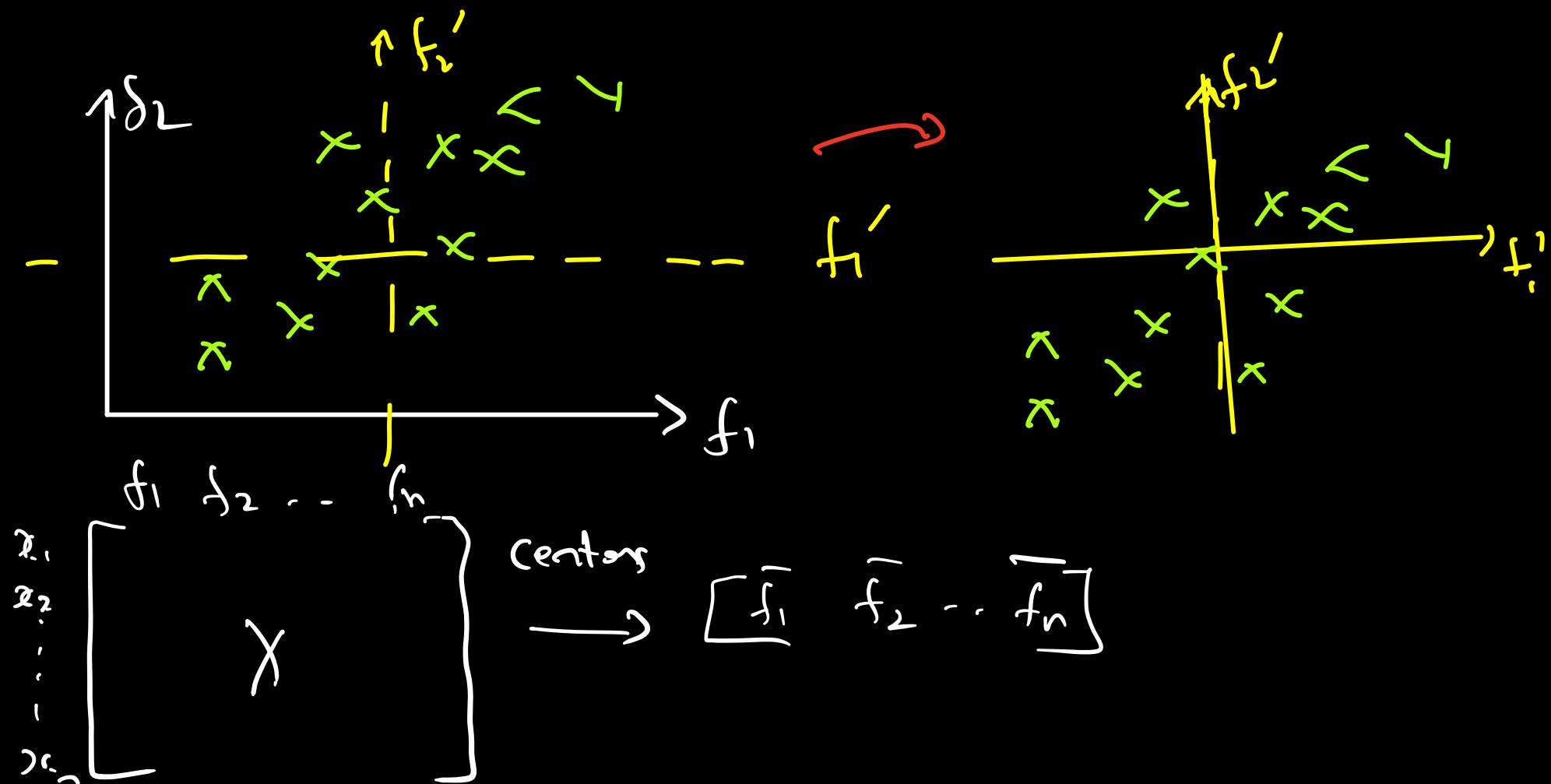


PCA Mathematics

→ Change of axis

→ Retain most information in few axis.

Step-1 : Centering



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

Step 2: Standardize / Scaling

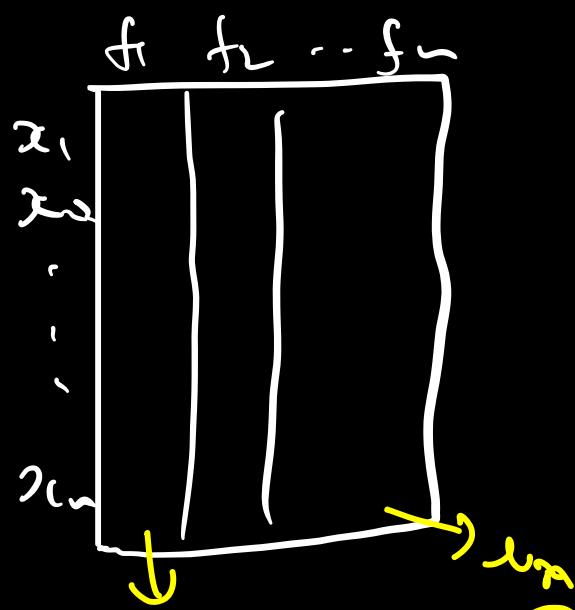
We need to make different features such as height (4-6) ft and weight (40 - 180) kg comparable.

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

Combining step 1 and 2

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

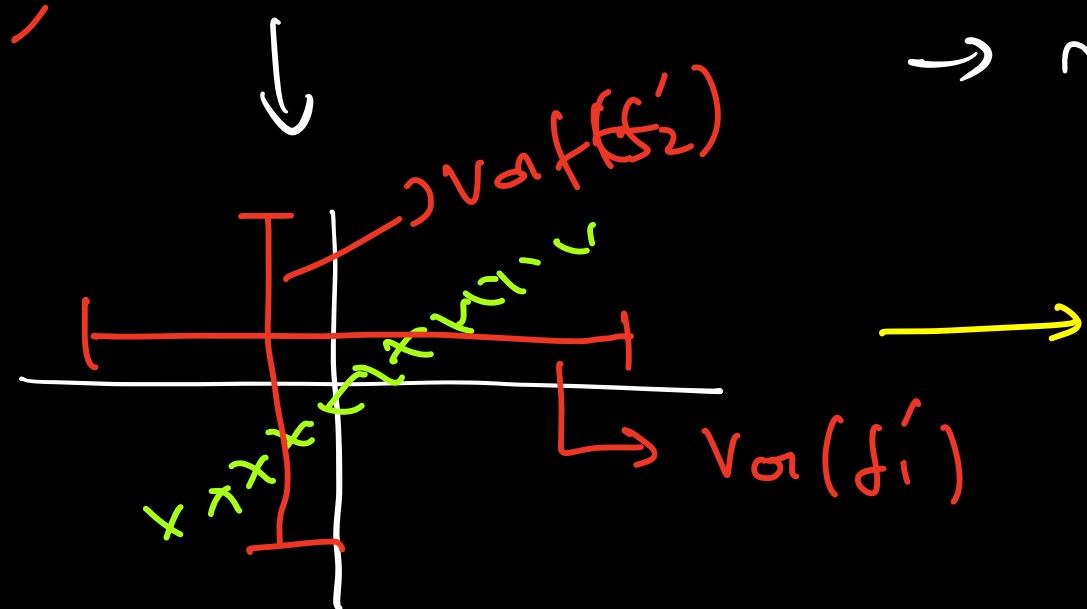
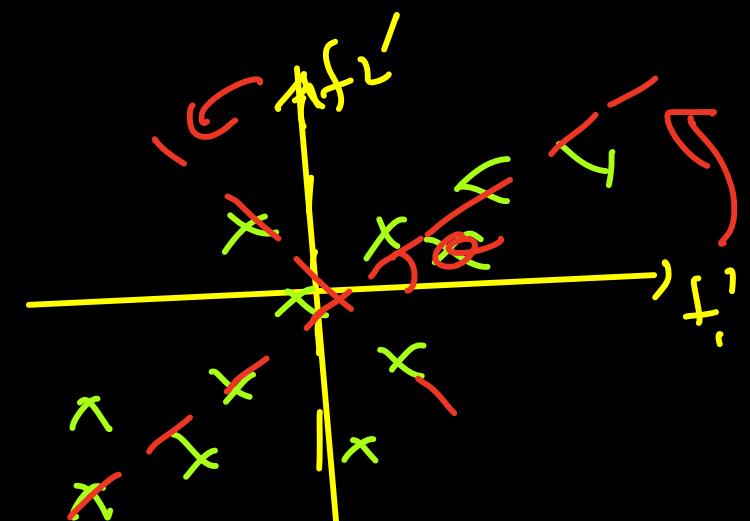
→ use sklearn



u_1, v_1

u_2, v_2

Step - 3: Rotate

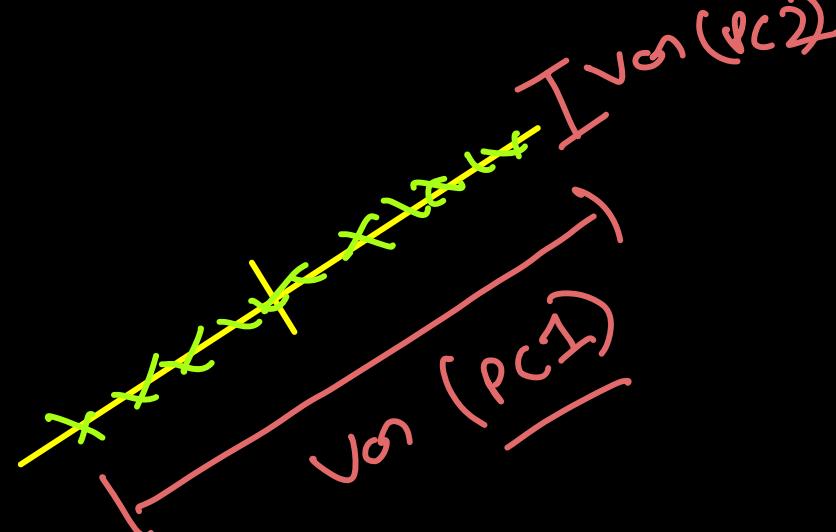


How to rotate?

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

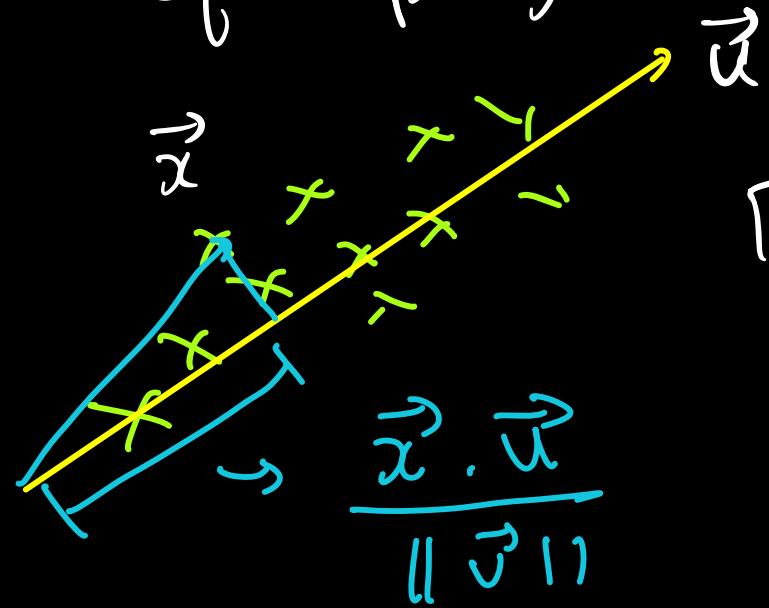
→ maximise variance

→ 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

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

However, projections can also be -ve

$$\max_{\vec{v}} \sum_{i=1}^n \left| \frac{\vec{v} \cdot \vec{x}_i}{\|\vec{v}\|} \right| \cdot \frac{1}{n} \rightarrow \text{But this is not differentially}$$

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

↓

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

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

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

\rightarrow We will use some identities w/o proof

\rightarrow We will learn a new concept

(dimension)

Variance \downarrow

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

~~\vec{v}~~

added square
↓ (homogeneous)

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

$\vdash \Rightarrow [\cdot] = \square \hookrightarrow \text{scalar}$

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

if A is invertible square matrix

$$\rightarrow \max_{\vec{u} \cdot \lambda} + \sum_{i=1}^n (x_i \cdot u)^T (x_i \cdot u) + \lambda (u^T u - 1)$$

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

$$\rightarrow \max_{\vec{u} \cdot \lambda} + \sum_{i=1}^n u^T x^T x \cdot u + \lambda (u^T u - 1)$$

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

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

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

matrix:

$$\frac{\vec{x}^T \vec{x}}{n} \perp \begin{bmatrix} \vec{x}^T \\ \vdots \\ f_1 \\ \vdots \\ f_d \end{bmatrix} \quad \begin{bmatrix} \vec{x} \\ f_1 \\ \vdots \\ f_d \end{bmatrix} \quad \begin{matrix} d \times n \\ \dots \\ n \times d \end{matrix}$$

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 ∇ , which is just $\text{ppcov}(x)$

→ \vec{v} must be such that when I multiply ∇ 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

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

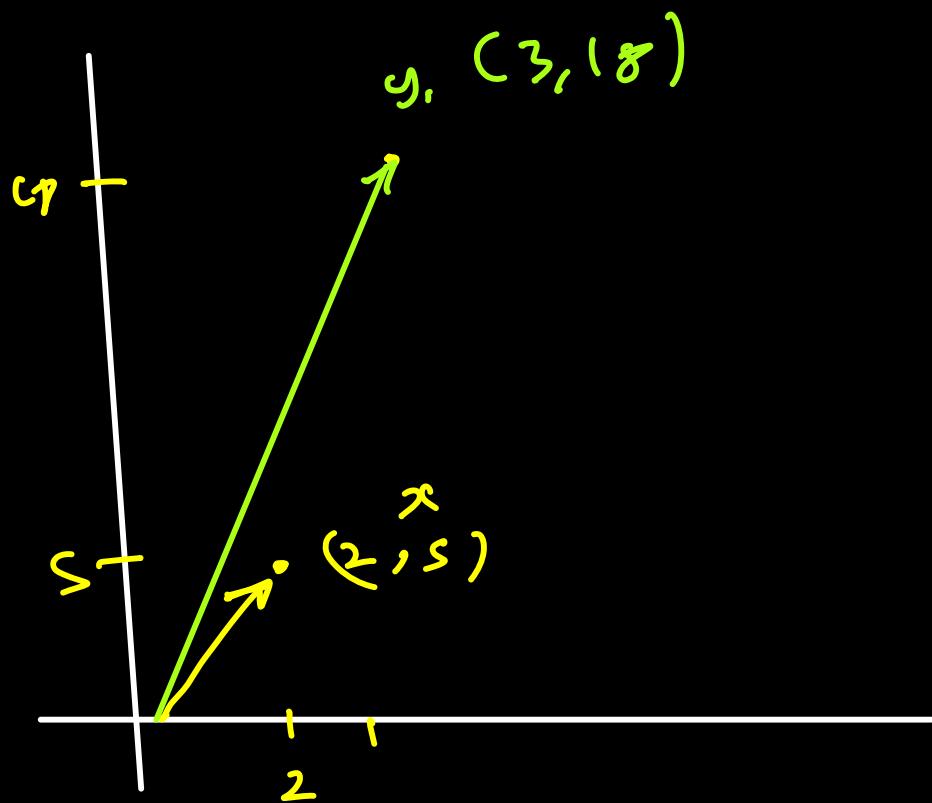
when we multiply $A \vec{x} = \vec{y}$, then

\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.

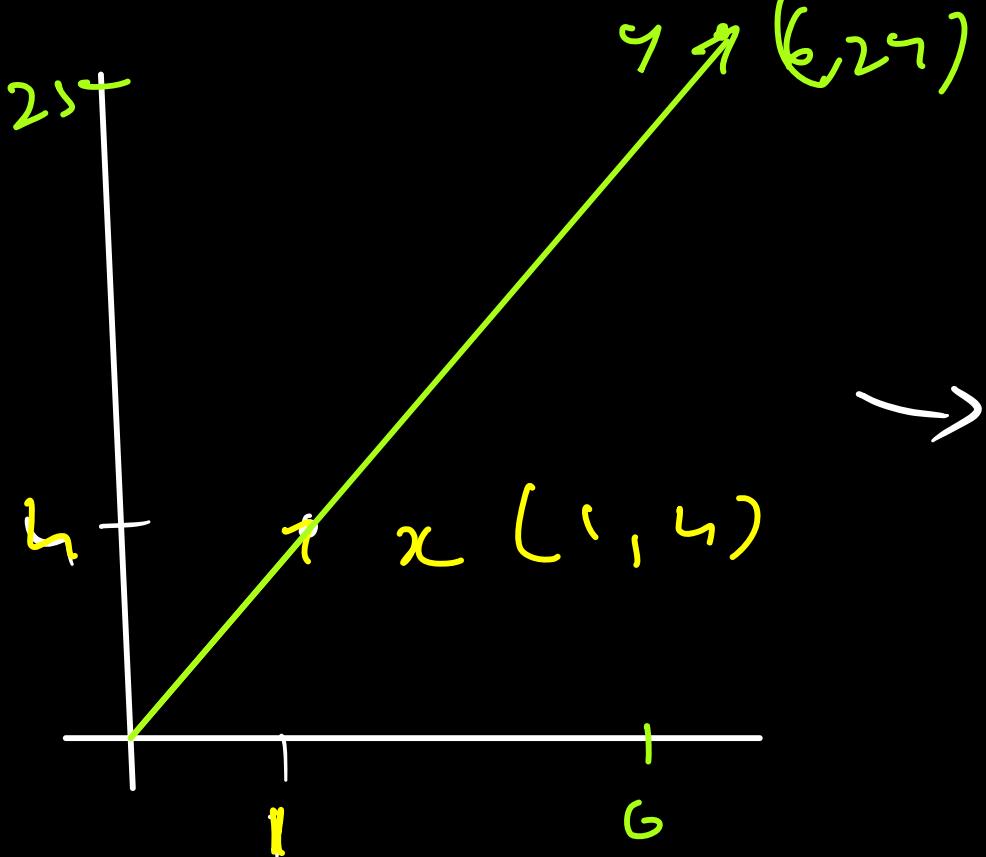
Eg: $\begin{bmatrix} -6 & 3 \\ 4 & 5 \end{bmatrix} \begin{bmatrix} 2 \\ 5 \end{bmatrix} = \begin{bmatrix} -12 + 15 \\ 8 + 10 \end{bmatrix} = \begin{bmatrix} 3 \\ 18 \end{bmatrix}$



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 $\underline{\underline{x}}$

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

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