

PCA

Have patience. All
things are difficult
before they become
easy



Wisdomward.in



Agenda:

→ Variants of G.D

SGD

Mini Batch G.D

Batch G.D

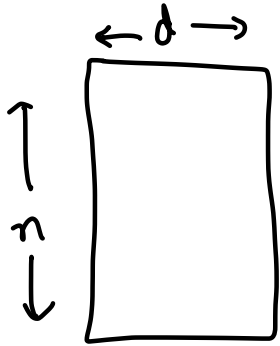
→ Dimⁿ. Reduction

→ PCA

Gradient Descent

$$w = w - \eta \cdot \frac{\partial L}{\partial w}$$

$$\frac{\partial L}{\partial w} = -\frac{1}{n} \sum_{i=1}^n y^i \cdot x^i + \frac{\lambda \bar{w}}{\|w\|}$$



$n = 1M$

→ slow updates in ' w '

Why?

$n = 1M$

→ Batch
Gradient
Descent

→ Calculate $\nabla_w L$ on subset of data

Batchsize

$m = 500$

mini-Batch
G.D

if $m = 1$

→ Stochastic G.D



- Batch gradient descent
- Mini-batch gradient Descent
- Stochastic gradient descent



epoch →



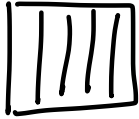
"10 epoch"



Dimⁿ Reduction

↓
features

#dim = #features



100 dimⁿ
(features)

Extract

10 dimⁿ
(create new features)

- max. info explained.
- feature extraction
- loss info. ↓

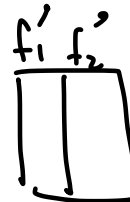
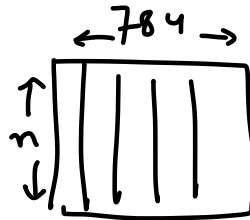
Why?

① Curse of Dimⁿ (d ↑ ↑ ↑ ↑)

→ training ML Slow

② Viz

2D/3D →



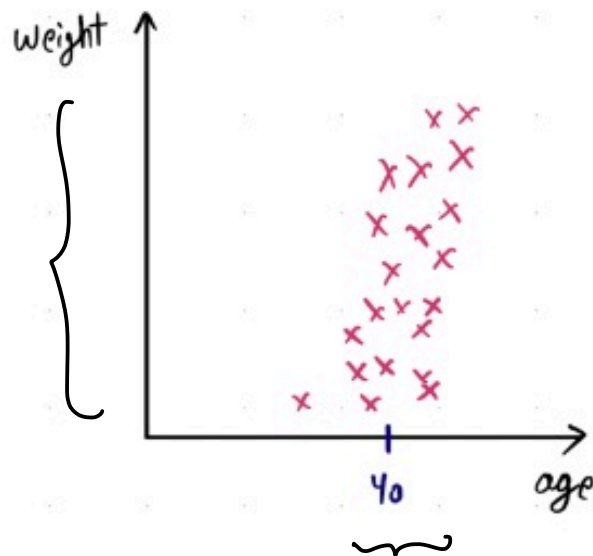
. Scatter()

Diabetes prediction

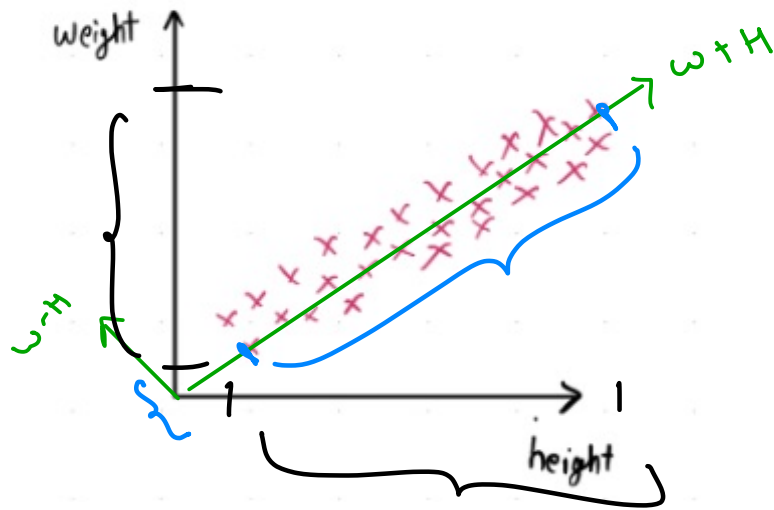
Weight	Age	Diabetic
...	...	Y
...	...	N
...	...	Y
...	...	Y
...

$f_1 = \text{weight}$

$f_2 = \text{Age}$



more spreadness = more info.



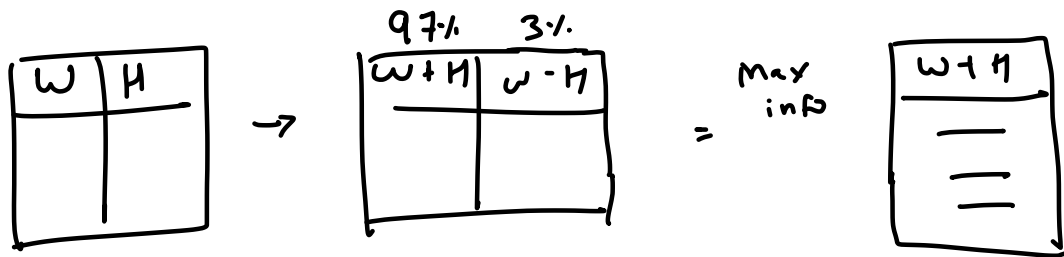
Which feature ?

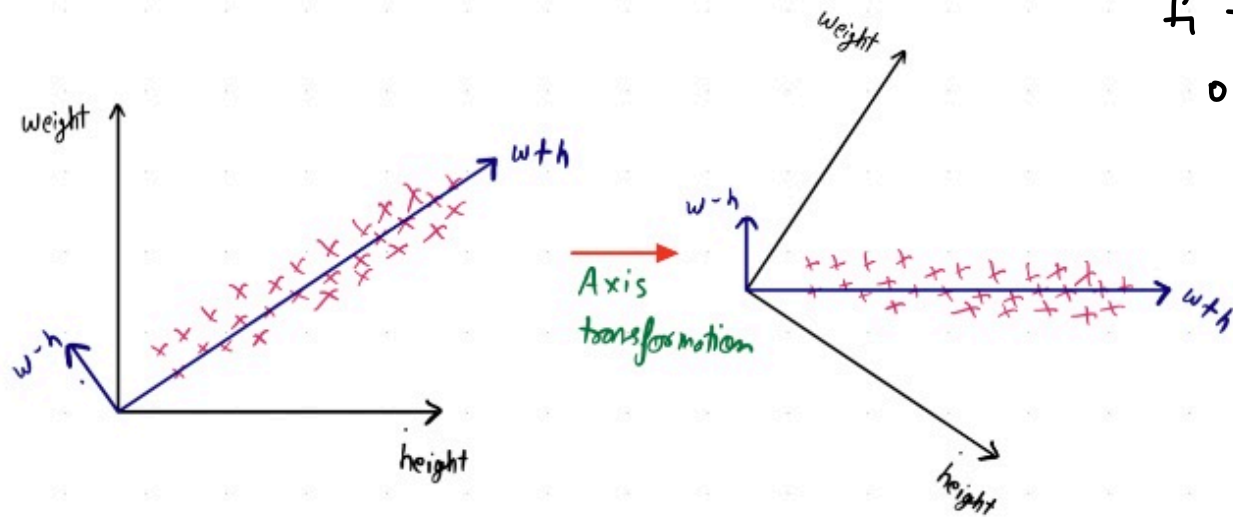
a) Weight

b) Height

☒ c) weight + weight

d) weight - height





$$f_1 \pm f_2$$

$$f_1 + 0.8f_2$$

$$0.5f_1 - 4f_2$$

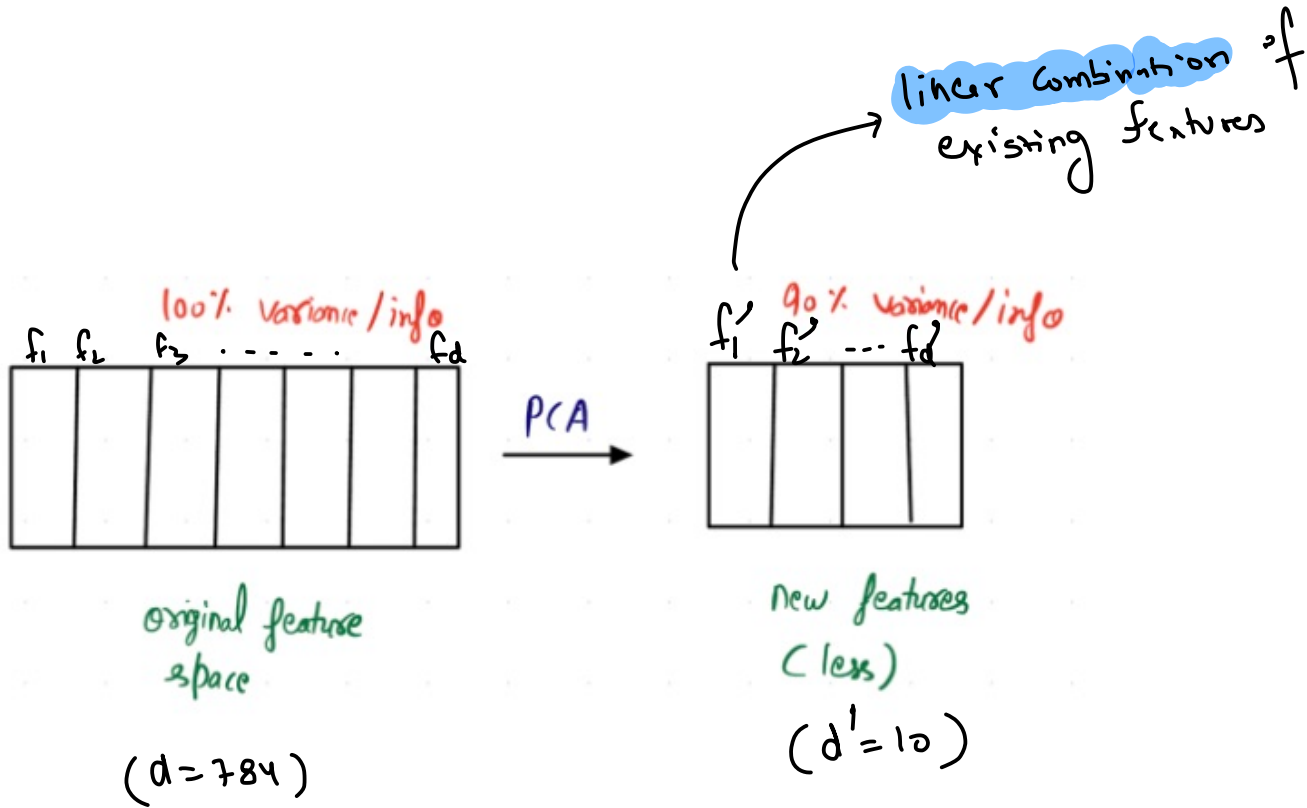
w	H



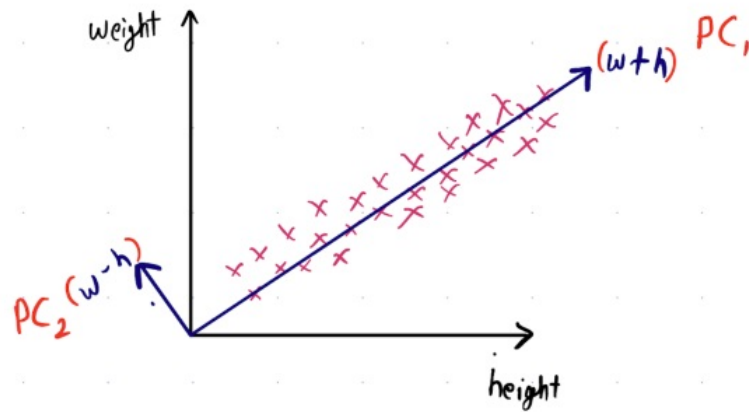
w+h	w-H

$$\xrightarrow[\text{Red}]{\text{Dim}^n}$$

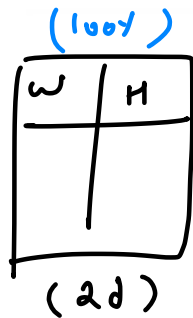
w+H



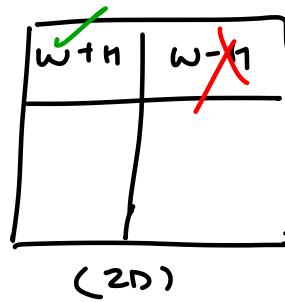
Principal Component Analysis [PCA]



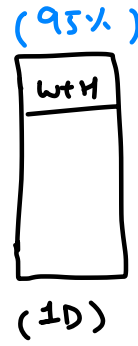
Principal
Components
→



PCA →



→



Goal: "95% info/variance should be preserved"

784 \rightarrow 1 [20%]

784 \rightarrow 32 [92%] ~~X~~

784 \rightarrow 45 [95.8%] ✓

x_1	x_2

Red

98% 2%

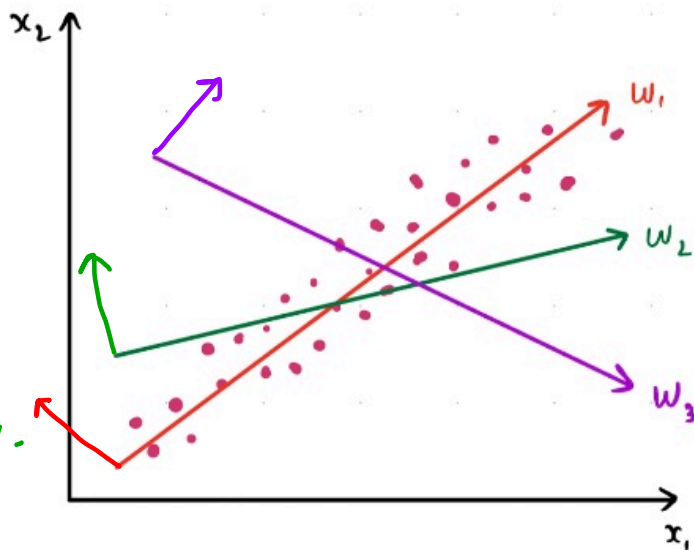
f_1'	f_2'

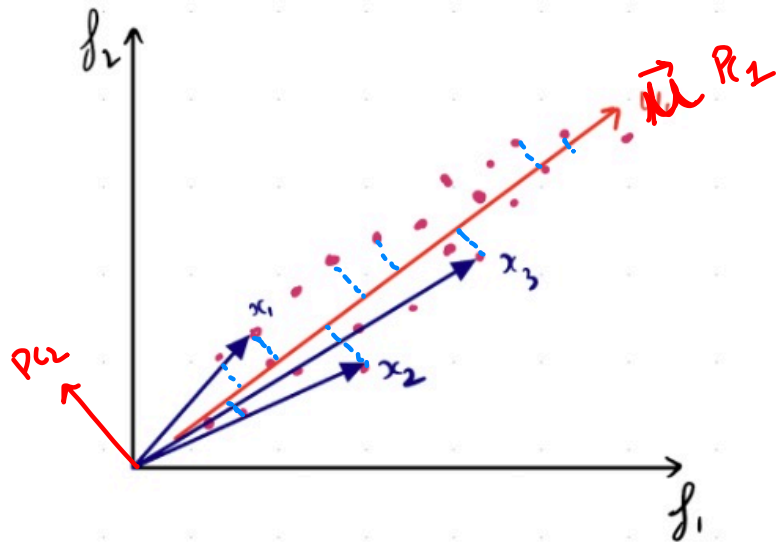
(Perfected)

Green

80% 20%

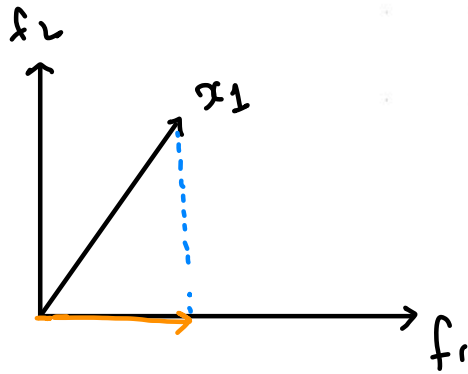
f_1'	f_2'





Projection of
Points to PC_1

find Variance
of Projected
Points



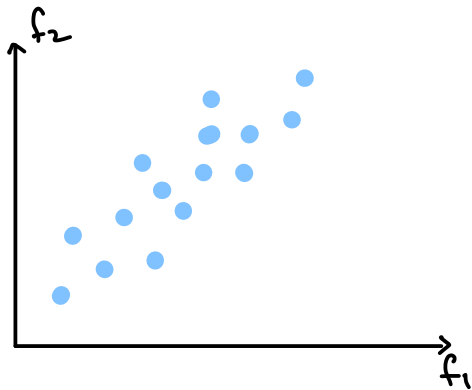
$$\text{Proj}_{f_1} x_1 = \frac{\vec{x}_1 \cdot \vec{f}_1}{\|\vec{f}_1\|}$$

In Principal Component Analysis (PCA), when rotating the axis to create new axis, which criterion is typically used to select the axis?

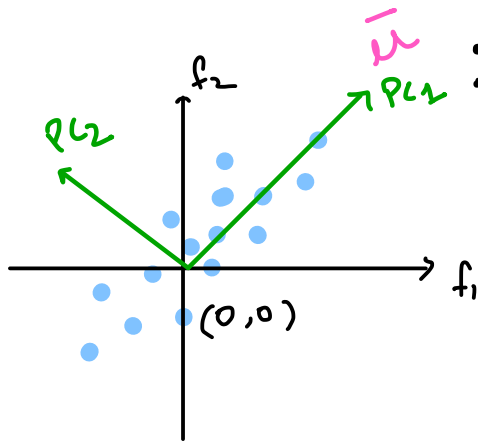
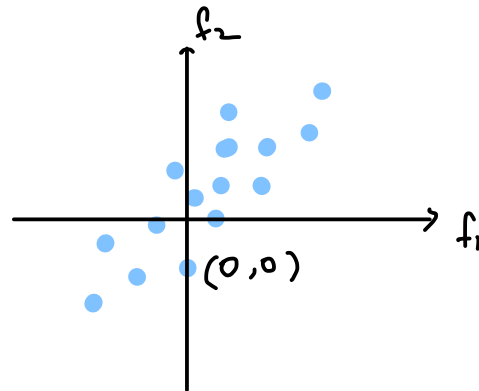
20 users have participated

A	Minimizing the variance along the new axis	25%
B	Selecting the axis with the least amount of data	0%
✓ C	Maximizing the variance along the new axis	75%
D	Choosing the axis randomly	0%

Maths



feature
Scaling
→
mean centering
 $X - \mu$



$\bar{\mu} : w^T x + w_0 = 0$

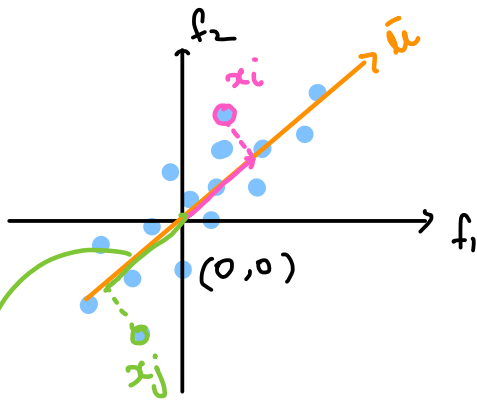
$$w_0 = 0$$

$$w^T x = 0$$

$$\mu^T x = 0$$

find $\bar{\mu}^*$!

$$\bar{\mu}^* = \begin{bmatrix} 2.5 \\ 3.1 \end{bmatrix}$$



$$\text{Proj}_{\vec{u}} x_j = \frac{x_j \cdot \vec{u}}{\|\vec{u}\|}$$

$$\text{Proj}_{\vec{\mu}} x_i = \frac{x_i \cdot \vec{\mu}}{\|\vec{\mu}\|}$$

Calculate Projection of all points on $\vec{\mu}$

$$\arg \max_{\vec{\mu}} \frac{1}{n} \sum_{i=1}^n \frac{x_i \cdot \vec{\mu}}{\|\vec{\mu}\|} \rightarrow \text{avg. Projection length on } \vec{\mu}$$

$$\underset{\vec{u}}{\operatorname{argmax}} \left[\frac{1}{n} \sum_{i=1}^n \frac{(x_i \cdot \vec{u})^2}{\|\vec{u}\|^2} \right]$$

→ differentiable
→ -ve & +ve
proj. length

$$\text{Constraint : } \|\vec{u}\| = 1$$

Constraint opt. Problem.



Convert into unconstrained opt.

$$\underset{\vec{u}, \lambda}{\operatorname{argmax}} \quad \frac{1}{n} \sum_{i=1}^n (x_i \cdot \vec{u})^2 + \lambda (\|\vec{u}\|^2 - 1)$$



Solve ? "Q.D"

$$A^2 = A A^T = \|A\|^2$$

$$\frac{(X \cdot u)^2}{n} + \lambda (\|u\|^2 - 1)$$

$$\frac{(X \cdot u)^T \cdot (X \cdot u)}{n} + \lambda (\|u\|^2 - 1)$$

$$(AB)^T = B^T A^T$$

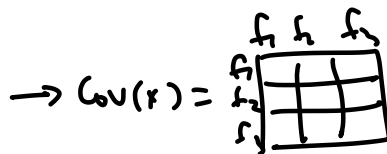
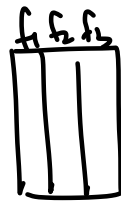
$\arg \max_{u, \lambda}$

$$u^T \frac{X^T \cdot X \cdot u}{n} + \lambda (u^T u - 1)$$

$\rightarrow \text{cov. Matrix}(X)$

$= V$

$n \times p \cdot \text{cov}(X)$



$$\operatorname{argmax}_{\mu, \lambda} \left[\mu^T V \mu + \lambda \mu^T \mu - \lambda \right]$$

Scalar derivative		Vector derivative	
$f(x)$	$\rightarrow \frac{df}{dx}$	$f(\mathbf{x})$	$\rightarrow \frac{df}{d\mathbf{x}}$
bx	$\rightarrow b$	$\mathbf{x}^T \mathbf{B}$	$\rightarrow \mathbf{B}$
bx	$\rightarrow b$	$\mathbf{x}^T \mathbf{b}$	$\rightarrow \mathbf{b}$
x^2	$\rightarrow 2x$	$\mathbf{x}^T \mathbf{x}$	$\rightarrow 2\mathbf{x}$
bx^2	$\rightarrow 2bx$	$\mathbf{x}^T \mathbf{B} \mathbf{x}$	$\rightarrow 2\mathbf{B} \mathbf{x}$

$$\frac{\partial L}{\partial \mu} = 0$$

$$\frac{\partial L}{\partial \lambda} = 0$$

$$2Vu + 2\mu\lambda = 0$$

$$2Vu = -2\mu\lambda$$

$$Vu = -\lambda \mu \Rightarrow$$

$$\mu^T \mu - 1 = 0$$

$$\mu^T \mu = 1$$

$$Vu = \lambda' \mu$$

μ - eigenvector
 λ' - eigenvalue

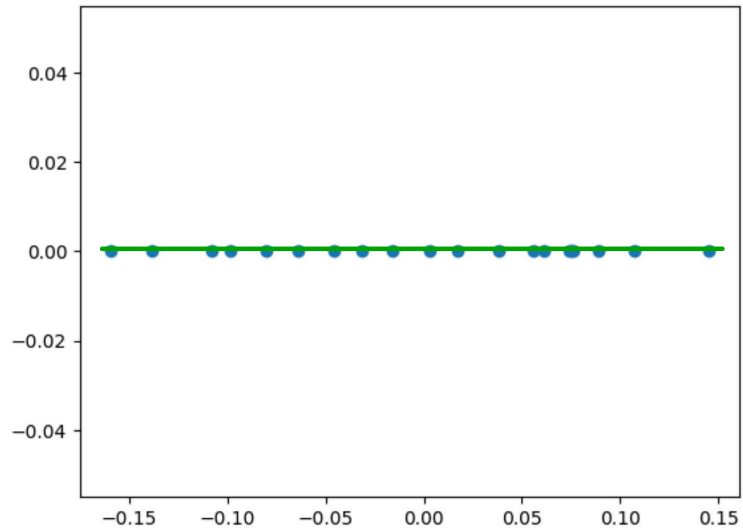
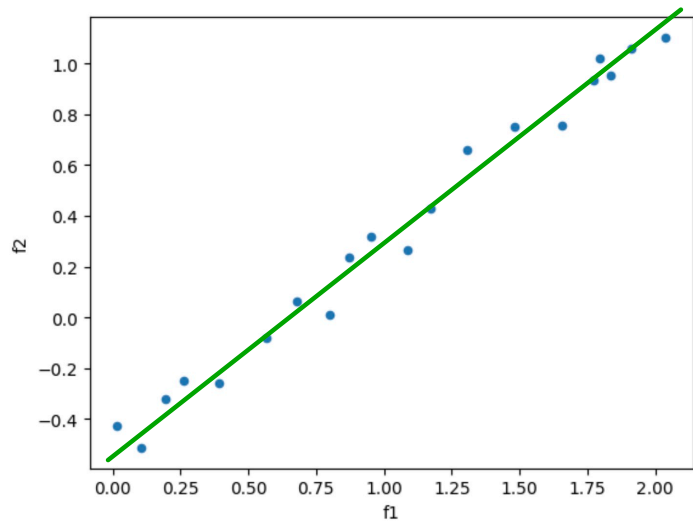
$$V = \text{np.cov}(X) = \frac{X^T X}{n}$$

↓
"Eigen decomposition"

<https://online.stat.psu.edu/statprogram/reviews/matrix-algebra/eigendecomposition>

$$\text{eigvec}, \text{eigval} = \text{np.linalg.eig}(V)$$

Higher eigenvalue \Rightarrow Better eigen vector / Principal component



$$P(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{\left(-\frac{(x-\mu)^2}{\sigma^2}\right)}$$

