

# Principal component analysis (PCA)

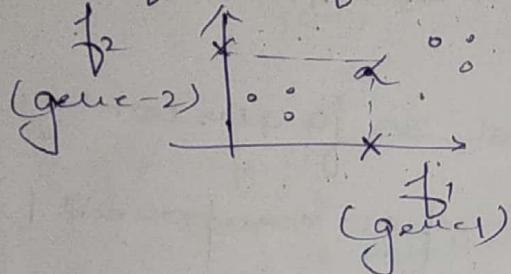
$m \rightarrow$  TE

$n \rightarrow$  features

(Unsupervised learning problem)

PCA finds the axes with maximum variance.

e.g. for 2 features



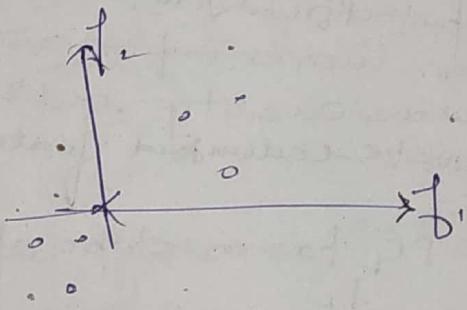
Step 1

Take the mean value with  $f_1$  and  $f_2$  denoted by  $\alpha$  and  $\beta$ .

Find  $\alpha$ .

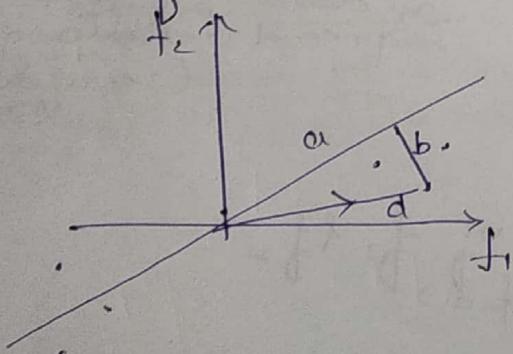
Step 2

Shift the origin to  $\alpha$ .



Step 3

Try to fit a line passing through origin



demands fixed  
for variable line

minimize  $\Sigma b^2$   
or maximize  $\Sigma a^2$

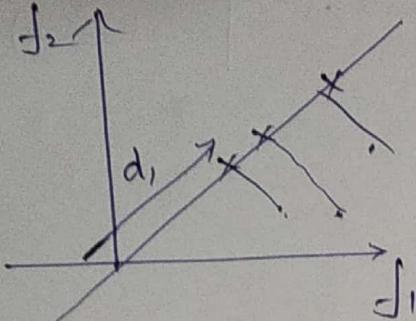
$$d^2 = a^2 + b^2$$

by pythagoras  $a, b$   
are oppositely  
related

and will give the axes with maximum variance

Step 4

project our those distances to that line.



$$SS = \underline{EV} \text{ (Eigen Value)}$$

$SS = \text{sum of square}$

$$= d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2 + d_6^2$$

∴ we just need the largest SS for eliminating noise.

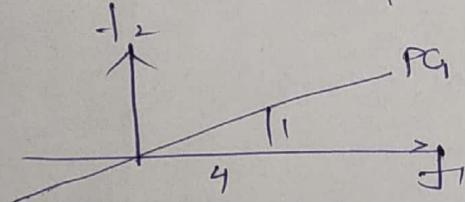
That line is called PC1 (principal component 1).

Advantages -

- ⇒ Avoid over fitting | cost
- ⇒ computational cost
- ⇒ gives more important features
- ⇒ dimensionality red.
- ⇒ eliminate redundant features

i.e. on which feature  
more contributing to PC1  
↳ can get b)  
the slope of PC1.

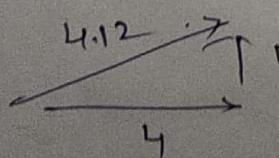
e.g. if  $PC_1$  has a slope of 0.25.



so, f1 is more contributing  
feature than f2  
or data is more  
spread out along f1  
(i.e. more variance)

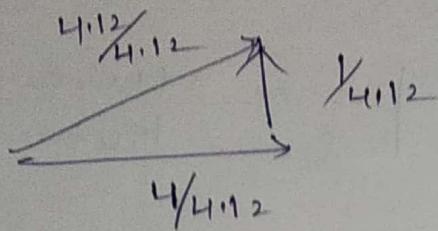
so, finally  
 $PC_1$  is a linear comb' of  $f_1$  &  $f_2$   
(cocktail recipe)

$PC_1$  is a linear comb' of variables



When we do PCA with SVD, we scale the length

$\frac{4.12}{\sqrt{4.12+0.1}}$



so, now it contributes 0.97

$$\text{and } f_r = \frac{0.242}{\sqrt{4.12^2 + 0.1^2}}$$

It follows is the same.

Eigen Vector / Singular Vector for PC<sub>1</sub>  
the unit vector that we just found along PC<sub>1</sub>

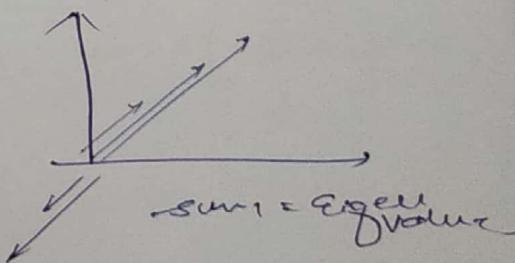
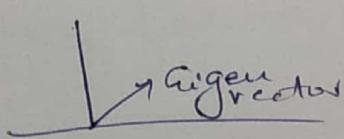
is known as Eigen / Singular vector for PC<sub>1</sub>.

### Loading Scores

The proportion of each features in linear combination  
are called loading scores.

Eigen value of PC<sub>1</sub> = SS

The sum of square of dist of proj of points on  
PC<sub>1</sub> (SS) = eigen val. of PC<sub>1</sub>



### Singular value of PC<sub>1</sub>

Singular Value of PC<sub>1</sub> =  $\sqrt{\text{Eigen value for PC}_1}$

$$= \sqrt{SS}$$

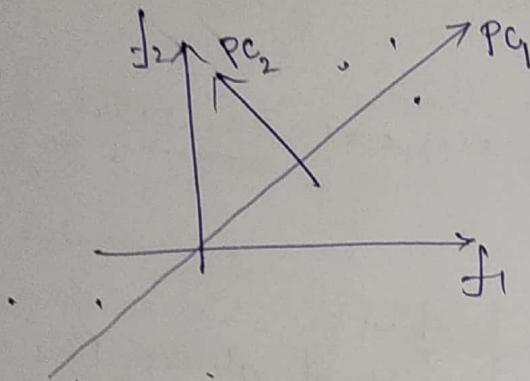
$\sqrt{SS}/N$

Note that  $SS/N$  is the variance along the PC<sub>1</sub>  
and we are maximizing it.  
For example using EV and variance

→ PC<sub>2</sub> will be directly the line to the PC<sub>1</sub> to the PG passing through the origin without any further optimiz.

Note :

Simply remember PG and PC<sub>2</sub> passes through the mean of all the data points & are linear comp. of f<sub>1</sub> and f<sub>2</sub>.



We know how to get the vector → loading scores

→ already snap controls of f<sub>1</sub> and f<sub>2</sub> so that f<sub>1</sub>, f<sub>2</sub> make any one of them -ve.

i.e. for PC<sub>1</sub>

$$\begin{cases} f_1 \rightarrow 0.97 \\ f_2 \rightarrow 0.242 \end{cases}$$

for PC<sub>2</sub>

$$\begin{cases} f_1 \rightarrow -0.242 \\ f_2 \rightarrow 0.97 \end{cases} \quad \left. \begin{array}{l} \text{loading scores of} \\ \text{PC}_2 \end{array} \right\}$$

To draw the final PCA plot

- mean on the origin

- later do make PC<sub>1</sub> as x-axis.

That's how PCA is done using SVD

(singular value  
Decomposition)

Note:-

for Eigenvalue distances are along the axes.

Variance

↑ since it is along the  
Variance      Eigenvalue  
                    ↓  
                    (Sample size - 1)

$$\text{Eigen var} = \frac{SS}{m-1}$$

Total variance = sum of var along both the PCs  
=  $\lambda_1 + \lambda_2$ .

Scree plot.

The graphical rep' of the % of variat that each PC accounts for.

\* Similar is the PCA for 3 variable  $f_1, f_2, f_3$   
we get PCA<sub>1</sub>, PCA<sub>2</sub>, PCA<sub>3</sub>, then-

No. of PCs = min (no. of features / no. of samples)

# Maths

## Variance of a matrix

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{N}$$

$$= \frac{1}{N} \sum x_i^2 - \bar{x}^2$$

$N$  = sample example

{ popu " variance }

actually

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

sample variance

degree of freedom

$$\text{or } \text{Var}(X) = E[(X - E(X))^2]$$

$$= E[X^2 - 2X E[X] + E[X]^2]$$

$$= E[X^2] - 2E[X]E[X] + E[X]^2$$

$$= E[X^2] - E[X]^2$$

## Covariance of a matrix

(egressally) the drn i.e how two variables are related to each other

highly  
or  
very  
not related if value very close to 0

} no. sig. of its value.

$$\text{cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N-1}$$

$$= \frac{\sum x_i y_i}{N-1} - \bar{x}\bar{y}$$

Note:-

If we consider a matrix as a transform then, a simple term Eigen value is the strength of that transform for a particular direction known as Eigen vector.

(see method to find Eigen value & Eigen vector)

Correlation [-1 to 1]

Coef of covariance gives the dir<sup>n</sup> or better mag. of how the variables are related

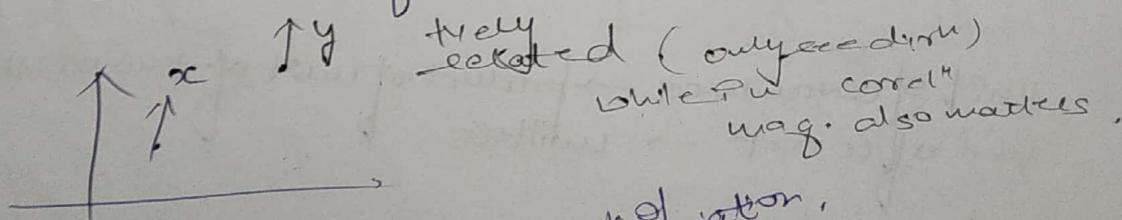
$$r, \text{ coed} = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y}$$

In layman's term,

cov matrix → physically / data structure view  
How similar the variances of features are.

Mathematically

Linear transform showing



gives dir<sup>n</sup> of variation,  
mag. also matches.

PCA physically.

projecting data onto Eigenvectors of covariance matrix

Mathematically PCS → Eigenvectors of (normalized empirical) covariance matrix.

Eigenvector + Eigenvalue of covariance matrix.

## Creation of covariance matrix

PCA is simply eigenvector of empirical covariance matrix

Eiger

Given a matrix  $X$

$\rightarrow$  matrix  $X$  features  
 $X =$  samples

steps

$\vec{d}_S \rightarrow g_{ij} - \vec{g}_{ij}$  column wise

$\hookrightarrow$  so gives  
most detailed  
feature and  
sample a

$$x = \boxed{\phantom{0}}$$

$$\cancel{S \cdot e^{b-2}} \quad X^T X$$

Mont

$$\Sigma = \text{cov}(X) = \frac{(X^T X)}{n-1} \curvearrowleft \text{for sample}$$

→ culture excavation  
→ arch. + muse.

*U*

n for popu

A diagram of a rectangle with vertices labeled as follows: the top-left vertex is  $f$ , the top-right vertex is  $k$ , the bottom-left vertex is  $l_1$ , the bottom-right vertex is  $l_2$ , the left edge midpoint is  $l_3$ , and the right edge midpoint is  $k'$ .

$$\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\sum_{i=1}^N (x_i - \bar{x}) y_i$$

$n \approx n_{\text{cr}}$  for very  
large value of  $\mu$ .

Unit of covariance  $\rightarrow$  product of unit of two variables  
unit of corr<sup>th</sup>  $\rightarrow$  unitless

A "corr" of covariance matrix are symmetric & each element in the matrix represent corr of covariance b/w these vari and varj diagonal = 1 in cov matrix

for 2d matrix I have 2 features which are different + which are different

$\begin{matrix} m \\ \text{features} \end{matrix}$

is the datapoints for each feature

$\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$

## Eigen Value

$$A = \begin{bmatrix} \quad & \end{bmatrix}$$

So,

$$A\vec{x} = \lambda \vec{x}$$

Eigen vector

$$\Rightarrow (A - \lambda I) \vec{x} = 0$$

$$\Rightarrow |A - \lambda I| = 0 \rightarrow \text{characteristic Eq}$$

$$\begin{vmatrix} -\lambda & & & \\ & -\lambda & & \\ & & -\lambda & \\ & & & -\lambda \end{vmatrix} = 0 \quad \text{and get the } \lambda.$$

$$\lambda = \dots, \text{equations of dimension}$$

Eigen space is the set of eigenvectors that correspond to some eigen value.

another set of vectors that satisfy this eqn-

$$(A - \lambda I) \vec{x} = 0$$

The eigen space is

$$\text{for } \lambda = \lambda_1, \lambda_2, \lambda_3$$

corresponding to each

$$(A - \lambda I) \vec{x} = 0$$

$$\begin{pmatrix} - & & & \\ & - & & \\ & & - & \\ & & & - \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = 0$$

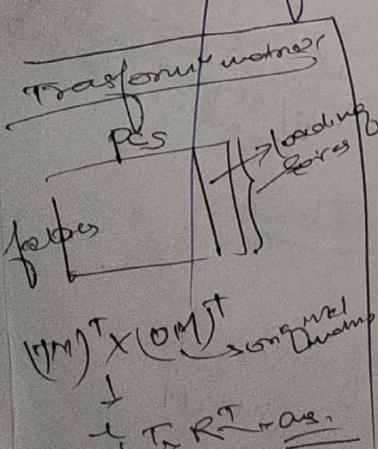
get the first eigenvector

$$x_1 - x_2 = 0$$

$$x_1 = x_2$$

$$\begin{bmatrix} 1 \\ 1 \end{bmatrix} v$$

Solving this we always get a drn





# Covariance vs Correl<sup>n</sup>

## creation

The covariance matrix

$$\therefore \text{cov}(x_i, y_j) = \frac{1}{N} \sum_{n=1}^N (x_i - \bar{x})(y_j - \bar{y}_j)$$

matrix for PCA

standardized vector  
of covar matrix  
correl matrix

$$\therefore \text{"Correl"}(x_i, y_j) = \frac{\text{cov}(x_i, y_j)}{\sigma_x \sigma_y}$$

$$= \frac{1}{N} \sum_{n=1}^N (x_i - \bar{x})(y_j - \bar{y}_j)$$

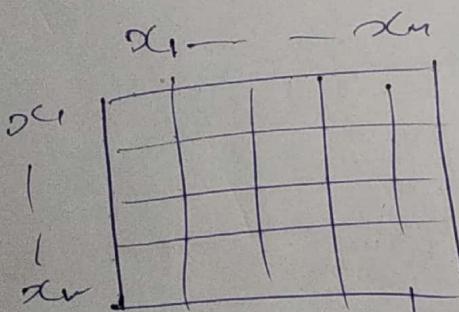
$$= \frac{1}{N} \sum_{n=1}^N (x_i - \bar{x})(y_j - \bar{y}_j) \quad H_{x_i y_j} \quad \text{corr} = 1$$

\* Covariance matrix measured when the features have same scale.

Correl when scale diff. like pressure & temperature depends on application

\* generally Correl<sup>n</sup> matrix contributes by PC1 & PC2 say 62% and 24%. Likewise covariance matrix 96%, 2% etc

Covariance matrix  
(Symmetric)



Correl matrix  
(Symmetric)

x_1	...	x_m
x_1	1	1
x_2	1	1
x_n	1	1

Standardized original data

EV will give a different data with max variances for height EV is EV roughly each cell in correlation matrix can't be computed many EVs or very large space needed, EV highest variance

## PCA

$$\Sigma = \frac{1}{m} \sum_{i=1}^m z^{(i)} z^{(i)T}$$

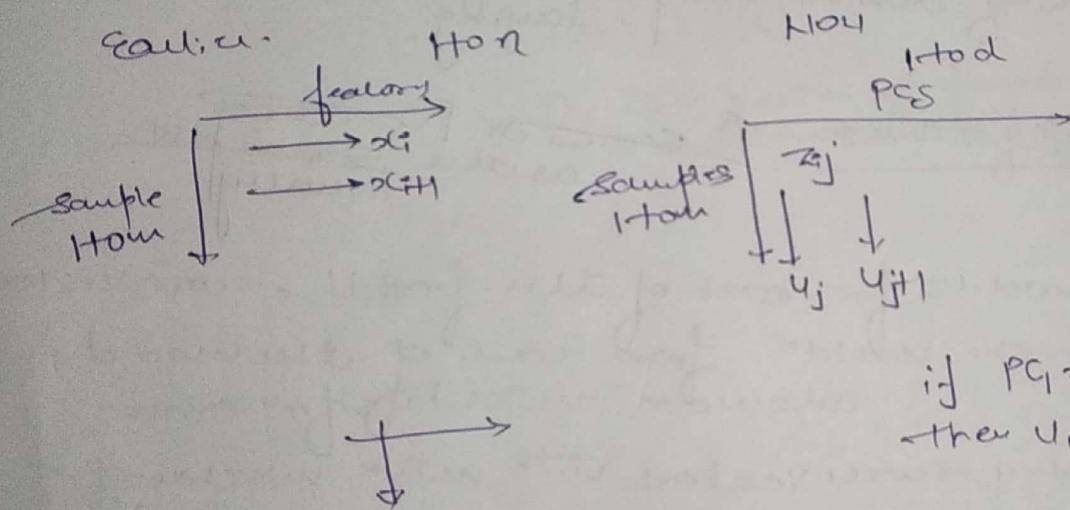
for a d-dimensional summary of the data that captures as much variance (as) possible.

We represent data on the basis of first d PCs.  
 $\{u_1, u_2, \dots, u_d\}$

and compute

$$z_{ij} = u_j^T z^{(i)} \quad \text{for } j=1 \text{ to } d \\ i=1 \text{ to } m$$

i.e. conversion to PC axis

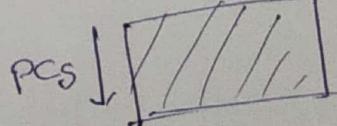
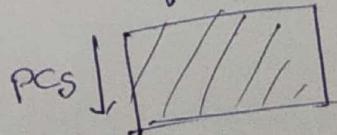
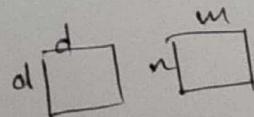
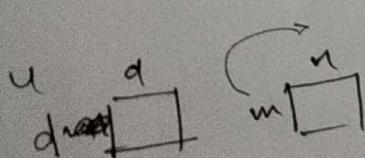


if  $PC_1 = f_1$  say  
 then  $u_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ .

so, 
$$z_{ij} = u_j^T z^{(i)}$$
 or  $z = u^T X$  simply

These entries are referred to as PC scores.

original features



Remember:

PCs are obtained only after ~~covariance~~ and scaling the original features.

Note:

Even if there are multiple PCs we will get some general classification by plotting first & second PCs of any defined features.

(it has most variance contribution & diff. dirn)

Since we have the one along which most of the variance change occurs.

Like

Walking  
Walking upstairs  
Walking downstairs

Standing  
Sitting  
Laying.

Classification Algorithms

steps →

of scores or PC scores  
as the score plot  
actually.

- ① Build models for each of the first  $H$  score vectors.
- ② Predict on validation from each of the model & calculate misclassification error.  
can be directly seen to estimate.
- ③ Select q score vectors with min misclassification error.
- ④ Use the q score vectors to build LDF on training data.
- ⑤ Predict on test data misclassification error.

## KNN

K Nearest Neighbours

$$P(y^{(i)}=j|x=x^{(i)}) = \frac{1}{K} \sum_{q=1}^K P(y^{(q)}=j)$$

we can't set to simultaneously select the no. of PCs or k PCs & the optimal value of k.  
We consider potential  $K \in \{1, 3, 5, 7, 9\}$

## B. Softmax Regression

$$P(y^{(i)}=j|x=x^{(i)}) = \frac{e^{\theta_j^T x^{(i)}}}{\sum_{l=1}^L e^{\theta_l^T x^{(i)}}}$$

$$j \in \{1, 2, \dots, L\}$$

$\theta_j$  is a vector of parameters associated with the single class j.