

Principal Component Analysis (PCA)

Steps involved in ml model building:

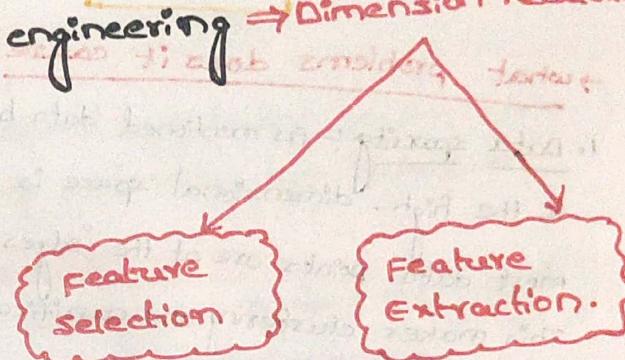
① Data ingestion

② EDA

③ preprocessing \rightarrow feature engineering \Rightarrow Dimension Reduction

④ Model Building

⑤ Model evaluation.



* Dimension - no.of columns or features

* Dimension reduction - process of Reducing the dimension or feature.

→ Before understanding FS and FE, we should know about curse of dimensionality.

How does the curse of dimensionality occur?

→ Let's say we have a dataset with $m \times n$ size.
↓
row column.

→ As we add more dimensions to our dataset, the volume of space

increases exponentially. This means that the data becomes sparse. Think of it this way: if you have line (1D), it's easy to fill it with few points. If you have a square (2D), you need more points to cover the area. Now imagine a cube (3D) - you'd need even more points to fill the space. This concept extends to higher dimensions, making the data extremely sparse!

Ex:- Data (Rowxcolumn)

$$f_1, f_2, \dots, f_{10} \rightarrow M_1$$

$$f_1, f_2, \dots, f_{100} \rightarrow M_2$$

$$f_1, f_2, \dots, f_{1000} \rightarrow M_3$$

independent features dependent feature

Accuracy

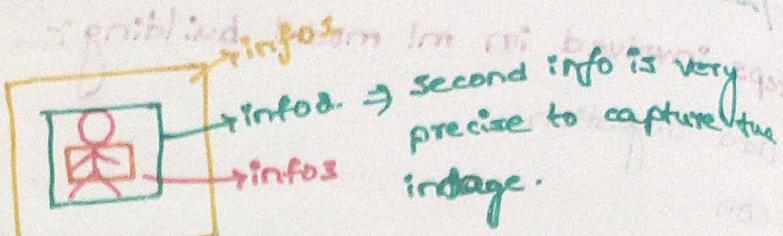
75% } \rightarrow less features

80% } \rightarrow optimal(having sufficient features)

72% } \rightarrow many
↓
overfitting \rightarrow sparsity

(lot's of 0's and less 1's.)

Ex:- which information is precise for our model to identify a image



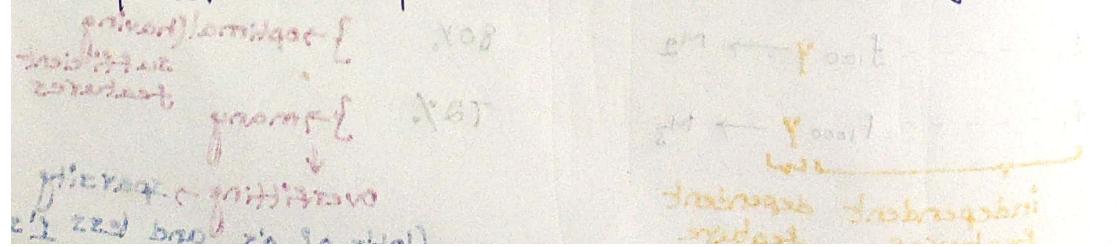
→ what problems does it cause?

1. Data sparsity: As mentioned data becomes sparse, meaning that most of the high-dimensional space is empty (in high dimensional space, most data points are at the "edges" or "corners" making data sparse). This makes clustering and classification tasks challenging.
2. Increased computation: More dimensions mean more computational resources and time to process the data.
3. overfitting: with Higher dimensions, models can become overly complex, fitting to the noise rather than the underlying pattern. This reduces the model's ability to generalize to new data.
4. Distance loose meaning: In high dimensions, the differences in distances between data points tend to become negligible, making measures like euclidean distance less meaningful.
5. performance degradation: Algorithms, especially those relying on distance measurements like k-nearest neighbours, can see a drop in performance.
6. visualization challenges: High dimensional data is hard to visualize, making exploratory data analysis more difficult.

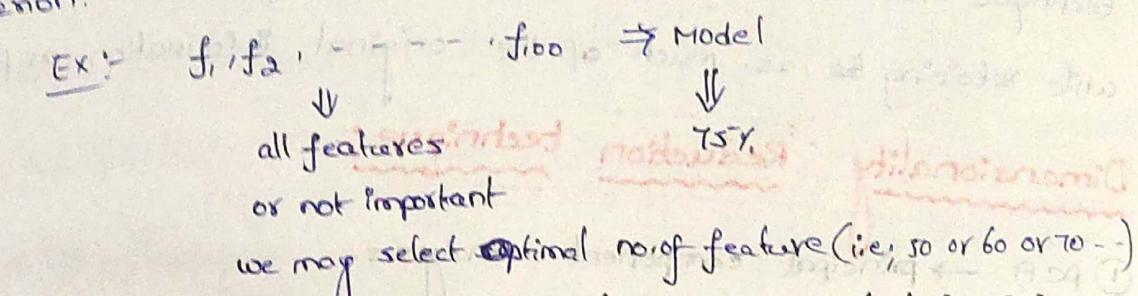
How to solve the curse of dimensionality:

Two different ways are there to resolve the curse of dimensionality

1. feature selection.
2. feature extraction / Dimension Reduction / Feature transformation.



feature selection: feature selection is a process that chooses a subset of features from original features so that the feature space is optimally reduced according to a certain criterion.



→ There are mainly two types of feature selection techniques:

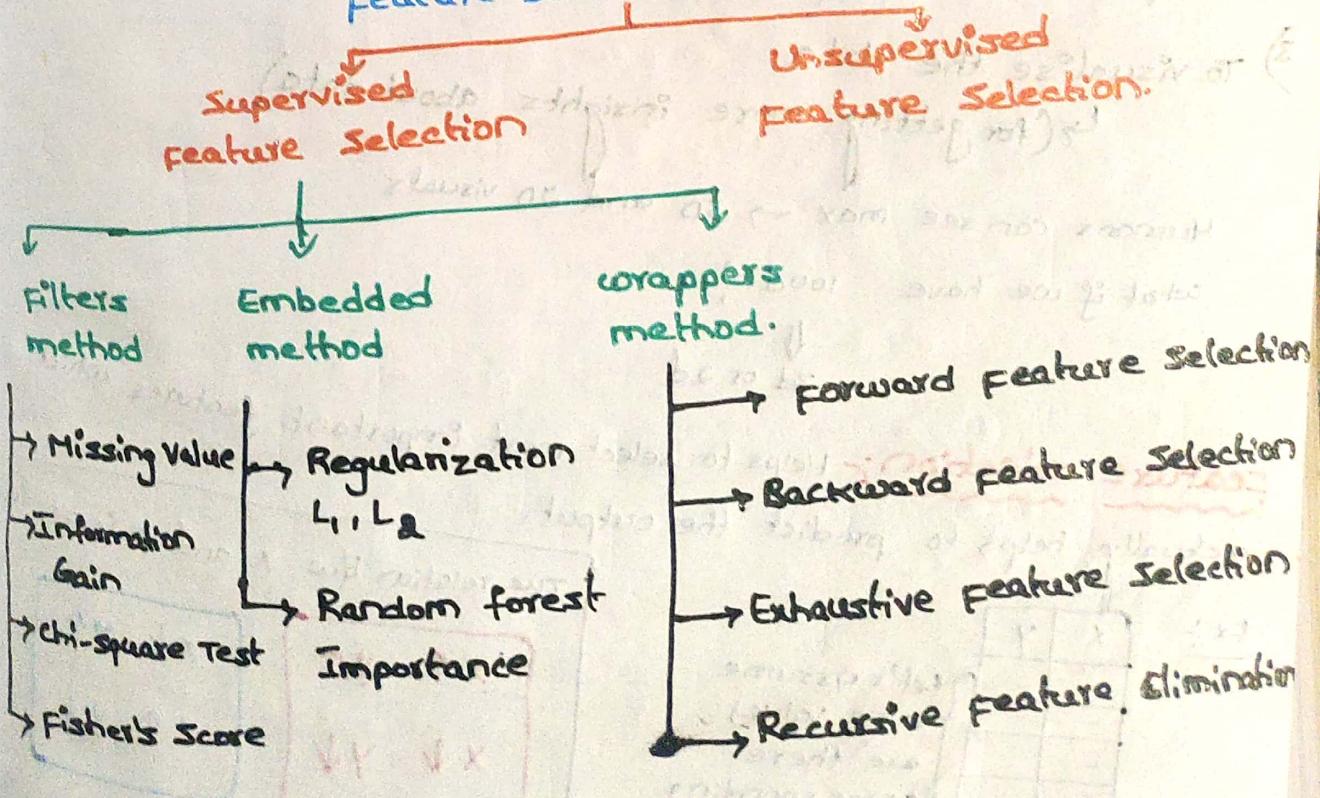
Supervised Feature Selection Techniques

→ Supervised feature selection techniques consider the target variable and can be used for the labelled dataset.

Unsupervised Feature Selection Techniques

→ Unsupervised feature selection techniques ignore the target variable and can be used for unlabelled dataset.

Feature Selection Techniques



2) Dimensionality reduction / Feature Extraction / Feature transformation

"Dimensionality Reduction" is a type of feature extraction technique that aims to reduce the number of input features while retaining as much of the original information as possible.

Dimensionality Reduction technique :-

- ① PCA → principal component Analysis.
- ② t-SNE → t distributed stochastic Neighbour embedding.
- ③ LDA → linear discriminant Analysis.

Why dimensionality Reduction?

- 1) prevent curse of dimensionality
- 2) To improve the performance of the model
(To train all the features it takes more time, but if we do the dimensionality reduction it will take less time and improves the performance of the model).
- 3) To visualise the data

↳ (for getting more insights about data)

Humans can see max \rightarrow 3D and 2D visuals

what if we have 1000 dimensions

↓ bottom

3d or 2d

bubbles

credit

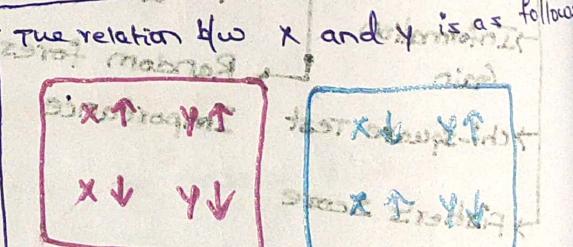
bottom

feature selection: Helps to select most important features which actually helps to predict the output.

Ex:-

x	y
1	1
2	2
3	3
4	4

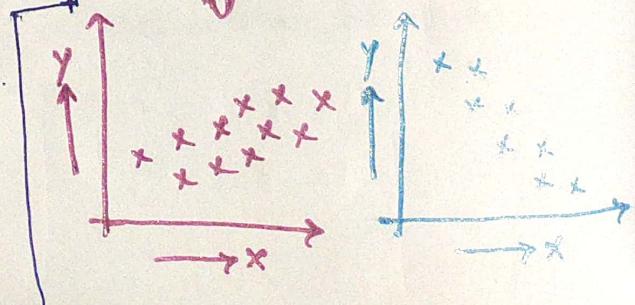
Let's assume some values are there corresponding to $x(1/p)$ and $y(0/p)$ features.

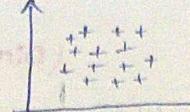


→ we can quantify the Relation in terms of mathematical terms.

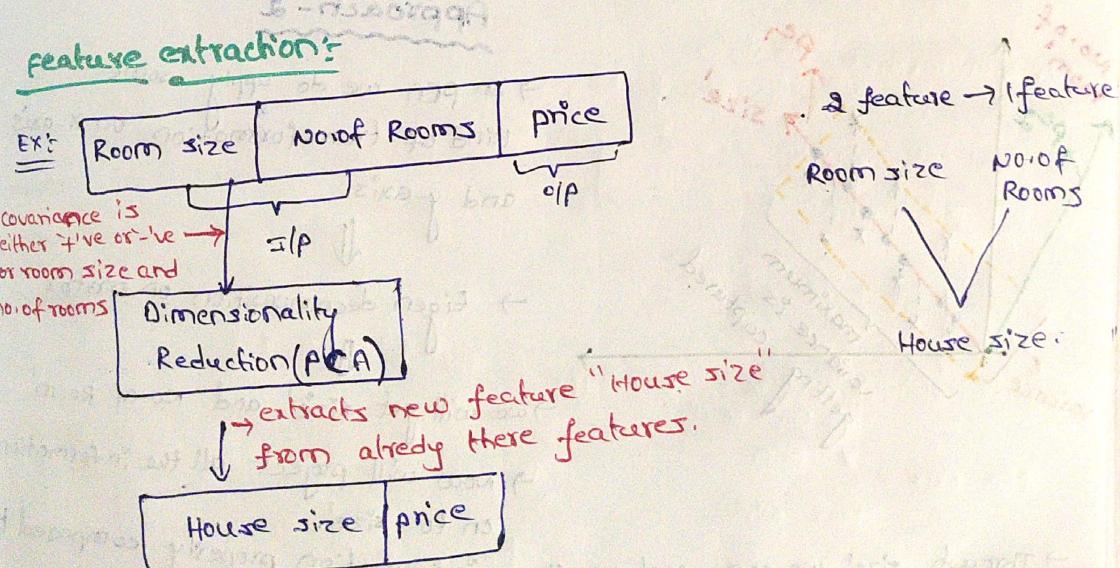
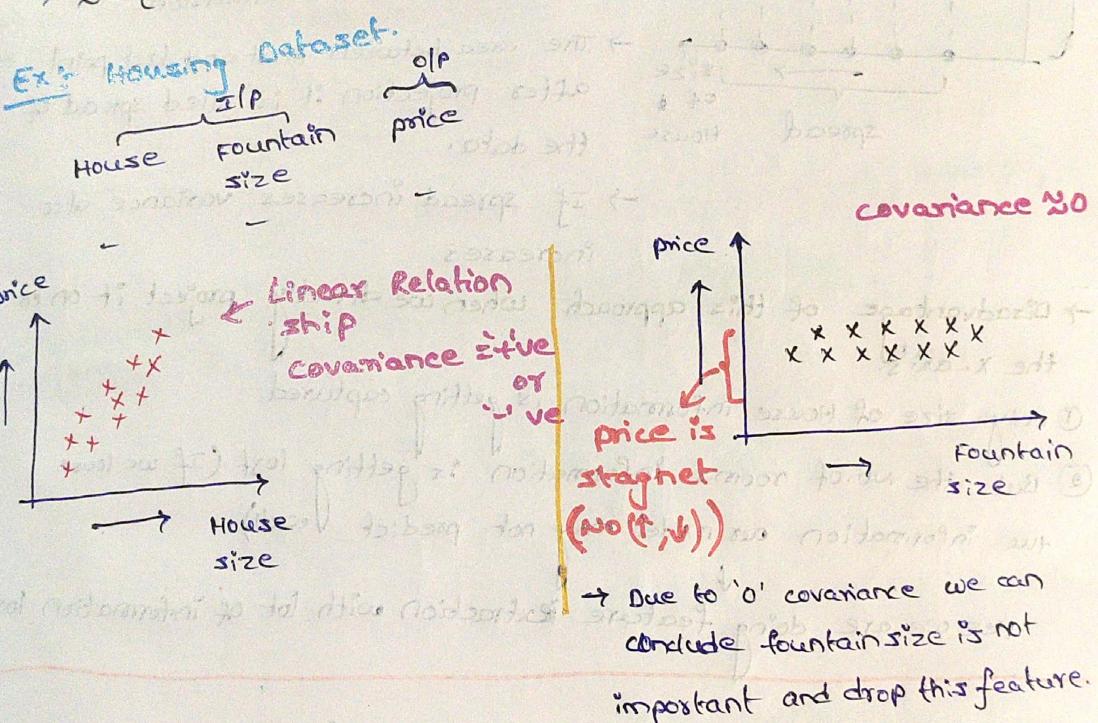
$$\text{cov}(x,y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{N-1}$$

= +ve
= -ve
= 0



$\text{cov}(x, y) = 0 \Rightarrow$ 

 → whichever features have this kind of highly positive covariance, those features are super important.
 x is important for predicting y .
 Pearson correlation coefficient is $= \frac{\text{cov}(x, y)}{\sigma_x \cdot \sigma_y} = -1 \text{ to } 1$
 → the more towards the value of +1, the more y 've correlated x & y
 → the more towards the value of -1, the more -ve correlated it is.
 → ≈ 0 {No relationship}



PCA Geometric Intuition

(Dimensionality Reduction)

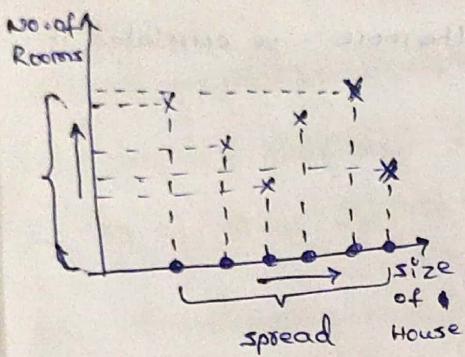
Size of House	No. of Room	price
-	-	↓ dp
-	-	-

with PCA



2 dimension \rightarrow 1 dimension

\rightarrow let's plot the above 3D features in the form of scatter plot.



Approach - 1

\rightarrow when we project all the data points on to the x-axis, we may get all the data points in one dimension.

\rightarrow The area between first and last point after projection it is called spread of the data.

\rightarrow If spread increases variance also increases.

\rightarrow Disadvantage of this approach when we directly project it on to the x-axis.

① only size of house information is getting captured.

② But the No. of rooms information is getting lost (If we lose the information our model may not predict well).

\rightarrow Here we are doing feature extraction with lot of information lost

Approach - 2

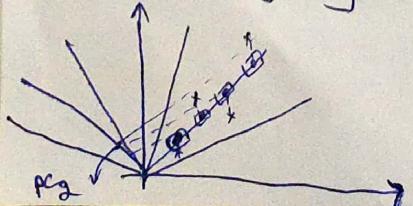
\rightarrow In PCA we do apply some kind of transformation on x-axis and y-axis.

\rightarrow Eigen decomposition on matrix

\rightarrow we will get size' and no. of room'

\rightarrow now will project all the information on to size'.

\rightarrow Through size' we can capture the information properly compared to approach 1 (projecting on to the x-axis).



\rightarrow PCA algorithm finds 2 principal components, when we project the points on to the line whichever is able to capture more variance.

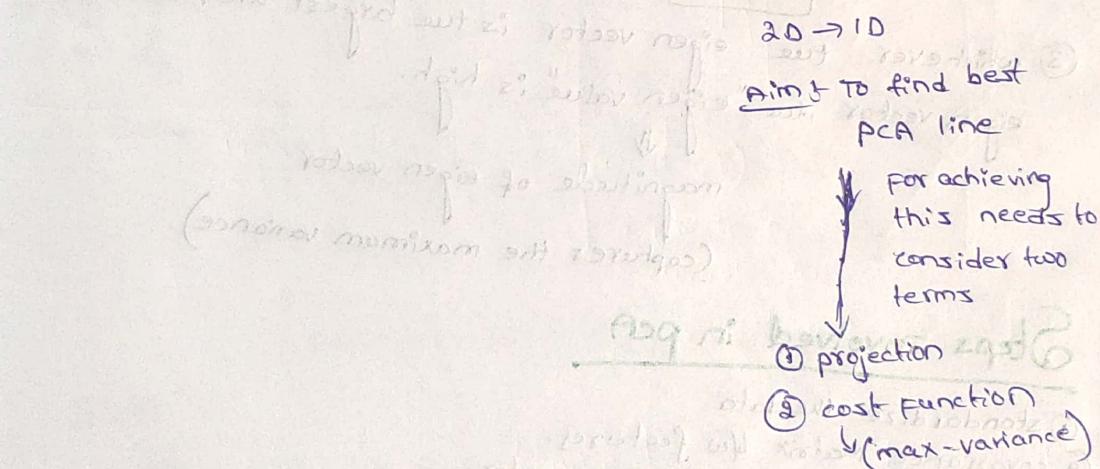
- PC_1 will capture the maximum amount of variance.
- PC_2 will capture the next maximum amount of variance after PC_1 .
- If we have three dimensions, PC_3 will capture the next maximum amount of information after PC_1 and PC_2 ($2D \rightarrow 1D$)
- The geometric intuition behind PCA is to find out principal components i.e., we need to find out one line that should be able to capture maximum amount of variance.
- Once we take this PC_1 , we can convert this ~~from~~ dataset of $2D$ from $1D$.

Why variance is important?

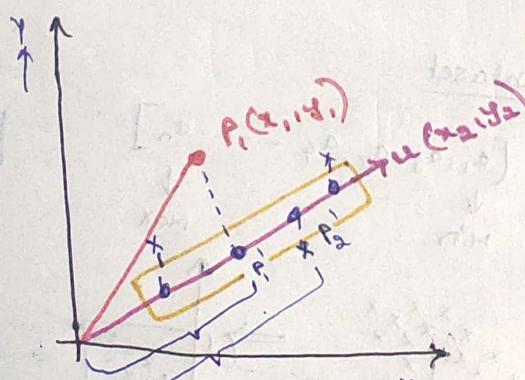
- Variance talks about the spread of the data.
- With the help of variance we can capture more information.

<u>2 Dimension</u>	<u>3 Dimension</u>
PC_1, PC_2 $\text{var}(PC_1) > \text{var}(PC_2)$	PC_1, PC_2, PC_3 $\text{var}(PC_1) > \text{var}(PC_2) > \text{var}(PC_3)$

Mathematical Intuition Behind PCA Algorithm.



- When we project points on to the line maximum variance is getting captured.
- Out of all these points let's take up one point.



These values talking about distance from origin.

projection of point P_i on

unit vector u

$$\text{proj}_{P_i} \cdot u = \frac{P_i \cdot u}{\|u\|}$$

$\|u\| = 1 \Rightarrow$ unit vector

$$\text{proj}_{P_i} \cdot u = P_i \cdot u \Rightarrow \text{scalar value}$$

$$P'_1, P'_2, P'_3, P'_4, \dots, P'_n$$

↳ scalar values

$P_1, P_1, P_2, P_3, P_4, \dots, P_n$ Goal :- find the best unit vector which captures maximum variance

$x_0, x_1, x_2, x_3, x_4, \dots, x_n$ (variance) (continues to training)

$$\text{Max-variance} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$$

cost function.

→ we should not go ahead and select different unit vectors to find out best vector.

→ In order to find that unit vector we have a method eigen decomposition

→ How do we find out which unit vector is best to capture maximum variance, so for that we will be using

"Eigen values and Eigen vectors"

→ ① First we need to find covariance matrix between features

→ ② Eigen vectors and Eigen values will find out from this covariance matrix

$$AV = \lambda V$$

→ ③ whichever the eigen vector is the largest one, for that eigen vector the "eigen value" is high.

↓
magnitude of eigen vector

(captures the maximum variance)

Steps involved in PCA

- ① standardize the data
- ② covariance matrix b/w features.
- ③ find out eigen value and eigen vector.
- ④ principal component.

scaling of the data:-

Min-Max standard scalar \rightarrow Data point \rightarrow Suppress the data [0-1]

$$\text{Min-Max} = \frac{x - \text{Min}}{\text{Max} - \text{Min}}$$

$$\text{Scaling factor} = \frac{\text{Max} - \text{Min}}{\text{Max} - \text{Min}} = 0$$

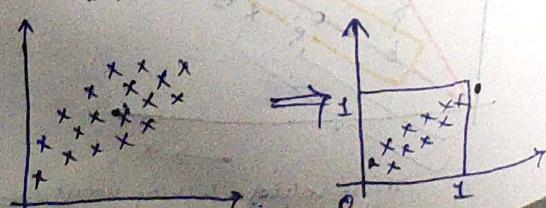
$$\begin{aligned} & \text{Max-Min} \\ & \therefore \frac{\text{max} - \text{min}}{\text{max} - \text{min}} = 1 \end{aligned}$$

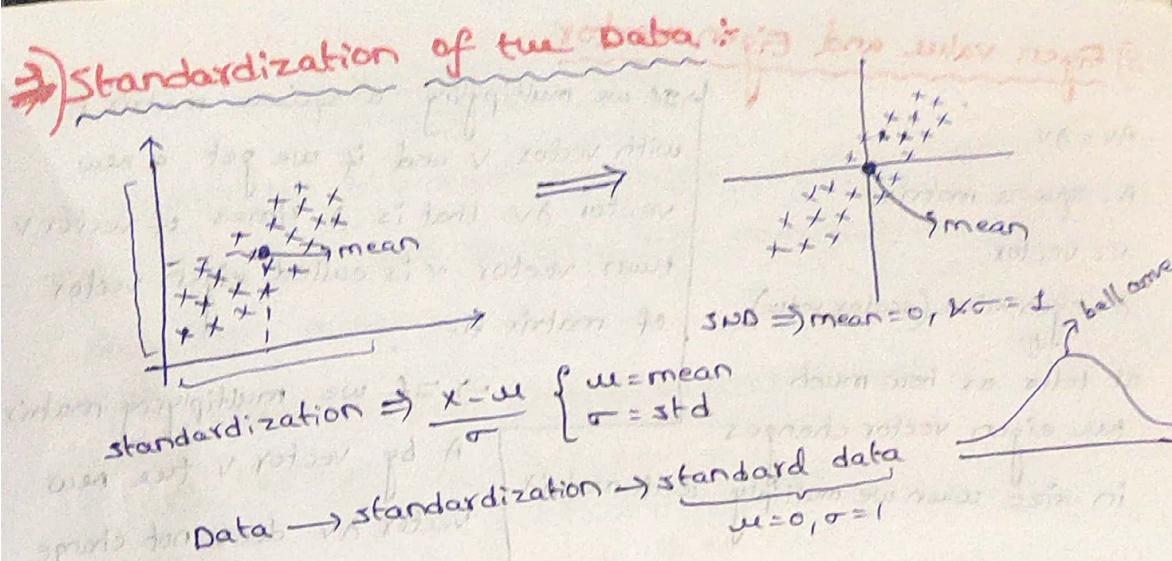
Dataset

$$[a_1, a_2, a_3, a_4, \dots, a_n]$$

↓ Min

↓ Max





Variance

$2, 8, 10, 12, 6, 5, 13, 14$

mean \downarrow

deviation $\Rightarrow (\text{mean} - \text{DP})^2$

mean of deviation \downarrow

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

$\downarrow \sigma^2 = \text{variance}$

variance \Rightarrow spread of the data around the mean.

2) covariance matrix of X & Y

$\begin{matrix} x & y \\ f_1 & f_2 \end{matrix} \Rightarrow \begin{matrix} \text{cov}(x,y) \\ \text{cov}(y,x) \end{matrix}$

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

$\text{cov}(x,x) = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 = \text{var}(x)$

$\text{cov}(y,y) = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2 = \text{var}(y)$

$\text{cov}(x,y) = \text{cov}(y,x)$

For 2D

$$\begin{bmatrix} x & y \\ x & \text{cov}(x,x) \quad \text{cov}(x,y) \\ y & \text{cov}(y,x) \quad \text{cov}(y,y) \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} \text{var}(x) & \text{cov}(x,y) \\ \text{cov}(y,x) & \text{var}(y) \end{bmatrix}$$

\Rightarrow covariance matrix of X & Y
 Features

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

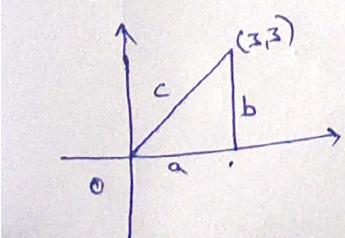
$\Rightarrow \text{std} = 1, \text{mean} = 0$

$\begin{bmatrix} x & y \\ s & E \end{bmatrix} = S^{-1} \text{Cov}^{-1} \begin{bmatrix} E \\ s \end{bmatrix}$

$\Rightarrow \frac{1}{N} \sum_{i=1}^N x_i y_i$

$x_1 y_1 + x_2 y_2 + x_3 y_3 + \dots + x_4 y_4$

\hookrightarrow dot product of two vectors



$$\begin{aligned} c^2 &= a^2 + b^2 \\ c &= \sqrt{a^2 + b^2} \\ c &= \sqrt{3^2 + 3^2} \\ &= \sqrt{18} = 4\sqrt{2} \end{aligned}$$

$$\begin{aligned} OC &= 4 \cdot 2 \begin{bmatrix} 3 \\ 3 \end{bmatrix} \\ &= 3\hat{i} + 3\hat{j} \end{aligned}$$

3) Eigen value and Eigen vector

$$Av = \lambda v$$

A = square matrix

v = vector

λ = eigen value (scalar value)

It tells us how much the eigen vector changes in size when we multiply with matrix.

If we multiplying a square matrix A with vector v and if we get a new vector λv that is λ times of vector v, then vector v is called eigen vector of matrix A.

If we multiplying matrix A by vector v the new vector λv does not change the direction.

$$Av - \lambda v = 0$$

$$(A - \lambda I)v = 0$$

$\det |A - \lambda I| = 0$ \Rightarrow identity matrix.

$\lambda_1, \lambda_2 \Rightarrow$ eigen values.

$$\downarrow \quad \downarrow$$

$$PC_1 \quad PC_2$$

$$3D \rightarrow 3D$$

$$\lambda_1, \lambda_2, \lambda_3$$

$$\downarrow \quad \downarrow \quad \downarrow$$

$$PC_1, PC_2, PC_3$$

$$\checkmark$$

$$1D$$

$$1D \Rightarrow 3D$$

$$3D \rightarrow 1D$$

$$\lambda_1, \lambda_2, \lambda_3$$

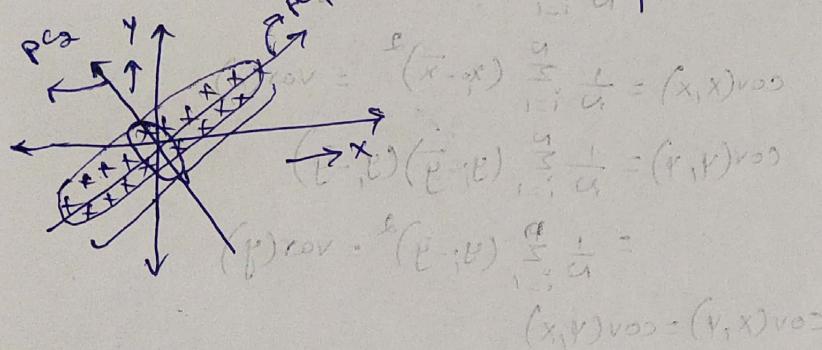
$$\downarrow \quad \downarrow \quad \downarrow$$

$$PC_1, PC_2$$

$$\checkmark$$

$$1D$$

$$1D \Rightarrow 1D$$



$$\begin{bmatrix} (v, x) v_{uvw} & (x, x) v_{uvw} \\ (v, y) v_{uvw} & (x, y) v_{uvw} \end{bmatrix}$$

For 2D column coordinates

$$\begin{bmatrix} (v, x) v_{uvw} & (x, x) v_{uvw} \\ (v, y) v_{uvw} & (x, y) v_{uvw} \end{bmatrix}$$

$$(v, x)(x, x) + (v, y)(x, y)$$

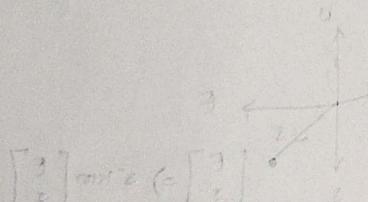
zeroes

$$(v, x)(x, x) + (v, y)(x, y)$$

zeroes

$$(v, x)(x, x) + (v, y)(x, y)$$

zeroes for subtracting terms



$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} x \\ y \end{bmatrix}$$

$$1+1=2$$

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} x \\ y \end{bmatrix}$$

