

PRINCIPAL COMPONENT ANALYSIS → Dimensionality Reduction

> CURSE OF DIMENSIONALITY:

> CURSE OF DIMENSIONALITY:

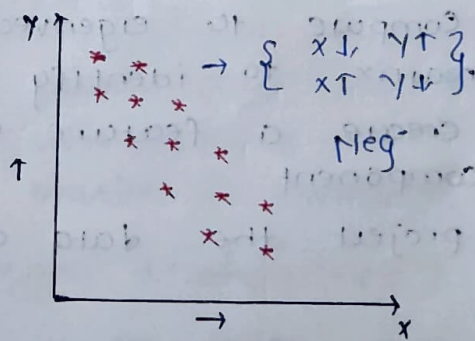
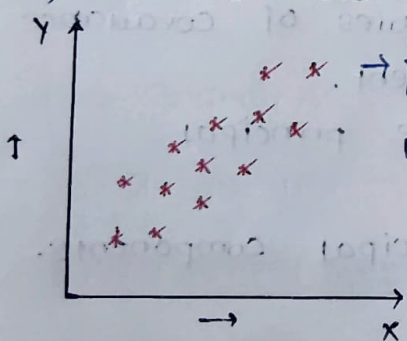
- Consider dataset with two feature x_1 and x_2 .
- When two independent variable are very strongly interacting with each other, that is the correlation coefficient is close to 1 then we are providing the same information to algorithms in two dimension, which is nothing but redundancy. it can lead poor performance in model.
- In short when we have too many dimension more than required, which is nothing but multicollinearity in data.
- PCA also helps to reduce dimension without any lose huge information.

> TWO DIFFERENT WAY OF REMOVE CURSE OF DIMENSIONALITY

1.) FEATURE SELECTION

2.) PCA - FEATURE EXTRACTION

1.) FEATURE SELECTION



$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n \frac{(x_i - \bar{x}) * (y_i - \bar{y})}{n-1} = \text{+ve or -ve} \rightarrow \text{Any value}$$

$$\text{pearson's corr} = \frac{\text{Cov}(x, y)}{\sigma_x * \sigma_y} = -1 \text{ to } 1$$

> In pearson correlation the more value towards 1 more positive correlated or more toward -1 more negative correlated or value towards 0 no correlation.

2.) FEATURE EXTRACTION - PCA

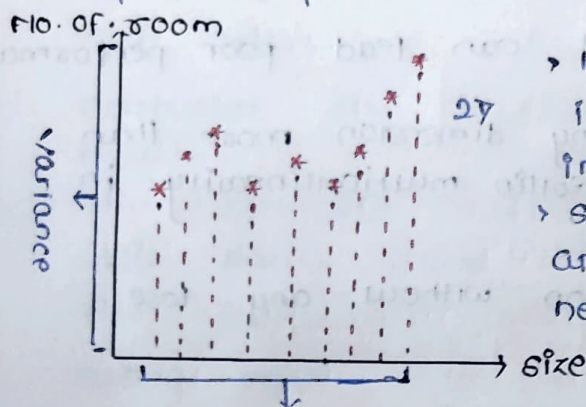
terminology:

- > variance: spread of datapoints from its mean
- > covariance: measure the relationship between x and y .

> GEOMETRIC INTUITION:

> consider housing dataset with no. of rooms and size.

Aim: predict price based on no. of rooms, and size.



> if you drop any one feature, it can make a huge loss of information.

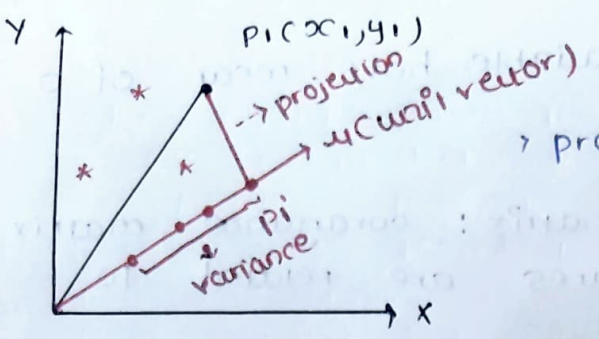
> so we use PCA. In PCA we apply some transformation that give new feature with maximum variance.

> In PCA we reduce dimension $2D \rightarrow 1D$

STEPS:-

- 1.) Compute covariance matrix.
- 2.) Compute the eigenvectors and eigenvalues of covariance matrix to identify principal component.
- 3.) Create a feature vector to decide principal component.
- 4.) project the data along the principal components.

* MATH INTUITION BEHIND PCA ALGORITHM * STEPS FOR PCA ALGORITHM



> projection of P_i on u

$$\text{proj}_{P_i} u = \frac{P_i \cdot u}{\|u\|^2}$$

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

$$\boxed{\text{proj}_{P_i} u = P_i \cdot u}$$

> when we do all the projection we will get

$\{P_0', P_1', P_2', P_3', P_4', \dots, P_n'\}$ → scalar value
 Talking about distance from origin,

different notation

$$\{x_0', x_1', x_2', x_3', x_4', \dots, x_n'\}$$

Max variance. or Cost function $\sum_{P=1}^n \frac{(x_i - \bar{x})^2}{n}$ } - Aim: find the best unit vector which capture maximum variance

> EIGEN VECTOR AND EIGEN VALUES : (eigen decomposition)

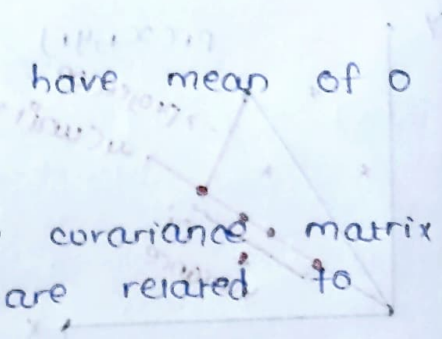
let say A is matrix and v is eigen vector and λ is eigenvalue

$$\therefore \{A v = \lambda v\}$$

* STEPS FOR PCA ALGORITHM *

1.) standardize the data : all variable have mean of 0 and standard deviation of 1.

2.) calculate the covariance matrix : covariance matrix gives you idea how features are related to each other.



3.) calculate the eigenvector and eigenvalue of covariance matrix.

-the eigenvector represent the direction in which the data varies most.

-the eigenvalue represent the amount of variation along each eigenvector.

4.) choose the principal component

5.) Transform the data into lower dimension.

EIGEN VECTOR AND EIGEN VALUE : (Eigen decomposition)

$$A \cdot v = \lambda \cdot v$$