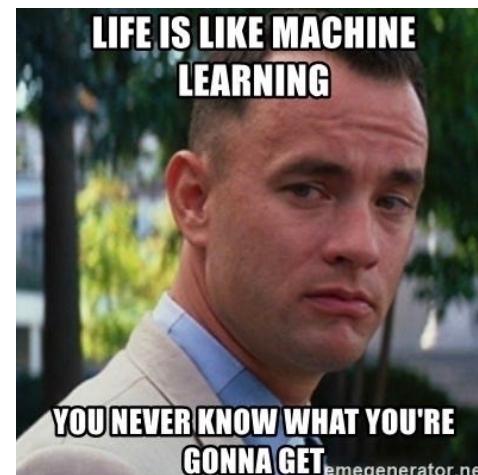
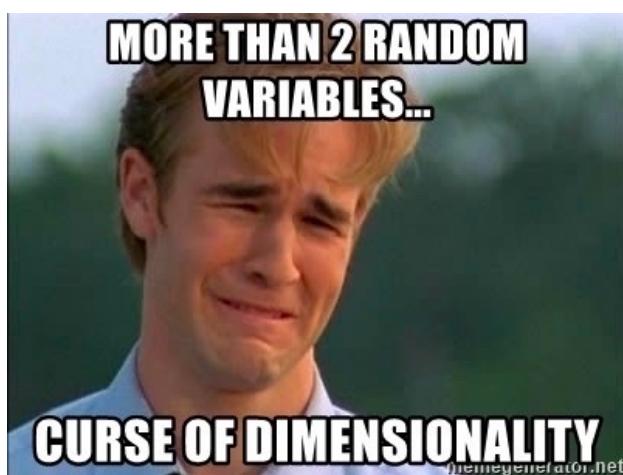
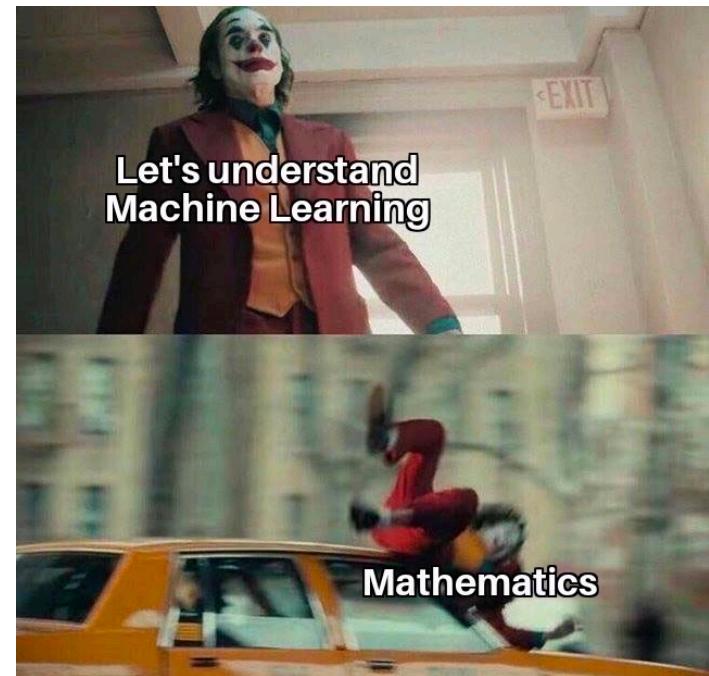
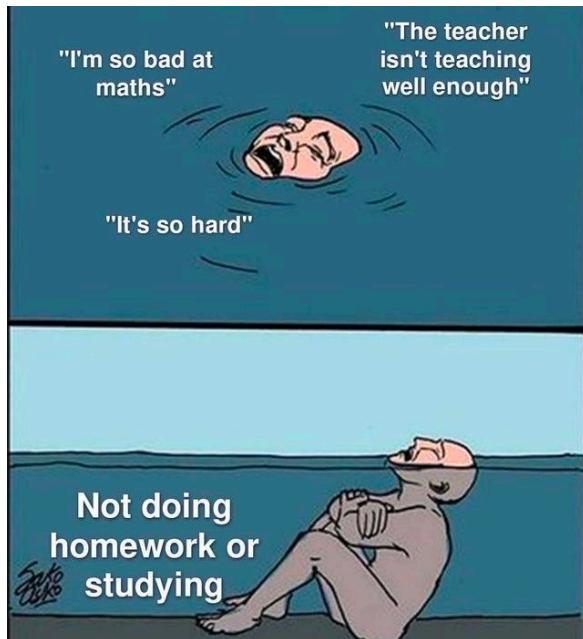


June 13, 2023.

DSML : Math for ML.

Optimization 5: Principal Component Analysis.



Class begins
@ 9:03 p.m.
sharp.

Dimensionality Reduction

$D = \{(\bar{x}_i, y_i) : \bar{x}_i \in \mathbb{R}^d\}_{i=1}^n$

Labelled Dataset

$\left\{ \begin{array}{c} \text{n rows} \\ \text{---} \\ \text{---} \end{array} \right\}$

$\left[\begin{array}{c|c|c|c|c|c} x_1 & x_2 & \dots & x_d & y \end{array} \right]$

$d + 1$ columns

$\bar{x}_i = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$

$d \gg 4$

target / label.

Examples:

- 1] Principal Component Analysis (PCA)
- 2] t-SNE.
- 3] UMAP.

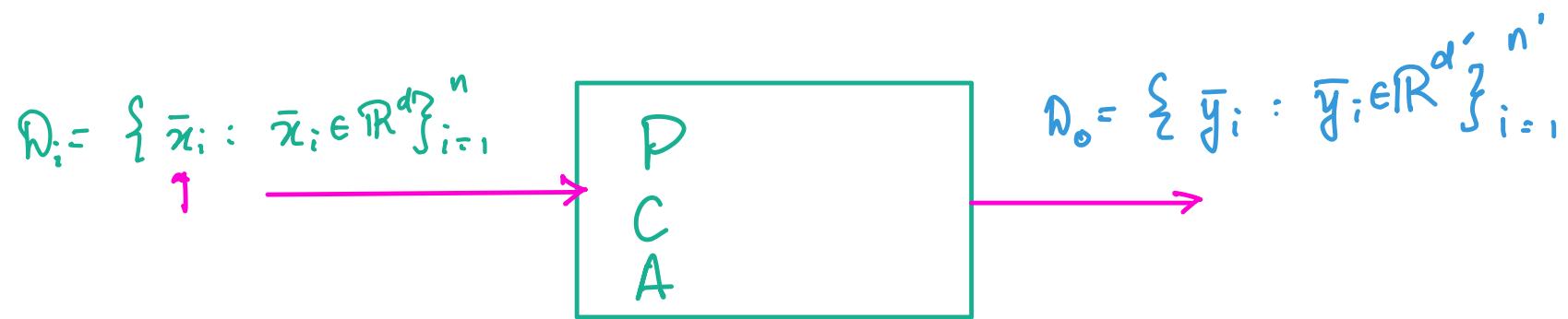
Objective: ① Visualize this dataset (Reduce dimension to 2 or 3).

② Keep as much information as possible.

Why is Dimensionality reduction needed?

- 1] To make it easier to visualize high dimensional data.
- 2] less dimensions \Rightarrow lesser memory to store data
 \Rightarrow less computation while doing ML.
- 3] Makes it easier to code / implement ML algos.

PCA as a black box:



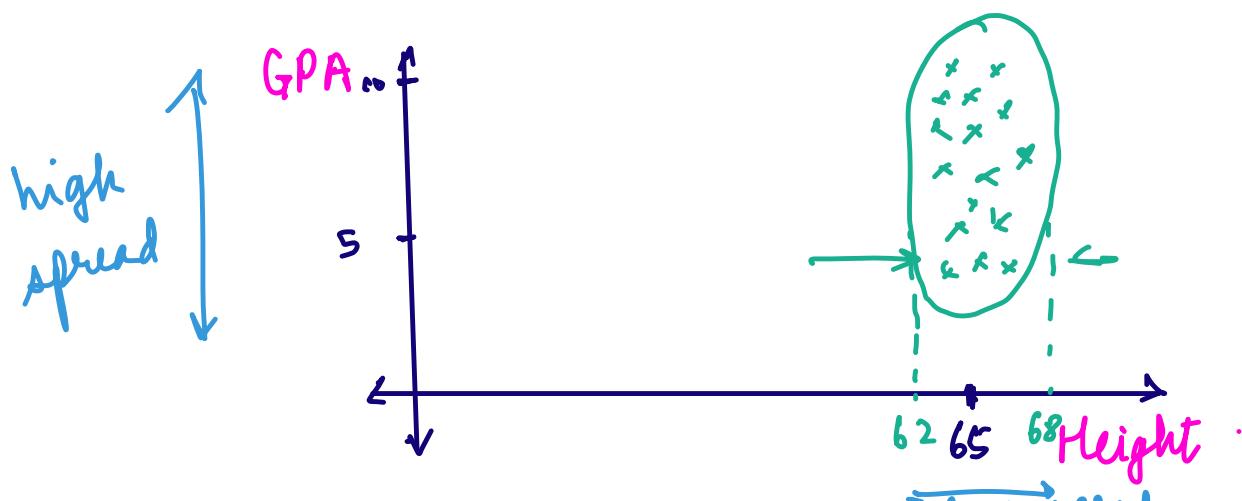
$$d \geq d'$$

$$n = n'$$

Examples:

1] GPA. (10) vs. Height.

spread of values
captures usefulness
of the feature.

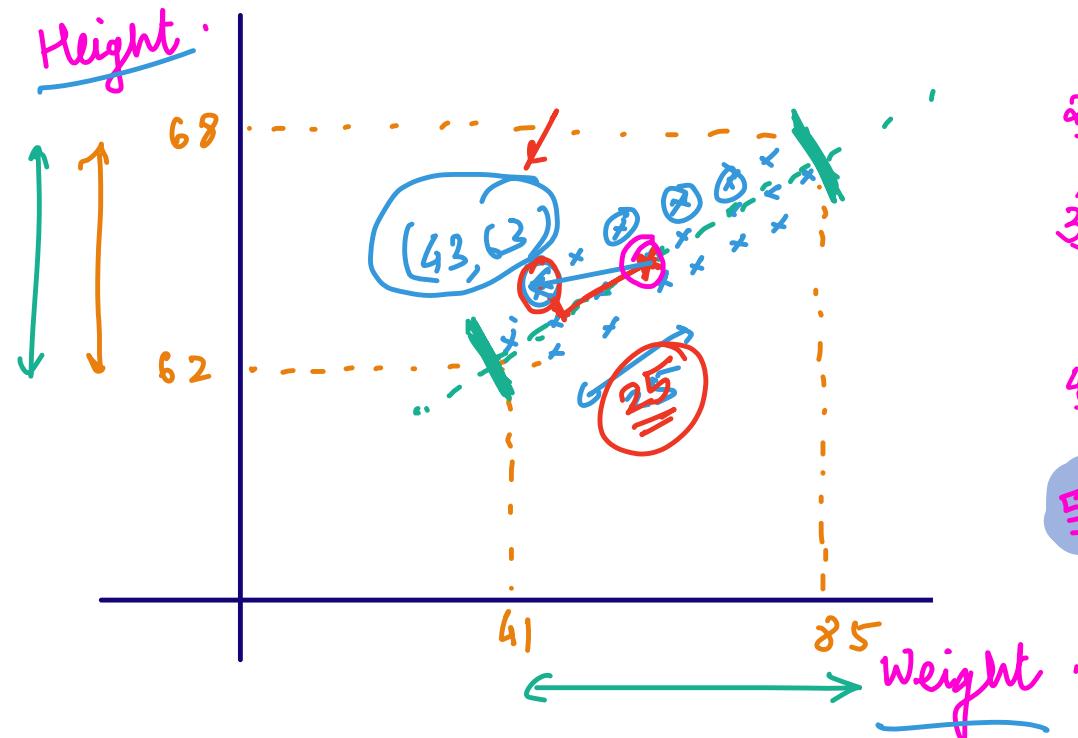


1] No correlation between height & GPA.

2] Range of values is low.

3] The standard deviation / variance / spread of the height feature is low.

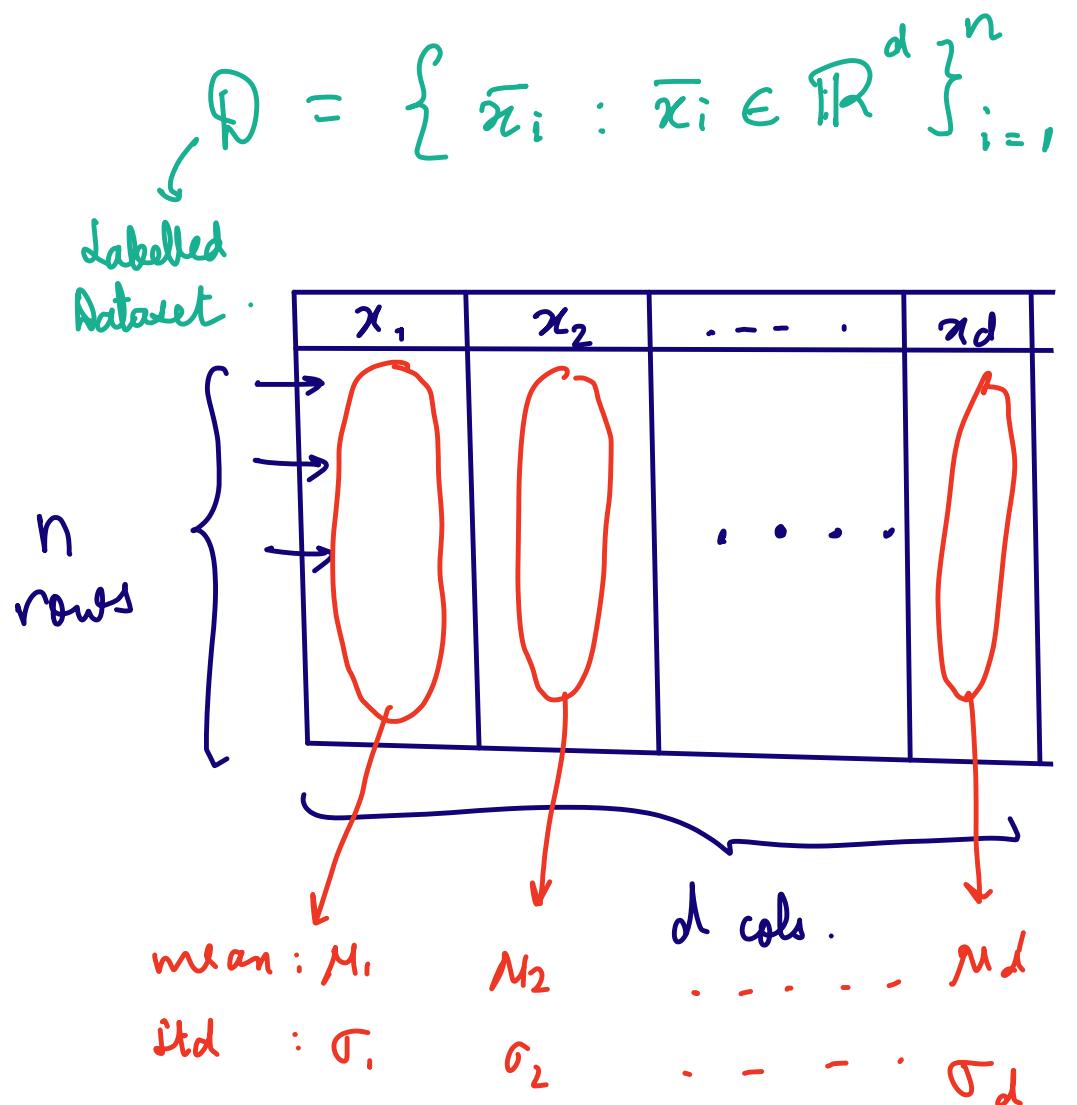
2) Height vs. Weight:



If we apply P.R. here, how to do it?

- 1] Calculate the centre of our data and change co-ordinate systems.
- 2] Standardize features
- 3] Find the direction of max. variance.
- 4] Project all points along this direction
- 5] Keep doing steps 3 & 4 for d' iterations.

Steps 1 & 2: Re-center & Standardize features. [column standardization.]



1] Recenter:

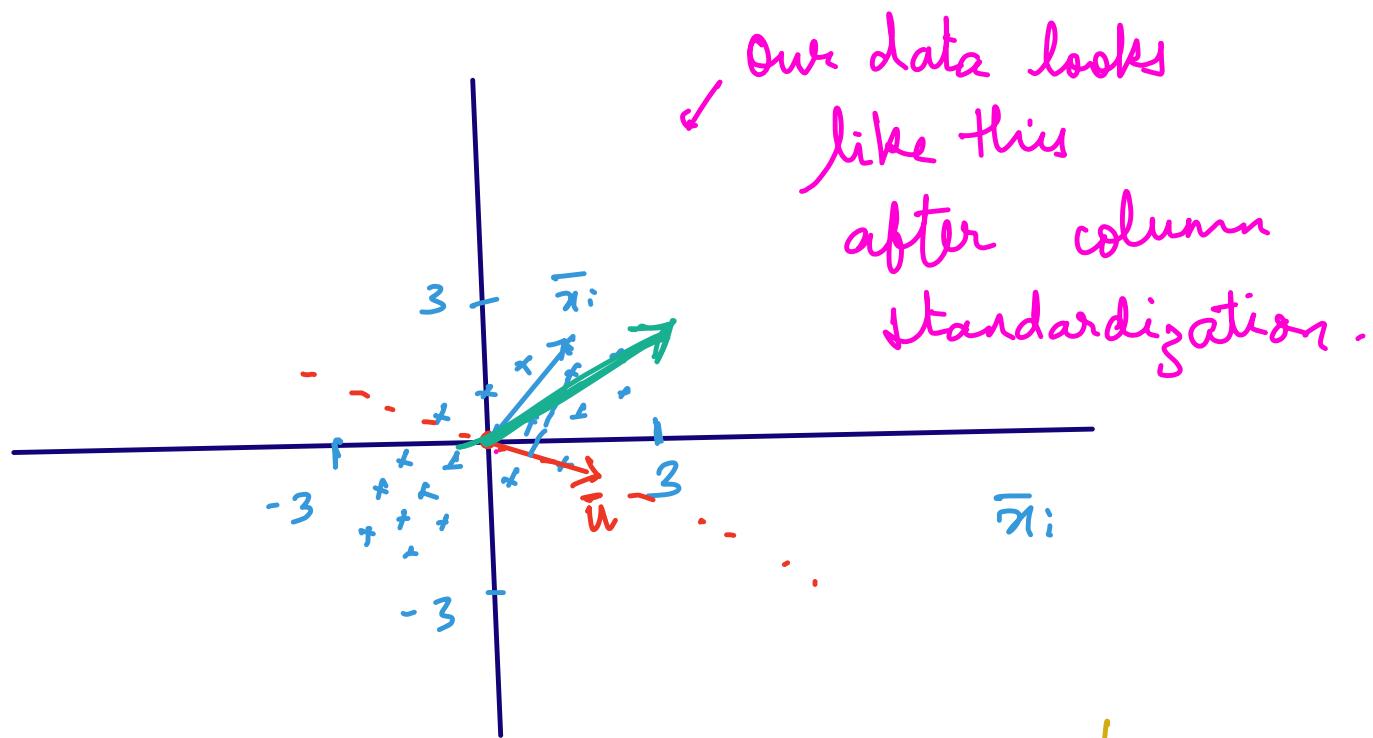
means: M_1, M_2, \dots, M_d .

stds: $\sigma_1, \sigma_2, \dots, \sigma_d$.

2] $\bar{x}_i' = \frac{x_i - M_i}{\sigma_i}$

Apply this to
all columns.

$$D' = \{ \bar{x}_i' : \bar{x}_i' \in \mathbb{R}^d \}_{i=1}^n$$



Next, we will find the direction of max. variance.

→ Design our cost function.

Sum of Projection of \bar{x}_i on \bar{u} :
Squares of

$$\sum_{i=1}^n \left(\frac{\bar{x}_i^\top \bar{u}}{\|\bar{u}\|} \right)^2$$

will be more
↓ for better \bar{u}

The optimization problem for PCA.

$$\max_{\bar{u}} \sum_{i=1}^n \left(\underbrace{\frac{\bar{x}_i^\top \bar{u}}{\|\bar{u}\|}}_{} \right)^2$$

Let's assume \bar{u} is a unit vector.



$$\max_{\bar{u}} \sum_{i=1}^n (\bar{x}_i^\top \bar{u})^2$$

$$\text{s.t. } \|\bar{u}\|^2 = 1.$$

Lagrange multipliers.

$$\max_{\bar{u}} \sum_{i=1}^n (\bar{x}_i^\top \bar{u})^2 + \lambda (\|\bar{u}\|^2 - 1)$$

$$\|\bar{v}\|^2 \text{ is same as} \\ \underset{\bar{u}}{\text{min}} \quad \sum_{i=1}^n (\bar{x}_i^\top \bar{u})^2 + \lambda (\|\bar{u}\|^2 - 1)$$

This is the problem we have to solve.

- Take the gradient of $f(\bar{x}_i; \bar{u})$ and set it to 0.

$$X \bar{u} = \bar{v}$$

$$\mathcal{D} = \left\{ \bar{x}_i : \bar{x}_i \in \mathbb{R}^d \right\}_{i=1}^n$$

$$X \rightarrow \begin{bmatrix} \bar{x}_1^\top \\ \bar{x}_2^\top \\ \vdots \\ \bar{x}_n^\top \end{bmatrix} \cdot \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_d \end{bmatrix}^{d \times 1} = \begin{bmatrix} \bar{x}_1^\top \cdot \bar{u} \\ \bar{x}_2^\top \cdot \bar{u} \\ \vdots \\ \bar{x}_n^\top \cdot \bar{u} \end{bmatrix}_{n \times 1}$$

$$\|\bar{v}\|^2 \text{ is same as } \max_{\bar{u}} \sum_{i=1}^n (\bar{x}_i^\top \bar{u})^2 + \lambda (\|\bar{u}\|^2 - 1)$$

f(\bar{x}_i ; \bar{u})

$$\Rightarrow \max_{\bar{u}} \bar{v}^\top \bar{v} + \lambda (\|\bar{u}\|^2 - 1).$$

$$\Rightarrow \max_{\bar{u}} (\bar{x}^\top \bar{u})^\top (\bar{x}^\top \bar{u}) + \lambda (\bar{u}^\top \bar{u} - 1)$$

$= \bar{u}^\top \bar{x}^\top \bar{x} \bar{u}$

$$(A \cdot B)^\top = B^\top A^\top$$

$$\max_{\bar{u}} \quad \boxed{\bar{u}^T x^T x \bar{u}} + \lambda (\bar{u}^T \bar{u} - 1).$$

$f(x; \bar{u})$

Now, take gradient ∇ set it to zero.

$$\text{let } f(\bar{x}) = \bar{x}^T A \bar{x}$$

$$\text{then } \nabla_{\bar{x}} f(\bar{x}) = \underbrace{(A + A^T) \bar{x}}_{\cdot}$$

$$\nabla_{\bar{u}} \bar{u}^T \bar{u} = 2 \bar{u} \quad A = x^T x \\ A^T = (x^T x)^T = x^T (x^T)^T = x^T x.$$

$$\nabla_{\bar{u}} f(x; \bar{u}) = 2 \cdot x^T x \bar{u} + \lambda 2 \bar{u}$$

$$x^T x \bar{u} + \lambda \bar{u} = 0$$

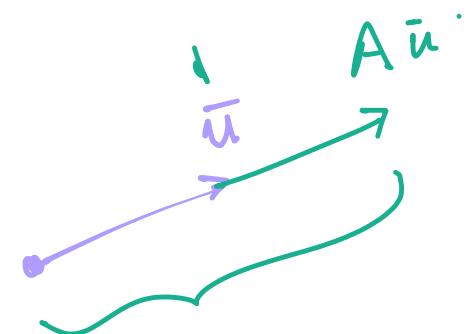
$x^T x \bar{u} = -\lambda \bar{u}$

$$X^T X \bar{u} = \underbrace{-\lambda}_{\lambda' = -\lambda} \bar{u}$$

$$\underbrace{X^T X}_{\text{matrix}} \bar{u} = \lambda' \bar{u}.$$

$$A \cdot \bar{u} = \lambda' \bar{u}$$

constant -



\bar{u} is an eigenvector of the matrix A with an eigenvalue of λ' .

Actually, \bar{u} is the eigenvector with the maximum eigenvalue.

In general, our dataset has d -dim vectors.

If I want to get d' dimensions for dimensionality reduction, then:

i] Calculate the covariance matrix $X^T X$.

ii] Next calculate the eigenvectors \vec{v} and eigenvalues of $X^T X$.

X ($n \times d$)

Here dimension of $X^T X = \underline{d \times d}$.

X^T ($d \times n$)

It is guaranteed that we will get d eigenvalues, eigenvector pairs.

Also, the d eigenvalues ≥ 0 .

After step 2, we will have:

$$\{ (\bar{u}_1, \lambda_1), (\bar{u}_2, \lambda_2), (\bar{u}_3, \lambda_3), \dots, (\bar{u}_d, \lambda_d) \}$$

3] Now, sort the above so that:

$$\lambda_1 \geq \lambda_2 \geq \lambda_3 \dots \geq \lambda_d.$$

4] Finally, my reduced dimension is obtained by selecting the top d' eigenvalue & eigenvector pairs.

$$\boxed{\bar{y}_i} =$$

$$\begin{bmatrix} (\bar{x}_i^\top \bar{u}_1) \\ (\bar{x}_i^\top \bar{u}_2) \\ \vdots \\ (\bar{x}_i^\top \bar{u}_{d'}) \end{bmatrix} \rightarrow \text{final output -}$$

Practical Considerations while Using PCA .

