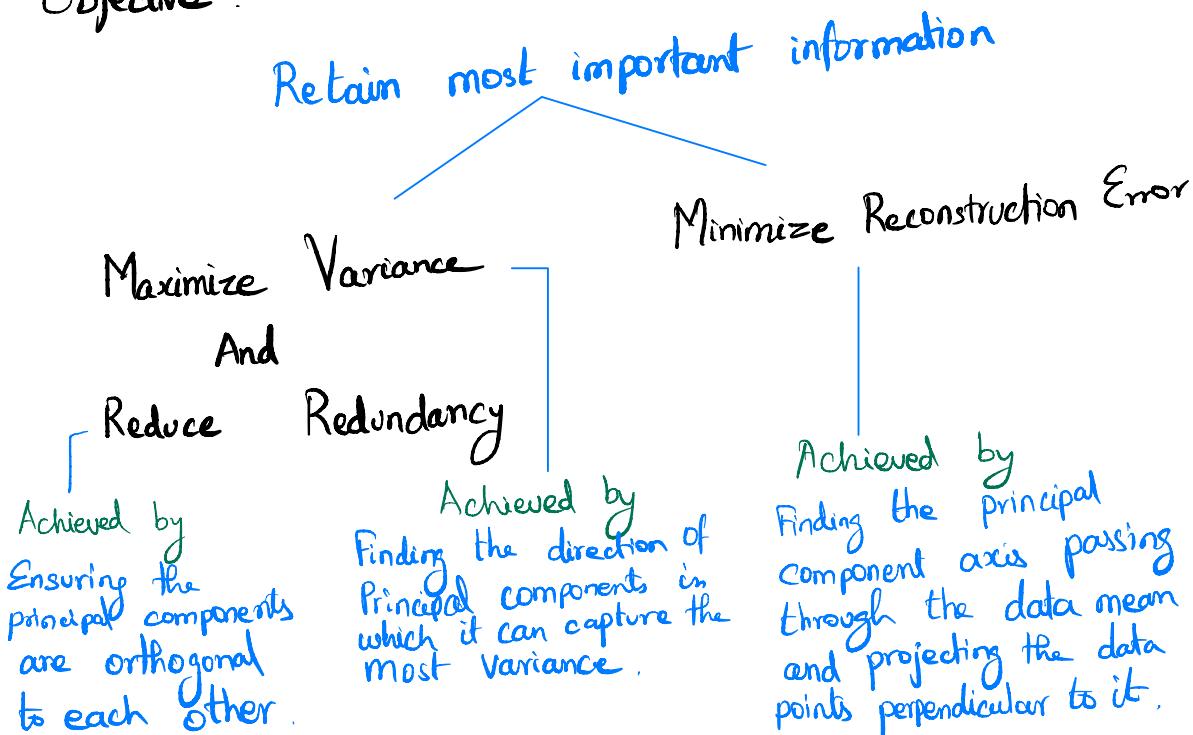




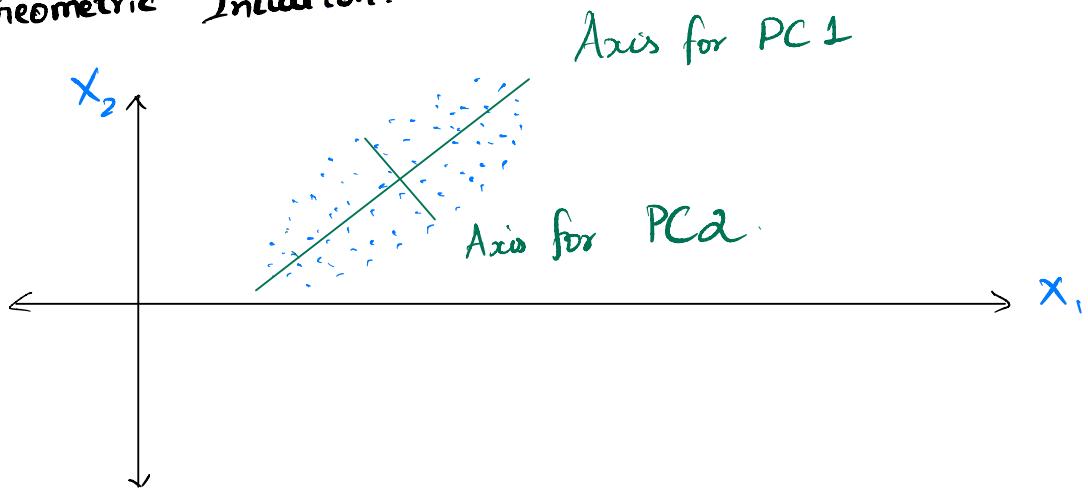
Principal Component Analysis: PCA

Aims to transform the original features of a dataset into a new set of Uncorrelated variables called Principal Components, with the objective of retaining most important information while discarding the least important.

Objective :



Geometric Intuition:



Dataset of d features can be transformed into k principal components, where $k \leq d$

Principal Components are ordered based on the amount of variance (information) they capture in the data.

Direction of PC1 \rightarrow Most variance

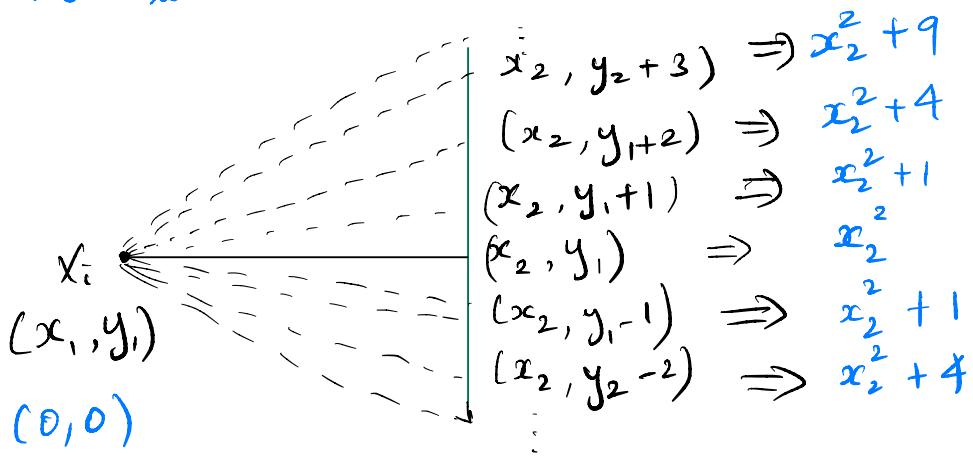
Direction of PCd \rightarrow Least variance.

Principal components are orthogonal to each other ensuring each PC captures unique and non-redundant aspect of the variability of the data.

Let's first solve how we are going to project our data points to these new principal component axes.

Assume we found our new axes

We need to project original data points onto this new line with minimal SSE



Triangle analogy: The straight "perpendicular" path is like one side of the triangle. It's always shorter than any other path you could take that forms a slanted side of the triangle

We need to find

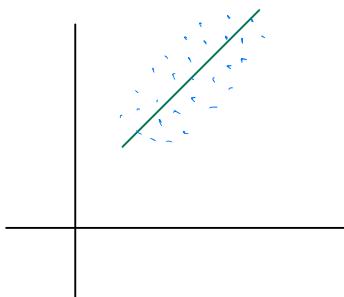
- 1) Direction of each principal component
- 2) New co-ordinates [w.r.t new axis original axis]

We can find these by solving our objectives.

② New co-ordinates:

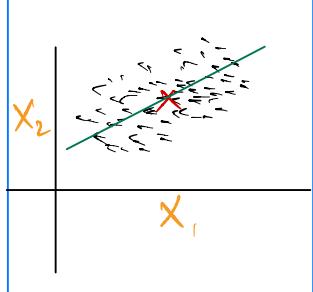
Assume we have the optimal direction for the principal component 1.

Let \vec{e} be the best direction in which we can capture most variance.

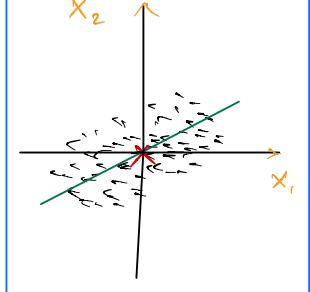


In order to find the new co-ordinates, we need a reference point, like the origin in our current co-ordinate system and then \vec{e} will be our new basis.

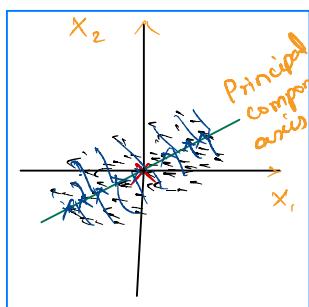
To minimize SSE, the new axis line should pass through the data mean.



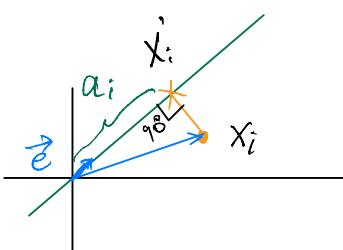
X
on original
co-ordinate system



$X_{\text{mean_centered}} = X - \bar{X}_{\text{mean}}$
on original
co-ordinate System



Project X on new axis



a_i - New co-ordinate
in the principal
component with
the basis \vec{e}

X_i - Co-ordinates of
data point of
interest in
original feature
space.

X'_i - Co-ordinates
of projected
data point
in original
feature space

To find a_i :

We know that angle b/w $(\overrightarrow{X_i - X'_i})$ & \vec{e} is 90° .
So, dot product is 0. $\Rightarrow (\overrightarrow{X_i - X'_i}) \cdot \vec{e} = 0$

Since, $X'_i = a_i \cdot \vec{e}$, $(\overrightarrow{X_i - a_i \cdot \vec{e}}) \cdot \vec{e} = 0$

$$X_i \cdot \vec{e} - a_i \cdot \vec{e} = 0$$

$$a_i = X_i \cdot \vec{e}$$

a_i gives the magnitude of X'_i

Optimal line that minimize the SSE & maximize the variance always go through the mean. So mean is the reference point or origin in new co-ord system. Mean center it to make MATH EASY

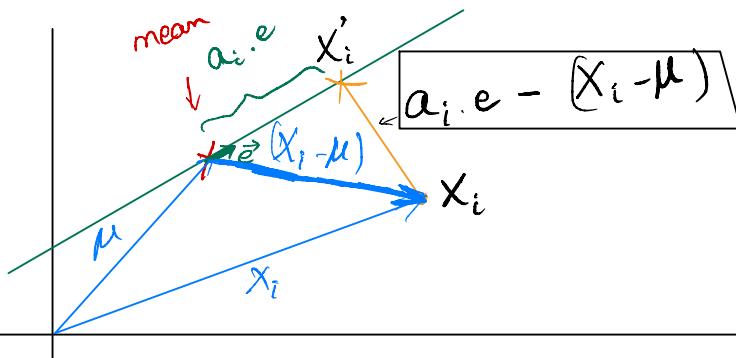
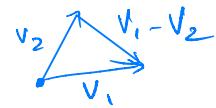
If not mean centered,
mean is our
reference point.

$$a_i = (X_i - \mu) \cdot \vec{e}$$

$$X'_i = \mu + a_i \cdot \vec{e}$$

If not mean centered

Reminder



Since, angle between vectors $\overrightarrow{a_i \cdot e - (x_i - \mu)}$ & \overrightarrow{e} is 90° , $(a_i \cdot e - (x_i - \mu)) \cdot e = 0$

$$a_i \cdot e - (x_i - \mu) \cdot e = 0$$

$$a_i \cdot e = (x_i - \mu) \cdot e$$

$$\text{And, } x'_i = \mu + a_i \cdot e$$

Two objectives: 1) Maximize Variance
2) Minimize SSE

② MINIMIZE SSE

$$\begin{aligned} \text{SSE} &= \sum_{i=1}^n \|(\mu + a_i \cdot e) - x_i\|^2 = \sum_{i=1}^n \|a_i \cdot e - (x_i - \mu)\|^2 \\ &= \sum_{i=1}^n a_i^2 \|e\|^2 - 2 \sum_{i=1}^n a_i e^T (x_i - \mu) + \underbrace{\sum_{i=1}^n \|x_i - \mu\|^2}_{\text{constant}} \end{aligned}$$

Minimize SSE

a_i

$$\boxed{SSE = \sum_i a_i^2 - 2 \sum_i a_i e^T (x_i - \mu) + \sum_i \|x_i - \mu\|^2}$$

$$\frac{\partial SSE}{\partial a_i} = 2a_i - 2 e^T (x_i - \mu) \stackrel{\text{want}}{=} 0$$

$$a_i - e^T (x_i - \mu) = 0$$

$$a_i = e^T (x_i - \mu)$$

[Same result as the one we got when solved for the line passes through the MEAN]

① MAXIMIZE VARIANCE AND FIND \vec{e}

Two ways to find \vec{e} \rightarrow By minimizing SSE \rightarrow By Maximizing Variance

$$\begin{aligned} (i) \text{ MINIMIZE SSE} \\ SSE &= \sum_i a_i^2 - 2 \sum_i a_i e^T (x_i - \mu) + \sum_i \|x_i - \mu\|^2 \\ &= \sum_i a_i^2 - 2 \sum_i a_i^2 + \sum_i \|x_i - \mu\|^2 \\ &= - \sum_i a_i^2 + \sum_i \|x_i - \mu\|^2 \\ &= - \sum_i (e^T (x_i - \mu))^2 + \sum_i \|x_i - \mu\|^2 \end{aligned}$$

$$\begin{aligned}
 SSE &= - \sum_i e^T (x_i - \mu)(x_i - \mu)^T e + \text{constant (e)} \\
 &= -e^T \left[\sum_i (x_i - \mu)(x_i - \mu)^T \right] e + \text{ignore} \\
 &= -e^T \underbrace{\sum_i}_{\substack{\Downarrow \\ \text{covariance Matrix} \& \text{it is constant w.r.t } e}} e + \text{ignore constant (e)} \\
 \text{To minimize } SSE, \text{ we have to minimize } -e^T \sum e \\
 \text{or maximize } e^T \sum e
 \end{aligned}$$

(ii) ^{MAXIMIZE VARIANCE}

$$\begin{aligned}
 \text{Variance } (X_i \text{ on } e) &= E \left[(\mu + a_i e - E[\mu + a_i e])^2 \right] \\
 &= E \left[(\mu + a_i e - \mu)^2 \right] \\
 &= E \left[(a_i e)^2 \right] \\
 &= E \left[e^T (x_i - \mu) \cdot e^T (x_i - \mu) \right] \\
 &= E \left[e^T (x_i - \mu)(x_i - \mu)^T e \right] \\
 &= e^T E \left[(x_i - \mu)(x_i - \mu)^T \right] e
 \end{aligned}$$

$$\text{Variance}_{(X_i \text{ on } e)} = e^T \Sigma e$$

Again, to maximize variance, we have to

$$\text{maximize } e^T \Sigma e$$

[Same as minimizing SSE ; it tells us both are same objectives]

Yet Another Intuition:

After projecting the data points X onto a principal component of direction vector \vec{e} , the new data points vector is given as $e^T X$

| we know, on point level
 $a_i = x_i \cdot e$

So, let's revisit two things Variance & Magnitude.

$\sigma^2 = \frac{1}{n-1} \sum_i (v_i - \bar{v})^2$	$\ v\ ^2 = \sum_i v_i^2$
---	--------------------------

Variance and Vector magnitude are monotonically related when mean-centered; so maximizing one maximizes the other.

Hence, Maximizing variance is equal to maximizing the vector magnitude $\|e^T X\|^2$ [note: $\|e\|=1$]

$$\begin{aligned}\|e^T X\|^2 &= (e^T X)(e^T X)^T \\ &= e^T X X^T e \\ &= e^T \Sigma e\end{aligned}$$

Once again, it is clear that we need to maximize $e^T \Sigma e$ in order to maximize variance and minimize SSE.

$e^T \Sigma e$ is the quadratic form
Let's solve \vec{e} that maximizes $e^T \Sigma e$
This is a constrained optimization. $\Rightarrow \|e\|=1$
 $\Rightarrow e^T e = 1$

We can use Lagrange Multiplier to solve this constrained optimization.

Lagrange Multiplier:

	$\underbrace{\text{Objective}}$ $e^T \Sigma e$	$\underbrace{\text{Constraint}}$ $\lambda(e^T e - 1)$
--	---	--

MAX L = $e^T \Sigma e - \lambda(e^T e - 1)$

idea: $\frac{\partial e^T \Sigma e}{\partial e}$ & $\frac{\partial (e^T e - 1)}{\partial e}$ have same tangent direction at the solution point.

Lagrangian multiplier = ratio of the tangents at solution point.

$$\frac{\partial L}{\partial e} = \frac{\partial e^T \Sigma e}{\partial e} - \lambda \frac{\partial e^T e}{\partial e}$$

$$= 2 \Sigma e - 2 \lambda e \stackrel{\text{want}}{=} 0$$

$$\Rightarrow \Sigma e = \lambda e$$

\downarrow Covariance Matrix $D \times D$ \downarrow scalar

"e is a special vector for Σ
 So, Σ does not rotate e by multiplication"

- ① Σ is symmetric
- ② $\Sigma = \mathbf{x}^T \mathbf{x} \Rightarrow$ Positive Definite
 $\Rightarrow \forall$ vector \mathbf{z} , $\mathbf{z}^T \Sigma \mathbf{z} \geq 0$
- ③ Eigen vectors of Σ : e_1, e_2, \dots, e_D are orthonormal : $e_i^T e_i = \|e_i\| = 1$; $e_i^T e_j = 0$
(perpendicular)

Spectral decomposition of Σ

$$\Sigma = \begin{bmatrix} & & & \\ | & | & | & | \\ e_1 & e_2 & \dots & e_D \\ | & | & & | \end{bmatrix} \cdot \begin{bmatrix} \lambda_1 & & & \\ & \text{Diag} & & \\ & & \lambda_2 & 0 \\ & & 0 & \ddots \\ & & & & \lambda_D \end{bmatrix} \cdot \begin{bmatrix} e_1^T \\ e_2^T \\ \vdots \\ e_D^T \end{bmatrix}$$

diag - eigen vals
 sorted
 $\lambda_1, \lambda_2, \dots, \lambda_D \geq 0$

Dimensionality Reduction:

want: $R < D$ dimensions for projections.

\Rightarrow Project X onto $e_1 | e_2 | \dots | e_R \Rightarrow R$ co-ordinate per data point

Covariance Matrix of new data (principal components) is R -dimensional APPROXIMATION of the original covariance matrix.

$$\sum_{D \times D \text{ original covar}} = \begin{bmatrix} D \times R \\ e_1 | e_2 | \dots | e_R | e_{R+1} | \dots | e_D \end{bmatrix} \quad \text{ignore}$$

$$R \times R \quad \begin{bmatrix} x_1 & x_2 & \dots & x_R \end{bmatrix} \quad \begin{bmatrix} x_1 & 0 & \dots & 0 \\ x_2 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ x_R & 0 & \dots & 0 \end{bmatrix} \quad \begin{bmatrix} R \times D \\ e_1 | e_2 | \dots | e_R | e_{R+1} | \dots | e_D \end{bmatrix} \quad \text{ignore}$$

$$\sum_{\text{new } R \text{-dim data } D \times D} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ e_1 | e_2 | \dots | e_R | 1 \end{bmatrix} \begin{bmatrix} x_1 & 0 & \dots & 0 \\ x_2 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ x_R & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} -e_1 | -e_2 | \dots | -e_R | e_{R+1} | \dots | e_D \end{bmatrix}$$