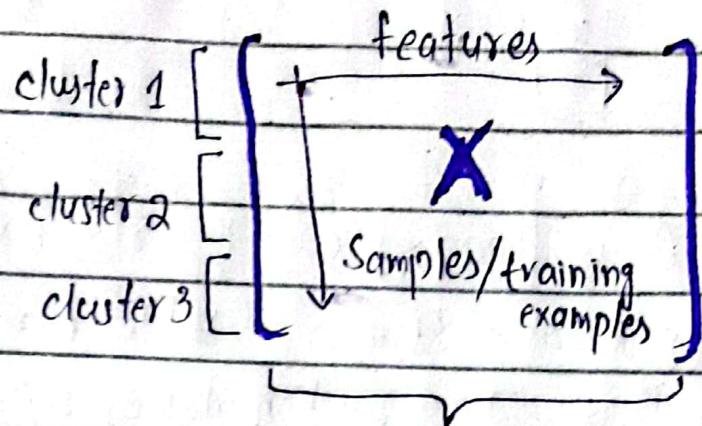


PCA

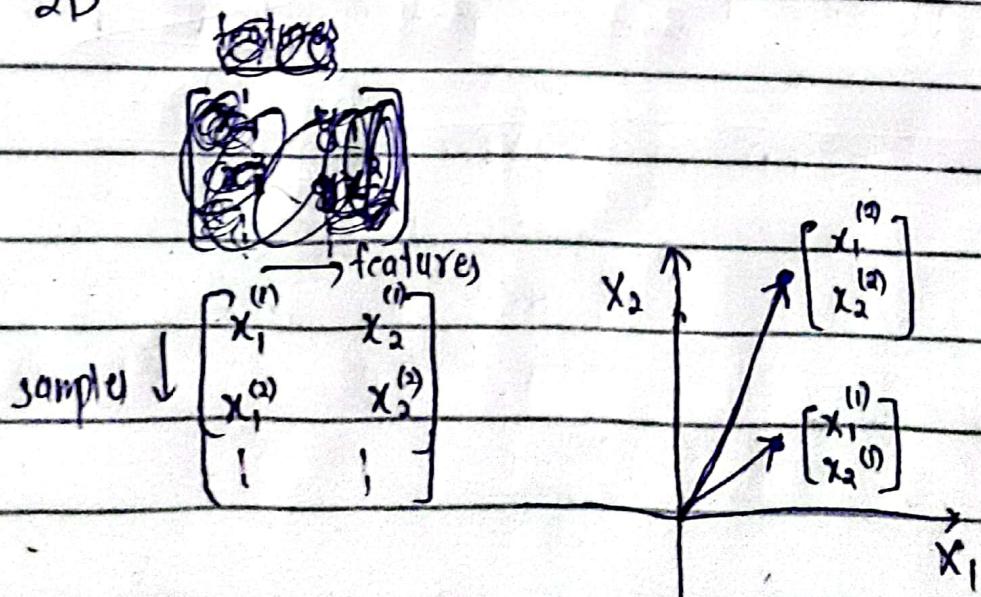
→ Method for dimensionality reduction



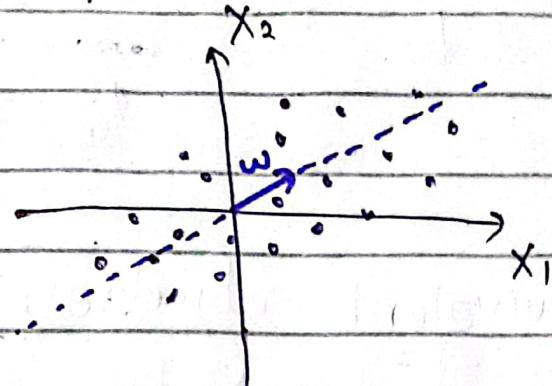
Reduce features by not discarding some features but transforming existing features to a smaller set of new features

Let's start with the example of reducing 2 dimensions to one dimension.

In 2D



training examples:
for many samples.

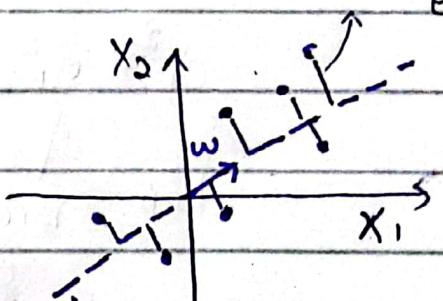


If w is a unit vector i.e $\|w\| = 1$ then
projections of sample vectors in X onto
this unit vector are given by:

$$Xw$$

How to choose w ?

→ To minimize the reconstruction error



OR

→ To maximize the variance

$$\dots \rightarrow w \rightarrow \dots$$

By maximizing variance, you ensure

variance of projected data

that variation in

original data is preserved when projected onto lower dimensional representation. If

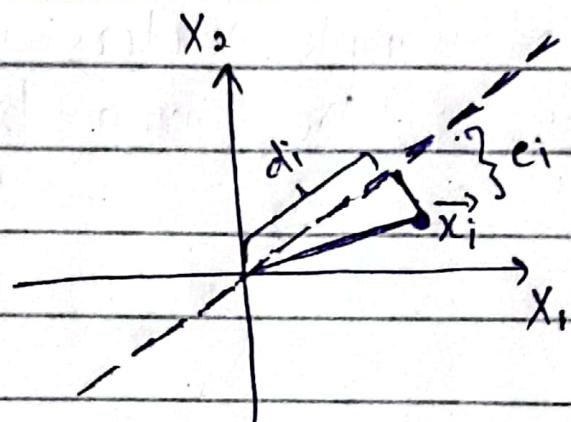
Variance is small then projected data points

are very close, losing useful structure
in original data. → ---

↑ small variance.

These are equivalent objectives:

Maximizing variance \Leftrightarrow minimizing error



$$d_i^2 + e_i^2 = \|x_i\|^2$$

projection reconstruction
error

$$\frac{1}{n} \sum_{i=1}^n d_i^2 + \frac{1}{n} \sum_{i=1}^n e_i^2 = \underbrace{\frac{1}{n} \sum_{i=1}^n \|x_i\|^2}_{\text{MSE}}$$

↑
const (doesn't depend on chosen axis)

if all features are centered (i.e. subtracted their mean)

then origin represent means of features and

$$\frac{1}{n} \sum_{i=1}^n d_i^2 = \frac{1}{n} \sum_{i=1}^n (d_i^2 - \text{mean}^2) = \text{Variance}$$

So

$$\text{Variance} + \text{Mean Squared Error} = \text{const}$$

Vari. of ~~data~~
projected
data

So either minimizing
MSE maximizes Variance
or vice versa.

So Do either one of them.

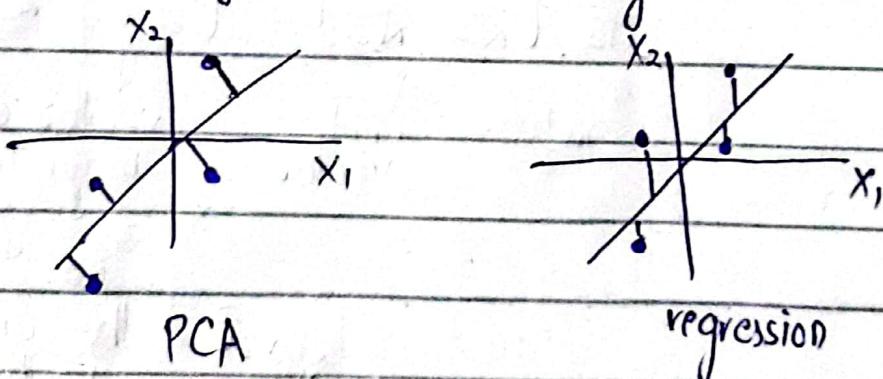
PCA vs regression:

In PCA, error is the perpendicular distance

b/w data point and projection axis while

In regression, error is the vertical distance

b/w predicted \hat{y} and actual y

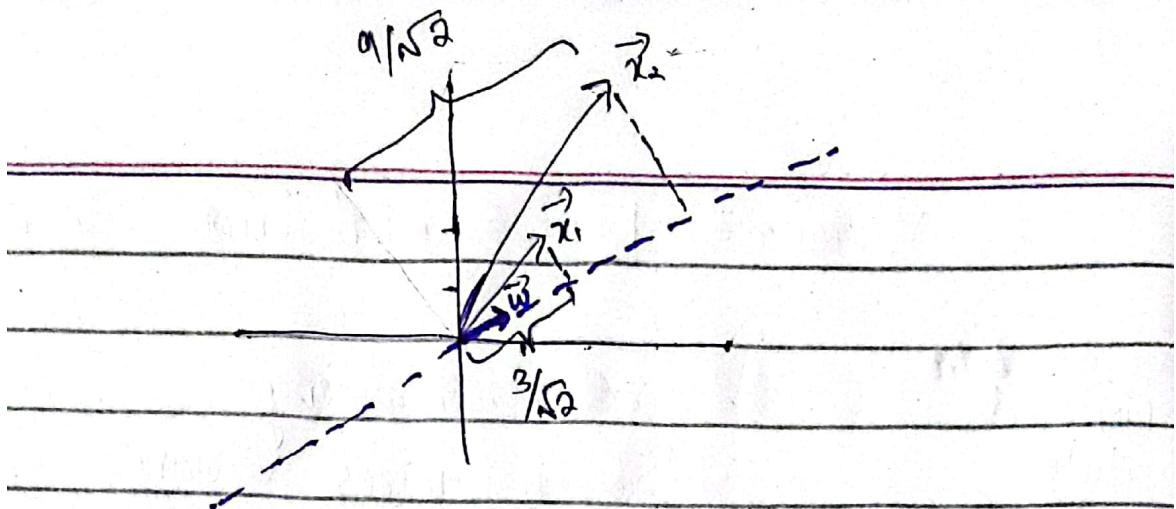


PCA loss function:

Minimizing reconstruction errors:

Xw gives projections of data vectors onto w :

$$\begin{bmatrix} 1 & 2 & 1 \\ 4 & 5 & 1 \end{bmatrix} \begin{bmatrix} \frac{3}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} = \boxed{\text{?}} \begin{bmatrix} \frac{3}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$$

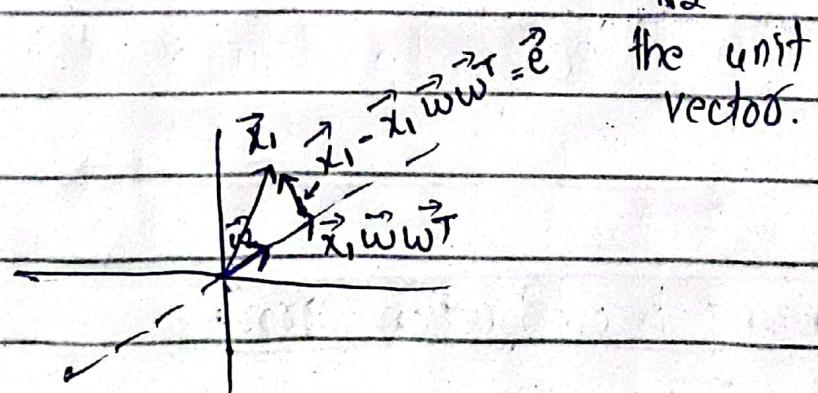


$Xw\omega^T$ will give vectors of these projections

$$\underbrace{\begin{bmatrix} \frac{3}{\sqrt{2}} \\ \frac{9}{\sqrt{2}} \end{bmatrix}}_{Xw} \underbrace{\begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}}_{\omega^T} = \underbrace{\begin{bmatrix} \frac{3}{2} & \frac{3}{2} \\ \frac{9}{2} & \frac{9}{2} \end{bmatrix}}_{Xw\omega^T}$$

$$\frac{3}{\sqrt{2}} \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} \frac{3}{2} & \frac{3}{2} \end{bmatrix}$$

\uparrow magnitude \uparrow unit vector \uparrow vector of magnitude
 $\vec{e}_i = \vec{x}_i - \vec{x}_i \vec{w} \vec{w}^T$ along the unit vector.



$$\vec{e}_i = \vec{x}_i - \vec{x}_i \vec{w} \vec{w}^T$$

$$\|\vec{e}_i\|_2 = \|\vec{x}_i - \vec{x}_i \vec{w} \vec{w}^T\|_2$$

gives length of error vector

$$\sum_i |\vec{e}_i|^2 = \sum_i \| \vec{x}_i - \vec{x}_i \vec{w} \vec{w}_T \|_2^2$$

Sum of squared errors = $\sum_i \| \vec{x}_i - \vec{x}_i \vec{w} \vec{w}_T \|_2^2$

$$= \left(\sqrt{\sum_i \| \vec{x}_i - \vec{x}_i \vec{w} \vec{w}_T \|_2^2} \right)^2$$

$$SSE = \left(\| X - X \vec{w} \vec{w}^T \|_F^2 \right)$$

$$L = \| X - X \vec{w} \vec{w}^T \|_F^2$$

ℓ_2 or
frobenius norm = ℓ_2 norm of each row vector or column vector in matrix

forbenius norm

$$= \sqrt{\sum_{i,j} |x_{ij}|}$$

sum of absolute values

$$= \sqrt{\sum_i \| r_i \|_2^2}$$

sum of square of length of row vectors

$$= \sqrt{\sum_j \| c_j \|_2^2}$$

sum of square of length of column vectors

2) Maximizing variance:

$$-L = \frac{1}{n} \underbrace{(\mathbf{X}\mathbf{w})^T (\mathbf{X}\mathbf{w})}_{\text{Sum of square of projections}}$$

Variance

$$-L = \frac{1}{n} \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}$$

if ~~\mathbf{X}~~ $\frac{1}{n} \mathbf{X}^T \mathbf{X} = \mathbf{C}$ ← covariance matrix

$$-L = \mathbf{w}^T \mathbf{C} \mathbf{w}, \text{ s.t. } \|\mathbf{w}\|^2 = 1$$

for example:

$$\mathbf{X} = \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \end{bmatrix}$$

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} x_1 & x_2 \\ y_1 & y_2 \end{bmatrix} \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \end{bmatrix}$$

→ These equations are true when 1) data is centered
with mean=0

2) \mathbf{w} is a unit vector

making $\mathbf{X}\mathbf{w}$ the projection of \mathbf{X} .

$$\vec{\mathbf{X}} \cdot \vec{\mathbf{w}} = |\vec{\mathbf{X}}| |\vec{\mathbf{w}}| \cos \theta$$

$$\vec{\mathbf{X}} \cdot \vec{\mathbf{w}} = |\vec{\mathbf{X}}| \cos \theta$$

$\vec{\mathbf{X}} \cdot \vec{\mathbf{w}}$ = projection.

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} x_1^2 + x_2^2 & x_1 y_1 + x_2 y_2 \\ x_1 y_1 + x_2 y_2 & y_1^2 + y_2^2 \end{bmatrix}$$

if mean of each feature ~~and~~ x, y is zero then

$$\frac{1}{n} \mathbf{X}^T \mathbf{X} = \begin{bmatrix} C_{xx} & C_{xy} \\ C_{xy} & C_{yy} \end{bmatrix}$$

Also otherwise:

$-L = \mathbf{w}^T \mathbf{C} \mathbf{w}$ can be increased to infinity using very large ($\rightarrow \infty$) \mathbf{w}

Solution:

Maximizing $\mathbf{w}^T \mathbf{C} \mathbf{w}$:

We will use Lagrange multiplier to solve this constrained optimization problem

$$\text{Constraint: } \|\mathbf{w}\|_2^2 = 1$$

$$\mathbf{w}^T \mathbf{w} = 1$$

$$-L = \mathbf{w}^T \mathbf{C} \mathbf{w} - \lambda (\mathbf{w}^T \mathbf{w} - 1)$$

Setting $\frac{\partial L}{\partial \mathbf{w}} = 0$,

$$\boxed{\mathbf{C} \mathbf{w} = \lambda \mathbf{w}}$$

\mathbf{w} should be the eigenvector of \mathbf{C}

To maximize variance. $\mathbf{w}^T \mathbf{C} \mathbf{w} = \lambda \mathbf{w}^T \mathbf{w} = \lambda$,

choose eigenvector with the largest eigenvalue λ .

Spectral theorem: for p features

$C = \frac{1}{n} X^T X$ is a symmetric $p \times p$ matrix.

It has p eigenvectors orthogonal to each other.

$$V = \begin{bmatrix} \vec{v}_1 & \vec{v}_2 & \dots & \vec{v}_p \end{bmatrix}$$

↑ $p \times p$ size

Projecting data on multiple eigenvectors of C :

Rotated data: XV

vector of
first data
point in
standard basis

$$\begin{bmatrix} \vec{x}_1 \\ \vdots \\ \vec{x}_n \end{bmatrix} \rightarrow \begin{bmatrix} \vec{w}_1 & \vec{w}_2 & \dots & \vec{w}_p \end{bmatrix}$$

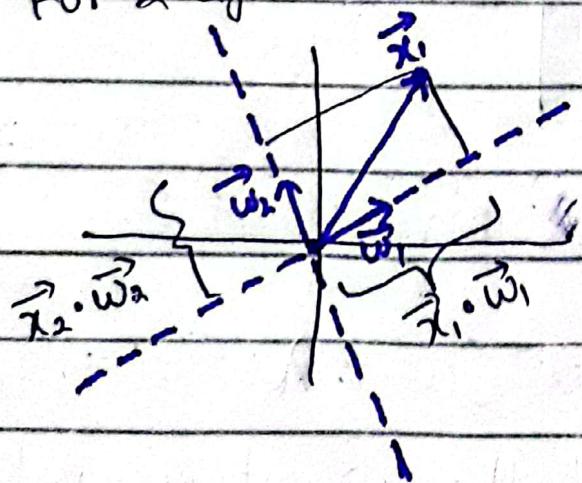
row representing

$$\begin{bmatrix} \vec{x}_1 \cdot \vec{w}_1 & \vec{x}_1 \cdot \vec{w}_2 & \dots & \vec{x}_1 \cdot \vec{w}_p \end{bmatrix}$$

vector
of first
data point

in eigenbasis.

For 2 eigenvectors:



The covariance (which was variance when using only one eigenvector) in the eigen basis is found as

Covariance after projection on eigen basis = $\frac{1}{n} (XV)^T XV$

$$= \frac{1}{n} VT X^T XV$$

$$= V^T CV$$

↑
Covariance of
original data
in standard basis

where $CV = V\Lambda V^T$

Diagonal matrix of eigenvalues

$$(I \text{ as } V \text{ is orthogonal}) \quad \Lambda = \begin{bmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \lambda_p \end{bmatrix}$$

$$= V^T V \Lambda \Rightarrow \Lambda$$

Covariance matrix after projection on eigen basis = Λ

New Data projected on eigenbasis has variance of λ_i along i th axis and covariance b/w all axis is 0

→ There is no correlation b/w the features of transformed data, avoiding any redundant feature.

As $VTCV = \Lambda$,

transformation

C in the perspective of eigen vectors of C is just scaling transformation Λ .

$$V(V^T C V)V^T = V(\Lambda)V^T$$

spectral theorem

$$C = V\Lambda V^T$$

eigen decomposition

scaling transformation Λ of

$$C = V\Lambda V^{-1}$$

for symmetric C ,

eigen vectors in the perspective of standard basis is

$$C = V\Lambda V^T$$

which is spectral theorem

just C

→ Each eigenvector in V is the principal component sorted in the descending order of λ in Λ .

Eigen decomposition is $C = V\Lambda V^{-1}$ which do not exist when V^{-1} do not exist. For any symmetric matrix, it simplifies to $C = V\Lambda V^T$ and exist. This is spectral theorem

Relationship to SVD:

SVD is given as

$$X = U S V^T$$

Diagonal matrix
of singular values

right singular
vectors of
 $X^T A \leftrightarrow$ eigenvectors
of $A^T A$

left singular
vectors of
 $X A \leftrightarrow$ eigenvectors
of $A A^T$

if X represent original data ~~then~~ which is centered then covariance matrix is

$$C = \frac{1}{n} X^T X$$

$$= \frac{1}{n} V S U^T U S V^T$$

$$= \frac{1}{n} V S^2 V^T$$

$$C = \frac{1}{n} V S^2 V^T$$

$$C = V \Lambda V^T$$

For diagonal matrix S :
 $\rightarrow S^T = S$
 $\rightarrow \underbrace{S S S \dots}_{n \text{ times}} = S^n$

This is the same form as found using eigen decomposition in previous section.

Orthogonal Diagonal Orthogonal

↓ ↓ ↓

$$C = \begin{bmatrix} \text{right singular vectors of } X \\ \text{of } X \end{bmatrix} \begin{bmatrix} \text{square of singular values of } X \text{ divided by } n \end{bmatrix} \begin{bmatrix} \text{right singular vectors of } X \end{bmatrix}^T \quad \left. \begin{array}{l} \text{From SVD} \\ \text{From eigen decomposition} \end{array} \right\}$$

which is equivalent to

$$C = \begin{bmatrix} \text{eigen vectors of } C \\ \text{of } C \end{bmatrix} \begin{bmatrix} \text{eigen values of } C \end{bmatrix} \begin{bmatrix} \text{eigen vectors of } C \end{bmatrix}^T \quad \left. \begin{array}{l} \text{From eigen decomposition} \end{array} \right\}$$

This means

eigen vectors of covariance matrix C = right singular vectors of data matrix X

eigen values of covariance matrix C = Square of singular values of data matrix X divided by sample size n

→ Principal components of data X is its right singular vectors with variance of square of their singular values divided by n .

Total variance:

We know from eigendecomposition of symmetric matrix \mathbf{A} (spectral theorem) that

$$\mathbf{A} = \mathbf{V}\Lambda\mathbf{V}^T$$

or

$$\text{transformation } \rightarrow \Lambda = \mathbf{V}^T \mathbf{A} \mathbf{V}$$

\mathbf{A} in the
perspective of
eigen basis

↑
eigen values
of \mathbf{A}

trace of matrix is sum of its
diagonal entries

From a property,

$$\text{tr}(\mathbf{A}) = \text{tr}(\Lambda)$$

For covariance of original data \mathbf{C} and covariance of transformed data (to its eigen basis) Λ

$$\text{tr}(\mathbf{C}) = \text{tr}(\Lambda)$$

↑
eigen values of \mathbf{C}

↔
Convariance of transformed data

→

Sum of variance of all the original features = Sum of variance explained by all the principal components

→ Explained variance ratio by PC i is $\lambda_i / \text{tr}(C)$

→ Using PCA, we reconstruct the data in few new axes, called principal components, which explain most of the variance while remaining axes are discarded as representing noise

→ Advantage of PCA:

1) Minimizing noise by discarding PCs with small variance

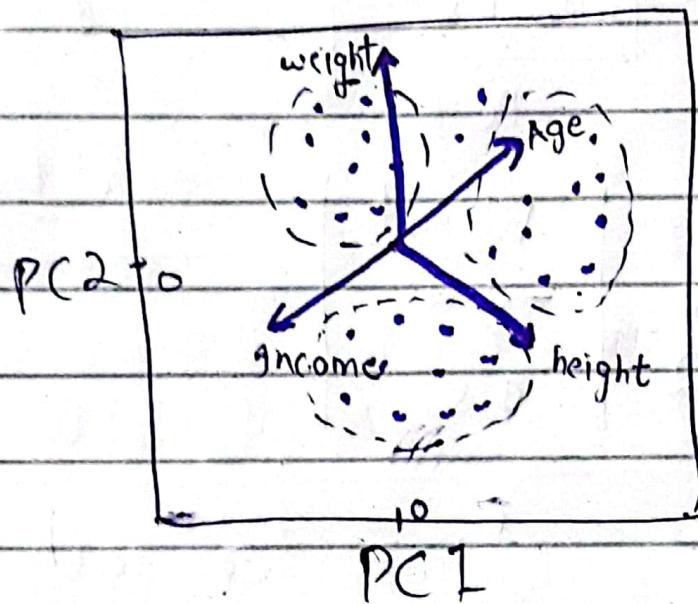
2) Removing redundant features by creating new uncorrelated features

PCA is used for:

- 1) Data Exploration
- 2) Data preprocessing

PCA for data exploration:

Biplot is created with datapoints in the coordinates of PC1 and PC2.



It contains feature vectors. The angle b/w the feature vectors represent the correlation b/w them.

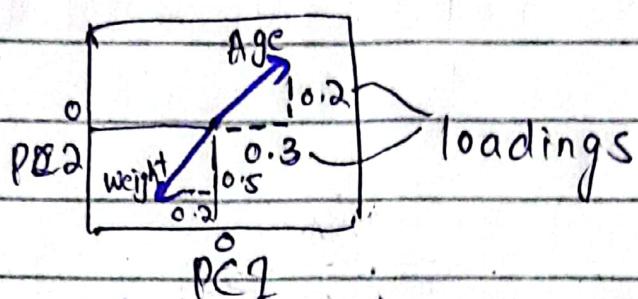
$90^\circ \rightarrow$ No correlation

greater than $90^\circ \rightarrow$ Negative correlation

less than $90^\circ \rightarrow$ Positive correlation.

Most importantly, the projection of each

feature vector on PC1 and PC2 axis gives the correlation of that feature with PC1 and PC2 respectively. These are called **loadings** for that feature.



Components:-

	Age	Weight	other features
PC1	0.3	-0.2	-----
PC2	0.2	0.5	-----

principal
components
vectors

loadings for each
feature

Intuition:

When projecting data point (x, y) on $\bar{PC} : [v_1 \ v_2]^T$, we have projection of

$$\text{projection on } \bar{PC} = x v_1 + y v_2$$

v_1 and v_2 , which are coordinates of PC vector, act as weights of feature x and y for the projection.

These weights are loadings and show correlation b/w feature and PC.

For example if

$$v_1 = 1$$

then $v_2 = 0$ as $\|\vec{PC}\|_2 = 1$

and

projection = x
on \vec{PC}

$$\vec{PC} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

which is along
the x axis.

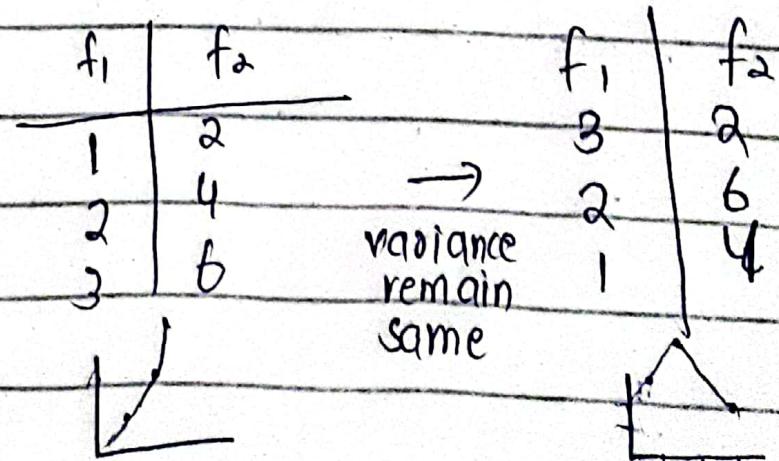
So PC and x
feature are highly
correlated.

→ Max value of loading is 1 ~~for~~ ^{when} using
unit eigenvector representing PC

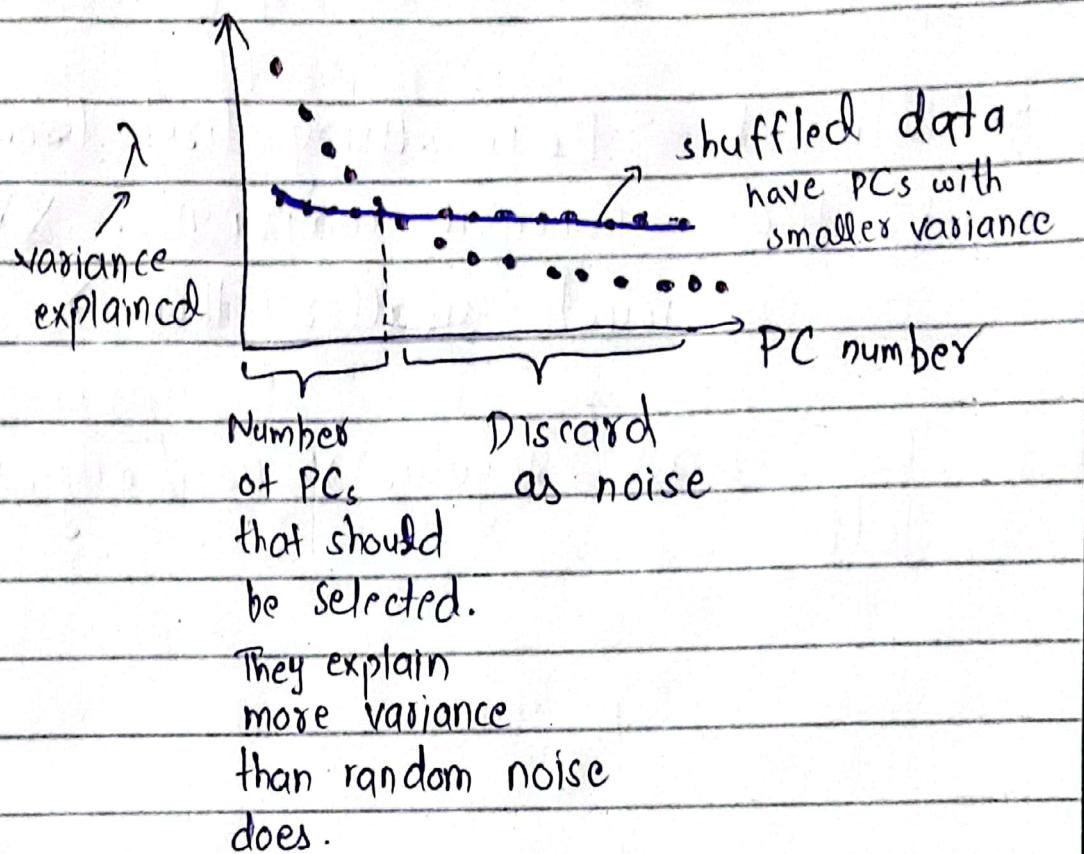
→ Before PCA, standardize the data to have same unit variance for each feature otherwise the features with high variance will dominate the PC. These will be selected as PC; the axes explaining max variance.

How many PCs to choose?

- 1) Heuristic of choosing number of PCs that will capture 90% of total variance.
- 2) Shuffling the features
It will remove any existing relationships b/w features and make data a random noise while preserving the variance of each feature.



- After shuffling
- Find the eigenvalues and compare them with eigenvalues of original unshuffled data.



PCA for preprocessing:

Reduce X with p features to a small number k of PCs using XV_k where V_k is a $p \times k$ matrix of unit norm eigenvectors with the largest eigenvalues. Then use XV_k for downstream processing.

$$\begin{bmatrix} X \\ n \times p \end{bmatrix} \begin{bmatrix} V_k \\ p \times k \end{bmatrix} = \begin{bmatrix} X V_k \\ n \times k \end{bmatrix} \text{ where } k \leq p$$

Advantage of XV_k over X :

- 1) XV_k has correlation b/w k features be equal to 0.
- 2) with $k \ll p$, i.e. using first few significant principal components, XV_k size is much smaller than X .
- 3) No small singular values / eigenvalues left

Summary of PCA

PCA is a dimensionality reduction technique used to transform data containing correlated features and noise to a smaller set of uncorrelated features, called principal components, which capture most of the variance in the original data while discarding remaining variance as caused by noise. PCA is used for data exploration and preprocessing.

- These principal components are found to be the eigenvectors of covariance matrix C of original data, sorted in the descending order of their eigenvalues.
- First principal component captures the maximum variance, the second principal component captures the second most, and so on.
- The covariance matrix of projected data on these eigenvectors (principal components) will be a diagonal matrix Λ containing eigenvalues in diagonal.

This shows that: 1) Principal components are uncorrelated due to 0s in off diagonal. 2) Variance explained by a principal component i is $\lambda_i / \text{eigen value}$

• Principal

- In relationship to SVD, principal components for data X is the right singular vectors of X and their variance is square of singular values of X divided by sample size n .
- From trace property; $\text{tr}(C) = \text{tr}(\Lambda)$ where Λ is the diagonal matrix of eigenvalues of covariance C of original data, which is also the covariance matrix of transformed data. It shows that sum of variance of all the original features is equal to the sum of variance explained by all the principal components.
- Explained variance ratio by PC i is $\lambda_i / \text{tr}(C)$

- Biplot is used to visualize the projection of each vector, representing a training example in n feature dimensions, onto the PC1 and PC2. The feature vectors are plotted with dimension of loadings. The loading of a feature f for PC_i is the weight given to that feature in finding the projections of training examples on PC_i . It shows the correlation b/w the features and principal components. The angle b/w feature vectors shows the correlation b/w them. greater than 90° , equal to 90° and less than 90° angle corresponds to negative, zero and positive correlation respectively.
- PCA can be used to preprocess data to a ~~unseen~~ transformed data of uncorrelated features and much smaller dimensionality.