

## Understanding Principal Component Analysis

### Table Of Contents

1. Idea Behind PCA
2. What are Principal Components
3. Eigen Decomposition Approach
4. Singular Value Decomposition Approach
5. Why do we maximize Variance
6. What is Explained Variance Ratio
7. How to select optimal no.of Prinicpal Components
8. Understanding Scree plot
9. Issues with PCA
10. Understanding Kernel PCA

## \* Principal Component Analysis

- The most popular dimensionality reduction technique is the principal component analysis, it is an unsupervised algorithm used for feature selection.

### Idea behind PCA

- PCA transforms and fits the data from a higher dimensional space to a lower dimensional subspace, this results into an entirely new coordinate system or the points where the first axis explains the most variance in the data (first principle component), this is the idea behind PCA.

### What are Principal Components

- Principal components are the axis, that explain the maximum variance in the data, the first principal component explains the most variance, the second a bit less and so on..
- Each new axis / Principal component found using PCA is a linear combination of the original features
- All the principal components are uncorrelated and orthogonal to each other.



## Math Behind PCA

→ Principal component analysis can be implemented in two approaches, each with its unique focus and method. They are as follows:

① Eigen decomposition method

② singular value decomposition method,  
lets derive both them

① Eigen decomposition Method

→ Eigen decomposition is a process that decomposes a square matrix  $A$  into a set of eigenvectors and eigenvalues.

→ For a given square matrix  $A$ , the eigenvectors are the non-zero vectors that satisfy the equation,

$$A = \lambda V$$

where,  $V$  is the eigen vector and  $\lambda$  is the corresponding eigenvalue.

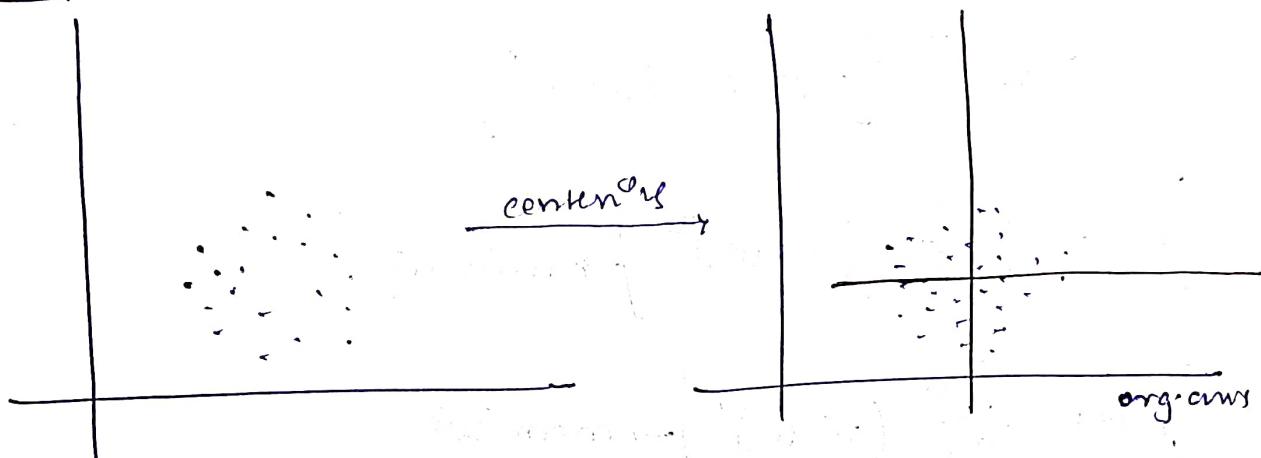
Steps for Eigen-decomposition approach

- ① Standardize the columns of  $A$ , so that each feature has a mean of zero, this ensures that each feature contributes equally to the analysis.

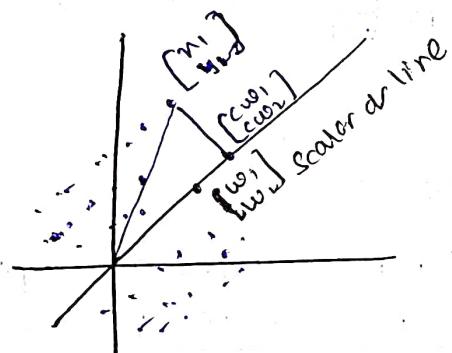
- ② compute co-variance matrix from the normalized data,  $\Sigma = A^T A / (m-1)$
- ③ Perform eigen decomposition of  $\Sigma$  (covariance matrix), to obtain eigenvalues & eigenvectors.
- ④ choose the top  $K$  eigen vectors associated with the largest eigenvalues to construct the principal components as  $A X_K$ ,  $X_K$   $\rightarrow$   $K$ th eigen vector.

### PROOF

Step 1



we have,  $D = \{n_1, n_2, \dots, n_m\}$   $n_i \in \mathbb{R}^d$



the projection point will be,  $(n^T w)w$ , considering  $w$  is a scalar or the best fit line/ the line that has least reconstruction error.

$$\text{Error}(\text{line}, \text{dataset}) = \sum_{i=1}^m \text{error}(\text{line}, n_i)$$

$$= \sum_{i=1}^m \|n_i - (n^T w)w\|^2$$

$$F(\omega) = \frac{1}{n} \sum_{i=1}^n \| n - (n^T \omega) \cdot \omega \|^2$$

This is the mean of the residuals,  
we can think of this as the avg. error.

$$= \frac{1}{n} \sum_{i=1}^n \left[ (n - (n^T \omega) \cdot \omega)^T \cdot (n - (n^T \omega) \cdot \omega) \right]$$

$$= \frac{1}{n} \sum_{i=1}^n \left[ n^T n_i - (n^T \omega)^2 - (n^T \omega)^2 + (n^T \omega)^2 \right]$$

$$= \frac{1}{n} \sum_{i=1}^n \left[ n_i^T n_i - (n^T \omega)^2 \right]$$

$$F(\omega) = \frac{1}{n} \sum_{i=1}^n - (n^T \omega)^2 \rightarrow \min \omega$$

$$F(\omega) = \frac{1}{n} \sum_{i=1}^n (n^T \omega)^2 \rightarrow \max \omega$$

$$F(\omega) = \frac{1}{n} \sum_{i=1}^n (n^T \omega) (n^T \omega) = \frac{1}{n} \sum_{i=1}^n \omega^T (n_i \cdot n_i^T) \omega$$

$$F(\omega) = \omega^T \left( \frac{1}{n} \sum_{i=1}^n n_i \cdot n_i^T \right) \omega$$

we know  
this is covariance matrix

Each column corresponds  
to one data point.

$$\therefore \max(\omega) \equiv \omega^T C \cdot \omega$$

\* where,  $C$  - covariance matrix  
 $\omega$  - eigen vector  
 corresponding to  
 max eigen value  
 of  $C$ ,

→ so, basically we can say that, the line which represents the data points or minimum reconstruction error is same as the line which minimizes the  $w^T C w$ , as  $C$ - covariance matrix,  
 and the line is given by eigen vector corresponding to max. eigen value of  $C$ ,  
 and this line is the principal component!!

★ First principal component ( $w_1$ ) =  $\arg \max_w (w^T C w)$

→ now, any line which minimizes the sum of errors/residuals, must also be orthogonal to  $w_1$ .

Hence, second principal component will be orthogonal to the first  $[w_2^T w_1 = 0]$  ★

→ By, continuing this procedure,  
 we get  $\{w_1, w_2 \dots w_d\}$

such that

$$\begin{aligned} \|w_k\| &= 1 \quad \forall k \text{ and} \\ w_i^T w_j &= 0 \quad \forall i \neq j \end{aligned}$$

Residuals after d rounds

$$\forall i \rightarrow n_i - [(n_i^T w_1) w_1 + (n_i^T w_2) w_2 + \dots + (n_i^T w_d) w_d] = r_i$$

$$\forall i \rightarrow n_i = (n_i^T w_1) w_1 + (n_i^T w_2) w_2 + \dots + (n_i^T w_d) w_d$$

→ If original dataset is of  $(d \times n)$ , then based on above  
 on, after  $K$  rounds of PCA, we get  $(d \times K) + (K \times n)$

Eg

+ suppose we have an initial dataset with 50 features and 100 samples.  $d=50$ ,  $n=100$ .

Data dimension =  $50 \times 100 = 1500$  data points to represent entire dataset.

After 5 PCA rounds,  $k=5$

Data dimension =  $(50 \times 5) + (5 \times 100) = 750$

So, now we only need half of datapoints to represent entire dataset.

+  $\{w_1, w_2, w_3, w_4, w_5\} \rightarrow$  Representation,

Common for  
dataset

+  $n_i \rightarrow [n_i^T w_1 \ n_i^T w_2 \ n_i^T w_3 \ n_i^T w_4 \ n_i^T w_5] \rightarrow$

Data  
Point  
Specific

+ Hence, to reconstruct  $n_i$  we don't need 1500 numbers like naive way features

[we only need 5 principal components]

+ And through Eigen decomposition approach, we calculate the covariance matrix, and find its eigenvalues & eigenvectors. And define the principal components from eigenvectors. we can reconstruct by

$A \times K$

## ② singular value decomposition method

→ SVD method factorizes the dataset in such a way that it becomes a product of the multiplication of 3 individual matrices.

$$X(\text{original data}) = U * \Sigma * V^T$$

where,  $U$  is the matrix that contains the principal components.

### Steps in SVD approach

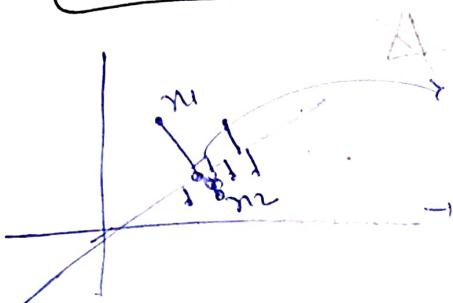
- ① Standardize columns of  $A$ , to ensure each feature contributes equally to analysis.
- ② Compute SVD of the centered data matrix to obtain matrices  $U$ ,  $\Sigma$ , and  $V^T$ , such that,  
$$X = U \Sigma V^T$$
 where  $U$  contains left singular vectors,  $\Sigma$  holds the singular values &  $V^T$  includes the right singular vectors.
- ③ Select principal components: The principal components are obtained by selecting the columns for  $U$  corresponding to the largest singular values in  $\Sigma$ , three principal components represent the directions of maximum variance in the data etc.

## Note:-

→ what does high variance mean?

when making a proxy,

what if two points get same proxy.



Now, How do we distinguish  
 $m_1 \& m_2$ .

→ And How do we reconstruct  
 $m_1 \& m_2$ .

→ So, If all the datapoints are crowded,  
then we can't reconstruct and the  
information is lost.

This happens due to less variance / small variance.

Hence, we want directions where projections not crowded.

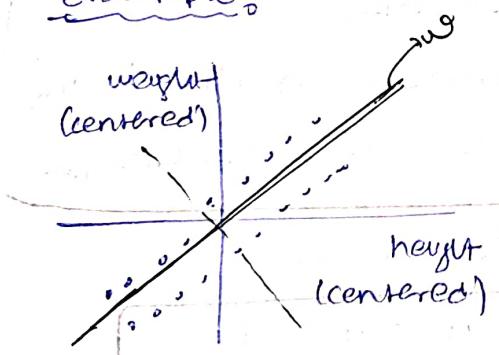
~~✓~~ Variance has to be high,

~~✓~~ so the datapoints will be spread out

and reconstruction of points can be done

properly.

## Example:



- height gives some  
information weight

→ we have two directions now  
 $w_1, w_2$  ..

& And these new directions are  
de-correlated, every direction is  
giving a new information  
from the other.

∴ Hence, we can say these  
directions are **de-correlated**.



thus's what the algorithm is  
generally doing. & it's called  
Principal Component Analysis.

\* Once we get principal components, project the original data onto the subspace spanned by the selected principal components, effectively reducing the dimensionality of the data.

(\*) What is the optimal no. of PCs needed?

\* Before we ask this qn, we need to ask how much of the information in a given dataset is lost by projecting the observations onto first few principal components, i.e,

the proportion of variance explained by each principal component

Explained variance ratio's variance by m<sup>th</sup> PC

total variance

$$\text{EVR} = \sum_{j=1}^n \left( \sum_{i=1}^p w_{ij} z_{ij} \right)^2$$

$$= \frac{\sum_{j=1}^n \sum_{i=1}^p w_{ij}^2 z_{ij}^2}{\sum_{j=1}^n \sum_{i=1}^p z_{ij}^2}$$

→ EVR for PC1 = Distance of PC1 point to centroid / Distance of PC1 point to PC1 + PC2 + ... + PCn

If EVR of PC1 = 0.91, then that means PC1 explains 91% of the variance of data.

→ Now that we have a metrics to get the variance explained, we can choose the no. of dimensions that add up to a sufficiently large proportion of the variance (eg- as %).

Eg:-

$$\text{EV}R \text{ or } PC1 = 0.88$$

$$\text{EV}R \text{ or } PC2 = 0.10$$

$$\text{EV}R \text{ or } PC3 = 0.05$$

→ Now, if we want to preserve (99%) variance, then we can just take PC1 & PC2, as they sum up to 0.94, 94%. 99%.

→ instead if we want to preserve (99.5%) of variance, then we should take PC1, PC2 & PC3. as the sum up to 0.99, 99.5%. 99.5%.

### Scree plots

→ Scree plots are the graphs that convey how much variance is explained by corresponding principal components.

→ So, by eyeballing the screeplot and looking for a point at which the proportion of variance explained by each subsequent PC drops off. This is often referred as



elbow in screeplots.

## Issue/concern with PCA

1 Time complexity

→ in particular,  
covariance matrix.

→ usually, it takes

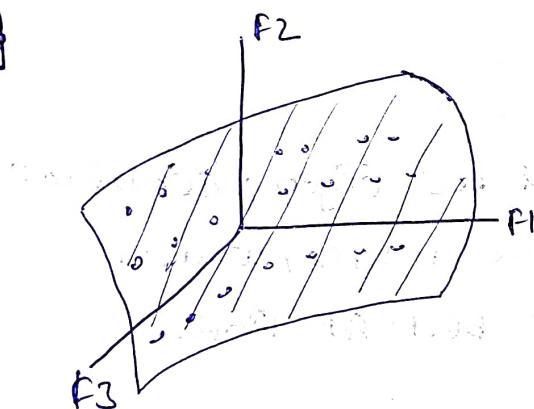
$$O(d^3)$$

Finding the eigen vectors  
and eigen values of  
covariance matrix is  
the time consuming step.

→ if no. of features were large, [like 32000 for  
image] then  $d^3$  is going to be  
very very large.

Eg Face recognition (Eigen faces)

2.

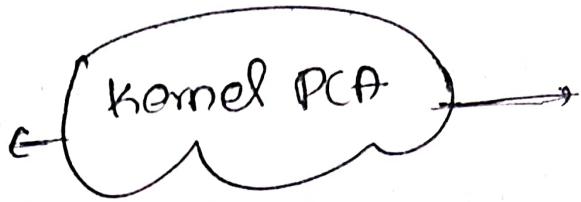


Data may not necessarily live in a low-dimensional linear subspace

→ features will not be linearly related.

### Notes

- we will have same solutions for both the above issue, solving one will solve the other.



• Input —  $\{n_1, \dots, n_n\}$   $n_i \in \mathbb{R}^d$ ,

kernel  $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$

Step ① compute  $K \in \mathbb{R}^{n \times n}$ ,

where  $K_{ij} = K(n_i, n_j)$   $\forall i, j$ .

And then we need to center the kernel

Step ② compute  $\beta_1, \dots, \beta_d$  eigen vectors &

$n_{1,1}, \dots, n_{d,d}$  eigen values.

Step and normalize to get,

$$\alpha_k = \frac{\beta_k}{\sqrt{n} \beta_k}$$

Step ③

$$w_k = \phi(n) \cdot \alpha_k$$

we cannot reconstruct the eigen vectors or the covariance matrix

Eigen computation of kernel matrix.  
Hence, thus defeats the purpose,  
as it needs  $\phi(n)$ , which we  
tried to avoid all long!

→ But we can still compute the compressed representation.

$$\phi(n_i)^T \cdot w_k = \phi(n_i)^T \cdot \left( \sum_{j=1}^n \phi(n_j) \alpha_{Bj} \right) = \sum_{j=1}^n \alpha_{Bj} \phi(n_i)^T \phi(n_j)$$

$$\phi(m_i)^T w_k = \sum_{j=1}^n \alpha_{kj} \cdot k_{pj}$$

→ for downstream task,  
this is usually  
good enough.

→ typically, the reason why we do kernel PCA  
is to get these projections and throw away  
original points and only retain these  
projections along each of these data points.

\* This is where dimensionality reduction happens.

(eg)

original datapoint was 100 dimensions,  
but then you only care about top 5.

→ So we map the 100 dimensional data into  
a 5 dimensional projection onto these  
5 most important directions.

→ that becomes a 5 dimensional representation  
of your 100 dimensional data

→ once you learn these representations  
in a low dimensional space, as put the  
projections onto these eigen directions  
you can use these projection values  
at your data point and work this for a  
downstream task (supervised learning task).

$$(\phi(m_i)^T w_k)^2 = (\phi(m_i) \cdot \frac{w_k}{\|w_k\|})^2 \cdot \frac{1}{\|w_k\|^2} = \text{proj}_{w_k}(m_i)$$

modified

Step (3)

compute  $\sum_{j=1}^n \alpha_{kj} \cdot k_{ij} + b$

$$x \in \mathbb{R}^d \rightarrow \left[ \sum_{j=1}^n \alpha_{1j} \cdot k_{1j}, \sum_{j=1}^n \alpha_{2j} \cdot k_{2j}, \dots, \sum_{j=1}^n \alpha_{Lj} \cdot k_{Lj} \right]$$

→ It might increase the dimension,  
nevertheless it will allow you to capture  
more complicated relationships.

Note

→ After step 2, we have to center the kernel.

→ Given,  $k \in \mathbb{R}^{n \times n}$

$$k_{ij} = k(n_i, n_j) + \epsilon_{ij}$$

create a new kernel  $k^c$  — centered.

$$k_{ij}^c = k_{ij} - \theta_i l_j^T - l_i \theta_j^T + p_i l_j$$

where,

$$\theta_i = \frac{1}{n} \sum_{k=1}^n k_i \cdot k_p$$

$$p_i = \frac{1}{n} \sum_{p,j} k_{ij}$$

→ Now, we can say we are able to  
solve both the issues by using kernel PCA.