

Principle Component Analysis (PCA)

DATE: _____

- Our discussions of unsupervised learning have focused on clustering techniques; which seek to "discover" labels on features data that has no historical labels -
- We will now shift toward unsupervised algorithms that focus on dimension reduction.

Motivation behind Dimension Reduction:

- Imagine a dataset with 30+ features how would you understand the key features?
- Visualization and Data Analysis have limitation when the number of features dimension increases.
- Dimension Reduction outcome:
 - Understand which features describe the most variance in the dataset.
 - Aid human understanding of large features sets, especially through visualization.

- Important Note:

- Dimensionality Reduction algorithm such as PCA do not simply choose a subset of the existing features-
- They create new dimensional components that are combination of proportions of the existing features-

• Section Overview

- Theory and Intuition of PCA.
- Manually create PCA Algorithm
- Utilize scikit learn to perform PCA-
- PCA Exercise project overview
- PLA Exercise Solution.

PCA - Model 1

Theory and Intuition

part - One

- 1901 - Karl Pearson publishes "On lines and planes of closest fit to systems of points in Space" based on the principal axis theorem in the field of geometry -
 - Pearson was the protege of Francis Galton and the Pearson correlation coefficient is named after it.
- 1933 - American mathematician and economist Harold Hotelling independently develops and names principal component analysis in the publication, "Analysis of a complex statistical variables into principal components".
- Hotelling's paper perfectly describes the purpose of PCA:
 - Analyzing a complex set of variable into principle components.
 - Let review the motivation behind and basic idea behind PCA

- principal components analysis outcomes:

- Reduce the number of data -
- Show which features explain the most variance in the data

- Dimension Reduction

- Helps visualize and understand complex dataset.
- can also act as a simpler data set for training data for machine learning algorithms
 - Reduce dimensions then train ML algorithm on smaller dataset
 - Help reduce N features to a desired smaller set of components through a transformation -

- It does not simply select subset of features

- Variance Explained

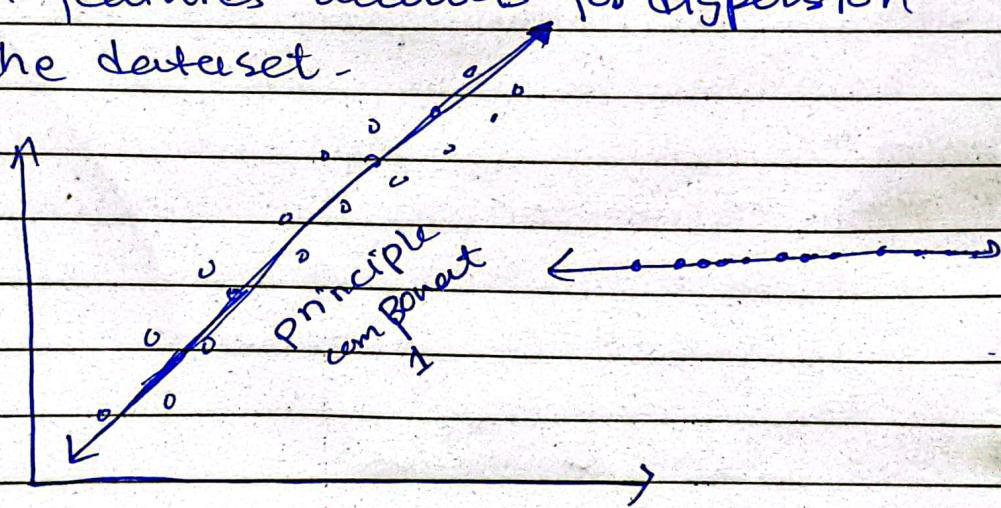
- We have often seen that certain features are more important or have more explanatory power than other features -
 - for example, size of house is probably much more important than the color of a house when explaining the price of house -

The idea of more important features is easy to understand when we can directly correlate feature to a known label. But what about unlabeled data?

- What measurement ^{we} can use to determine feature "importance"?

variance

Measure the proportion to which each features accounts for dispersion in the dataset.



- principle components is a linear combination of original features account for the more influence it has over the principle components-

- Here we went from 2 features down to 1 principal component-

- This single principal component can "explain" some percentage of the original data, for example 90% of the variance explained by

principal component-

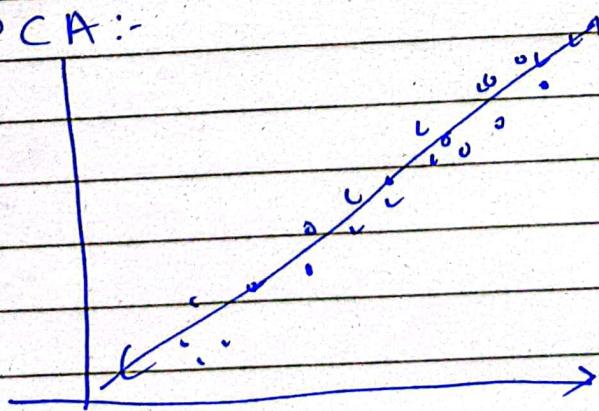
- 100% variance in the data is explained by all the original features.
- We trade off some of the explained variance for less dimensions.
- This can be significant saving for data set ~~by~~ with many dimensions, but only a few strong features.

Theory and Intuition part - two

- principle Component Analysis operates by creating a new set of dimensions (the principal components) that are normalized linear combinations of the original features.

$$z_1 = \phi_{11}x_1 + \phi_{12}x_2 + \dots + \phi_{1p}x_p$$

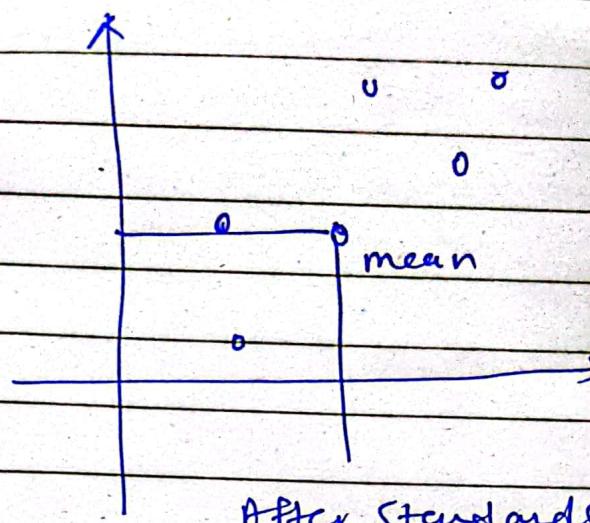
PCA:-



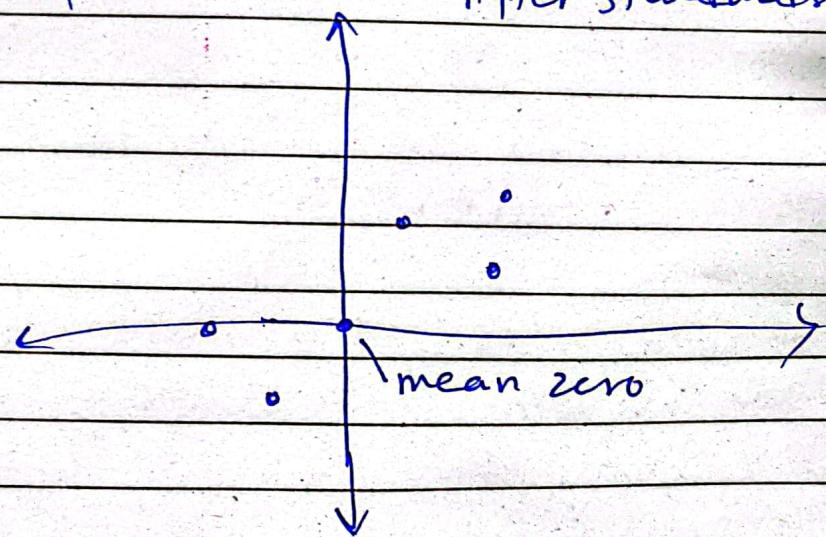
$$z_1 = \phi_{11}x_1 + \phi_{12}x_2$$

Begin with two dimension dataset

x_1	x_2
1	2
2	1
3	2
4	4
5	3
6	5



After Standardization



- 1 step standardize the data
- calculate the covariance matrix data:

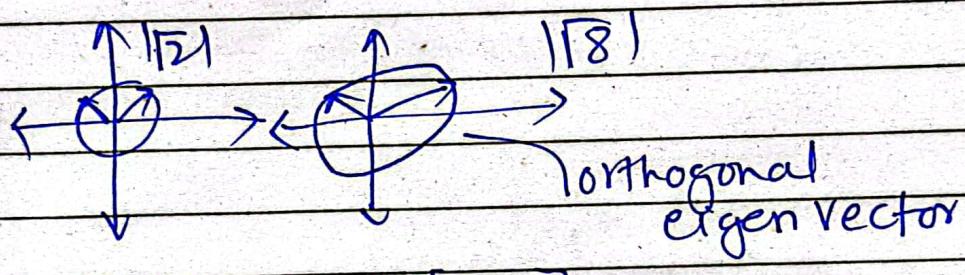
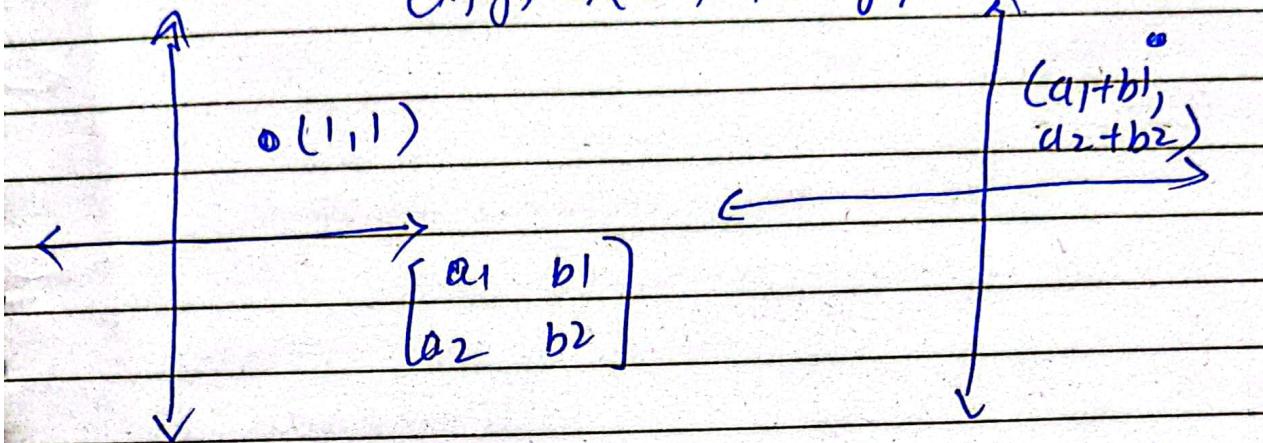
Linear Transformation

$(x, y) \rightarrow (ax + by, cx + dy)$

$\begin{pmatrix} a & b \\ c & d \end{pmatrix}$
 $\begin{pmatrix} x \\ y \end{pmatrix} \rightarrow \begin{pmatrix} ax + by \\ cx + dy \end{pmatrix}$

- Linear Transformation of Data:

$$(x, y) \rightarrow (a_1 x + b_1 y, a_2 x + b_2 y)$$

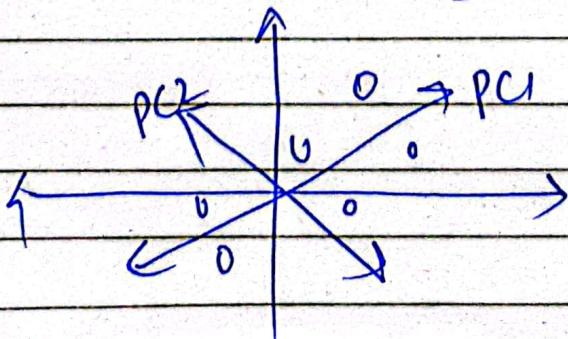


$$\begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \begin{bmatrix} 2 \\ 2 \end{bmatrix} \quad \circ \text{ Eigen}$$

$$\begin{bmatrix} a_1 & b_1 \\ a_2 & b_2 \end{bmatrix}$$

- Apply Linear Transformation:

$$\begin{bmatrix} \text{var}(X) & \text{cov}(X, Y) \\ \text{cov}(X, Y) & \text{var}(Y) \end{bmatrix}$$



- EigenValue measure variance explained:

- PCA Steps

- Get original data
- Calculate Covariance Matrix
- Calculate EigenVector
- sort EigenVectors by EigenValues
- Choose N largest EigenValues
- project original data onto Eigen vector

PCA - Manual Implementation in python

```
df = pd.read_csv('cancer_tumors')
```

~~df~~
scaled_X = scaled · fit · transform(df)
scaled_X = scaled_X · mean(axis=0)

covariance_matrix = np.cov(scaled_X, rowvar=False)

eigen_values, eigen_vectors = np.linalg.eig(covariance_matrix)

N W F S

DATE: _____

PCA N features Spaces \rightarrow N PC space

num_components = 2

sorted_idx = np.argsort(eigen_values)[::-1]

eigen_value[eigen_idx] = eigen_values[sorted_idx]
[^{num of}
^{combined}]

eigen_vectors

[:, sorted_idx]

principal_cmp = np.dot(scaled_x, eigen_vectors)

principle_components

plt.scatterplot(principal_components[:, 0],
principal_components[:, 1])

Scikit-learn Implementation

df = pd.read_csv(" " "

scaler = scaler.fit_transform()

from sklearn.decomposition import PCA

help(PCA)

pca_model = PCA(n_components=2)

pca_model.fit(scaled_x)

new_model.transform(scaled_x)

M T W T F S

DATE: _____

pc = pca_model.fit_transform(scaled_X)

plt.scatter(pc_result[:, 0], pc_result[:, 1])

C = cancer.dictionary
{target}

pca_model.components_

df_temp = pd.DataFrame(pca_model.components_,
index=['PC1', 'PC2'],
columns=df.columns)

pca_model.explained_variance_ratio

pca_30 = PCA(n_components=30)

pca_30.fit(scaled_X)

pca_30.explained_variance_ratio

explained_variance = []

for n in range(1, 30):

pca = PCA(n_components=n)

pca.fit(scaled_X)

Model deployment

REVIEWERS

DATE: _____

Section Overview

- General theory and concepts
- Model persistence Basic
- Model deployment as an API

- General theory and concept

- When is a model ready for deployment?
- How often to retain your model?
- When to revisit model algorithm choice and assumption.

Model Deployment Consideration

- choosing a model
- purpose of deployment
- performance expectation
- Retraining Intervals