# Principal Component Analysis

**Tushar B. Kute,**
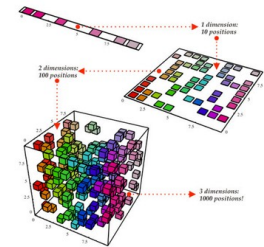http://tusharkute.com

# Dimensionality Reduction

- Dimensionality reduction or dimension reduction is the process of reducing the number of random variables under consideration by obtaining a set of principal variables.

- It can be divided into feature selection and feature extraction.

  – Feature selection approaches try to find a subset of the original variables (also called features or attributes).

  – Feature projection or Feature extraction transforms the data in the high-dimensional space to a space of fewer dimensions.

# Large Dimensions

- Large number of features in the dataset is one of the factors that affect both the training time as well as accuracy of machine learning models. You have different options to deal with huge number of features in a dataset.

  - Try to train the models on original number of features, which take days or weeks if the number of features is too high.

  - Reduce the number of variables by merging correlated variables.

  - Extract the most important features from the dataset that are responsible for maximum variance in the output. Different statistical techniques are used for this purpose e.g. linear discriminant analysis, factor analysis, and principal component analysis.

# Principal Component Analysis

- Principal component analysis, or PCA, is a statistical technique to convert high dimensional data to low dimensional data by selecting the most important features that capture maximum information about the dataset.

- The features are selected on the basis of variance that they cause in the output.

- The feature that causes highest variance is the first principal component. The feature that is responsible for second highest variance is considered the second principal component, and so on.

- It is important to mention that principal components do not have any correlation with each other.

# Steps in PCA

- Standardization

- Covariance Matrix Computation

- Computer Eigen vector and eigen values

- Feature vector

- Recast the Data Along the Principal Components Axes

- The aim of this step is to standardize the range of the continuous initial variables so that each one of them contributes equally to the analysis.

- More specifically, the reason why it is critical to perform standardization prior to PCA, is that the latter is quite sensitive regarding the variances of the initial variables.

- That is, if there are large differences between the ranges of initial variables, those variables with larger ranges will dominate over those with small ranges

- Mathematically, this can be done by subtracting the mean and dividing by the standard deviation for each value of each variable.

$$z = \frac{value - mean}{standard\ deviation}$$

- Once the standardization is done, all the variables will be transformed to the same scale.

# Covariance Matrix Computation

- The aim of this step is to understand how the variables of the input data set are varying from the mean with respect to each other, or in other words, to see if there is any relationship between them.

- Because sometimes, variables are highly correlated in such a way that they contain redundant information.

- So, in order to identify these correlations, we compute the covariance matrix.

tusharkute
.com

# Covariance Matrix Computation

$$\begin{bmatrix} Cov(x,x) & Cov(x,y) & Cov(x,z) \\ Cov(y,x) & Cov(y,y) & Cov(y,z) \\ Cov(z,x) & Cov(z,y) & Cov(z,z) \end{bmatrix}$$

# Covariance Matrix Computation

- What do the covariances that we have as entries of the matrix tell us about the correlations between the variables?

- It's actually the sign of the covariance that matters :

  - if positive then : the two variables increase or decrease together (correlated)

  - if negative then : One increases when the other decreases (Inversely correlated)

- Now, that we know that the covariance matrix is not more than a table that summaries the correlations between all the possible pairs of variables

# Eigenvector and eigenvalues

- Eigenvectors and eigenvalues are the linear algebra concepts that we need to compute from the covariance matrix in order to determine the principal components of the data.

- Before getting to the explanation of these concepts, let's first understand what do we mean by principal components.

# Eigenvector and eigenvalues

- Principal components are new variables that are constructed as linear combinations or mixtures of the initial variables.

- These combinations are done in such a way that the new variables (i.e., principal components) are uncorrelated and most of the information within the initial variables is squeezed or compressed into the first components.

- So, the idea is 10-dimensional data gives you 10 principal components, but PCA tries to put maximum possible information in the first component, then maximum remaining information in the second and so on, until having something like shown in the scree plot below.

# Eigenvector and eigenvalues

# Principal Components

- Organizing information in principal components this way, will allow you to reduce dimensionality without losing much information, and this by discarding the components with low information and considering the remaining components as your new variables.

- An important thing to realize here is that, the principal components are less interpretable and don't have any real meaning since they are constructed as linear combinations of the initial variables.

- Geometrically speaking, principal components represent the directions of the data that explain a maximal amount of variance, that is to say, the lines that capture most information of the data.

# Principal Components

- As there are as many principal components as there are variables in the data, principal components are constructed in such a manner that the first principal component accounts for the largest possible variance in the data set.

# Principal Components

- For example, let's assume that the scatter plot of our data set is as shown below, can we guess the first principal component ?

- Yes, it's approximately the line that matches the purple marks because it goes through the origin and it's the line in which the projection of the points (red dots) is the most spread out.

- Or mathematically speaking, it's the line that maximizes the variance (the average of the squared distances from the projected points (red dots) to the origin).

# Principal Components

# Example

- Let's suppose that our data set is 2-dimensional with 2 variables x,y and that the eigenvectors and eigenvalues of the covariance matrix are as follows:

$$v1 = \begin{bmatrix} 0.6778736 \\ 0.7351785 \end{bmatrix} \qquad \lambda_1 = 1.284028$$

$$v2 = \begin{bmatrix} -0.7351785 \\ 0.6778736 \end{bmatrix} \qquad \lambda_2 = 0.04908323$$

- 

If we rank the eigenvalues in descending order, we get λ1>λ2, which means that the eigenvector that corresponds to the first principal component (PC1) is v1 and the one that corresponds to the second component (PC2) isv2.

# Feature Vector

- As we saw in the previous step, computing the eigenvectors and ordering them by their eigenvalues in descending order, allow us to find the principal components in order of significance.

- In this step, what we do is, to choose whether to keep all these components or discard those of lesser significance (of low eigenvalues), and form with the remaining ones a matrix of vectors that we call Feature vector.

- So, the feature vector is simply a matrix that has as columns the eigenvectors of the components that we decide to keep.

- This makes it the first step towards dimensionality reduction, because if we choose to keep only p eigenvectors (components) out of n, the final data set will have only p dimensions.

# Example:

- Continuing with the example from the previous step, we can either form a feature vector with both of the eigenvectors v1 and v2:

$$\begin{bmatrix} 0.6778736 & -0.7351785 \\ 0.7351785 & 0.6778736 \end{bmatrix}$$

- Or discard the eigenvector v2, which is the one of lesser significance, and form a feature vector with v1 only:

$$\begin{bmatrix} 0.6778736 \\ 0.7351785 \end{bmatrix}$$

- Discarding the eigenvector v2 will reduce dimensionality by 1, and will consequently cause a loss of information in the final data set.

- The aim is to use the feature vector formed using the eigenvectors of the covariance matrix, to reorient the data from the original axes to the ones represented by the principal components (hence the name Principal Components Analysis).

- This can be done by multiplying the transpose of the original data set by the transpose of the feature vector.

$$FinalDataSet = FeatureVector^T * StandardizedOriginalDataSet^T$$

# Advantages of PCA

- The training time of the algorithms reduces significantly with less number of features.

- It is not always possible to analyze data in high dimensions. For instance if there are 100 features in a dataset. Total number of scatter plots required to visualize the data would be 100(100-1)2 = 4950. Practically it is not possible to analyze data this way.

# Normalization of features

- It is imperative to mention that a feature set must be normalized before applying PCA. For instance if a feature set has data expressed in units of Kilograms, Light years, or Millions, the variance scale is huge in the training set. If PCA is applied on such a feature set, the resultant loadings for features with high variance will also be large. Hence, principal components will be biased towards features with high variance, leading to false results.

- Finally, the last point to remember before we start coding is that PCA is a statistical technique and can only be applied to numeric data. Therefore, categorical features are required to be converted into numerical features before PCA can be applied.

# Example:

# Reading the dataset

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

dataset = pd.read_csv('wisc.csv')
print dataset.head()

# Input features
X = dataset.drop(['id','diagnosis'], 1)

# Output feature
y = dataset['diagnosis']
```

tusharkute
.com

# Normalize

```python
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)

from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

tusharkute
.com

# Apply PCA

```python
from sklearn.decomposition import PCA
from sklearn.ensemble import RandomForestClassifier

# Create any calssifier
classifier = RandomForestClassifier(n_estimators=10,
                random_state=0)

# Object of PCA
pca = PCA()

X_train0 = pca.fit_transform(X_train)
X_test0 = pca.transform(X_test)
classifier.fit(X_train0, y_train)
y_pred = classifier.predict(X_test0)
```
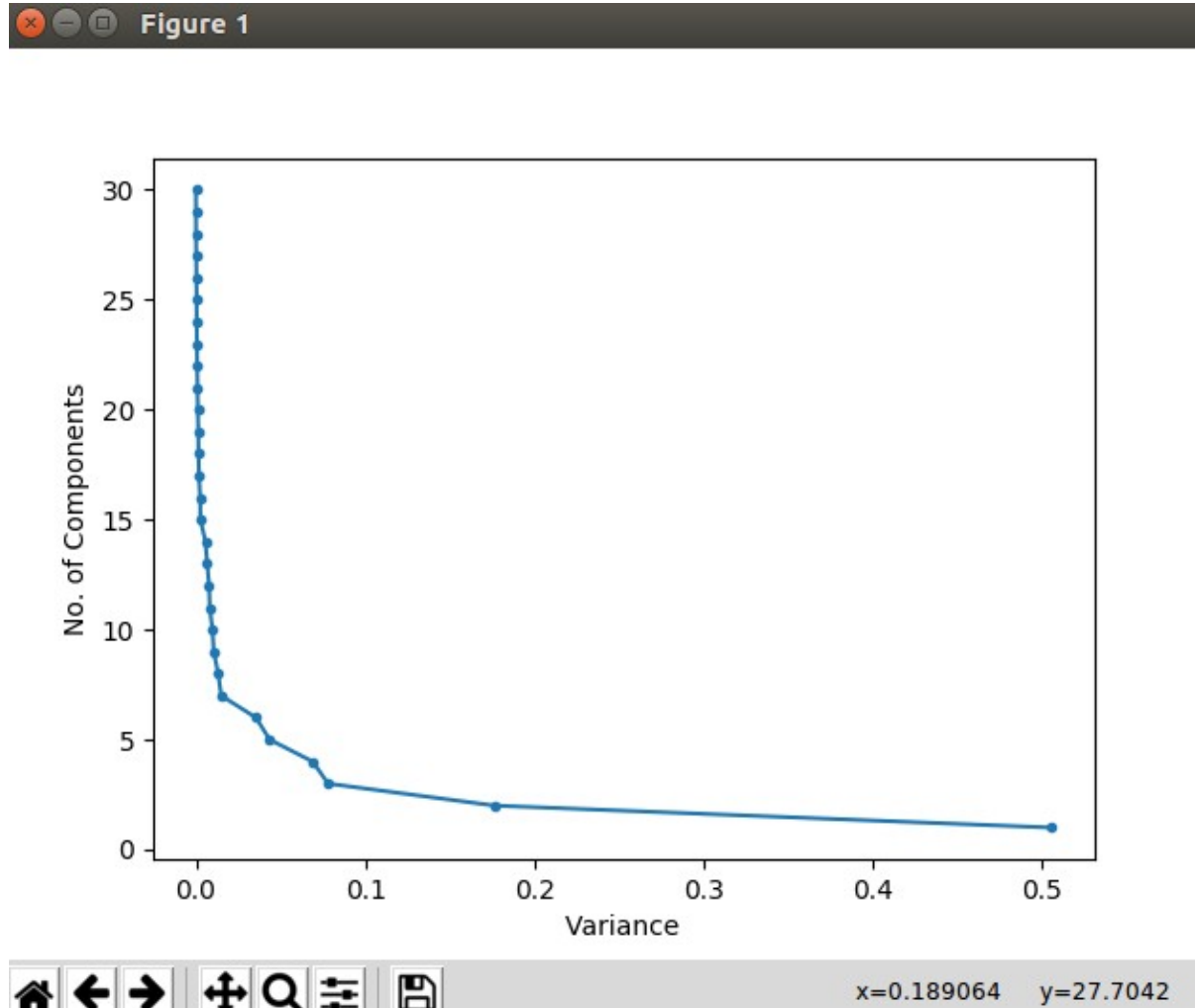
# Calculate variance

```python
from sklearn.metrics import confusion_matrix,
                             accuracy_score
print 'Confusion Matrix (default):\n',
            confusion_matrix(y_test, y_pred)
print 'Accuracy (default):',
            accuracy_score(y_test, y_pred) * 100
variance = pca.explained_variance_ratio_
print variance

# Plot variance vs. no. of components
plt.plot(variance, range(1,31), '.-')
plt.xlabel('Variance')
plt.ylabel('No. of Components')
plt.show()
```

# Variance plot

# Variance Ratio

- The PCA class contains explained_variance_ratio_ which returns the variance caused by each of the principal components.

# Principal Components = 1

```python
# Principal components = 1
pca = PCA(n_components=1)
X_train1 = pca.fit_transform(X_train)
X_test1 = pca.transform(X_test)

classifier.fit(X_train1, y_train)
y_pred = classifier.predict(X_test1)

print 'Confusion Matrix (1):\n', confusion_matrix(y_test, y_pred)
print 'Accuracy (1):', accuracy_score(y_test, y_pred) * 100
```

# Principal Components = 2

```python
# Principal components = 2
pca = PCA(n_components=2)
X_train2 = pca.fit_transform(X_train)
X_test2 = pca.transform(X_test)

classifier.fit(X_train2, y_train)
y_pred = classifier.predict(X_test2)

print 'Confusion Matrix (2):\n', confusion_matrix(y_test, y_pred)
print 'Accuracy (2):', accuracy_score(y_test, y_pred) * 100
```

tusharkute
.com

# Principal Components = 3

```python
# Principal components = 3
pca = PCA(n_components=3)
X_train3 = pca.fit_transform(X_train)
X_test3 = pca.transform(X_test)
classifier.fit(X_train3, y_train)
y_pred = classifier.predict(X_test3)

print 'Confusion Matrix (3):\n', confusion_matrix(y_test, y_pred)
print 'Accuracy (3):', accuracy_score(y_test, y_pred) * 100
```

# Linear Discriminant Analysis

- Linear discriminant analysis (LDA), normal discriminant analysis (NDA), or discriminant function analysis is a generalization of Fisher's linear discriminant, a method used in statistics, pattern recognition, and machine learning to find a linear combination of features that characterizes or separates two or more classes of objects or events.

- The resulting combination may be used as a linear classifier, or, more commonly, for dimensionality reduction before later classification.
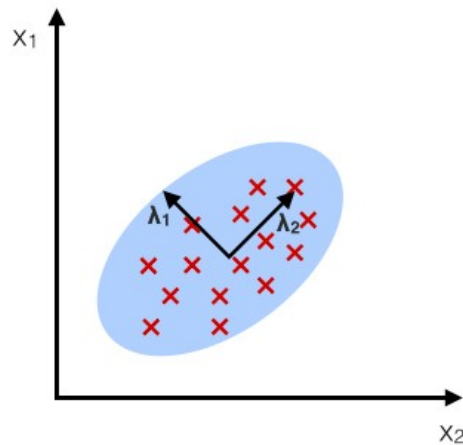
- Compute d-dimensional mean vectors for different classes from the dataset, where d is the dimension of feature space.

- Compute in-between class and with-in class scatter matrices.

- Compute eigen vectors and corresponding eigen values for the scatter matrices.

- Choose k eigen vectors corresponding to top k eigen values to form a transformation matrix of dimension d x k.

- Transform the d-dimensional feature space X to k-dimensional feature space X_lda via the transformation matrix.
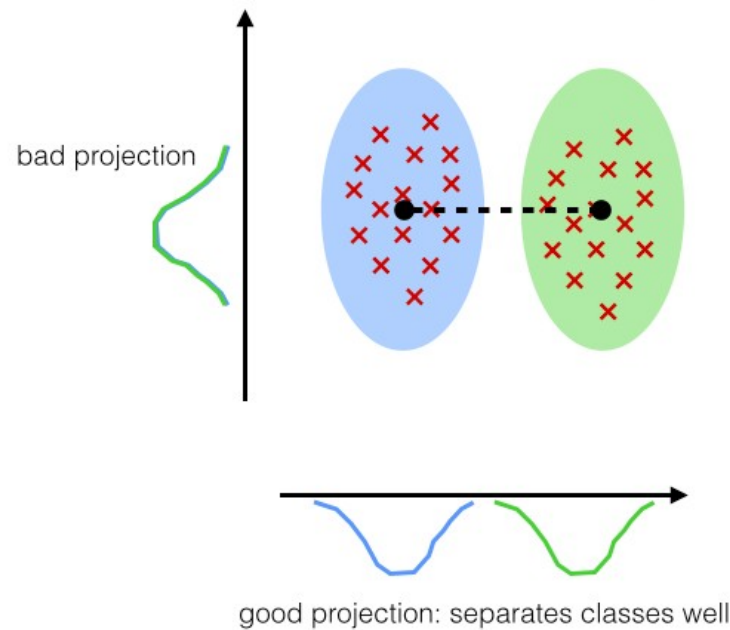
**PCA:**
component axes that maximize the variance

**LDA:**
maximizing the component axes for class-separation

Reference: sebastianraschka.com

# Useful resources

- https://stackabuse.com
- http://archive.ics.uci.edu/ml/index.php
- https://scikit-learn.org
- https://en.wikipedia.org
- www.towardsdatascience.com
- www.analyticsvidhya.com
- www.kaggle.com
- www.github.com

# Thank you

@mitu_skillologies

/mITuSkillologies

@mitu_group

/company/mitu-skillologies

c/MITUSkillologies

**Web Resources**
http://mitu.co.in
http://tusharkute.com

contact@mitu.co.in

tushar@tusharkute.com