

Feature Selection

Agenda

In this lesson, we will cover the following concepts with the help of a business use case:

- Feature Selection
 - Dimensionality Reduction:
 - Dimensionality Reduction Technique
 - Pros and Cons of Dimensionality Reduction
- Prominent Methods Used for Feature Selection:
 - Factor Analysis
 - Principal Component Analysis
 - LDA

What Is Feature Selection?

Feature selection is a method that helps in the inclusion of the significant variables that help form a model with good predictive power.

Features or variables that are redundant or irrelevant can negatively impact the performance of the model, thus it becomes necessary to remove them.

Benefits of Feature Selection

- It reduces overfitting as the unwanted variables are removed, and the focus is on the significant variables.
- It removes irrelevant information, which helps to improve the accuracy of the model's predictions.
- It reduces the computation time involved to get the model's predictions.

Dimensionality Reduction

Dimensionality reduction is the method of transforming a collection of data having large dimensions into data of smaller dimensions while ensuring that identical information is conveyed concisely.

Dimensionality Reduction Techniques

Some of the techniques used for dimensionality reduction are:

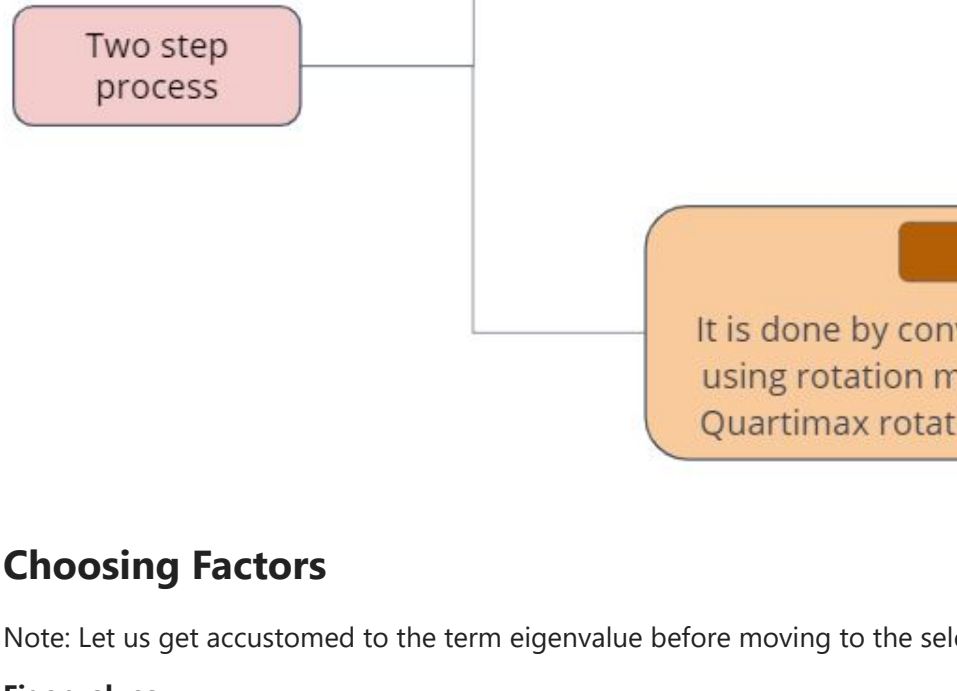
1. Imputing missing values.
2. Dropping low-variance variables
3. Decision trees (DT)
4. Random forest (RF)
5. Reducing highly correlated variables
6. Backward feature elimination
7. Factor analysis
8. Principal component analysis (PCA)

Pros of Dimensionality Reduction

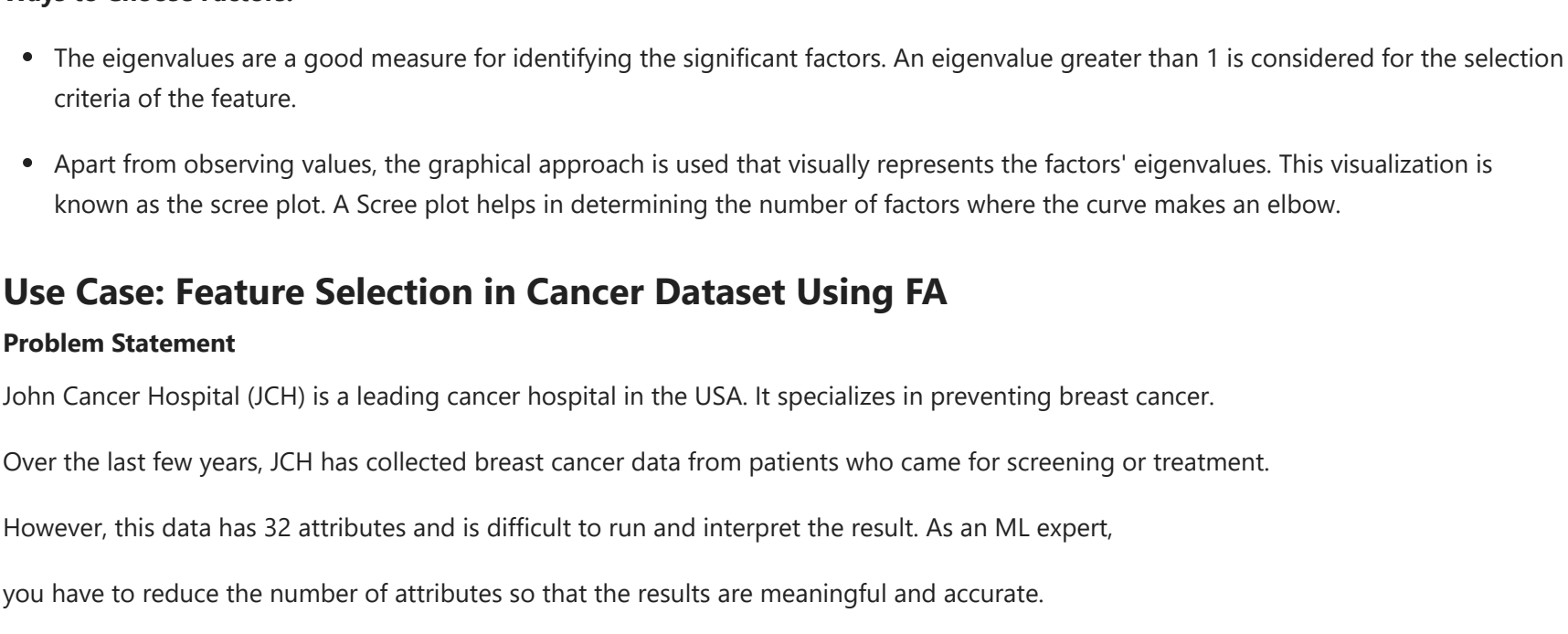
- It helps to compress data, reducing the storage space needed.
 - It cuts down on computing time.
 - It also aids in the removal of redundant features.
- Cons of Dimensionality Reduction**
- Some data may will be lost as a result.
 - PCA tends to find linear relationships between variables, which is not always ideal.
 - When mean and covariance are insufficient to describe datasets, PCA fails.
 - We use certain thumb rules when we do not know how many principal components to keep in practice.

Gist of Factor Analysis

- Factor analysis is used to:
 - Explain variance among the observed variables
 - Condense the set of observed variables into the factors
- $$Y_i = \beta_{i0} + \beta_{i1}F_1 + \beta_{i2}F_2 + \dots + (1)\epsilon_i$$
- Factor explains the amount of variance in the observed variables.
- In other words, factor analysis is a method that investigates linear relation of a number of variables of interest V1, V2,....., V_i with a smaller number of unobservable factors F1, F2,....., F_k.



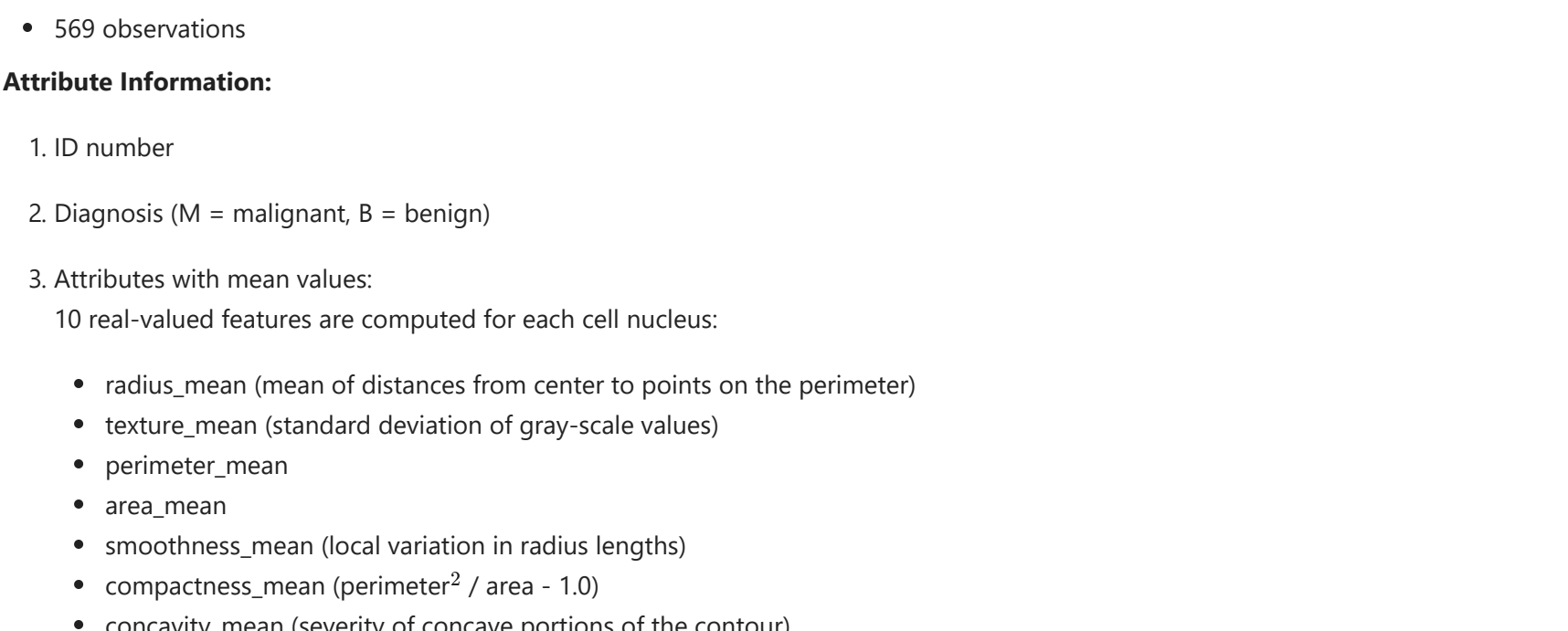
Types of FA



Work Process of FA

The objective of the factor analysis is the reduction of the number of observed variables and find the unobservable variables.

It is a two-step process.



Choosing Factors

Note: Let us get accustomed to the term eigenvalue before moving to the selecting the number of factors.

Eigenvalues:

It represents the explained variance of each factor from the total variance, and is also known as the characteristic roots.

Ways to Choose Factors:

- The eigenvalues are a good measure for identifying the significant factors. An eigenvalue greater than 1 is considered for the selection criteria of the factor.
- Apart from observing values, the graphical approach is used that visually represents the factors' eigenvalues. This visualization is known as the scree plot. A scree plot helps in determining the number of factors where the curve makes an elbow.

Use Case: Feature Selection in Cancer Dataset Using FA

Problem Statement

John Cancer Hospital (JCH) is a leading cancer hospital in the USA. It specializes in preventing breast cancer.

Over the last few years, JCH has collected breast cancer data from patients who came for screening or treatment.

However, this data has 32 attributes and is difficult to run and interpret the result. As an ML expert,

you have to reduce the number of attributes so that the results are meaningful and accurate.

Use FA for feature selection.

Dataset

Features of the dataset are computed from a digitized image of a Fine-Needle Aspirate (FNA) of a breast mass.

They describe the characteristics of the cell nuclei present in the image.

Data Dictionary

Dimensions:

- 32 variables
- 569 observations

Attribute Information:

1. ID number
2. Diagnosis (M = malignant, B = benign)
3. Attributes with mean values:
 - 10 real-valued features are computed for each cell nucleus:
 - radius_mean (mean of distances from center to points on the perimeter)
 - texture_mean (standard deviation of gray-scale values)
 - perimeter_mean
 - area_mean
 - smoothness_mean (local variation in radius lengths)
 - compactness_mean (perimeter² / area - 1.0)
 - concavity_mean (severity of concave portions of the contour)
 - concave points_mean (number of concave portions of the contour)
 - symmetry_mean
 - fractal dimension_mean ("coastline approximation" - 1)
4. Attributes with standard error and worst/largest:
 - radius_se
 - texture_se
 - perimeter_se
 - area_se
 - smoothness_se
 - compactness_se
 - concavity_se
 - concave points_se
 - symmetry_se
 - fractal dimension_se
 - radius_worst
 - texture_worst
 - perimeter_worst
 - area_worst
 - smoothness_worst
 - compactness_worst
 - concavity_worst
 - concave points_worst
 - symmetry_worst
 - fractal dimension_worst

Solution

Import Libraries

In Python, Numpy is a package that includes multidimensional array objects as well as a number of derived objects. Matplotlib is an amazing visualization library in Python for 2D plots of arrays. Pandas is used for data manipulation and analysis. So these are the core libraries that are used for the EDA process.

These libraries are written with an import keyword.

```
In [1]: import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
import matplotlib inline
```

Import and Check the Data

Before reading data from a csv file, you need to download the "breast-cancer-data.csv" dataset from the resource section and upload it into the Lab. We will use the Up arrow icon, which is shown on the left side under the View icon. Click on the Up arrow icon and upload the file from wherever it has downloaded into your system.

After this, you will see the downloaded file will be visible on the left side of your lab along with all the .pynb files.

```
In [2]: df = pd.read_csv('breast-cancer-data.csv')
df.head()
```

The p-value is 0, and this indicates that the test is statistically significant and highlights that the correlation matrix is not an identity matrix.

Kaiser-Meyer-Olkin Test

The Kaiser-Meyer-Olkin (KMO) test determines if data is suitable for factor analysis.

- The Kaiser-Meyer-Olkin (KMO) test determines if data is suitable for factor analysis.
- It assesses the suitability of each observed variable as well as the entire model.

```
from_factor_analyzer.factor_analyzer import calculate_kmo
kmo_all, kmo_model=calculate_kmo(df)

C:\Users\vaipika.gupta\AppData\Local\anaconda3\lib\site-packages\factor_analyzer\utils.py:249: UserWarning: The inverse of the
variance-covariance matrix could not be calculated using the Moore-Penrose generalised linearisation, due to its de-
terminant being at or very close to zero.
warnings.warn('The inverse of the variance-covariance matrix '
```

```
kmo_model
```

```
0.2599517032832366
```

pd.read_csv function is used to read the "breast-cancer-data.csv" file and df.head() will show the top 5 rows of the dataset.

- dataframe or df is a variable that will store the data read by the csv file.
- head will show the rows and () default take the 5 top rows as output.
- one more example - df.head(3) will show the top 3 rows.

shape function

```
In [3]: df.shape
```

Out[3]: (569, 32)

df.shape will show the number of rows and columns in the dataframe.

info Function

```
In [4]: # Check the data , there should be no missing values
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 569 entries, 0 to 568
Data columns (total 32 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   id                     569 non-null    int64
 1   diagnosis              569 non-null    object
 2   radius_mean            569 non-null    float64
 3   texture_mean           569 non-null    float64
 4   perimeter_mean         569 non-null    float64
 5   area_mean              569 non-null    float64
 6   smoothness_mean        569 non-null    float64
 7   compactness_mean       569 non-null    float64
 8   concavity_mean         569 non-null    float64
 9   concave points_mean    569 non-null    float64
10  symmetry_mean          569 non-null    float64
11  fractal_dimension_mean  569 non-null    float64
12  radius_se              569 non-null    float64
13  texture_se              569 non-null    float64
14  perimeter_se           569 non-null    float64
15  area_se                569 non-null    float64
16  smoothness_se          569 non-null    float64
17  compactness_se         569 non-null    float64
18  concavity_se           569 non-null    float64
19  concave points_se      569 non-null    float64
20  symmetry_se            569 non-null    float64
21  fractal_dimension_se    569 non-null    float64
22  radius_worst           569 non-null    float64
23  texture_worst           569 non-null    float64
24  perimeter_worst        569 non-null    float64
25  area_worst             569 non-null    float64
26  smoothness_worst       569 non-null    float64
27  compactness_worst      569 non-null    float64
28  concavity_worst        569 non-null    float64
29  concave points_worst   569 non-null    float64
30  symmetry_worst         569 non-null    float64
31  fractal_dimension_worst 569 non-null    float64
dtypes: float64(30), int64(1), object(1)
memory usage: 142.4 KB
```

- The dataframe's information is printed using the info() function.
- The number of columns, column labels, column data types, memory use, range index, and the number of cells in each column are all included in the data (non-null values).

```
In [5]: # Defining the array as np.array
feature_names = np.array(['mean radius', 'mean texture', 'mean perimeter', 'mean area',
                           'mean smoothness', 'mean compactness', 'mean concavity',
                           'mean concave points', 'mean symmetry', 'mean fractal dimension',
                           'radius error', 'texture error', 'perimeter error', 'area error',
                           'smoothness error', 'compactness error', 'concavity error',
                           'concave points error', 'symmetry error', 'fractal dimension error',
                           'worst radius', 'worst texture', 'worst perimeter', 'worst area',
                           'worst smoothness', 'worst compactness', 'worst concavity',
                           'worst concave points', 'worst symmetry', 'worst fractal dimension'])
```

```
In [6]: ## Convert diagnosis column to 1/0 and store in new column target
from sklearn.preprocessing import LabelEncoder
```

- The sklearn.preprocessing package contains a number of useful utility methods and transformer classes for converting raw feature vectors into an format that is suitable for downstream estimators.
- LabelEncoder encodes labels with a value between 0 and n_classes-1 where n is the number of distinct labels.
- These libraries are written with an import keyword.

```
In [7]: # Encode label diagnosis
#M -> 1
#B -> 0

#Converting diagnosis to numerical variable in df
df['diagnosis'] = df['diagnosis'].map({'M': 1, 'B': 0}).astype(int)
```

In the above code, we are encoding the column diagnosis in which we are encoding M as 1 and B as 0.

Factor Analysis

- A linear statistical model is a factor analysis.
- It is used to condense a group of observable variables into an unobserved variable termed factors and to explain the variance among the observed variables.

Adequacy Test

Before you perform factor analysis, you need to evaluate the "factorability" of our dataset. Can we find the factors in the dataset? Checking factorability or sampling adequacy can be done in two ways:

1- The Bartlett's Test

2- Test of Kaiser-Meyer-Olkin

```
In [8]: #Install factor analyzer
!pip install factor_analyzer
```

Bartlett's Test

Bartlett's test of sphericity checks whether or not the observed variables intercorrelate at all using the observed correlation matrix against the identity matrix. If the test is found to be statistically insignificant, you should not employ a factor analysis.

Note: This test checks for the intercorrelation of observed variables by comparing the observed correlation matrix against the identity matrix.

```
In [9]: !pip install factor_analyzer

Requirement already satisfied: factor_analyzer in c:\users\alpika.gupta\anaconda3\lib\site-packages (0.4.0)
Requirement already satisfied: scikit-learn in c:\users\alpika.gupta\anaconda3\lib\site-packages (from factor_analyzer) (1.1.2)
Requirement already satisfied: scipy in c:\users\alpika.gupta\anaconda3\lib\site-packages (from factor_analyzer) (1.7.1)
Requirement already satisfied: pandas in c:\users\alpika.gupta\anaconda3\lib\site-packages (from factor_analyzer) (1.3.4)
Requirement already satisfied: pytz>=2017.3 in c:\users\alpika.gupta\anaconda3\lib\site-packages (from pandas->factor_analyzer) (2021.3)
Requirement already satisfied: python-dateutil>=2.7.3 in c:\users\alpika.gupta\anaconda3\lib\site-packages (from pandas->factor_analyzer) (2.8.2)
Requirement already satisfied: six>=1.5 in c:\users\alpika.gupta\anaconda3\lib\site-packages (from python-dateutil>=2.7.3->pandas->factor_analyzer) (1.16.0)
Requirement already satisfied: joblib>=1.0.0 in c:\users\alpika.gupta\anaconda3\lib\site-packages (from scikit-learn->factor_analyzer) (1.1.0)
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\alpika.gupta\anaconda3\lib\site-packages (from scikit-learn->factor_analyzer) (2.2.0)
```

- In the above code, we are installing the factor analyzer.
- Factor analysis is an exploratory data analysis method used to search for influential underlying factors or latent variables from a set of observed variables.

Now, you are trying to perform factor analysis by using the factor analyzer module. Use the below code for calculating bartlett_sphericity.

```
In [10]: from factor_analyzer import FactorAnalyzer
from factor_analyzer.factor_analyzer import calculate_bartlett_sphericity
chi_square_value, p_value, df = calculate_bartlett_sphericity(df)

Out[10]: (40227.4857034462, 0.0)
```

- In the above code, we are importing the factor analyzer and calculate bartlett_sphericity.
- In this Bartlett's test, the p-value is 0. The test was statistically significant, indicating that the observed correlation matrix against the identity matrix.

Inference:

The p-value is 0, and this indicates that the test is statistically significant and highlights that the correlation matrix is not an identity matrix.

Kaiser-Meyer-Olkin Test

- The Kaiser-Meyer-Olkin (KMO) test determines if data is suitable for factor analysis.
- It assesses the suitability of each observed variable as well as the entire model.

```
In [11]: from factor_analyzer.factor_analyzer import calculate_kmo
kmo_all, kmo_model = calculate_kmo(df)
```

C:\Users\alpika.gupta\anaconda3\lib\site-packages\factor_analyzer\utils.py:249: UserWarning: The inverse of the variance-covariance matrix was calculated using the Moore-Penrose generalized matrix inversion, due to its determinant being at or very close to zero.
warnings.warn("The inverse of the variance-covariance matrix ")

```
In [12]: kmo_model

Out[12]: 0.2599517032832366
```

- In the above code, we are calculating the KMO. KMO estimates the proportion of variance among all the observed variables.
- The overall KMO for the data is 0.25, which is excellent. This value indicates that you can proceed with your planned factor analysis.

```
In [13]: df.head(2)
```

should be the feature names. If this method is used to create the "ModelSpecification", then factor names and variable names will be added as properties to that object.

Parameters

Y: array-like

2 rows x 32 columns

Inference:

The KMO value is less than 0.5, and this indicates that we need to delete the insignificant variables.

Finding Significant Variables

```
In [14]: corr = df.corr()
corr.style.background_gradient(cmap='coolwarm')
```

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	convexity_mean	fractal_dimension_mean
id	1.000000		0.074626		0.099770	0.096893		-0.012968		0.000096	
diagnosis	0.039760	1.000000		0.491785	0.740636	0.708984		0.358560		0.596534	
radius_mean	0.074626	0.730309	1.000000		0.323782	0.997855		0.987357		0.506124	
texture_mean	0.099770	0.415185	0.323782	1.000000		0.329533		0.321086		0.236702	
perimeter_mean	0.097319	0.742636	0.997855	0.329533	1.000000	0.986507		0.207278		0.556936	
area_mean	0.068890	0.708984	0.987357	0.321086	0.986507	1.000000		0.177028		0.498502	
smoothness_mean	-0.012968	0.358560	0.170581	-0.023389	0.207278	0.177028	1.000000			0.659123	
compactness_mean	0.000096	0.596534	0.506124	0.236702	0.556936	0.498502		0.659123		1.000000	
concavity_mean	0.044159	0.696360	0.506124	0.302418	0.716136	0.685983		0.521984		0.881321	
convexity_mean	-0.022114	-0.330498	0.147441	-0.076437	0.183027	0.151291		0.557775		0.602641	
fractal_dimension_mean	0.005251	0.596534	0.506124	0.236702	0.556936	0.498502	1.000000			0.659123	
radius_worst	0.143040	0.567134	0.679090	0.275869	0.691765	0.732542		0.301467		0.494743	
texture_worst	0.064720	0.730309	0.323782	0.386358	0.912045	0.987357		0.236702		0.556936	
perimeter_worst	0.079986	0.742636	0.997855	0.329533	1.000000	0.986507		0.207278		0.556936	
area_worst	0.091789	0.730309	0.997855	0.329533	1.000000	0.986507		0.207278		0.556936	
smoothness_worst	0.044159	0.696360	0.506124	0.302418	0.716136	0.685983		0.521984		0.881321	
compactness_worst	0.000096	0.596534	0.506124	0.236702	0.556936	0.498502		0.659123		1.000000	
concavity_worst	0.044159	0.696360	0.506124	0.302418	0.716136	0.685983		0.521984		0.881321	
convexity_worst	-0.022114	-0.330498	0.147441	-0.076437	0.183027	0.151291		0.557775		0.602641	
fractal_dimension_worst	0.005251	0.596534	0.506124	0.236702	0.556936	0.498502	1.000000			0.659123	
radius_worst	0.143040	0.567134	0.679090	0.275869	0.691765	0.732542		0.301467		0.494743	
texture_worst	0.064720	0.730309	0.323782	0.386358	0.912045	0.987357		0.236702		0.556936	
perimeter_worst	0.079986	0.742636	0.997855	0.329533	1.000000	0.986507		0.207278		0.556936	
area_worst	0.091789	0.730309	0.997855	0.329533	1.000000	0.986507		0.207278		0.556936	
smoothness_worst	0.044159	0.696360	0.506124	0.302418	0.716136	0.685983		0.521984		0.881321	
compactness_worst	0.000096	0.596534	0.506124	0.236702	0.556936	0.498502		0.659123		1.000000	
concavity_worst	0.044159	0.696360	0.506124	0.302418	0.716136	0.685983		0.521984		0.881321	
convexity_worst	-0.022114	-0.330498	0.147441	-0.076437	0.183027	0.151291		0.557775		0.602641	
fractal_dimension_worst	0.005251	0.596534	0.506124	0.236702	0.556936	0.498502	1.000000			0.659123	

If your main goal is to visualize the correlation matrix rather than creating a plot per se, the convenient pandas styling options are a viable built-in solution, as shown in the above code.


```
[28]: model = ModelSpecificationParser.parse_model_specification_from_dict(df_corr,model_dict)
df_corr = Pandas.DataFrame and model_specification dictionary with factors

def fit(self, X, y=None):
    """
    Perform confirmatory factor analysis.

    Parameters
    -----
    X : array-like
        The data to use for confirmatory
        factor analysis. If this is just a
        covariance matrix, make sure 'is_cov_matrix'
        was set to True.
    y : ignored

    Raises
    -----
    ValueError: If the specification is not None or a "ModelSpecification" object
    AssertionError: If 'is_cov_matrix=True' and the matrix is not square.
    AssertionError: If len(bounds) != len(x0)

In [29]: # Performs confirmatory factor analysis
cfa = ConfirmatoryFactorAnalyzer(model_spec, disp=False)
cfa.fit(df_corr.values)

C:\Users\alpika.gupta\Anaconda3\lib\site-packages\factor_analyzer\confirmatory_factor_analyzer.py:1732: UserWarning: The optimization routine failed to converge: ABRNORMAL_TERMINATION_IN_INSRCH
warnings.warn("The optimization routine failed.")

Out[29]: ConfirmatoryFactorAnalyzer
ConfirmatoryFactorAnalyzer(disp=False, n_obs=569, specification=ConfirmatoryFactorAnalyzer.ModelSpecification object at 0x000016076764A680)

In [30]: # cfa.loadings will give you the factor loading matrix
# The factor loading is a matrix which shows the relationship of each variable to the underlying factor.
# It shows the correlation coefficient for observed variable and factor.
# It shows the variance explained by the observed variables.
cfa.loadings

Out[30]: array([[ 9.5238985e+00,  0.0000000e+00],
       [ 3.47315903e+02,  0.0000000e+00],
       [ 7.86489367e+03,  0.0000000e+00],
       [ 4.9873855e+02,  0.0000000e+00],
       [ 9.11336000e+02,  0.0000000e+00],
       [ 1.32123059e+04,  0.0000000e+00],
       [ 0.00000000e+00,  0.18813073e+03],
       [ 0.00000000e+00, -1.45499332e+03],
       [ 0.00000000e+00, -3.13842115e+02],
       [ 0.00000000e+00,  4.12273057e+03],
       [ 0.00000000e+00,  1.11725801e+02]])

In [31]: # This will give you the factor covariance matrix and the type of this will be numpy array
cfa.factor_covariance_

Out[31]: array([[1. ,  0.33443073],
       [0.33443073,  1. ]])

In [32]: # transform(x) used to get the factor scores for new data set.

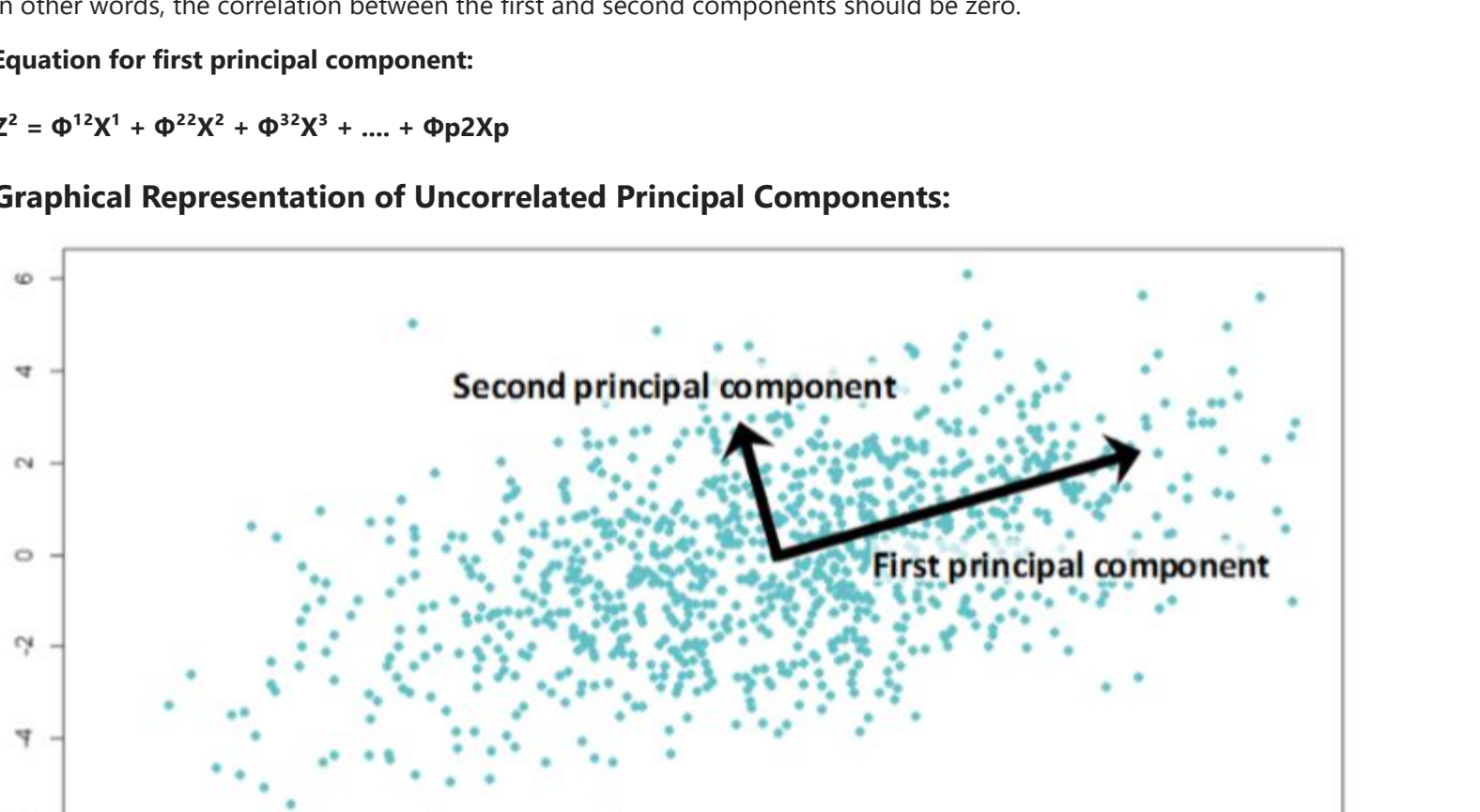
# Parameters: X (array-like, shape (n_samples, n_features)) - The data to score using the fitted factor model.
# Returns: scores - The latent variables of X.
# Return type: numpy array, shape (n_samples, n_components)

In [33]: cfa.transform(df_corr.values)

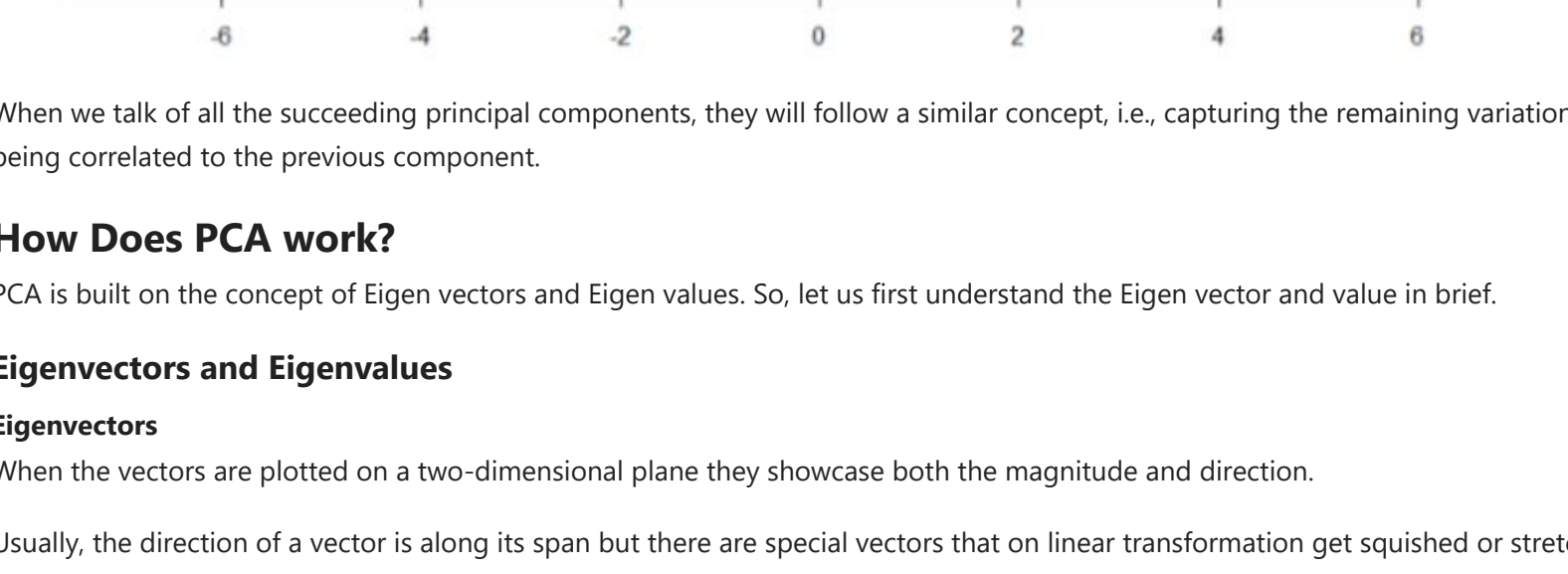
Out[33]: array([[ 0.08030784,  0.00103084],
       [ 0.08139898,  0.00108067],
       [ 0.06382796,  0.0010036 ],
       ...,
       [ 0.01961916,  0.0012 ],
       [ 0.0725346 ,  0.00141418],
       [-0.04865471, -0.00067236]])
```

Now, let us take a look at PCA and LDA-specific benefits.

PCA and LDA # NO Review



Gist of PCA



What Is Principal Component?

A principal component is a normalized linear combination of the original predictors in a dataset.

First Principal Component (Z¹)

- The aim of PCA is to find components that account for **maximum variance in the data** that includes the error and within-variable variance.
 - It finds the direction of the highest variability in the data. Greater the variability captured in the first component, the greater the information captured by the component.
 - The first principal component develops a line that is closest to the data points.
- In other words, it minimizes the sum of squared distance between a data point and the line.

Equation for first principal component:

$$Z^1 = \phi^1 X^1 + \phi^2 X^2 + \phi^3 X^3 + \dots + \phi^p X^p$$

- Z¹ = Principal component
- φ¹ = Loading vector
- X¹ to X^p = Normalized predictors that have mean equals to 0 and standard deviation equals to 1.

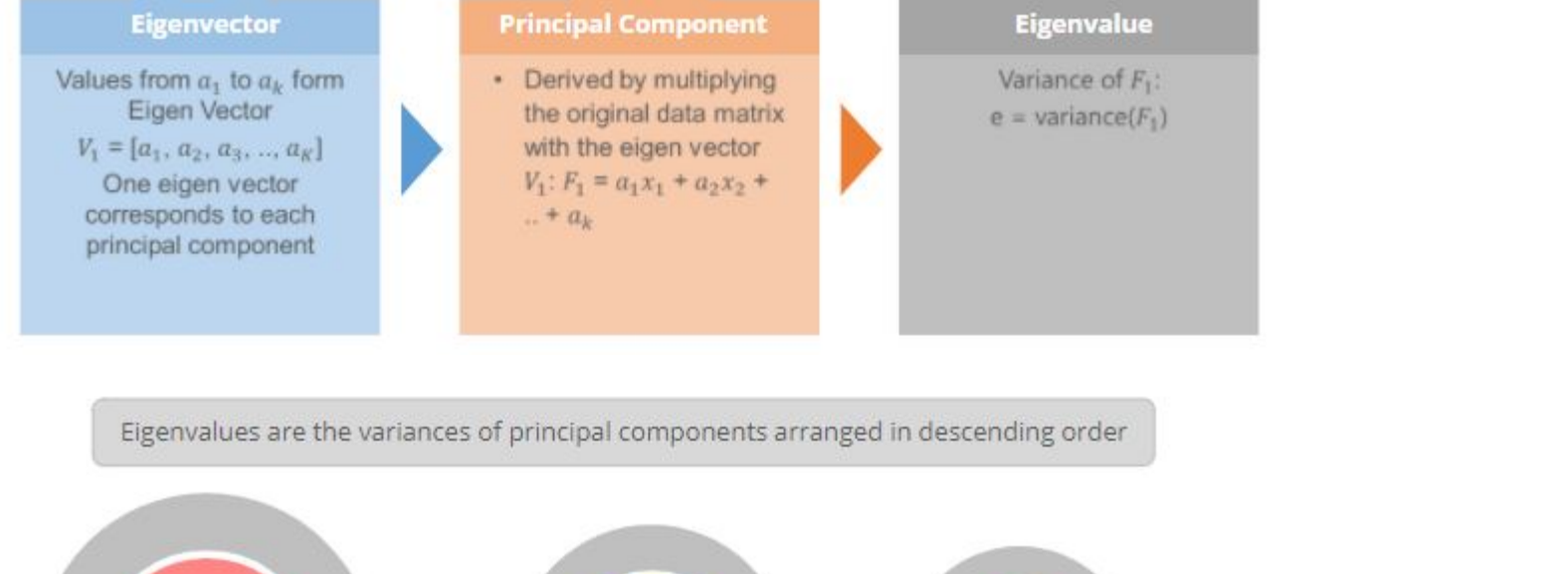
Second Principal Component (Z²)

The aim of the linear combination of original predictors is to capture the remaining variance in the dataset where Z² is uncorrelated to Z¹. In other words, the correlation between the first and second components should be zero.

Equation for first principal component:

$$Z^2 = \phi^1 X^1 + \phi^2 X^2 + \phi^3 X^3 + \dots + \phi^p X^p$$

Graphical Representation of Uncorrelated Principal Components:



When we talk of all the succeeding principal components, they will follow a similar concept, i.e., capturing the remaining variation without being correlated to the previous component.

How Does PCA work?

PCA is built on the concept of Eigen vectors and Eigen values. So, let us first understand the Eigen vector and value in brief.

Eigen vectors and Eigen values

Eigenvectors

When the vectors are plotted on a two-dimensional plane they showcase both the magnitude and direction.

Usually, the direction of a vector is along its span but there are special vectors that on linear transformation get squished or stretched instead of falling off their span.

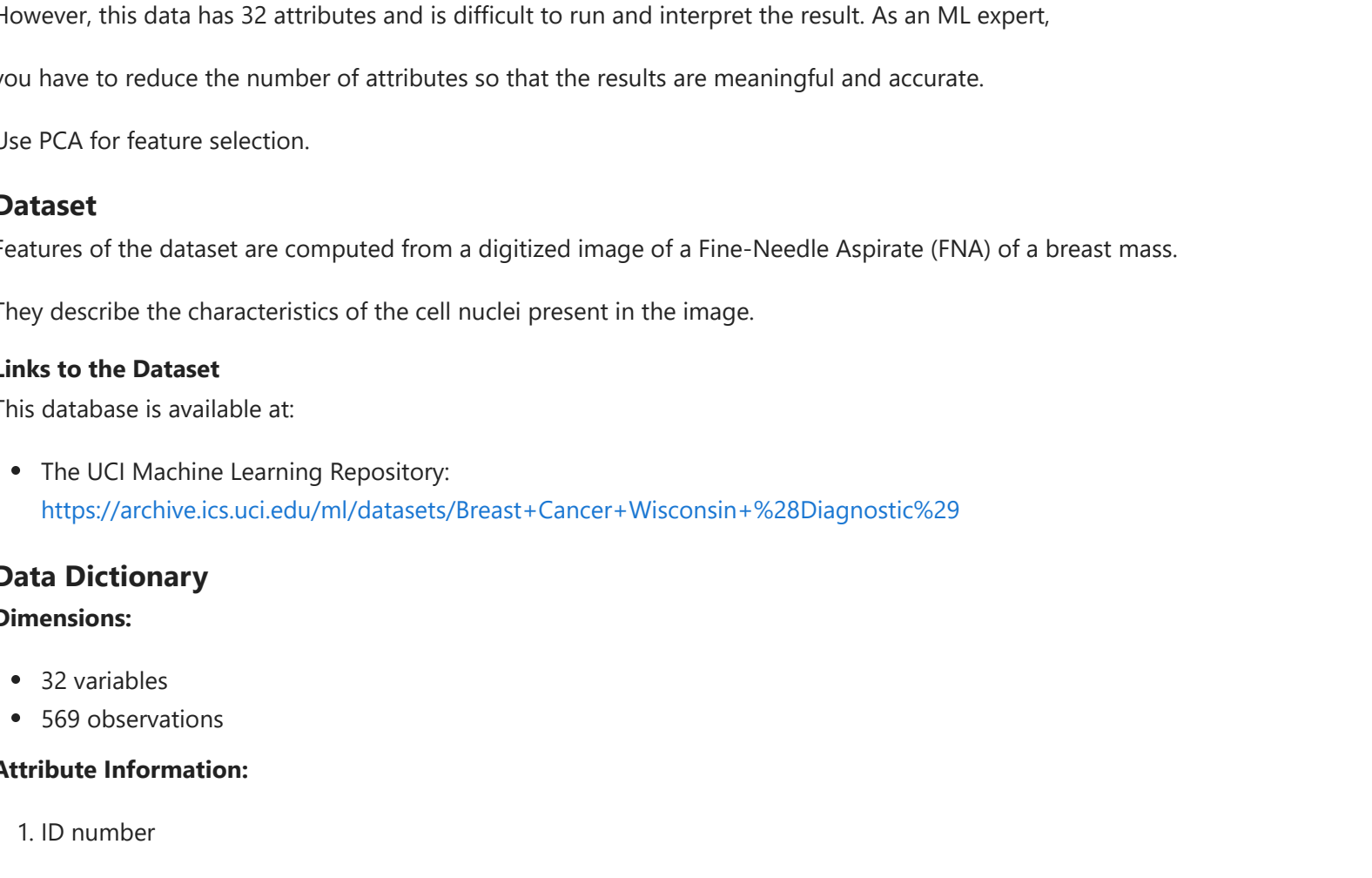
Eigenvalues

These are the constant values that increase or decrease the Eigenvectors along their span when transformed linearly.

PCA Process

Finding PC1

The first principal component is the direction of maximum variance and is obtained by solving eigenvectors



Use Case: Feature Selection in Cancer Dataset Using PCA

Problem Statement

John Cancer Hospital (JCH) is a leading cancer hospital in the USA. It specializes in preventing breast cancer.

Over the last few years, JCH has collected breast cancer data from patients who came for screening or treatment.

However, this data has 32 attributes and is difficult to run and interpret the result. As an ML expert,

you have to reduce the number of attributes so that the results are meaningful and accurate.

Use PCA for feature selection.

Dataset

Features of the dataset are computed from a digitized image of a Fine-Needle Aspirate (FNA) of a breast mass.

They describe the characteristics of the cell nuclei present in the image.

Links to the Dataset

This database is available at:

- The UCI Machine Learning Repository:
<https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+%28Diagnostic%29>

Data Dictionary

Dimensions:

- 32 variables
- 569 observations

Attribute Information:

- ID number
- Diagnosis (M = malignant, B = benign)
- Attributes with mean values:
Ten real-valued features are computed for each cell nucleus:
 - radius_mean (mean of distances from center to points on the perimeter)
 - texture_mean (standard deviation of gray-scale values)
 - perimeter_mean
 - area_mean
 - smoothness_mean (local variation in radius lengths)
 - compactness_mean (severity of concave portions of the contour)
 - concave points_mean (number of concave portions of the contour)
 - symmetry_mean
 - fractal dimension_mean ("coastline approximation" - 1)
- Attributes with standard error and worst/largest:
 - radius_se
 - texture_se
 - perimeter_se
 - area_se
 - smoothness_se
 - compactness_se
 - concavity_se
 - concave points_se
 - symmetry_se
 - fractal dimension_se
 - radius_worst
 - texture_worst
 - perimeter_worst
 - area_worst
 - smoothness_worst
 - compactness_worst
 - concavity_worst
 - concave points_worst
 - symmetry_worst
 - fractal dimension_worst

Solution

Import Libraries

```
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib inline
```

Import and Check the Data

Before reading data from a csv file, you need to download the "breast-cancer-data.csv" dataset from the resource section and upload it into the lab.

We will use the Up arrow icon which is shown on the left side under the View icon. Click on the Up arrow icon and upload the file wherever it has been downloaded into your system.

After this, you will see the downloaded file will be visible on the left side of your lab along with all the .pynb files.

```
df = pd.read_csv('breast-cancer-data.csv')
df.head()
```

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	point
0	842302	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.08601	
1	842517	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.06891	
2	8430093	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	
3	8434801	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	
4	84358402	M	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	

5 rows x 32 columns

pd.read_csv function is used to read the "breast-cancer-data.csv" file and df.head() will show the top 5 rows of the dataset.

- dataframe or df is a variable that will store the data read by the csv file.
- head will show the rows and (j default take the top 5 rows as output.

```
df.shape
```

```
(569, 32)
```

```
# Check the data , there should be no missing values
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 569 entries, 0 to 568
Data columns (total 32 columns):
# Column Non-Null Count Dtype
---
0 id 569 non-null int64
1 diagnosis 569 non-null object
2 radius_mean 569 non-null float64
3 texture_mean 569 non-null float64
4 perimeter_mean 569 non-null float64
5 area_mean 569 non-null float64
6 smoothness_mean 569 non-null float64
7 compactness_mean 569 non-null float64
8 concavity_mean 569 non-null float64
9 concave points_mean 569 non-null float64
10 symmetry_mean 569 non-null float64
11 fractal_dimension_mean 569 non-null float64
12 radius_se 569 non-null float64
13 texture_se 569 non-null float64
14 perimeter_se 569 non-null float64
15 area_se 569 non-null float64
16 smoothness_se 569 non-null float64
17 compactness_se 569 non-null float64
18 concavity_se 569 non-null float64
19 concave points_se 569 non-null float64
20 symmetry_worst 569 non-null float64
21 fractal_dimension_worst 569 non-null float64
22 radius_worst 569 non-null float64
23 texture_worst 569 non-null float64
24 perimeter_worst 569 non-null float64
25 area_worst 569 non-null float64
26 smoothness_worst 569 non-null float64
27 compactness_worst 569 non-null float64
28 concavity_worst 569 non-null float64
29 concave points_worst 569 non-null float64
30 symmetry_worst 569 non-null float64
31 fractal_dimension_worst 569 non-null float64
dtypes: float64(30), int64(1), object(1)
memory usage: 142.4+ KB
```

- The dataframe's information is printed using the info() function.
- The number of columns, column labels, column data types, memory use, range index, and the number of cells in each column are all included in the data (non-null values).

```
feature_names = np.array(['mean radius', 'mean texture', 'mean perimeter', 'mean area',
                           'mean smoothness', 'mean compactness', 'mean concavity',
                           'mean concave points', 'mean symmetry', 'mean fractal dimension',
                           'radius error', 'texture error', 'perimeter error', 'area error',
                           'smoothness error', 'compactness error', 'concavity error',
                           'concave points error', 'symmetry error', 'fractal dimension error',
                           'worst radius', 'worst texture', 'worst perimeter', 'worst area',
                           'worst smoothness', 'worst compactness', 'worst concavity',
                           'worst concave points', 'worst symmetry', 'worst fractal dimension'])
```

```
### Convert diagnosis column to 1/0 and store in new column target
```

```
from sklearn.preprocessing import LabelEncoder
```

```
# Encode label diagnosis
M -> 1
B -> 0
```

```
target_data = df["diagnosis"]
encoder = LabelEncoder()
target_data = encoder.fit_transform(target_data)
```

In the above code, in output we are getting all the rows, but only with the last column.

```
df.drop(["diagnosis"],axis = 1, inplace = True)
```

In the above code, you store the encoded column in a dataframe and drop the diagnosis column for simplicity.

Principal Component Analysis

Let's use PCA to find the first two principal components and visualize the data in this new, two-dimensional space with a single scatter-plot.

```
from sklearn.preprocessing import StandardScaler
```

In the above code, you will scale data so that each feature has a single unit variance.

```
scaler = StandardScaler()
scaler.fit(df)
```

```
StandardScaler()
```

```
scaled_data = scaler.transform(df)
```

In the above code, you will initially create an object of the StandardScaler() function. Further, you will use fit() along with the object assigned to df and standardize it.

Two Principal Components

```
from sklearn.decomposition import PCA
```

In the above code, you can transform this data into its first 2 principal components.

```
pca = PCA(n_components=2)
```

Now, you will pass the number of components (n_components=2).

```
pca.fit(scaled_data)
```

```
PCA(n_components=2)
```

Finally, call fit function to aggregate the data. Here, several components represent the lower dimension in which you will project your higher dimension data.

```
#shape of data
scaled_data.shape
```

```
(569, 31)
```

```
x_pca.shape
```

```
(569, 2)
```

Plot # NO Review

```
# Reduced 30 dimensions to just 2! Let's plot these two dimensions out!
# Draw inference from the plot
plt.figure(figsize=(9,6))
plt.scatter(x_pca[:,0],x_pca[:,1],c=target_data,cmap='viridis')
plt.xlabel('First Principal Component')
plt.ylabel('Second Principal Component')
```

```
Text(0, 0.5, 'Second Principal Component')
```



Interpreting the components:

Unfortunately, with this great power of dimensionality reduction, comes the cost of being able to easily understand what these components represent.

The components correspond to combinations of the original features. The components themselves are stored as an attribute of the fitted PCA object.

```
pca.components_
```

```
array([[ 0.02291216,  0.21891302,  0.10384388,  0.22753491,  0.22104577,
        -0.14241471,  0.2230073 ,  0.25828025,  0.26073811,  0.13797774,
        0.06414779,  0.20611747,  0.01741339,  0.2144652,  0.20307642,
        0.01467821,  0.1702884 ,  0.15354367,  0.18340675,  0.04241502,
        0.10249607,  0.22809335,  0.12651545,  0.23683734,  0.22493224,
        0.12782441,  0.20988456,  0.22860218,  0.2507482,  0.12267993,
        0.13156024,
        -0.03406469, -0.2332714 ,  0.0600442 ,  0.214549 ,  -0.23066802,
        0.18642221,  0.15245473,  0.06054163,  0.03416739,  0.19068498,
        0.36653106, -0.1059357 ,  0.08954779, -0.08980704, -0.15277129,
        0.20318988,  0.22250336,  0.13684609,  0.129946518,  0.18355863,
        0.27958414, -0.219293604, -0.04550122, -0.19929359, -0.21898546,
        0.17256296,  0.14425364,  0.09852625, -0.00753437,  0.14261944,
        0.27570208]])
```

Explained Variance:

The explained variance tells you how much information (variance) can be attributed to each of the principal components. This is important as you can convert n-dimensional space to 2-dimensional space, you lose some of the variance (information).

```
pca.explained_variance_ratio_
```

```
array([0.42864701,  0.18376792])
```

Three Principal Components

```
pca_3 = PCA(n_components=3)
pca_3.fit(scaled_data)
x_pca_3 = pca_3.transform(scaled_data)
```

In this Numpy matrix array, each row represents a principal component, and each column relates back to the original features. You can visualize this relationship with a heatmap:

```
x_pca_3.shape
```

```
(569, 3)
```

What is the total variance attributed by three Components?

```
pca_3.explained_variance_ratio_
```

```
array([0.42864701,  0.18376792,  0.09146436])
```

Gist of LDA

LDA identifies the linear combination of the observed variables that maximize the class separation on the availability of the prior information of the classes.

For example:

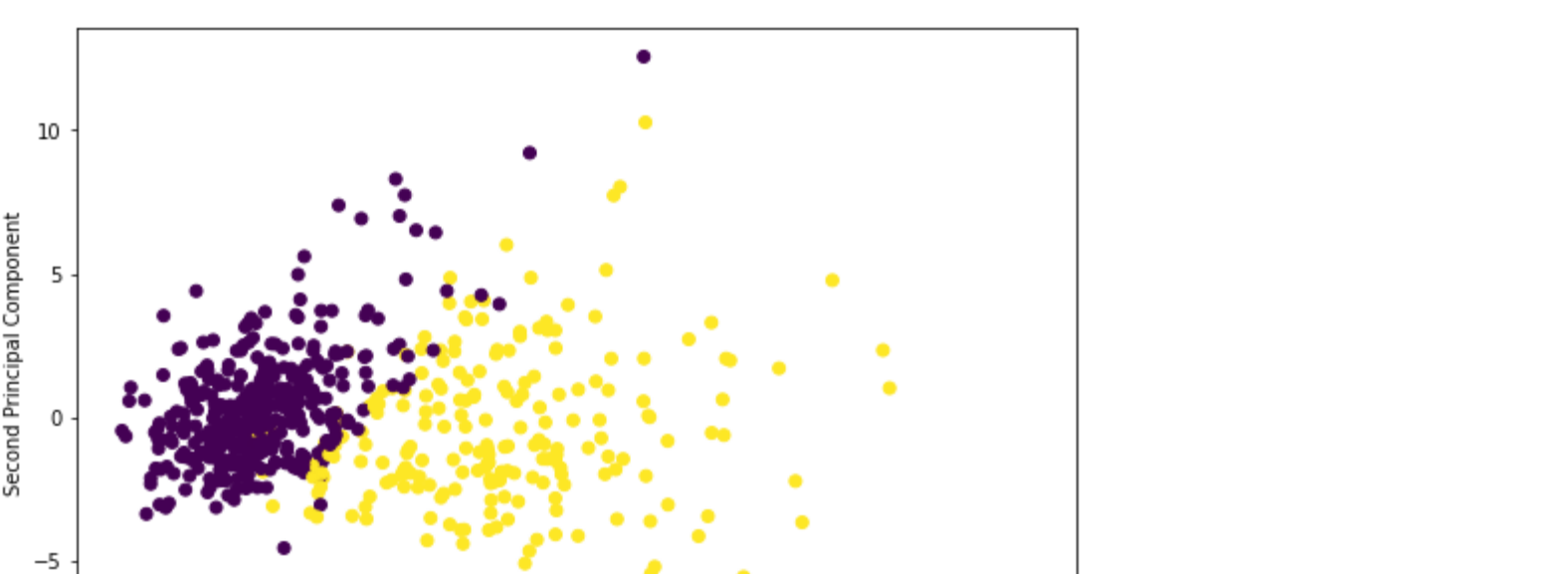
A variable in the training dataset that specifies the class of each observation

LDA Process

- Assume a set of D-dimensional samples $X(1, X(2), \dots, X(N)$, N1 of which belong to class w_1 and N2 to class w_2
- Obtain a scalar y by projecting the samples x onto a line: $Y = W^T X$
- Of all the possible lines, select the one that maximizes the separability of the scalars

Line of Maximum Separability

The maximum separable line finds out the feature subspace such that class separability is also optimized.



Finding Maximum Separable Line

1. Measure of Separation

- Better class separability
- Larger difference between the means
- Mean vector within each class of x and y is:

$$\mu_1 = \frac{1}{N_1} \sum_{x \in \mu_1} x \quad \text{and} \quad \mu_2 = \frac{1}{N_2} \sum_{y \in \mu_2} y = \frac{1}{N} \sum_{x \in \mu_2} W^T x = W^T \mu_2$$

- Objective function is the distance between the projected means:

$$J(W) = \|\mu_1 - \mu_2\| = \|W^T (\mu_1 - \mu_2)\|$$

2. Linear Discriminant

- Function of difference between the means normalized by measure of scatter (an equivalent variance):

$$J(W) = \frac{\|\mu_1 - \mu_2\|^2}{S_1^2 + S_2^2}$$

Variables from same class are projected very close to each other and at the same time, the projected means are as far apart as possible

Note: Variables from same class are projected very close to each other and at the same time, the projected means are as far apart as possible

3. Optimum Projection

- A measure of the scatter in multivariate feature space x :

$$S_1 = \sum (x - \mu_1)(x - \mu_1)^T$$

$$S_1 + S_2 = S_W$$

S_W is the within class scatter matrix

List-Group">

- Scatter of the projection y can be expressed as a function of the scatter matrix in feature space x

$$\tilde{S}_1 = \sum (y - \tilde{\mu}_1)(y - \tilde{\mu}_1)^T = \sum (W^T x - \tilde{\mu}_1)(W^T x - \tilde{\mu}_1)^T = W^T (x - \mu_1)(x - \mu_1)^T W = W^T S_1 W$$

$$S_W = \sum (y - \tilde{\mu}_1)(y - \tilde{\mu}_1)^T = W^T S_1 W$$

List-Group">

- Difference between the projected means

$$\|\tilde{\mu}_1 - \tilde{\mu}_2\|^2 = (W^T \mu_1 - W^T \mu_2)^T (W^T \mu_1 - W^T \mu_2) = W^T (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T W = W^T S_B W$$

4. Obtain the Maxima

List-Group">

- Express the linear discriminant in terms of S_W and S_B :

$$J(W) = \frac{W^T S_B W}{W^T S_W W}$$

List-Group">

- Find the maxima of $J(W)$ by differentiating and equating to zero

$$W^* = \arg \max_W \frac{W^T S_B W}{W^T S_W W} = S_W^{-1} (\mu_$$

Attribute Information:

1. ID number
2. Diagnosis (M = malignant, B = benign)
3. Attributes with mean values:
- Ten real-valued features are computed for each cell nucleus:
- radius_mean (mean of distances from center to points on the perimeter)
 - texture_mean (standard deviation of gray-scale values)
 - perimeter_mean
 - area_mean
 - smoothness_mean (local variation in radius lengths)
 - compactness_mean (perimeter² / area - 1.0)
 - concavity_mean (severity of concave portions of the contour)
 - concave points_mean (number of concave portions of the contour)
 - symmetry_mean
 - fractal_dimension_mean ("coastline approximation" - 1)
4. Attributes with standard error and worst/largest:

- radius_se
- texture_se
- perimeter_se
- area_se
- smoothness_se
- compactness_se
- concavity_se
- concave points_se
- symmetry_se
- fractal_dimension_se
- radius_worst
- texture_worst
- perimeter_worst
- area_worst
- smoothness_worst
- compactness_worst
- concavity_worst
- concave points_worst
- symmetry_worst
- fractal_dimension_worst

Solution

Import and Check the Data # CE Review

```
In [58]: import pandas as pd
import numpy as np
from matplotlib import pyplot as plt
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
from sklearn.model_selection import KFold, cross_val_score, GridSearchCV, train_test_split
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
import seaborn as sns
```

In the above code, you are importing different files: pandas, numpy, pyplot, LinearDiscriminantAnalysis, KFold, cross_val_score, GridSearchCV, train_test_split, accuracy_score, classification_report, confusion_matrix and so to know them better, you can start the notebook from the beginning.

```
In [59]: # Load data
df = pd.read_csv('breast-cancer-data.csv')
```

Before reading data from a csv file, you need to download the "breast-cancer-data.csv" dataset from the resource section and upload it into the lab. We will use the Up arrow icon, which is shown on the left side under the View icon. Click on the Up arrow icon and upload the file wherever it has been downloaded into your system.

After this, you will see the downloaded file will be visible on the left side of your lab with all the .pynb files.

```
In [60]: df = df.head(559)
```

```
In [61]: df.shape
```

Out[61]: (559, 32)

df.shape will show the number of rows and columns in the dataframe.

```
In [62]: # map categorical variable 'diagnosis' into numeric
df.diagnosis = df.diagnosis.map({'M': 1, 'B': 0})
```

In the above code, you are encoding the column diagnosis in which we are encoding M as 1 and B as 0.

```
In [63]: # head will show the rows and () default take 5 top rows as output.
df.head(5)
```

Out[63]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	point
0	842302	1	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	
1	842517	1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0669	
2	84300903	1	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	
3	84348301	1	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	
4	84158402	1	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	

5 rows × 32 columns

```
In [64]: df.columns
```

Out[64]: Index(['id', 'diagnosis', 'radius_mean', 'texture_mean', 'perimeter_mean', 'area_mean', 'smoothness_mean', 'compactness_mean', 'concavity_mean', 'concave points_mean', 'symmetry_mean', 'fractal_dimension_mean', 'radius_se', 'texture_se', 'perimeter_se', 'area_se', 'smoothness_se', 'compactness_se', 'concavity_se', 'concave points_se', 'symmetry_se', 'fractal_dimension_se', 'radius_worst', 'texture_worst', 'perimeter_worst', 'area_worst', 'smoothness_worst', 'compactness_worst', 'concavity_worst', 'concave points_worst', 'symmetry_worst', 'fractal_dimension_worst'], dtype='object')

```
In [65]: # Drop redundant column 'id'
df.drop('id', axis=1, inplace=True)
```

```
In [66]: # Check for NA values
df.isna().any()
```

Out[66]:

	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	fractal_dimension_mean	radius_se	texture_se	perimeter_se	area_se	smoothness_se	compactness_se	concavity_se	concave points_se	symmetry_se	fractal_dimension_se	radius_worst	texture_worst	perimeter_worst	area_worst	smoothness_worst	compactness_worst	concavity_worst	concave points_worst	symmetry_worst	fractal_dimension_worst	dtype: bool
	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False

```
In [68]: X_train, X_val, y_train, y_val = train_test_split(df.iloc[:,2-1], df['diagnosis'],
train_size=0.8, test_size=0.2, random_state=120)
```

In the above code, you split the data into training and validation sets (since we already have a separate test set).

```
In [69]: from sklearn.preprocessing import Normalizer
norm = Normalizer()
norm.fit(X_train)
X_train_norm = norm.transform(X_train)
X_val_norm = norm.transform(X_val)
```

In above code, you are Normalizing the features.

LDA

Let's throw in linear classifiers, i.e., Linear Discriminant Analysis for good measure and because our dataset is small.

```
In [70]: lda = LinearDiscriminantAnalysis()
lda.fit(X_train_norm, y_train)
lda_predicted = lda.predict(X_val_norm)

print('LDA Accuracy is: {}'.format(accuracy_score(y_val, lda_predicted)))
print('LDA Classification Report')
print(classification_report(y_val, lda_predicted))
```

LDA Accuracy is: 1.0
LDA Classification Report

	precision	recall	f1-score	support
0	1.00	1.00	1.00	76
1	1.00	1.00	1.00	36
accuracy			1.00	112
macro avg	1.00	1.00	1.00	112
weighted avg	1.00	1.00	1.00	112

- In the above code, you are doing the Linear Discriminant Analysis, in which you find the LDA Accuracy and LDA Classification Report.
- Analogous to the principle of beam search, let us stick with the three best algorithms, LDA, RFC, and GBC, and tune these models further, the one that gives the highest accuracy after tuning. Instantiate a new LDA model and check its performance.

```
In [71]: confusion_matrix_lda = pd.DataFrame(confusion_matrix(y_val, lda_predicted),
index = ['Actual Negative', 'Actual Positive'],
columns = ['Predicted Negative', 'Predicted Positive'] )
confusion_matrix_lda
```

Out[71]:

	Predicted Negative	Predicted Positive
Actual Negative	76	0
Actual Positive	0	36