

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 Techniques
 - Pros and Cons of Dimensionality Reduction
 - Factor Analysis

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:

- Imputing missing values
- Dropping low-variance variables
- Decision trees (DT)
- Random forest (RF)
- Reducing highly correlated variables
- Backward feature elimination
- Factor analysis

Pros of Dimensionality Reduction

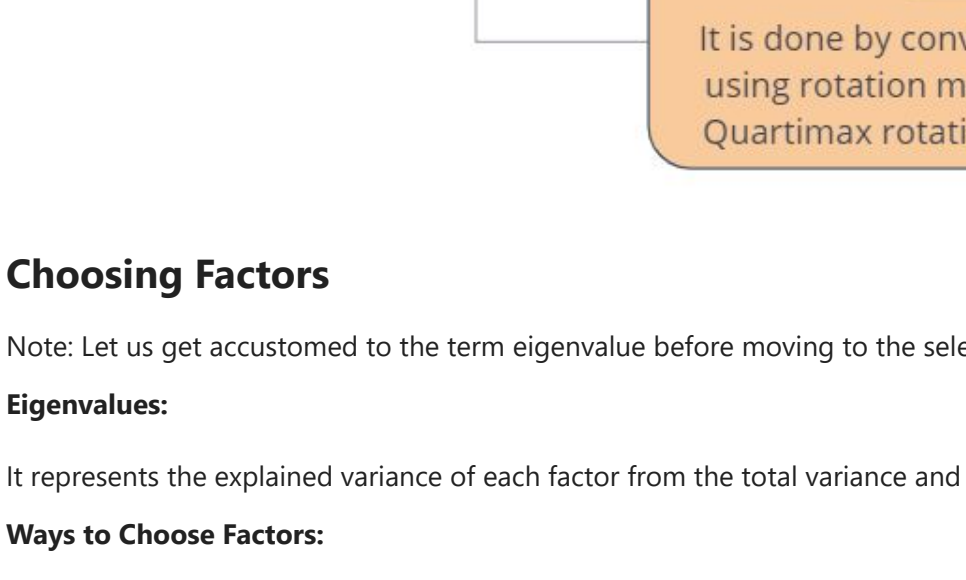
- It helps to compress data, reducing the storage space needed.
- It cuts down on computing time.
- It also aids in the removal of redundant features.

Cons of Dimensionality Reduction

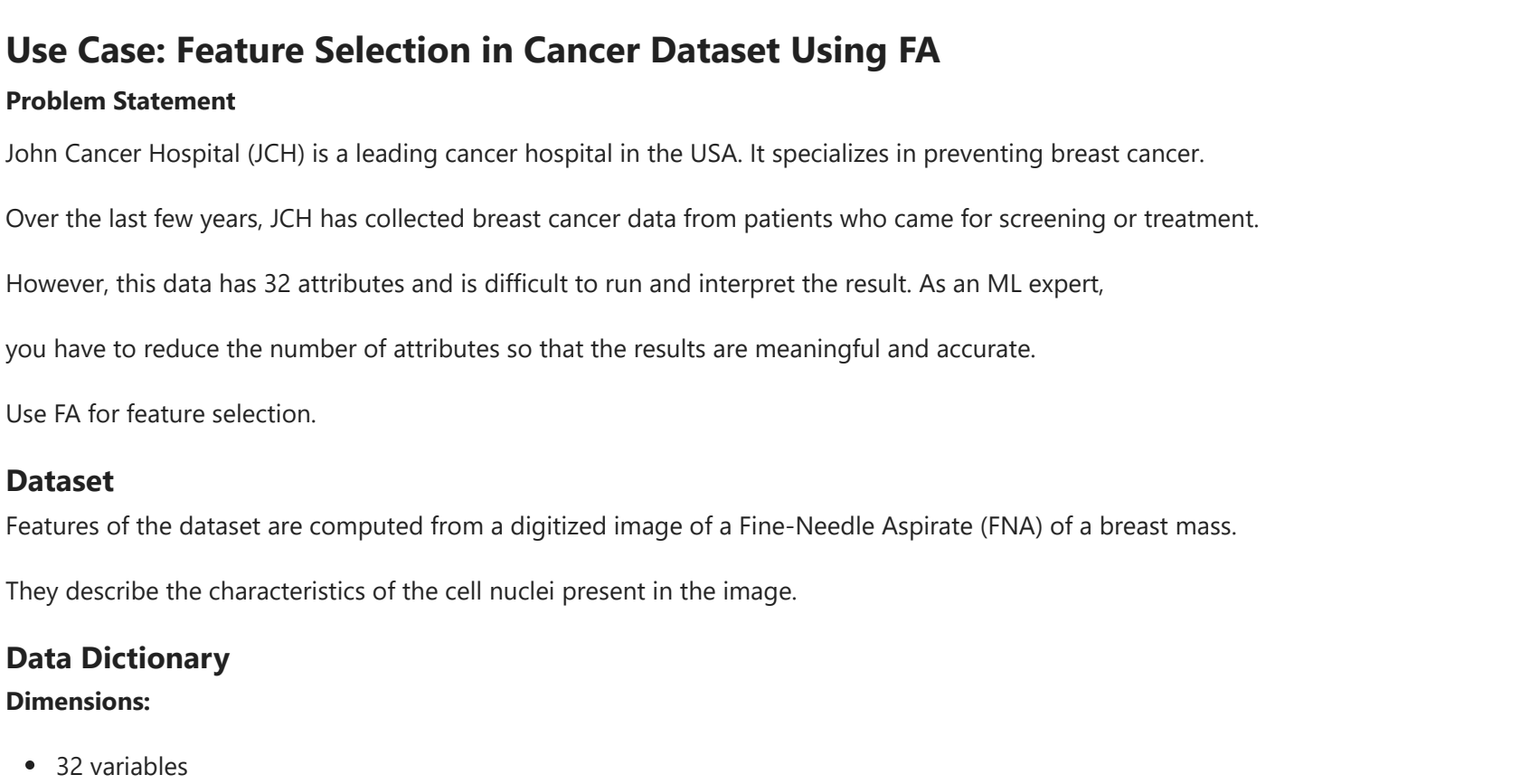
- Some data may be lost as a result.
- We use certain thumb rules when we do not know how many principal components to keep in practice.

Gist of Factor Analysis

- Factor analysis is used to:
 - Explain variance among the observed variables
 - Condense the set of observed variables into the factors
$$Y_i = \beta_{i0} + \beta_{i1}F_1 + \beta_{i2}F_2 + (1)e_i$$
- 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,....., VL with a smaller number of unobservable factors F1, F2,....., Fk.



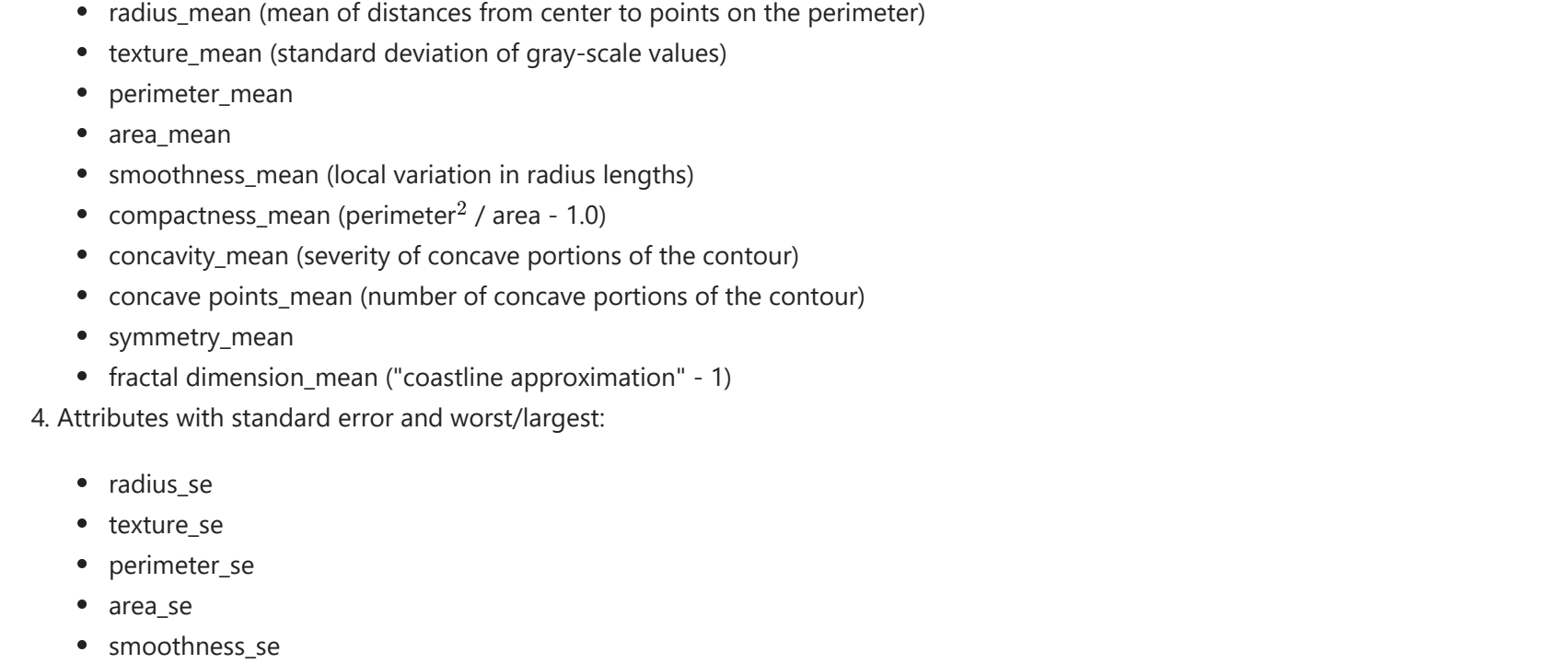
Types of FA



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 feature.
- Apart from observing plots, 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:

- ID number
- Diagnosis (M = malignant, B = benign)
- 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 (perimeter^2 / 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)
- 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 source 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 .pybn files.

```
In [2]: 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	0.3001	
1	842517	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	
2	84300903	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	
3	84348301	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 × 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 () default take the 5 top rows as output.
- one more example - df.head() will show the top 5 rows.

shape function

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

(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    object
 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_worst          569 non-null    float64
13  texture_worst         569 non-null    float64
14  perimeter_worst       569 non-null    float64
15  area_worst            569 non-null    float64
16  smoothness_worst      569 non-null    float64
17  compactness_worst     569 non-null    float64
18  concavity_worst       569 non-null    float64
19  concave points_worst  569 non-null    float64
20  symmetry_worst        569 non-null    float64
21  fractal_dimension_worst 569 non-null    float64
22  radius_se             569 non-null    float64
23  texture_se            569 non-null    float64
24  perimeter_se          569 non-null    float64
25  area_se               569 non-null    float64
26  smoothness_se        569 non-null    float64
27  compactness_se       569 non-null    float64
28  concavity_se         569 non-null    float64
29  concave points_se     569 non-null    float64
30  symmetry_se          569 non-null    float64
31  fractal_dimension_se  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',
'radius_se', 'texture_se', 'perimeter_se', 'area_se',
'compactness', 'mean compactness', 'mean concavity',
'concave points', 'mean symmetry', 'mean fractal dimension',
'radius_worst', 'texture_worst', 'perimeter_worst', 'area_worst',
'smoothness_worst', 'compactness_worst', 'concavity_worst',
'concave points_worst', 'symmetry_worst', 'fractal dimension_worst',
'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'])
```

```
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 a format that is suitable for downstream estimators.
- LabelEncoder encodes labels with an integer 0 and n_classes-1 where n is the number of distinct labels.
- These libraries are written with a value 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

-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: numpy in c:\users\alpika.gupta\anaconda3\lib\site-packages (from factor_analyzer) (1.20.3)
Requirement already satisfied: pandas in c:\users\alpika.gupta\anaconda3\lib\site-packages (from factor_analyzer) (1.20.3)
Requirement already satisfied: python-dateutil in c:\users\alpika.gupta\anaconda3\lib\site-packages (from pandas->factor_analyzer) (2.8.2)
Requirement already satisfied: six>=1.3 in c:\users\alpika.gupta\anaconda3\lib\site-packages (from python-dateutil->factor_analyzer) (1.10.0)
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.
- pip is used to install the packages.
- 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, we 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.factor_analyzer import calculate_bartlett_sphericity
from factor_analyzer.factor_analyzer import calculate_bartlett_sphericity
chi_square_value, p_value=calculate_bartlett_sphericity(df)
kmo_model=calculate_kmo(df)
```

(40227.48570934462, 0.0)

- In the above code, you are installing 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 is not an 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_corr)
```

(0.2599517032832666, 0.2599517032832666)

- 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)
```

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	point	cor
0	842302	1	17.99	10.38	122.8	1001.0	0.11840	0.27760	0.3001		0.1
1	842517	1	20.57	17.77	132.9	1326.0	0.08474	0.07864	0.0869		0.0

2 rows × 32 columns

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')
```

```

The model specification object.
Type: ModelSpecification

loadings_

The factor loadings matrix.
Type: numpy array

errorvars

The error variance matrix
Type: numpy array

factorvarcovs

The factor covariance matrix.
Type: numpy array

loglikelihood

The log likelihood from the optimization routine.
Type: float

```



```
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
cfal = ConfirmatoryFactorAnalyzer(model_spec, disp=False)
cfal.fit(df_corr.values)

C:\Users\alpika.gupta\Anaconda3\lib\site-packages\factor_analyzer\confirmatory_factor_analyzer.py:732: UserWarn
ing: The optimization routine failed to converge: ABNORMAL_TERMINATION_IN_LNSRCH
warnings.warn('The optimisation routine failed')

Out[29]:
ConfirmatoryFactorAnalyzer
ConfirmatoryFactorAnalyzer(disp=False, n_obs=569,
specification=ConfirmatoryFactorAnalyzer.confirmatory_factor_analyzer.ModelSpecification obj
ect at 0x00000160F76C4A68)

In [30]:
# cfal.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.
cfal.loadings_

Out[30]:
array([[ 9.52388865e+00,  0.00000000e+00],
       [ 5.47131803e+00,  0.00000000e+00],
       [ 7.86488967e+03,  0.00000000e+00],
       [-1.49873855e+02,  0.00000000e+00],
       [ 9.11336000e+00,  0.00000000e+00],
       [ 1.32123059e+04,  0.00000000e+00],
       [ 0.00000000e+00, -1.98153073e+03],
       [ 0.00000000e+00, -1.45489332e+03],
       [ 0.00000000e+00, -3.13842115e+02],
       [ 0.00000000e+00,  4.12173570e+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
cfal.factor_varcovs_

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]:
cfal.transform(df_corr.values)

Out[33]:
array([[ 0.08030754,  0.00103094],
       [ 0.08199898,  0.00108067],
       [ 0.06382796,  0.0010036 ],
       ...,
       [ 0.01961916,  0.0012       ],
       [ 0.07253346 ,  0.00114118],
       [-0.04865471, -0.00067256]])
```

Note: In this lesson, we saw the use of the feature selection methods, but in the next lesson we are going to use one of these methods as a sub-component of “Supervised Learning - Regression and Classification”.