

Problem1

```
In [62]: import seaborn as sns
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import scipy
from sklearn import Input
from sklearn.decomposition import PCA
from scipy import spatial
from numpy.linalg import norm
import math
from sklearn.metrics import jaccard_score

In [2]: data=sns.load_dataset("titanic")

In [3]: data.head()
```

survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	adult	male	deck	embark_town	alive	alone
0	0	3	male	22.0	1	0	7.2500	S	Third	man	True	NaN	Southampton	no	False
1	1	1	female	38.0	1	0	71.2833	C	First	woman	False	C	Cherbourg	yes	False
2	1	1	woman	38.0	0	0	53.0000	S	Third	woman	False	NaN	Southampton	yes	True
3	1	1	male	35.0	1	0	53.0000	S	First	woman	False	C	Southampton	yes	False
4	0	3	male	35.0	0	0	53.0000	S	Third	man	True	NaN	Southampton	no	True

```
In [4]: data.shape
Out[4]: (891, 145)

In [5]: data['survived'].value_counts()
Out[5]: 0    549
      1    342
      Name: survived, dtype: int64

In [6]: survived=data[data["survived"]==1]
survived.data.shape
Out[6]: (342, 145)

In [7]: sns.set_style('darkgrid')
survived.data.head()
```

```
Out[7]: <matplotlib.axes._subplots.AxesSubplot at 0x24afdc3088>
```

Of the 342 people survived mannum number of the survivors are from the age group of 20 to 40 followed by age group of 0 to 10

```
In [8]: sns.catplot(x="sex", kind="count", data=survived.data)
Out[8]: <seaborn.axisgrid.FacetGrid at 0x24afdc0ebc8>
```

Of the people who survived females are greater in number compared to males

```
In [9]: survived.data.describe()
Out[9]:
```

	survived	pclass	age	sibsp	parch	fare
count	342.0	342.000000	296.000000	342.000000	342.000000	342.000000
mean	0.0	1.968902	26.100000	0.475944	0.648012	46.990000
std	1.0	0.963321	14.960902	0.709888	0.771312	66.569986
min	0.0	1.000000	0.420000	0.000000	0.000000	0.000000
25%	1.0	1.000000	19.000000	0.000000	0.000000	12.470000
50%	1.0	2.000000	28.000000	0.000000	0.000000	26.000000
75%	1.0	3.000000	36.000000	1.000000	1.000000	57.000000
max	1.0	3.000000	80.000000	4.000000	5.000000	512.920000

Of the people who survived the middle 50 % of the survivors are of the age group 19 to 36 which confirms our earlier analysis through histogram

Problem2

```
In [30]: autodata=pd.read_csv("https://archive.ics.uci.edu/ml/machine-learning-databases/auto-mpg/auto-mpg.data",delim_whitespace=True,names=["mpg","cylinders","displacement","horsepower","weight","acceleration","model_year","origin","car_name"],na_values=["n"])

In [31]: autodata.head()
Out[31]:
```

mpg	cylinders	displacement	horsepower	weight	acceleration	model_year	origin	car_name	
0	18.0	8	307.0	1300	300.0	11.0	70	1	chevrolet chevelle malibu
1	15.0	8	300.0	1650	300.0	12.0	70	1	buick skylark 320
2	18.0	8	300.0	1600	300.0	11.0	70	1	plymouth satellite
3	16.0	8	344.0	1500	340.0	12.0	70	1	amc rebel sst
4	17.0	8	302.0	1400	340.0	10.5	70	1	ford torino

```
In [32]: autodata["horsepower"].isna().value_counts()
Out[32]: False    392
      True     6
      Name: horsepower, dtype: int64

In [33]: autodata.describe()
Out[33]:
```

	mpg	cylinders	displacement	horsepower	weight	acceleration	model_year	origin
count	398.000000	398.000000	398.000000	392.000000	398.000000	398.000000	398.000000	398.000000
mean	21.514573	5.454774	102.425879	104.469388	2970.424623	15.568000	76.010050	1.572864
std	7.815884	1.701004	104.269338	38.491150	846.841774	2.757659	3.697267	0.820255
min	9.000000	3.000000	61.000000	46.000000	1613.000000	8.000000	70.000000	1.000000
25%	17.500000	4.000000	104.000000	70.000000	2225.000000	13.500000	73.000000	1.000000
50%	23.000000	4.000000	146.000000	93.500000	2963.000000	15.500000	76.000000	1.000000
75%	29.000000	8.000000	262.000000	126.000000	3608.000000	17.170000	79.000000	2.000000
max	46.000000	8.000000	455.000000	230.000000	5140.000000	24.800000	82.000000	3.000000

```
In [34]: print("The VARIANCE before applying imputing is {}").format(autodata["horsepower"].var())
The VARIANCE before applying imputing is 1481.569329745862

In [35]: df=autodata.copy()
simpleImputer=SimpleImputer(missing_values=np.nan,strategy="mean")
df["horsepower"]=simpleImputer.fit_transform(df[["horsepower"]])
print("The VARIANCE after MEAN imputing is {}").format(df["horsepower"].var())
The VARIANCE after MEAN imputing is 1459.177916862676
```

```
In [36]: df=autodata.copy()
simpleImputer=SimpleImputer(missing_values=np.nan,strategy="median")
df["horsepower"]=simpleImputer.fit_transform(df[["horsepower"]])
print("The VARIANCE after MEDIAN imputing is {}").format(df["horsepower"].var())
The VARIANCE after MEDIAN imputing is 1466.96989389816

In [37]: df=autodata.copy()
simpleImputer=SimpleImputer(missing_values=np.nan,strategy="most_frequent")
df["horsepower"]=simpleImputer.fit_transform(df[["horsepower"]])
print("The VARIANCE after MODE imputing is {}").format(df["horsepower"].var())
The VARIANCE after MODE imputing is 1499.6361252164324
```

Mean imputing results in lowest variance as it replaces the value with the mean so the distribution of data doesn't change significantly

The other strategy is to drop the missing values and building a model to predict the missing values by giving them as query points

Problem 3

```
In [38]: irisdata=sns.load_dataset("iris")

In [39]: type(irisdata)
Out[39]: pandas.core.frame.DataFrame

In [20]: irisdata.shape
Out[20]: (150, 5)

In [21]: irisdata.head()
Out[21]:
```

sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2 setosa
1	4.9	3.0	1.4	0.2 setosa
2	4.7	3.2	1.3	0.2 setosa
3	4.6	3.1	1.5	0.2 setosa
4	5.0	3.6	1.4	0.2 setosa

```
In [22]: irisdata.dtypes
Out[22]: sepal_length    float64
      sepal_width    float64
      petal_length    float64
      petal_width    float64
      species        object
      dtype: object

In [23]: sns.pairplot(data=irisdata,hue="species")
plt.show()
```

```
In [24]: pca=PCA(n_components=3)
pca.fit_transform(irisdata.iloc[:,1:4])
print(pca.explained_variance_ratio_)
[0.92461872 0.05306648 0.01718261]

In [25]: allvariance=list()
for column in range(len(irisdata.columns)-1):
    allvariance.append(irisdata.iloc[:,column].var())
    print("Variance of the {} is {}".format(irisdata.columns[column],irisdata.iloc[:,column].var()))
    print("Total Variance is {}".format(sum(allvariance)))
Variance of the sepal_length is 0.685693123642505
Variance of the sepal_width is 0.4899796184545186
Variance of the petal_length is 0.162778523489942
Variance of the petal_width is 0.0019062909821829
Total Variance is 4.572957946979867
```

```
In [26]: for column in range(len(allvariance)):
    print("percentage of variance explained by FEATURE {} is {}".format(irisdata.columns[column],allvariance[column]/sum(allvariance)))
percentage of variance explained by FEATURE sepal_length is 0.14994532899467353
percentage of variance explained by FEATURE sepal_width is 0.04154418723232823
percentage of variance explained by FEATURE petal_length is 0.00414579319797653
percentage of variance explained by FEATURE petal_width is 0.12785263964348513

From the features the petal_length has maximum variance followed by petal_width so these two are the important features to classify the iris flowers this can be seen through the pair plot

In [27]: for i in range(len(pca.explained_variance_ratio_)):
    print("percentage variance explained by principal component {} is {}".format(i+1,pca.explained_variance_ratio_[i]))
percentage variance explained by principal component 1 is 0.9246187232817271
percentage variance explained by principal component 2 is 0.053066483117667825
percentage variance explained by principal component 3 is 0.01718260988792976
```

The Principal component 1 has maximum variance i.e. 92.4% followed by principal component 2 i.e. 5.3% and the principal component 3 has least variance and explains 1% of the total data.

The max variance explained by the features is 92.14 % by the Petal_length where as the Principal component 1 explains 92.4%. So from this we can say that the Principal component 1 explains more about the data than the petal_length.

Problem 4

```
In [28]: pcadataPoints=pd.DataFrame(pca.fit_transform(irisdata.iloc[:,1:4]))

In [29]: pcadataPoints[0]
Out[29]: 0    -2.684126
      1    -2.743442
      2    -2.888991
      3    -2.785348
      4    -2.728717
      ..
      145    1.944519
      146    1.527157
      147    1.784586
      148    1.980942
      149    1.390109
      Name: 0, Length: 150, dtype: float64

In [30]: for column in range(len(irisdata.columns)-1):
    plt.scatter(x=irisdata.iloc[:,column],y=pcadataPoints[0])
    plt.xlabel(irisdata.columns[column])
    plt.ylabel("PC1")
    plt.show()
```

```
In [31]: for column in range(len(irisdata.columns)-1):
    print("The correlation B/W {} and Principal Component 1 is {}".format(irisdata.columns[column],np.corrcoef(irisdata.iloc[:,column],pcadataPoints[0])[0][1]))
The correlation B/W sepal_length & Principal Component 1 is 0.8974817619562883
The correlation B/W sepal_width & Principal Component 1 is -0.3987484724537002
The correlation B/W petal_length & Principal Component 1 is 0.9978739422413107
The correlation B/W petal_width & Principal Component 1 is 0.8664733079330669
```

We can see the petal_length has highest correlation coefficient followed by petal_width,sepal_width,sepal_length.

We can observe that the sepal_width has inverse relation.

The above scatter plot show the same i.e. Higher the correlation shows the points

Problem 5

```
In [32]: pca1=PCA(n_components=4)
pca1.fit_transform(irisdata.iloc[:,1:4])
pcadataPoints=pd.DataFrame(pca1.fit_transform(irisdata.iloc[:,1:4]))
pca1.components_
Out[32]: array([[ 0.36138659,  0.08492251,  0.85667061,  0.3582892 ],
      [ 0.69058877,  0.73910143,  0.17377256, -0.07183302 ],
      [ 0.58202985,  0.59791983,  0.07623608,  0.54833431 ],
      [ 0.33487129,  0.33374281,  0.47983995, -0.73005741]])

In [33]: for column in range(len(allvariance)):
    print("percentage of variance explained by FEATURE {} is {}".format(irisdata.columns[column],allvariance[column]/sum(allvariance)))
percentage of variance explained by FEATURE sepal_length is 0.14994532899467353
percentage of variance explained by FEATURE sepal_width is 0.04154418723232823
percentage of variance explained by FEATURE petal_length is 0.00414579319797653
percentage of variance explained by FEATURE petal_width is 0.12785263964348513

In [34]: allvariancePCA=list()
for column in range(len(pcadataPoints.columns)-1):
    allvariancePCA.append(pcadataPoints.iloc[:,column].var())
    print("Variance of the {} is {}".format(pcadataPoints.columns[column],pcadataPoints.iloc[:,column].var()))
    print("Total Variance is {}".format(sum(allvariancePCA)))
Variance of the 0 is 4.2282417690324864
Variance of the 1 is 4.242674729865339
Variance of the 2 is 0.07829960094291939
Total Variance is 4.548232504664165

In [35]: allvariance=list()
for column in range(len(irisdata.columns)-1):
    allvariance.append(irisdata.iloc[:,column].var())
    print("Variance of the {} is {}".format(irisdata.columns[column],irisdata.iloc[:,column].var()))
    print("Total Variance is {}".format(sum(allvariance)))
Variance of the sepal_length is 0.685693123642505
Variance of the sepal_width is 0.4899796184545186
Variance of the petal_length is 0.162778523489942
Variance of the petal_width is 0.0019062909821829
Total Variance is 4.572957946979867

In [36]: for i in range(len(pca1.explained_variance_ratio_)):
    print("percentage variance explained by principal component {} is {}".format(i+1,pca1.explained_variance_ratio_[i]))
percentage variance explained by principal component 1 is 0.9246187232817271
percentage variance explained by principal component 2 is 0.053066483117667825
percentage variance explained by principal component 3 is 0.01718260988792976
percentage variance explained by principal component 4 is 0.0052123837275377
```

If we use 4 principal components it explains almost same amount of variance as that of the original features.

To explain 95 % of the variance of data we need only Principal component 1 & Principal component 2.

As Principal component 1 & Principal component 2 explains 0.9776 i.e. 97.76% of the data.

This implies we only need two principal components.

Recitation Exercises

Chapter 1

Question 1

Data mining is the process of automatically discovering useful information in large data. Data mining techniques allows us to find useful patterns that might otherwise remain unknown if mined manually.

- 1) No it doesn't need Data Mining techniques as it can be found querying database.
- 2) No it doesn't need Data Mining techniques as it can be found querying database.
- 3) No it doesn't need Data Mining techniques as it can be found querying database.
- 4) No it doesn't need Data Mining techniques as it can be found querying database.
- 5) No it doesn't need Data Mining techniques as it can be found using mathematical calculations.
- 6) Yes the outcomes of future stock prices needs substantial amount of application of Data Mining Techniques
- 7) Yes the predictions of abnormalities needs substantial amount of application of Data Mining Techniques on seismic waves data.
- 8) Yes the predictions of earthquakes needs substantial amount of application of Data Mining Techniques on seismic waves data.
- 9) No it does not need data mining techniques as it can be done using basic signal processing

Chapter 2

Question 2

- 1) Time AM PM: Binary, qualitative, ordinal.
- 2) Brightness light meter: Continuous quantitative ratio.
- 3) Brightness people: Discrete qualitative ordinal.
- 4) Angles 0 to 360: Continuous quantitative ratio.
- 5) Medalist Olympic: Discrete qualitative ordinal.
- 6) Medalist Sea level: Continuous quantitative ratio.
- 7) patients: Discrete qualitative ordinal.
- 8) ISBN: Discrete qualitative ordinal.
- 9) pass light: Discrete qualitative ordinal.
- 10) Military rank: Discrete qualitative ordinal.
- 11) Distance from campus: Continuous quantitative ratio.
- 12) Density: Continuous quantitative ratio.
- 13) Coat check number: Discrete qualitative ordinal.

Question 7

Daily temperature has more autocorrelation & deals uses the concept of temporal autocorrelation which says, if two measurements are close in time then the values of those measurements are often very similar. In realistic scenario temperature across closer locations are very similar where as rainfall is not.

Question 15

- 1) We get equal number of elements from each group.
- 2) We don't guarantee any specific number from each group the elements are picked randomly.

Question 16

- 1) Term that occurs in every document has very less weight or 0' weight
- 2) Term that occurs only once has very high weight or maximum weight

Question 17

- 1) A, B) = (a^2+b^2)
- 2) y=x^2

Question 18

```
1)
x=0101010001
y=0001001000
Jaccard Distance 2/(10-5)=2/5
Hamming Distance 3
2)
Simple Matching Coefficient and Hamming distance are equivalent.As the hamming distance gives how many different points are present.Where as Simple Matching Coefficient gives the number of same attributes divided by total number of data points.
Cosine is similar to Jaccard as both look one type of attributes in the entire data
3)
jaccard distance seems more appropriate as it can be shown the similarity or common genes between the points
4)
Hamming distance we want to know the difference between two genes. Hamming distance is 2
```

```
In [86]: points=[[[1,1,1],[0,1,0,1],[0,-1,0,1],[1,1,0,1,0,1],[2,-1,0,2,0,-3]]]
points=[[[2,2,2],[1,0,0,1],[0,-0,-1,0],[1,1,1,0,0,0,1],[1,-1,-1,-0,0,-1]]]
disjoints=[['E','E','E']]
for i in range(len(points)):
    array=np.array(points[i])
    print("calculating the below metrics for points {} and {}".format(array,array2))
    print("un")
    if distancetype=="E":
        print("The euclidean Distance is {}".format(np.linalg.norm(array-array2))
    elif distancetype=="J":
        print("The jaccard value is {}".format(jaccard_score(array,array2)))
    elif distancetype=="C":
        print("The cosine Distance is {}".format(np.linalg.norm(array-array2))
    print("The Jaccard value is {}".format(jaccard_score(array,array2)))
    print("The Cosine Coefficient is {}".format(np.corrcoef(array,array2)[0][1]))
    print("The Correlation Coefficient is {}".format(np.corrcoef(array,array2)[0][1]))
    print("The cosine Coefficient is {}".format(np.corrcoef(array,array2)[0][1]))
    univector2 = points[i][1] / np.linalg.norm(points[i][1])
    univector1 = points[i][0] / np.linalg.norm(points[i][0])
    dotProduct=np.dot(univector1,univector2)
    x=np.arccos(dotProduct)
    print("The cosine Similarity is {} and angle between the points is {}".format(scipy.spatial.distance.cosine(a
rray,array2),
    math.cos(x)))
    print("un")
calculating the Below metrics for points [1 1 1 1] and [2 2 2 2]
```

The euclidean Distance is 2.0
The Correlation Coefficient is nan
The Cosine Similarity is 0.8 and angle between the points is 1.8

calculating the Below metrics for points [0 1 0 1] and [1 0 -1 0]

The euclidean Distance is 2.0
The jaccard value is 0.0
The Correlation Coefficient is -1.0
The cosine Similarity is 1.0 and angle between the points is 6.123233995736766e-17

calculating the Below metrics for points [0 -1 0 1] and [1 0 -1 0]

The euclidean Distance is 2.0
The Correlation Coefficient is 0.8
The cosine Similarity is 1.0 and angle between the points is 6.123233995736766e-17

calculating the Below metrics for points [1 1 0 1 0 1] and [1 1 0 1 0 1]

The jaccard value is 0.6
The Correlation Coefficient is 0.25000000000000000
The cosine Similarity is 0.25 and angle between the points is 0.75

calculating the Below metrics for points [2 -1 0 2 0 -3] and [-1 -1 -1 0 0 -1]

The Correlation Coefficient is -2.866583523299565e-17
The cosine Similarity is 1.0 and angle between the points is 6.123233995736766e-17

Question 19

From the above we can see for points

```
1) x=(1,1,1),y=(2,2,2)
cosine=1.0
Euclidean=2.0
2) x=(0,1,0,1),y=(0,0,1,0)
cosine=0.5
cosine=0.5
cosine=0.5
Euclidean=2.0
Jaccard=0.0
3) x=(0,-1,0,1,0,1),y=(1,0,-1,0,-1,0)
cosine=0.0
cosine=0.0
cosine=0.0
Euclidean=2.0
4) x=(1,0,1,0,1,0),y=(1,1,0,1,0,1)
cosine=0.75
correlation=0.25
Jaccard=0.0
5) x=(2,-1,0,2,0,-3),y=(-1,-1,0,0,-1,0)
cosine=0.5
cosine=0.5
cosine=0.5
Euclidean=2.0
Jaccard=0.0
6) x=(2,-1,0,2,0,-3),y=(-1,-1,0,0,-1,0)
cosine=0.5
cosine=0.5
cosine=0.5
Euclidean=2.0
Jaccard=0.0
```