



## Data Mining

### Lab-5 - Data Preprocessing

137 | Vishal Baraiya |  
23010101014

**1) First, you need to read the titanic dataset from local disk and display Last five records**

```
In [3]: import pandas as pd  
import numpy as np
```

```
In [4]: # Try reading the CSV with ISO-8859-1 encoding  
df = pd.read_csv("titanic.csv")  
df
```

Out[4]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.25
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...)	female	38.0	1	0	PC 17599	71.28
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.92
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.10
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.05
...	...	...	...	...	...	...	...	...	...	...
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.00
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.00
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.45
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.00
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.73

891 rows × 12 columns



In [5]: `df.tail(5)`

Out[5]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.00
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.00
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.45
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.00
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.75

## 2) Handle Missing Values in data set [use dropna(), fillna(), and interpolate]

In [6]: `data_dropna = df.dropna()  
data_dropna`

Out[6]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Far
<b>1</b>	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.283
<b>3</b>	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.100
<b>6</b>	7	0	1	McCarthy, Mr. Timothy J	male	54.0	0	0	17463	51.862
<b>10</b>	11	1	3	Sandstrom, Miss. Marguerite Rut	female	4.0	1	1	PP 9549	16.700
<b>11</b>	12	1	1	Bonnell, Miss. Elizabeth	female	58.0	0	0	113783	26.550
...	...	...	...	...	...	...	...	...	...	...
<b>871</b>	872	1	1	Beckwith, Mrs. Richard Leonard (Sallie Monypeny)	female	47.0	1	1	11751	52.554
<b>872</b>	873	0	1	Carlsson, Mr. Frans Olof	male	33.0	0	0	695	5.000
<b>879</b>	880	1	1	Potter, Mrs. Thomas Jr (Lily Alexenia Wilson)	female	56.0	0	1	11767	83.158
<b>887</b>	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.000
<b>889</b>	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.000

183 rows × 12 columns



```
In [7]: # delete row in which every column value is null
# data_dropna = df.dropna(how='all')

# delete row in which any
# data_dropna = df.dropna(how='any',axis = 1)
```

```
In [11]: # using fillna
data_fillna = df.fillna(30)
data_fillna
```

Out[11]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
--	-------------	----------	--------	------	-----	-----	-------	-------	--------	------

0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.25
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.28
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.92
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.10
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.05
...	...	...	...	...	...	...	...	...	...	...
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.00
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.00
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	30.0	1	2	W./C. 6607	23.45
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.00
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.75

891 rows × 12 columns



```
In [12]: data_fillna = df.fillna({'Age' : 35, 'Cabin' : 'Not Available'})  
data_fillna
```

Out[12]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.25
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.28
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.92
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.10
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.05
...	...	...	...	...	...	...	...	...	...	...
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.00
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.00
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	35.0	1	2	W./C. 6607	23.45
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.00
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.75

891 rows × 12 columns





```
In [14]: age_mean = df.Age.mean()  
data_fillna = df.fillna({'Age':age_mean})  
data_fillna
```

Out[14]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket
<b>0</b>	1	0	3	Braund, Mr. Owen Harris	male	22.000000	1	0	A/5 21171
<b>1</b>	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.000000	1	0	PC 17599
<b>2</b>	3	1	3	Heikkinen, Miss. Laina	female	26.000000	0	0	STON/O2. 3101282
<b>3</b>	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.000000	1	0	113803
<b>4</b>	5	0	3	Allen, Mr. William Henry	male	35.000000	0	0	373450
...	...	...	...	...	...	...	...	...	...
<b>886</b>	887	0	2	Montvila, Rev. Juozas	male	27.000000	0	0	211536
<b>887</b>	888	1	1	Graham, Miss. Margaret Edith	female	19.000000	0	0	112053
<b>888</b>	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	29.699118	1	2	W./C. 6607
<b>889</b>	890	1	1	Behr, Mr. Karl Howell	male	26.000000	0	0	111369
<b>890</b>	891	0	3	Dooley, Mr. Patrick	male	32.000000	0	0	370376

891 rows × 12 columns



```
In [16]: data_interpolate = df.interpolate()  
data_interpolate
```

C:\Users\ASUS\AppData\Local\Temp\ipykernel\_1768\2280711911.py:1: FutureWarning: Data Frame.interpolate with object dtype is deprecated and will raise in a future version. Call obj.infer\_objects(copy=False) before interpolating instead.  
data\_interpolate = df.interpolate()

Out[16]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.25
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.28
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.92
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.10
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.05
...	...	...	...	...	...	...	...	...	...	...
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.00
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.00
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	22.5	1	2	W./C. 6607	23.45
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.00
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.75

891 rows × 12 columns



### 3) Apply Scaling to AGE attribute with min max, decimal scaling and z score.

```
In [17]: df.fillna(df.Age.mean(), inplace=True)

data2 = df.copy()
minAge = df.Age.min()
maxAge = df.Age.max()
data2['MinMaxAge'] = (data2['Age']-minAge)/(maxAge-minAge)
data2
```

Out[17]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket
<b>0</b>	1	0	3	Braund, Mr. Owen Harris	male	22.000000	1	0	A/5 21171
<b>1</b>	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.000000	1	0	PC 17599
<b>2</b>	3	1	3	Heikkinen, Miss. Laina	female	26.000000	0	0	STON/O2. 3101282
<b>3</b>	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.000000	1	0	113803
<b>4</b>	5	0	3	Allen, Mr. William Henry	male	35.000000	0	0	373450
...	...	...	...	...	...	...	...	...	...
<b>886</b>	887	0	2	Montvila, Rev. Juozas	male	27.000000	0	0	211536
<b>887</b>	888	1	1	Graham, Miss. Margaret Edith	female	19.000000	0	0	112053
<b>888</b>	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	29.699118	1	2	W./C. 6607
<b>889</b>	890	1	1	Behr, Mr. Karl Howell	male	26.000000	0	0	111369
<b>890</b>	891	0	3	Dooley, Mr. Patrick	male	32.000000	0	0	370376

891 rows × 13 columns



```
In [19]: data3 = df.copy()
maxAge = df.Age.max()
noOfDigits = len(str(int(maxAge)))

data2['AgeDS'] = data2['Age'] / ( 10 ** noOfDigits )
data2
```

Out[19]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket
<b>0</b>	1	0	3	Braund, Mr. Owen Harris	male	22.000000	1	0	A/5 21171
<b>1</b>	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.000000	1	0	PC 17599
<b>2</b>	3	1	3	Heikkinen, Miss. Laina	female	26.000000	0	0	STON/O2. 3101282
<b>3</b>	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.000000	1	0	113803
<b>4</b>	5	0	3	Allen, Mr. William Henry	male	35.000000	0	0	373450
...	...	...	...	...	...	...	...	...	...
<b>886</b>	887	0	2	Montvila, Rev. Juozas	male	27.000000	0	0	211536
<b>887</b>	888	1	1	Graham, Miss. Margaret Edith	female	19.000000	0	0	112053
<b>888</b>	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	29.699118	1	2	W./C. 6607
<b>889</b>	890	1	1	Behr, Mr. Karl Howell	male	26.000000	0	0	111369
<b>890</b>	891	0	3	Dooley, Mr. Patrick	male	32.000000	0	0	370376

891 rows × 14 columns





```
In [21]: meanAge = df.Age.mean()
stdAge = df.Age.std()
data3['AgeZScore'] = (data3['Age']-minAge)/stdAge
data3
```

Out[21]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket
<b>0</b>	1	0	3	Braund, Mr. Owen Harris	male	22.000000	1	0	A/5 21171
<b>1</b>	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.000000	1	0	PC 17599
<b>2</b>	3	1	3	Heikkinen, Miss. Laina	female	26.000000	0	0	STON/O2. 3101282
<b>3</b>	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.000000	1	0	113803
<b>4</b>	5	0	3	Allen, Mr. William Henry	male	35.000000	0	0	373450
...	...	...	...	...	...	...	...	...	...
<b>886</b>	887	0	2	Montvila, Rev. Juozas	male	27.000000	0	0	211536
<b>887</b>	888	1	1	Graham, Miss. Margaret Edith	female	19.000000	0	0	112053
<b>888</b>	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	29.699118	1	2	W./C. 6607
<b>889</b>	890	1	1	Behr, Mr. Karl Howell	male	26.000000	0	0	111369
<b>890</b>	891	0	3	Dooley, Mr. Patrick	male	32.000000	0	0	370376

891 rows × 13 columns



In [ ]: