

Lab - 5 - Data Preprocessing

Mohil Parmar | 23010101192

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

```
In [4]: import pandas as pd
```

```
In [6]: df = pd.read_csv("pre_titanic.csv")
```

```
In [8]: df.head()
```

```
Out[8]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S

```
In [10]: df.tail()
```

Out[10]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.00	NaN	S
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.00	B42	S
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.45	NaN	S
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.00	C148	C
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.75	NaN	Q

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

```
In [13]: df_row = df.dropna(how="any",axis=0)
df_row
```

Out[13]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
6	7	0	1	McCarthy, Mr. Timothy J	male	54.0	0	0	17463	51.8625	E46	S
10	11	1	3	Sandstrom, Miss. Marguerite Rut	female	4.0	1	1	PP 9549	16.7000	G6	S
11	12	1	1	Bonnell, Miss. Elizabeth	female	58.0	0	0	113783	26.5500	C103	S
...
871	872	1	1	Beckwith, Mrs. Richard Leonard (Sallie Monypeny)	female	47.0	1	1	11751	52.5542	D35	S
872	873	0	1	Carlsson, Mr. Frans Olof	male	33.0	0	0	695	5.0000	B51 B53 B55	S
879	880	1	1	Potter, Mrs. Thomas Jr (Lily Alexenia Wilson)	female	56.0	0	1	11767	83.1583	C50	C
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000	B42	S
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000	C148	C

183 rows × 12 columns

In [15]:

```
df.isna().sum()
```

```
Out[15]: PassengerId      0
         Survived        0
         Pclass          0
         Name            0
         Sex             0
         Age            177
         SibSp           0
         Parch           0
         Ticket          0
         Fare            0
         Cabin          687
         Embarked        2
         dtype: int64
```

```
In [17]: df_col = df.dropna(how="any",axis=1)
         df_col
```

Out[17]:

	PassengerId	Survived	Pclass	Name	Sex	SibSp	Parch	Ticket	Fare
0	1	0	3	Braund, Mr. Owen Harris	male	1	0	A/5 21171	7.2500
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	1	0	PC 17599	71.2833
2	3	1	3	Heikkinen, Miss. Laina	female	0	0	STON/O2. 3101282	7.9250
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	1	0	113803	53.1000
4	5	0	3	Allen, Mr. William Henry	male	0	0	373450	8.0500
...
886	887	0	2	Montvila, Rev. Juozas	male	0	0	211536	13.0000
887	888	1	1	Graham, Miss. Margaret Edith	female	0	0	112053	30.0000
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	1	2	W./C. 6607	23.4500
889	890	1	1	Behr, Mr. Karl Howell	male	0	0	111369	30.0000
890	891	0	3	Dooley, Mr. Patrick	male	0	0	370376	7.7500

891 rows × 9 columns

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

In [20]: `meanAge=df.Age.mean()`

In [22]: `meanAge`

Out[22]: 29.69911764705882

In [24]: `data_withfillna=df.fillna({'Age':meanAge , 'Cabin':"Not aval."})`
`data_withfillna`

Out[24]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.000000	1	0	A/5 21171	7.2500	Not aval.	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.000000	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.000000	0	0	STON/O2. 3101282	7.9250	Not aval.	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.000000	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.000000	0	0	373450	8.0500	Not aval.	S
...
886	887	0	2	Montvila, Rev. Juozas	male	27.000000	0	0	211536	13.0000	Not aval.	S
887	888	1	1	Graham, Miss. Margaret Edith	female	19.000000	0	0	112053	30.0000	B42	S
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	29.699118	1	2	W./C. 6607	23.4500	Not aval.	S
889	890	1	1	Behr, Mr. Karl Howell	male	26.000000	0	0	111369	30.0000	C148	C
890	891	0	3	Dooley, Mr. Patrick	male	32.000000	0	0	370376	7.7500	Not aval.	Q

891 rows × 12 columns

In [28]:

```
# Step 1: Min-Max Scaling
data_withfillna['Age_MinMax'] = (data_withfillna['Age'] - data_withfillna['Age'].min()) / (data_withfillna['Age'].max() - data_withfillna['Age'].min())
```

Out[28]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked	Age_MinMax
0	1	0	3	Braund, Mr. Owen Harris	male	22.000000	1	0	A/5 21171	7.2500	Not aval.	S	0.271174
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.000000	1	0	PC 17599	71.2833	C85	C	0.472229
2	3	1	3	Heikkinen, Miss. Laina	female	26.000000	0	0	STON/O2. 3101282	7.9250	Not aval.	S	0.321438
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.000000	1	0	113803	53.1000	C123	S	0.434531
4	5	0	3	Allen, Mr. William Henry	male	35.000000	0	0	373450	8.0500	Not aval.	S	0.434531
...
886	887	0	2	Montvila, Rev. Juozas	male	27.000000	0	0	211536	13.0000	Not aval.	S	0.334004
887	888	1	1	Graham, Miss. Margaret Edith	female	19.000000	0	0	112053	30.0000	B42	S	0.233476
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	29.699118	1	2	W./C. 6607	23.4500	Not aval.	S	0.367921

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked	Age_MinMax
889	890	1	1	Behr, Mr. Karl Howell	male	26.000000	0	0	111369	30.0000	C148	C	0.321438
890	891	0	3	Dooley, Mr. Patrick	male	32.000000	0	0	370376	7.7500	Not aval.	Q	0.396833

891 rows × 13 columns

```
In [32]: # Step 2: Decimal Scaling
d = len(str(int(abs(data_withfillna['Age'].max())))) # get number of digits
data_withfillna['Age_DecimalScaling'] = data_withfillna['Age'] / (10**d)
data_withfillna
```


Out[32]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked	Age_MinMax	Age_
0	1	0	3	Braund, Mr. Owen Harris	male	22.000000	1	0	A/5 21171	7.2500	Not aval.	S	0.271174	
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.000000	1	0	PC 17599	71.2833	C85	C	0.472229	
2	3	1	3	Heikkinen, Miss. Laina	female	26.000000	0	0	STON/O2. 3101282	7.9250	Not aval.	S	0.321438	
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.000000	1	0	113803	53.1000	C123	S	0.434531	
4	5	0	3	Allen, Mr. William Henry	male	35.000000	0	0	373450	8.0500	Not aval.	S	0.434531	
...
886	887	0	2	Montvila, Rev. Juozas	male	27.000000	0	0	211536	13.0000	Not aval.	S	0.334004	
887	888	1	1	Graham, Miss. Margaret Edith	female	19.000000	0	0	112053	30.0000	B42	S	0.233476	
888	889	0	3	Johnston, Miss. Catherine	female	29.699118	1	2	W./C. 6607	23.4500	Not aval.	S	0.367921	

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked	Age_MinMax	Age_
				Helen "Carrie"										
889	890	1	1	Behr, Mr. Karl Howell	male	26.000000	0	0	111369	30.0000	C148	C	0.321438	
890	891	0	3	Dooley, Mr. Patrick	male	32.000000	0	0	370376	7.7500	Not aval.	Q	0.396833	

891 rows × 14 columns

```
In [34]: # Step 3: Z-Score Normalization
data_withfillna['Age_Zscore'] = (data_withfillna['Age'] - data_withfillna['Age'].mean()) / data_withfillna['Age'].std()
data_withfillna
```

Out[34]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked	Age_MinMax	Age_
0	1	0	3	Braund, Mr. Owen Harris	male	22.000000	1	0	A/5 21171	7.2500	Not aval.	S	0.271174	
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.000000	1	0	PC 17599	71.2833	C85	C	0.472229	
2	3	1	3	Heikkinen, Miss. Laina	female	26.000000	0	0	STON/O2. 3101282	7.9250	Not aval.	S	0.321438	
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.000000	1	0	113803	53.1000	C123	S	0.434531	
4	5	0	3	Allen, Mr. William Henry	male	35.000000	0	0	373450	8.0500	Not aval.	S	0.434531	
...
886	887	0	2	Montvila, Rev. Juozas	male	27.000000	0	0	211536	13.0000	Not aval.	S	0.334004	
887	888	1	1	Graham, Miss. Margaret Edith	female	19.000000	0	0	112053	30.0000	B42	S	0.233476	
888	889	0	3	Johnston, Miss. Catherine	female	29.699118	1	2	W./C. 6607	23.4500	Not aval.	S	0.367921	

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked	Age_MinMax	Age_
				Helen "Carrie"										
889	890	1	1	Behr, Mr. Karl Howell	male	26.000000	0	0	111369	30.0000	C148	C	0.321438	
890	891	0	3	Dooley, Mr. Patrick	male	32.000000	0	0	370376	7.7500	Not aval.	Q	0.396833	

891 rows × 15 columns

In []: