

# DATA CLEANING,MISSING VALUE TREATMENT

```
#Name : Ankita Gulde
#Roll no. : 44
#section : 3A
```

```
#Aim : TO Perform Data Processing, Data cleaning,Missing value treatment
```

```
import pandas as pd
```

```
import os
```

```
os.getcwd()
```

```
os.chdir("C:\\Users\\HP\\Desktop")
```

```
df=pd.read_csv("titanic.csv")
```

```
df
```

	pclass	survived	name	sex	age	sibsp	parch	ticket	fare	cabin	embarked	boat	body	home.dest
0	1.0	1.0	Allen, Miss. Elisabeth Walton	female	29.0000	0.0	0.0	24160	211.3375	B5	S	2	NaN	St Louis, MO
1	1.0	1.0	Allison, Master.			1.0	2.0	113781	151.5500	C22 C26	S	11	NaN	Montreal, PQ / Chesterville, ON
2	1.0	0.0	Allison, Miss. Helen Loraine	female	2.0000	1.0	2.0	113781	151.5500	C22 C26	S	NaN	NaN	Montreal, PQ / Chesterville, ON
3	1.0	0.0	Allison, Mr. Hudson Joshua Creighton	male	30.0000	1.0	2.0	113781	151.5500	C22 C26	S	NaN	135.0	Montreal, PQ / Chesterville, ON
4	1.0	0.0	Allison, Mrs. Hudson J C (Bessie Waldo Daniels)	female	25.0000	1.0	2.0	113781	151.5500	C22 C26	S	NaN	NaN	Montreal, PQ / Chesterville, ON
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
1306	3.0	0.0	Zakarian, Mr. Mapriededer	male	26.5000	0.0	0.0	2656	7.2250	NaN	C	NaN	304.0	NaN
1307	3.0	0.0	Zakarian, Mr. Ortin	male	27.0000	0.0	0.0	2670	7.2250	NaN	C	NaN	NaN	NaN
1308	3.0	0.0	Zimmerman, Mr. Leo	male	29.0000	0.0	0.0	315082	7.8750	NaN	S	NaN	NaN	NaN
1309	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

1310 rows × 14 columns

```
df.head(40)
```

	pclass	survived	name	sex	age	sibsp	parch	ticket	fare	cabin	embarked	boat	body	home.dest
0	1.0	1.0	Walton			0.0	0.0	24160	211.3375	B5	S	2	NaN	St Louis, MO
1	1.0	1.0	Allison, Master. Hudson Trevor	male	0.9167	1.0	2.0	113781	151.5500	C22 C26				Montreal, PQ / Chesterville, ON
2	1.0	0.0	Loraine	female	2.0000	1.0	2.0	113781	151.5500	C22 C26	S	NaN	NaN	Montreal, PQ / Chesterville, ON
3	1.0	0.0	Allison, Mr. Hudson Joshua Creighton	male	30.0000	1.0	2.0	113781	151.5500	C22 C26	S	NaN	135.0	Montreal, PQ / Chesterville, ON
4	1.0	0.0	Allison, Mrs. Hudson J C (Bessie Waldo Daniels)	female	25.0000	1.0	2.0	113781	151.5500	C22 C26	S	NaN	NaN	Montreal, PQ / Chesterville, ON
5	1.0	1.0	Anderson, Mr. Harry	male	48.0000	0.0	0.0	19952	26.5500	E12	S	3	NaN	New York, NY



6 1. 1  
0 .  
7 1. 0  
0 .  
0  
8 1. 1  
0 .  
9 1. 0  
0 .  
11 1. 0  
0 0 .  
0  
11 1. 1  
1 0 .  
11 1. 1  
2 0 .  
11 1. 1  
3 0 .  
11 1. 1  
4 0 .  
11 1. 0  
5 0 .  
11 1. 0  
6 0 .  
11 1. 1  
7 0 .  
11 1. 1  
8 0 .  
11 1. 0  
9 0 .  
21 1. 1  
0 0 .  
0  
21 1. 1  
1 0 .  
21 1. 1  
2 0 .  
21 1. 1  
3 0 .  
21 1. 1  
4 0 .  
0  
21 1. 0  
5 0 .  
21 1. 1  
6 0 .  
21 1. 1  
7 0 .  
21 1. 1  
8 0 .  
21 1. 1  
9 0 .  
0  
31 1. 0  
0 0 .  
31 1. 1  
1 0 .  
31 1. 1  
2 0 .  
33  
1.0  
1.0

d:17) N  
r:37' a  
e:11 L  
s:11 L  
M:11 N  
r:11 a  
t:11 R  
M  
r:11  
s:11 N  
e:11 a  
t:11 N  
An



[illegible]

Borebank, Mr. John

Bowen, Miss. Grace female Scott	45.00 00	0.0	0.0	PC 1760 8	262.37 50	NaN	C	4	NaN Coburnstown, NY
werman, Miss. Elsie Edith female	22.00 00	0.0	1.0	11350 5	55.0000	E33	S	6	NaN St Leonards-on-Sea, England Ohio
Bradley, Mr. George ("George Arthur male	NaN	0.0	0.0	11142 7	26.5500	NaN	S	9	NaN Los Angeles, CA

Bo

Brayton")

38	1.0	0.0	Brady, Mr. John											
			Bertram	male	41.0000	0.0	0.0	113054	30.5000	A21	S	NaN	NaN	Pomeroy, WA
39	1.0	0.0	Brandeis, Mr. Emil	male	48.0000	0.0	0.0	PC 17591	50.4958	B10	C	NaN	208.0	Omaha, NE

```
In [12]: df.tail(10)
```

pclass survived			name	sex	age	sibsp	parch	ticket	fare	cabin	embarked	boat	body	home.dest
1300	3.0	1.0	Yasbeck, Mrs. Antoni (Selini Alexander)	female	15.0	1.0	0.0	2659	14.4542	NaN	C	NaN	NaN	NaN
1301	3.0	0.0	Youseff, Mr. Gerious	male	45.5	0.0	0.0	2628	7.2250	NaN	C	NaN	312.0	NaN
1302	3.0	0.0	Yousif, Mr. Wazli	male	NaN	0.0	0.0	2647	7.2250	NaN	C	NaN	NaN	NaN
1303	3.0	0.0	Youseff, Mr. Gerious	male	NaN	0.0	0.0	2627	14.4583	NaN	C	NaN	NaN	NaN
1304	3.0	0.0	Zabour, Miss. Hileni	female	14.5	1.0	0.0	2665	14.4542	NaN	C	NaN	328.0	NaN
1305	3.0	0.0	Zabour, Miss. Thamine	female	NaN	1.0	0.0	2665	14.4542	NaN	C	NaN	NaN	NaN
1306	3.0	0.0	Zakarian, Mr. Mapriededer	male	26.5	0.0	0.0	2656	7.2250	NaN	C	NaN	304.0	NaN
1307	3.0	0.0	Zakarian, Mr. Ortin	male	27.0	0.0	0.0	2670	7.2250	NaN	C	NaN	NaN	NaN
1308	3.0	0.0	Zimmerman, Mr. Leo	male	29.0	0.0	0.0	315082	7.8750	NaN	S	NaN	NaN	NaN
1309	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

```
In [13]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1310 entries, 0 to 1309
Data columns (total 14 columns):
#   Column      Non-Null Count  Dtype
---  -
1.  pclass      1309 non-null   float64
2.  survived    1309 non-null   float64
3.  name        1309 non-null   object
4.  sex         1309 non-null   object
5.  age         1046 non-null   float64
6.  sibsp       1309 non-null   float64
7.  parch       1309 non-null   float64
8.  ticket      1309 non-null   object
9.  fare        1308 non-null   float64
10. cabin      295 non-null    object
11. embarked   1307 non-null   object
12. boat       486 non-null    object
13. body       121 non-null    float64
14. home.dest  745 non-null    object
dtypes: float64(7), object(7)
memory usage: 143.4+ KB
```

	pclass	survived	age	sibsp	parch	fare	body
count	1309.000000	1309.000000	1046.000000	1309.000000	1309.000000	1308.000000	121.000000
mean	2.294882	0.381971	29.881135	0.498854	0.385027	33.295479	160.809917
std	0.837836	0.486055	14.413500	1.041658	0.865560	51.758668	97.696922
min	1.000000	0.000000	0.166700	0.000000	0.000000	0.000000	1.000000
25%	2.000000	0.000000	21.000000	0.000000	0.000000	7.895800	72.000000
50%	3.000000	0.000000	28.000000	0.000000	0.000000	14.454200	155.000000
75%	3.000000	1.000000	39.000000	1.000000	0.000000	31.275000	256.000000
max	3.000000	1.000000	80.000000	8.000000	9.000000	512.329200	328.000000

```
In [18]: df.isna()
```

```
Out[18]: 2
```



Out[18]:

	pclas s	surviv ed	nam e	sex	age	sibsp	parc h	ticke t	fare	cabi n	embark ed	boat	bod y	home.de st
0	False	False	Fals e	Fals e	Fals e	Fals e	False	Fals e	Fals e	Fals e	False	Fals e	True	False
1	False	False	Fals e	Fals e	Fals e	Fals e	False	Fals e	Fals e	Fals e	False	Fals e	True	False
2	False	False	Fals e	Fals e	Fals e	Fals e	False	Fals e	Fals e	Fals e	False	True	True	False
3	False	False	Fals e	Fals e	Fals e	Fals e	False	Fals e	Fals e	Fals e	False	True	Fals e	False
4	False	False	Fals e	Fals e	Fals e	Fals e	False	Fals e	Fals e	Fals e	False	True	True	False
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
130 5	False	False	Fals e	Fals e	True	Fals e	False	Fals e	Fals e	True	False	True	True	True
130 6	False	False	Fals e	Fals e	Fals e	Fals e	False	Fals e	Fals e	True	False	True	Fals e	True
130 7	False	False	Fals e	Fals e	Fals e	Fals e	False	Fals e	Fals e	True	False	True	True	True
130 8	False	False	Fals e	Fals e	Fals e	Fals e	False	Fals e	Fals e	True	False	True	True	True
130 9	True	True	True	True	True	True	True	True	True	True	True	True	True	True

1310 rows × 14 columns

In [19]:

```
df.isna().any()
```

Out[19]:

pclass	True
survived	True
name	True
sex	True
age	True
sibsp	True
parch	True
ticket	True
fare	True
cabin	True
embarked	True
boat	True
body	True
home.dest	True

In [20]:

pclass	1
survived	1

In [21]:

```
df["age"].fillna(29.699118)
```

Out[21]:

sex	1
age	264
sibsp	1
parch	1
ticket	1
fare	2
name	1310, dtype: float64
embarked	3
boat	824
body	118 9
home.dest	565
dtype:	int64

1	0.916700
2	2.000000
3	30.000000
4	25.000000
...	
1305	29.699118
1306	26.500000
1307	27.000000
1308	29.000000
1309	29.699118

```
In [22]: df.isna().sum()
```

```
Out[22]: pclass      1
survived    1
name        1
sex         1
age        264
sibsp       1
parch       1
ticket      1
fare        2
cabin      1015
embarked    3
boat        824
body        1189
home.dest   565
dtype: int64
```

```
In [23]: df.any()
```

```
Out[23]: pclass      True
survived    True
name        True
sex         True
age         True
sibsp       True
parch       True
ticket      True
fare        True
cabin       True
embarked    True
boat        True
body        True
home.dest   True
```

```
In [24]: df=df.dropna()
```

```
In [25]: df.any()
```

```
Out[25]: pclass      False
survived    False
name        False
sex         False
age         False
sibsp       False
parch       False
ticket      False
fare        False
cabin       False
embarked    False
boat        False
body        False
home.dest   False
```

```
In [26]: df.isna().sum()
```

```
Out[26]:
```

```
pclass      0.0
survived    0.0
name        0.0
sex         0.0
age         0.0
sibsp       0.0
dtype: float64
parch       0.0
ticket      0.0
fare        0.0
cabin       0.0
embarked    0.0
boat        0.0
body        0.0
home.dest   0.0
```

```
In [ ]:
```