

S3 Bucket to Sage Maker-AWS

```
In [6]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline

In [2]: data=pd.read_csv('C:\Users\Arun\Desktop\train-1.csv')

In [3]: data

Out[3]:
```

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Utilities	...	PoolArea	PoolQC	Fence	MiscFeature	MiscVal	MoSold	YrSold	SaleType	SaleCondition
0	1	60	RL	65.0	8450	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	NaN	NaN	0	2	2008	WD	Normal
1	2	20	RL	80.0	9600	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	NaN	NaN	0	5	2007	WD	Normal
2	3	60	RL	68.0	11250	Pave	NaN	IR1	Lvl	AllPub	...	0	NaN	NaN	NaN	0	9	2008	WD	Normal
3	4	70	RL	60.0	9550	Pave	NaN	IR1	Lvl	AllPub	...	0	NaN	NaN	NaN	0	2	2006	WD	Abnrm
4	5	60	RL	84.0	14260	Pave	NaN	IR1	Lvl	AllPub	...	0	NaN	NaN	NaN	0	12	2008	WD	Normal
...
1455	1456	60	RL	62.0	7917	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	NaN	NaN	0	8	2007	WD	Normal
1458	1457	20	RL	85.0	13175	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	MGrV	NaN	0	2	2010	WD	Normal
1457	1458	70	RL	66.0	9042	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	GrV	Shed	2500	5	2010	WD	Normal
1458	1459	20	RL	68.0	9717	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	NaN	NaN	0	4	2010	WD	Normal
1459	1460	20	RL	75.0	9937	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	NaN	NaN	0	6	2008	WD	Normal

1460 rows × 21 columns

Data Cleaning

```
In [7]: data.shape
Out[7]: (1460, 81)

In [4]: data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1460 entries, 0 to 1459
Data columns (total 81 columns):
0      dtype: object
1      dtype: object
2      dtype: object
3      dtype: object
4      dtype: object
5      dtype: object
6      dtype: object
7      dtype: object
8      dtype: object
9      dtype: object
10     dtype: object
11     dtype: object
12     dtype: object
13     dtype: object
14     dtype: object
15     dtype: object
16     dtype: object
17     dtype: object
18     dtype: object
19     dtype: object
20     dtype: object
21     dtype: object
22     dtype: object
23     dtype: object
24     dtype: object
25     dtype: object
26     dtype: object
27     dtype: object
28     dtype: object
29     dtype: object
30     dtype: object
31     dtype: object
32     dtype: object
33     dtype: object
34     dtype: object
35     dtype: object
36     dtype: object
37     dtype: object
38     dtype: object
39     dtype: object
40     dtype: object
41     dtype: object
42     dtype: object
43     dtype: object
44     dtype: object
45     dtype: object
46     dtype: object
47     dtype: object
48     dtype: object
49     dtype: object
50     dtype: object
51     dtype: object
52     dtype: object
53     dtype: object
54     dtype: object
55     dtype: object
56     dtype: object
57     dtype: object
58     dtype: object
59     dtype: object
60     dtype: object
61     dtype: object
62     dtype: object
63     dtype: object
64     dtype: object
65     dtype: object
66     dtype: object
67     dtype: object
68     dtype: object
69     dtype: object
70     dtype: object
71     dtype: object
72     dtype: object
73     dtype: object
74     dtype: object
75     dtype: object
76     dtype: object
77     dtype: object
78     dtype: object
79     dtype: object
80     dtype: object
81     dtype: object
dtypes: float64(3), int64(35), object(43)
memory usage: 924.0+ KB

In [5]: data.describe()

Out[5]:
```

	Id	MSSubClass	LotFrontage	LotArea	OverallQual	OverallCond	YearBuilt	YearRemodAdd	MasVnrArea	BsmntFinSF1	...	WoodDeckSF	OpenPorchSF	EnclosedPorch	3S	
count	1460.000000	1460.000000	1201.000000	1460.000000	1460.000000	1460.000000	1460.000000	1460.000000	1452.000000	1460.000000	...	1460.000000	1460.000000	1460.000000	1460.000000	1460.000000
mean	72.000000	56.897260	70.049989	10516.828082	6.099315	5.575342	1971.267808	1984.865753	1452.000000	181.066207	...	44.639716	66.660274	21.954110	66.660274	21.954110
std	42.610093	42.300751	24.284752	9981.264932	1.382997	1.112799	30.202904	20.645407	181.066207	456.080991	...	125.338794	66.256028	61.119149	66.256028	61.119149
min	1.000000	20.000000	21.000000	1300.000000	1.000000	1.000000	1872.000000	1950.000000	0.000000	0.000000	...	0.000000	0.000000	0.000000	0.000000	0.000000
25%	365.750000	20.000000	58.000000	7553.500000	5.000000	5.000000	1954.000000	1967.000000	0.000000	0.000000	...	0.000000	0.000000	0.000000	0.000000	0.000000
50%	730.500000	50.000000	69.000000	9478.500000	6.000000	5.000000	1973.000000	1994.000000	0.000000	383.500000	...	0.000000	25.000000	0.000000	0.000000	0.000000
75%	1095.250000	70.000000	80.000000	11601.500000	7.000000	6.000000	2000.000000	2004.000000	166.000000	712.250000	...	168.000000	68.000000	0.000000	68.000000	0.000000
max	1460.000000	190.000000	313.000000	215245.000000	10.000000	9.000000	2010.000000	2010.000000	1600.000000	5644.000000	...	857.000000	547.000000	552.000000	552.000000	552.000000

8 rows × 38 columns

```
In [9]: print(data.SalePrice.value_counts())
data['SalePrice'].value_counts(normalize=True)

1490880    20
1490880    17
1559880    14
1459880    14
1980360    13
...
2026560    1
1649880    1
2083800    1
2815800    1
1475880    1

Name: SalePrice, Length: 663, dtype: int64

Out[9]:
```

	Id	MSSubClass	LotFrontage	LotArea	OverallQual	OverallCond	YearBuilt	YearRemodAdd	MasVnrArea	BsmntFinSF1	...	WoodDeckSF	OpenPorchSF	EnclosedPorch	3S	
count	1460.000000	1460.000000	1201.000000	1460.000000	1460.000000	1460.000000	1460.000000	1460.000000	1452.000000	1460.000000	...	1460.000000	1460.000000	1460.000000	1460.000000	1460.000000
mean	72.000000	56.897260	70.049989	10516.828082	6.099315	5.575342	1971.267808	1984.865753	1452.000000	181.066207	...	44.639716	66.660274	21.954110	66.660274	21.954110
std	42.610093	42.300751	24.284752	9981.264932	1.382997	1.112799	30.202904	20.645407	181.066207	456.080991	...	125.338794	66.256028	61.119149	66.256028	61.119149
min	1.000000	20.000000	21.000000	1300.000000	1.000000	1.000000	1872.000000	1950.000000	0.000000	0.000000	...	0.000000	0.000000	0.000000	0.000000	0.000000
25%	365.750000	20.000000	58.000000	7553.500000	5.000000	5.000000	1954.000000	1967.000000	0.000000	0.000000	...	0.000000	0.000000	0.000000	0.000000	0.000000
50%	730.500000	50.000000	69.000000	9478.500000	6.000000	5.000000	1973.000000	1994.000000	0.000000	383.500000	...	0.000000	25.000000	0.000000	0.000000	0.000000
75%	1095.250000	70.000000	80.000000	11601.500000	7.000000	6.000000	2000.000000	2004.000000	166.000000	712.250000	...	168.000000	68.000000	0.000000	68.000000	0.000000
max	1460.000000	190.000000	313.000000	215245.000000	10.000000	9.000000	2010.000000	2010.000000	1600.000000	5644.000000	...	857.000000	547.000000	552.000000	552.000000	552.000000

8 rows × 38 columns

```
In [10]: print(data.SalePrice.value_counts())
data['SalePrice'].value_counts(normalize=True)

1490880    20
1490880    17
1559880    14
1459880    14
1980360    13
...
2026560    1
1649880    1
2083800    1
2815800    1
1475880    1

Name: SalePrice, Length: 663, dtype: float64

In [11]: data.isnull().sum()

Out[11]:
```

	Id	MSSubClass	LotFrontage	LotArea	OverallQual	OverallCond	YearBuilt	YearRemodAdd	MasVnrArea	BsmntFinSF1	...	WoodDeckSF	OpenPorchSF	EnclosedPorch	3S	
count	1460	1460	1201	1460	1460	1460	1460	1460	1452	1460	...	1460	1460	1460	1460	1460
mean	72.0	56.89726	70.049989	10516.828082	6.099315	5.575342	1971.267808	1984.865753	1452.0	181.066207	...	44.639716	66.660274	21.954110	66.660274	21.954110
std	42.610093	42.300751	24.284752	9981.264932	1.382997	1.112799	30.202904	20.645407	181.066207	456.080991	...	125.338794	66.256028	61.119149	66.256028	61.119149
min	1	20	21	1300	1	1	1872	1950	0	0	...	0	0	0	0	0
25%	365.75	20	58	7553.5	5	5	1954	1967	0	0	...	0	0	0	0	0
50%	730.5	50	69	9478.5	6	5	1973	1994	0	383.5	...	0	25	0	0	0
75%	1095.25	70	80	11601.5	7	6	2000	2004	166	712.25	...	168	68	0	68	0
max	1460	190	313	215245	10	9	2010	2010	1600	5644	...	857	547	552	552	552

8 rows × 38 columns

```
In [12]: data.isnull().sum()

Out[12]:
```

	Id	MSSubClass	LotFrontage	LotArea	OverallQual	OverallCond	YearBuilt	YearRemodAdd	MasVnrArea	BsmntFinSF1	...	WoodDeckSF	OpenPorchSF	EnclosedPorch	3S	
count	1460	1460	1201	1460	1460	1460	1460	1460	1452	1460	...	1460	1460	1460	1460	1460
mean	72.0	56.89726	70.049989	10516.828082	6.099315	5.575342	1971.267808	1984.865753	1452.0	181.066207	...	44.639716	66.660274	21.954110	66.660274	21.954110
std	42.610093	42.300751	24.284752	9981.264932	1.382997	1.112799	30.202904	20.645407	181.066207	456.080991	...	125.338794	66.256028	61.119149	66.256028	61.119149
min	1	20	21	1300	1	1	1872	1950	0	0	...	0	0	0	0	0
25%	365.75	20	58	7553.5	5	5	1954	1967	0	0	...	0	0	0	0	0
50%	730.5	50	69	9478.5	6	5	1973	1994	0	383.5	...	0	25	0	0	0
75%	1095.25	70	80	11601.5	7	6	2000	2004	166	712.25	...	168	68	0	68	0
max	1460	190	313	215245	10	9	2010	2010	1600	5644	...	857	547	552	552	552

8 rows × 38 columns

```
In [13]: data['Alley'].isnull().sum()

Out[13]: 1369

In [15]: #heatmap to show the correlation between various variables of the dataset

plt.figure(figsize=(20, 15))
cor = data.corr()
sns.heatmap(cor, annot=True,
            bottom, top ax.get_ylim(),
            ax.set_ylim(bottom = 0.5, top = 0.5))
plt.show()
```

```
In [17]: sns.histplot(data[data.BsmntCond=='No'],
                    data[data.BsmntCond=='Yes'])

Out[17]: <AxesSubplot: xlabel='Count'>
```

```
In [21]: fig, (ax1, ax2) = plt.subplots(1, 2, figsize = (15, 5))

#scatter plot 1
ax1.scatter(x=data['YrSold'], y=data['SalePrice'])
ax1.set_title('Year Sold v/s SalePrice')

#scatter plot 2
ax2.scatter(x=data['SaleCondition'], y=data['SalePrice'])
ax2.set_title('SaleCondition v/s SalePrice')
plt.draw()
```

```
In [23]: sns.countplot(data['LotShape'])

Out[23]:
```

```
In [24]: data.isna().any()[lambda x: x]

Out[24]:
```

	LotFrontage	Alley	MSSubClass	MSZoning	LotFrontage	LotArea	Street	LotShape	LandContour	Utilities	LotConfig	...	EnclosedPorch	3SsnPorch	ScreenPorch	PoolArea	MiscVal	MoSold	YrSold	SaleType	SaleCondition
0	1	60	RL	65.0	8450	Pave	NaN	Reg	Lvl	AllPub	Inside	...	0	0	0	0	0	0	2	2008	WD
1	2	20	RL	80.0	9600	Pave	NaN	Reg	Lvl	AllPub	Inside	...	0	0	0	0	0	0	5	2007	WD
2	3	60	RL	68.0	11250	Pave	NaN	IR1	Lvl	AllPub	FR2	...	0	0	0	0	0	0	9	2008	WD
3	4	70	RL	60.0	9550	Pave	NaN	IR1	Lvl	AllPub	FR2	...	272	0	0	0	0	0	2	2006	WD
4	5	60	RL	84.0	14260	Pave	NaN	IR1	Lvl	AllPub	Corner	...	0	0	0	0	0	0	12	2008	WD

5 rows × 27 columns

```
In [31]: data=pd.get_dummies(data)
data1=

Out[31]:
```

	Id	MSSubClass	LotFrontage	LotArea	OverallQual	OverallCond	YearBuilt	YearRemodAdd	MasVnrArea	BsmntFinSF1	...	SaleType	ContLw	SaleType	New	SaleType	Oth	SaleType	WD	Sale
0	1	60	60.0	8450	6	5	2003	2003	196.0	708	...	0	0	0	0	0	0	0	0	1
1	2	20	80.0	9600	7	5	2001	2002	162.0	486	...	0	0	0	0	0	0	0	1	1
2	3	60	68.0	11250	7	5	1915	1970	0.0	216	...	0	0	0	0	0	0	0	1	1
3	4	70	60.0	9550	7	5	2000	2000	350.0	655	...	0	0	0	0	0	0	0	1	1
4	5	60	84.0	14260	8	5

1460 rows × 277 columns

```
In [34]: X=data1.dropna(axis=1)
y=data1['SalePrice']

Train and test
```

```
In [36]: from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state=101)
print(X_train.shape)
print(X_test.shape)
print(y_train.shape)
print(y_test.shape)

(1022, 274)
(438, 274)
(
```