

```
In [3]: import pandas as pd
import numpy as np
from matplotlib import pyplot as plt

In [4]: df = pd.read_csv("Dataset.csv")

In [5]: df.head()

Out[5]:
   Id  MSSubClass  MSZoning  LotFrontage  LotArea  Street  Alley  LotShape  LandContour  Utilities  ...  PoolArea  PoolQC  Fence  MiscFeature  MiscVal  MoSold  YrSold  SaleType  SaleCondition  SalePrice
0    1           60        RL          65.0    8450   Pave   NaN      Reg        Lvl     AllPub  ...      0      NaN    NaN      NaN          0         2    2008        WD          Normal      146000000
1    2           20        RL          80.0    9600   Pave   NaN      Reg        Lvl     AllPub  ...      0      NaN    NaN      NaN          0         5    2007        WD          Normal      180921195
2    3           60        RL          68.0   11250   Pave   NaN     IR1        Lvl     AllPub  ...      0      NaN    NaN      NaN          0         9    2008        WD          Normal      79442502
3    4           70        RL          60.0   9550   Pave   NaN     IR1        Lvl     AllPub  ...      0      NaN    NaN      NaN          0         2    2006        WD          AbnormalMort  34900000
4    5           60        RL          84.0   14260   Pave   NaN     IR1        Lvl     AllPub  ...      0      NaN    NaN      NaN          0        12    2008        WD          Normal      129975000

5 rows x 81 columns

In [6]: df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1460 entries, 0 to 1459
Data columns (total 81 columns):
#   Column              Non-Null Count  Dtype
---  --
0    Id                  1460 non-null   int64
1    MSSubClass          1460 non-null   int64
2    MSZoning            1460 non-null   object
3    LotFrontage         1201 non-null   float64
4    LotArea             1460 non-null   int64
5    Street              1460 non-null   object
6    Alley              91 non-null     object
7    LotShape            1460 non-null   object
8    LandContour         1460 non-null   object
9    Utilities           1460 non-null   object
10   LotConfig           1460 non-null   object
11   LandSlope           1460 non-null   object
12   Neighborhood        1460 non-null   object
13   Condition1          1460 non-null   object
14   Condition2          1460 non-null   object
15   BldgType            1460 non-null   object
16   HouseStyle          1460 non-null   object
17   OverallQual          1460 non-null   int64
18   OverallCond         1460 non-null   int64
19   YearBuilt            1460 non-null   int64
20   YearRemodAdd        1460 non-null   int64
21   RoofStyle           1460 non-null   object
22   RoofMatl            1460 non-null   object
23   Exterior1st         1460 non-null   object
24   Exterior2nd         1460 non-null   object
25   MasVnrType          1452 non-null   object
26   MasVnrArea          1452 non-null   float64
27   ExterQual            1460 non-null   object
28   ExterCond           1460 non-null   object
29   Foundation          1460 non-null   object
30   BsmtQual            1423 non-null   object
31   BsmtCond            1423 non-null   object
32   BsmtExposure        1422 non-null   object
33   BsmtFinType1        1423 non-null   object
34   BsmtFinSF1          1460 non-null   int64
35   BsmtFinType2        1422 non-null   object
36   BsmtFinSF2          1460 non-null   int64
37   BsmtUnfSF           1460 non-null   int64
38   TotalBsmtSF         1460 non-null   int64
39   Heating             1460 non-null   object
40   HeatingQC           1460 non-null   object
41   CentralAir          1460 non-null   object
42   Electrical           1459 non-null   object
43   1stFlrSF            1460 non-null   int64
44   2ndFlrSF            1460 non-null   int64
45   LowQualFinSF        1460 non-null   int64
46   GrLivArea           1460 non-null   int64
47   BsmtFullBath        1460 non-null   int64
48   BsmtHalfBath        1460 non-null   int64
49   FullBath            1460 non-null   int64
50   HalfBath            1460 non-null   int64
51   BedroomAbvGr        1460 non-null   int64
52   KitchenAbvGr        1460 non-null   int64
53   KitchenQual         1460 non-null   object
54   TotRmsAbvGrd        1460 non-null   int64
55   Functional          1460 non-null   object
56   Fireplaces          1460 non-null   int64
57   FireplaceQu         770 non-null   object
58   GarageType          1379 non-null   object
59   GarageYrBlt         1379 non-null   float64
60   GarageFinish        1379 non-null   object
61   GarageCars          1460 non-null   int64
62   GarageArea          1460 non-null   int64
63   GarageQual          1379 non-null   object
64   GarageCond          1379 non-null   object
65   PavedDrive          1460 non-null   object
66   WoodDeckSF          1460 non-null   int64
67   OpenPorchSF         1460 non-null   int64
68   EnclosedPorch       1460 non-null   int64
69   3SsnPorch           1460 non-null   int64
70   ScreenPorch         1460 non-null   int64
71   PoolArea            1460 non-null   int64
72   PoolQC              7 non-null     object
73   Fence              281 non-null   object
74   MiscFeature         54 non-null   object
75   MiscVal             1460 non-null   int64
76   MoSold              1460 non-null   int64
77   YrSold              1460 non-null   int64
78   SaleType            1460 non-null   object
79   SaleCondition        1460 non-null   object
80   SalePrice           1460 non-null   int64
dtypes: float64(3), int64(35), object(43)
memory usage: 924.0+ KB

In [7]: df.describe()

Out[7]:
   Id  MSSubClass  LotFrontage  LotArea  OverallQual  OverallCond  YearBuilt  YearRemodAdd  MasVnrArea  BsmtFinSF1  ...  WoodDeckSF  OpenPorchSF  EnclosedPorch  3Ssn
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  1460.0
mean    730.500000   56.897260    70.049958   10516.828082    6.099315    5.575342   1971.267808   1984.865753   103.685262  443.639726  ...    94.244521    46.660274    21.954110    3.4
std    421.610009   42.300571   24.284752   9981.264932    1.382997    1.112799    30.202904    20.645407   181.066207   456.098091  ...   125.338794    66.256028    61.119149   29.3
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    0.0
25%    365.750000   20.000000   59.000000   7553.500000    5.000000    5.000000   1954.000000   1967.000000    0.000000    0.000000  ...    0.000000    0.000000    0.000000    0.0
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    0.0
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    0.0
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   508.0

8 rows x 38 columns

In [8]: df.duplicated().sum()

Out[8]:
0

In [9]: df.isna().sum()

Out[9]:
Id                  0
MSSubClass          0
MSZoning            0
LotFrontage        259
LotArea             0
...
MoSold              0
YrSold              0
SaleType            0
SaleCondition       0
SalePrice           0
Length: 81, dtype: int64

In [10]: df.drop('Id', axis= 1, inplace = True)

In [11]: threshold = len(df)*0.20
column = df.isna().sum() > threshold
columnToBeDropped = df.columns[column]

In [12]: columnToBeDropped.shape

Out[12]:
(5,)

In [13]: df.drop(columnToBeDropped, axis= 1, inplace = True)
df.shape

Out[13]:
(1460, 75)

In [14]: # To check how many null values exist in the dataframe
df.isna().sum()

Out[14]:
MSSubClass          0
MSZoning            0
LotFrontage        259
LotArea             0
Street              0
...
MoSold              0
YrSold              0
SaleType            0
SaleCondition       0
SalePrice           0
Length: 75, dtype: int64

In [15]: df.duplicated().sum()

Out[15]:
0

In [16]: df['SalePrice'].describe()

Out[16]:
count      1460.000000
mean    180921.195890
std      79442.502883
min       34900.000000
25%     129975.000000
50%     163000.000000
75%     214000.000000
max      755000.000000
Name: SalePrice, dtype: float64

In [17]: df.shape

Out[17]:
(1460, 75)

In [18]: for col in df.select_dtypes(include=['int64', 'float64']):
    if df[col].isnull().sum() > 0:
        median_value = df[col].median()
        df[col].fillna(median_value, inplace=True)

In [19]: df.isnull().sum()

Out[19]:
MSSubClass          0
MSZoning            0
LotFrontage         0
LotArea             0
Street              0
..
MoSold              0
YrSold              0
SaleType            0
SaleCondition       0
SalePrice           0
Length: 75, dtype: int64

In [20]: for col in df.select_dtypes(include=['object']):
    if df[col].isnull().sum() > 0:
        mode_value = df[col].mode()[0]
        df[col].fillna(mode_value, inplace=True)

In [21]: df.isnull().sum()

Out[21]:
MSSubClass          0
MSZoning            0
LotFrontage         0
LotArea             0
Street              0
..
MoSold              0
YrSold              0
SaleType            0
SaleCondition       0
SalePrice           0
Length: 75, dtype: int64

In [22]: df.dtypes.value_counts()

Out[22]:
object      38
int64       34
float64      3
dtype: int64

In [23]: df.head(5)

Out[23]:
   MSSubClass  MSZoning  LotFrontage  LotArea  Street  LotShape  LandContour  Utilities  LotConfig  LandSlope  ...  EnclosedPorch  3SsnPorch  ScreenPorch  PoolArea  MiscVal  MoSold  YrSold  SaleType  SaleCondition  SalePrice
0           60        RL          65.0    8450   Pave      Reg        Lvl     AllPub    Inside      Gtl  ...           0           0           0           0           0         2    2008        WD          Normal      146000000
1           20        RL          80.0    9600   Pave      Reg        Lvl     AllPub    FR2       Gtl  ...           0           0           0           0           0         5    2007        WD          Normal      180921195
2           60        RL          68.0   11250   Pave     IR1       Lvl     AllPub    Inside      Gtl  ...           0           0           0           0           0         9    2008        WD          Normal      79442502
3           70        RL          60.0   9550   Pave     IR1       Lvl     AllPub    Corner      Gtl  ...          272           0           0           0           0         2    2006        WD          AbnormalMort  34900000
4           60        RL          84.0   14260   Pave     IR1       Lvl     AllPub    FR2       Gtl  ...           0           0           0           0           0        12    2008        WD          Normal      129975000

5 rows x 75 columns

In [24]: from sklearn.preprocessing import LabelEncoder

In [25]: label_encoder = LabelEncoder()
for col in df.select_dtypes(include=['object']).columns:
    df[col] = label_encoder.fit_transform(df[col])

In [26]: df.head(5)

Out[26]:
   MSSubClass  MSZoning  LotFrontage  LotArea  Street  LotShape  LandContour  Utilities  LotConfig  LandSlope  ...  EnclosedPorch  3SsnPorch  ScreenPorch  PoolArea  MiscVal  MoSold  YrSold  SaleType  SaleCondition  SalePrice
0           60         3         65.0    8450         1         3         3         0         4         0  ...           0           0           0           0           0         2    2008        WD          Normal      146000000
1           20         3         80.0    9600         1         3         3         0         2         0  ...           0           0           0           0           0         5    2007        WD          Normal      180921195
2           60         3         68.0   11250         1         0         3         0         4         0  ...           0           0           0           0           0         9    2008        WD          Normal      79442502
3           70         3         60.0   9550         1         0         3         0         0         0  ...          272           0           0           0           0         2    2006        WD          AbnormalMort  34900000
4           60         3         84.0   14260         1         0         3         0         2         0  ...           0           0           0           0           0        12    2008        WD          Normal      129975000

5 rows x 75 columns

In [27]: df.dtypes.value_counts()

Out[27]:
int32      38
int64      34
float64      3
dtype: int64
```