

Data Cleaning Techniques



deleting-col-rows

May 18, 2023

1 Methods to handle missing value

- 1) Deleting rows and columns that contains missing value
- 2) Fill missing value manually
- 3) Global Constant
- 4) Measure of central tendency (Mean, Median, Mode)
- 5) Measure of central tendency for each class
- 6) Most probable value

1.0.1 1) Deleting rows and columns that contains missing value

Importing necessary libraries

```
[1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
```

Loading dataset

```
[2]: df= pd.read_csv('train.csv')
df.head()
```

```
[2]:
```

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	\
0	1	60	RL	65.0	8450	Pave	NaN	Reg	
1	2	20	RL	80.0	9600	Pave	NaN	Reg	
2	3	60	RL	68.0	11250	Pave	NaN	IR1	
3	4	70	RL	60.0	9550	Pave	NaN	IR1	
4	5	60	RL	84.0	14260	Pave	NaN	IR1	

	LandContour	Utilities	...	PoolArea	PoolQC	Fence	MiscFeature	MiscVal	MoSold	\
0	Lvl	AllPub	...	0	NaN	NaN	NaN	0	2	
1	Lvl	AllPub	...	0	NaN	NaN	NaN	0	5	
2	Lvl	AllPub	...	0	NaN	NaN	NaN	0	9	
3	Lvl	AllPub	...	0	NaN	NaN	NaN	0	2	
4	Lvl	AllPub	...	0	NaN	NaN	NaN	0	12	

	YrSold	SaleType	SaleCondition	SalePrice
0	2008	WD	Normal	208500
1	2007	WD	Normal	181500
2	2008	WD	Normal	223500
3	2006	WD	Abnorml	140000
4	2008	WD	Normal	250000

[5 rows x 81 columns]

```
[3]: df.shape # Checking the dimesion of the dataset
```

```
[3]: (1460, 81)
```

```
[4]: # Dataset has 81 colmns but it shows few column, to see every column we use pd.
      ↪ set_option
```

```
pd.set_option('display.max_columns',None)
pd.set_option('display.max_rows',None)
```

```
[5]: df.head(6)
```

```
[5]:
```

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	\
0	1	60	RL	65.0	8450	Pave	NaN	Reg	
1	2	20	RL	80.0	9600	Pave	NaN	Reg	
2	3	60	RL	68.0	11250	Pave	NaN	IR1	
3	4	70	RL	60.0	9550	Pave	NaN	IR1	
4	5	60	RL	84.0	14260	Pave	NaN	IR1	
5	6	50	RL	85.0	14115	Pave	NaN	IR1	

	LandContour	Utilities	LotConfig	LandSlope	Neighborhood	Condition1	\
0	Lvl	AllPub	Inside	Gtl	CollgCr	Norm	
1	Lvl	AllPub	FR2	Gtl	Veenker	Feedr	
2	Lvl	AllPub	Inside	Gtl	CollgCr	Norm	
3	Lvl	AllPub	Corner	Gtl	Crawfor	Norm	
4	Lvl	AllPub	FR2	Gtl	NoRidge	Norm	
5	Lvl	AllPub	Inside	Gtl	Mitchel	Norm	

	Condition2	BldgType	HouseStyle	OverallQual	OverallCond	YearBuilt	\
0	Norm	1Fam	2Story	7	5	2003	
1	Norm	1Fam	1Story	6	8	1976	
2	Norm	1Fam	2Story	7	5	2001	
3	Norm	1Fam	2Story	7	5	1915	
4	Norm	1Fam	2Story	8	5	2000	
5	Norm	1Fam	1.5Fin	5	5	1993	

	YearRemodAdd	RoofStyle	RoofMatl	Exterior1st	Exterior2nd	MasVnrType	\
0	2003	Gable	CompShg	VinylSd	VinylSd	BrkFace	

1	1976	Gable	CompShg	MetalSd	MetalSd	None
2	2002	Gable	CompShg	VinylSd	VinylSd	BrkFace
3	1970	Gable	CompShg	Wd Sdng	Wd Shng	None
4	2000	Gable	CompShg	VinylSd	VinylSd	BrkFace
5	1995	Gable	CompShg	VinylSd	VinylSd	None

	MasVnrArea	ExterQual	ExterCond	Foundation	BsmtQual	BsmtCond	BsmtExposure	\
0	196.0	Gd	TA	PConc	Gd	TA	No	
1	0.0	TA	TA	CBlock	Gd	TA	Gd	
2	162.0	Gd	TA	PConc	Gd	TA	Mn	
3	0.0	TA	TA	BrkTil	TA	Gd	No	
4	350.0	Gd	TA	PConc	Gd	TA	Av	
5	0.0	TA	TA	Wood	Gd	TA	No	

	BsmtFinType1	BsmtFinSF1	BsmtFinType2	BsmtFinSF2	BsmtUnfSF	TotalBsmtSF	\
0	GLQ	706	Unf	0	150	856	
1	ALQ	978	Unf	0	284	1262	
2	GLQ	486	Unf	0	434	920	
3	ALQ	216	Unf	0	540	756	
4	GLQ	655	Unf	0	490	1145	
5	GLQ	732	Unf	0	64	796	

	Heating	HeatingQC	CentralAir	Electrical	1stFlrSF	2ndFlrSF	LowQualFinSF	\
0	GasA	Ex	Y	SBrkr	856	854	0	
1	GasA	Ex	Y	SBrkr	1262	0	0	
2	GasA	Ex	Y	SBrkr	920	866	0	
3	GasA	Gd	Y	SBrkr	961	756	0	
4	GasA	Ex	Y	SBrkr	1145	1053	0	
5	GasA	Ex	Y	SBrkr	796	566	0	

	GrLivArea	BsmtFullBath	BsmtHalfBath	FullBath	HalfBath	BedroomAbvGr	\
0	1710	1	0	2	1	3	
1	1262	0	1	2	0	3	
2	1786	1	0	2	1	3	
3	1717	1	0	1	0	3	
4	2198	1	0	2	1	4	
5	1362	1	0	1	1	1	

	KitchenAbvGr	KitchenQual	TotRmsAbvGrd	Functional	Fireplaces	FireplaceQu	\
0	1	Gd	8	Typ	0	NaN	
1	1	TA	6	Typ	1	TA	
2	1	Gd	6	Typ	1	TA	
3	1	Gd	7	Typ	1	Gd	
4	1	Gd	9	Typ	1	TA	
5	1	TA	5	Typ	0	NaN	

GarageType	GarageYrBlt	GarageFinish	GarageCars	GarageArea	GarageQual	\
------------	-------------	--------------	------------	------------	------------	---

0	Attchd	2003.0	RFn	2	548	TA
1	Attchd	1976.0	RFn	2	460	TA
2	Attchd	2001.0	RFn	2	608	TA
3	Detchd	1998.0	Unf	3	642	TA
4	Attchd	2000.0	RFn	3	836	TA
5	Attchd	1993.0	Unf	2	480	TA

	GarageCond	PavedDrive	WoodDeckSF	OpenPorchSF	EnclosedPorch	3SsnPorch	\
0	TA	Y	0	61	0	0	
1	TA	Y	298	0	0	0	
2	TA	Y	0	42	0	0	
3	TA	Y	0	35	272	0	
4	TA	Y	192	84	0	0	
5	TA	Y	40	30	0	320	

	ScreenPorch	PoolArea	PoolQC	Fence	MiscFeature	MiscVal	MoSold	YrSold	\
0	0	0	NaN	NaN	NaN	0	2	2008	
1	0	0	NaN	NaN	NaN	0	5	2007	
2	0	0	NaN	NaN	NaN	0	9	2008	
3	0	0	NaN	NaN	NaN	0	2	2006	
4	0	0	NaN	NaN	NaN	0	12	2008	
5	0	0	NaN	MnPrv	Shed	700	10	2009	

	SaleType	SaleCondition	SalePrice
0	WD	Normal	208500
1	WD	Normal	181500
2	WD	Normal	223500
3	WD	Abnorml	140000
4	WD	Normal	250000
5	WD	Normal	143000

```
[27]: df.info() # info() is to find the information about the dataset like columns,
        ↪non_null values, data type etc
```

```
<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

```

```
[7]: df.isnull().sum() # To check number of null values each column contains
```

```

[7]: Id                0
     MSSubClass        0
     MSZoning          0
     LotFrontage      259
     LotArea          0
     Street           0
     Alley           1369
     LotShape         0
     LandContour      0
     Utilities        0
     LotConfig        0
     LandSlope        0
     Neighborhood     0
     Condition1       0
     Condition2       0
     BldgType         0
     HouseStyle       0

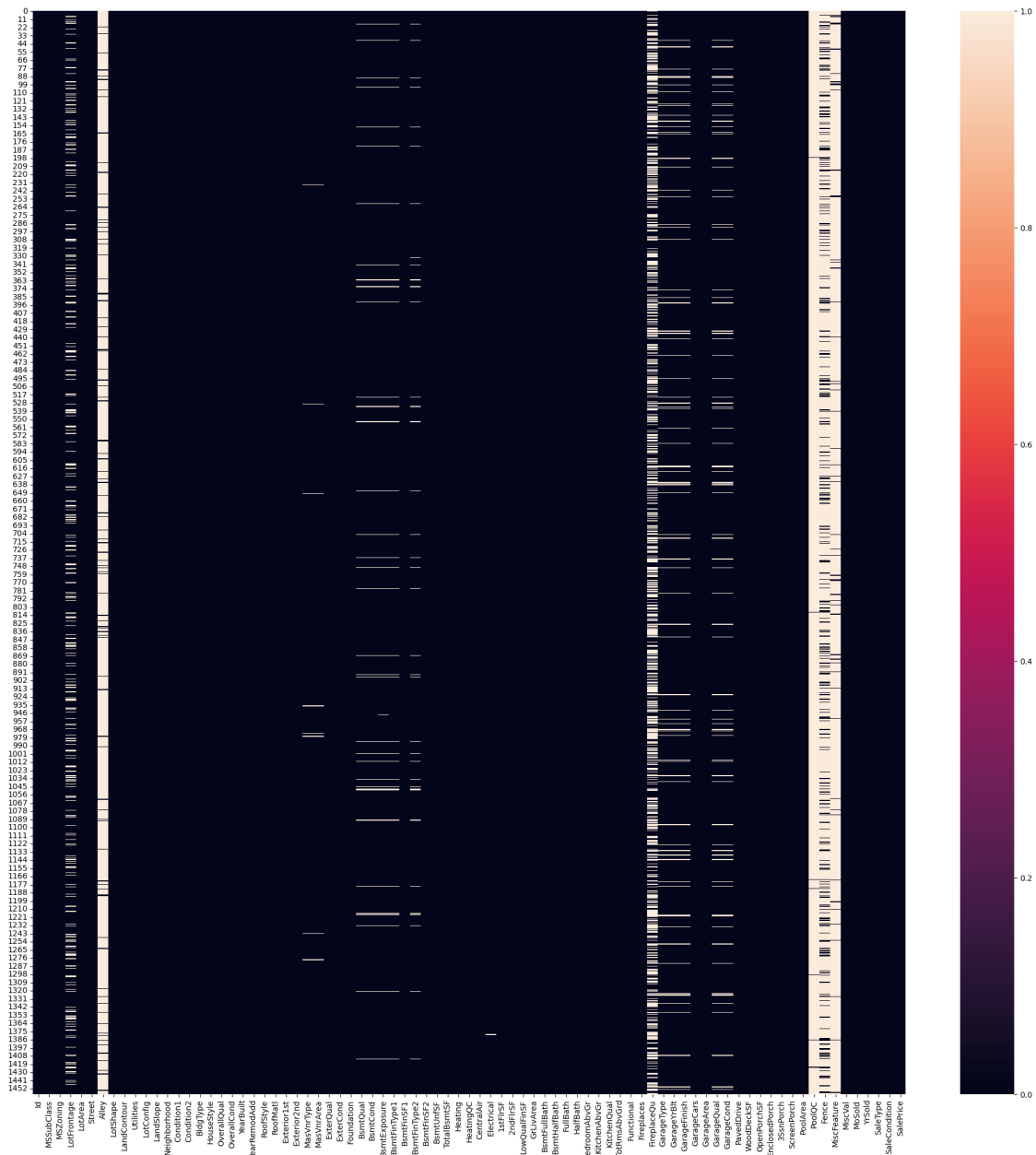
```

OverallQual	0
OverallCond	0
YearBuilt	0
YearRemodAdd	0
RoofStyle	0
RoofMatl	0
Exterior1st	0
Exterior2nd	0
MasVnrType	8
MasVnrArea	8
ExterQual	0
ExterCond	0
Foundation	0
BsmtQual	37
BsmtCond	37
BsmtExposure	38
BsmtFinType1	37
BsmtFinSF1	0
BsmtFinType2	38
BsmtFinSF2	0
BsmtUnfSF	0
TotalBsmtSF	0
Heating	0
HeatingQC	0
CentralAir	0
Electrical	1
1stFlrSF	0
2ndFlrSF	0
LowQualFinSF	0
GrLivArea	0
BsmtFullBath	0
BsmtHalfBath	0
FullBath	0
HalfBath	0
BedroomAbvGr	0
KitchenAbvGr	0
KitchenQual	0
TotRmsAbvGrd	0
Functional	0
Fireplaces	0
FireplaceQu	690
GarageType	81
GarageYrBltn	81
GarageFinish	81
GarageCars	0
GarageArea	0
GarageQual	81


```
GarageCond      81
PavedDrive      0
WoodDeckSF      0
OpenPorchSF     0
EnclosedPorch   0
3SsnPorch       0
ScreenPorch     0
PoolArea        0
PoolQC         1453
Fence           1179
MiscFeature     1406
MiscVal         0
MoSold          0
YrSold          0
SaleType        0
SaleCondition    0
SalePrice       0
dtype: int64
```

```
[8]: # To check distribution of null values in the dataset we use heat map. Here, we
      ↪ need to clean the dataset.
      plt.figure(figsize=(25,25))
      sns.heatmap(df.isnull())
```

```
[8]: <AxesSubplot:>
```



```
[9]: null_var= df.isnull().sum()/df.shape[0]*100 # To find the percentage of null
      ↪values
      null_var
```

```
[9]: Id                0.000000
      MSSubClass        0.000000
      MSZoning          0.000000
      LotFrontage      17.739726
      LotArea           0.000000
```

Street	0.000000
Alley	93.767123
LotShape	0.000000
LandContour	0.000000
Utilities	0.000000
LotConfig	0.000000
LandSlope	0.000000
Neighborhood	0.000000
Condition1	0.000000
Condition2	0.000000
BldgType	0.000000
HouseStyle	0.000000
OverallQual	0.000000
OverallCond	0.000000
YearBuilt	0.000000
YearRemodAdd	0.000000
RoofStyle	0.000000
RoofMatl	0.000000
Exterior1st	0.000000
Exterior2nd	0.000000
MasVnrType	0.547945
MasVnrArea	0.547945
ExterQual	0.000000
ExterCond	0.000000
Foundation	0.000000
BsmtQual	2.534247
BsmtCond	2.534247
BsmtExposure	2.602740
BsmtFinType1	2.534247
BsmtFinSF1	0.000000
BsmtFinType2	2.602740
BsmtFinSF2	0.000000
BsmtUnfSF	0.000000
TotalBsmtSF	0.000000
Heating	0.000000
HeatingQC	0.000000
CentralAir	0.000000
Electrical	0.068493
1stFlrSF	0.000000
2ndFlrSF	0.000000
LowQualFinSF	0.000000
GrLivArea	0.000000
BsmtFullBath	0.000000
BsmtHalfBath	0.000000
FullBath	0.000000
HalfBath	0.000000
BedroomAbvGr	0.000000

```

KitchenAbvGr      0.000000
KitchenQual       0.000000
TotRmsAbvGrd     0.000000
Functional        0.000000
Fireplaces        0.000000
FireplaceQu      47.260274
GarageType        5.547945
GarageYrBltd      5.547945
GarageFinish      5.547945
GarageCars        0.000000
GarageArea        0.000000
GarageQual        5.547945
GarageCond        5.547945
PavedDrive        0.000000
WoodDeckSF        0.000000
OpenPorchSF       0.000000
EnclosedPorch     0.000000
3SsnPorch         0.000000
ScreenPorch       0.000000
PoolArea          0.000000
PoolQC            99.520548
Fence              80.753425
MiscFeature       96.301370
MiscVal           0.000000
MoSold            0.000000
YrSold            0.000000
SaleType          0.000000
SaleCondition     0.000000
SalePrice         0.000000
dtype: float64

```

```
[10]: drop_column = null_var[null_var>17].keys() #Checking for the columns which
      ↪contain null values more than 17%
```

```
[10]: Index(['LotFrontage', 'Alley', 'FireplaceQu', 'PoolQC', 'Fence',
      ↪      'MiscFeature'],
      ↪      dtype='object')
```

```
[11]: df2=df.drop(columns=drop_column) # We are dropping all the columns that contain
      ↪null values more than 17%
      ↪drop_column
```

```
[12]: df2.shape
```

```
[12]: (1460, 75)
```

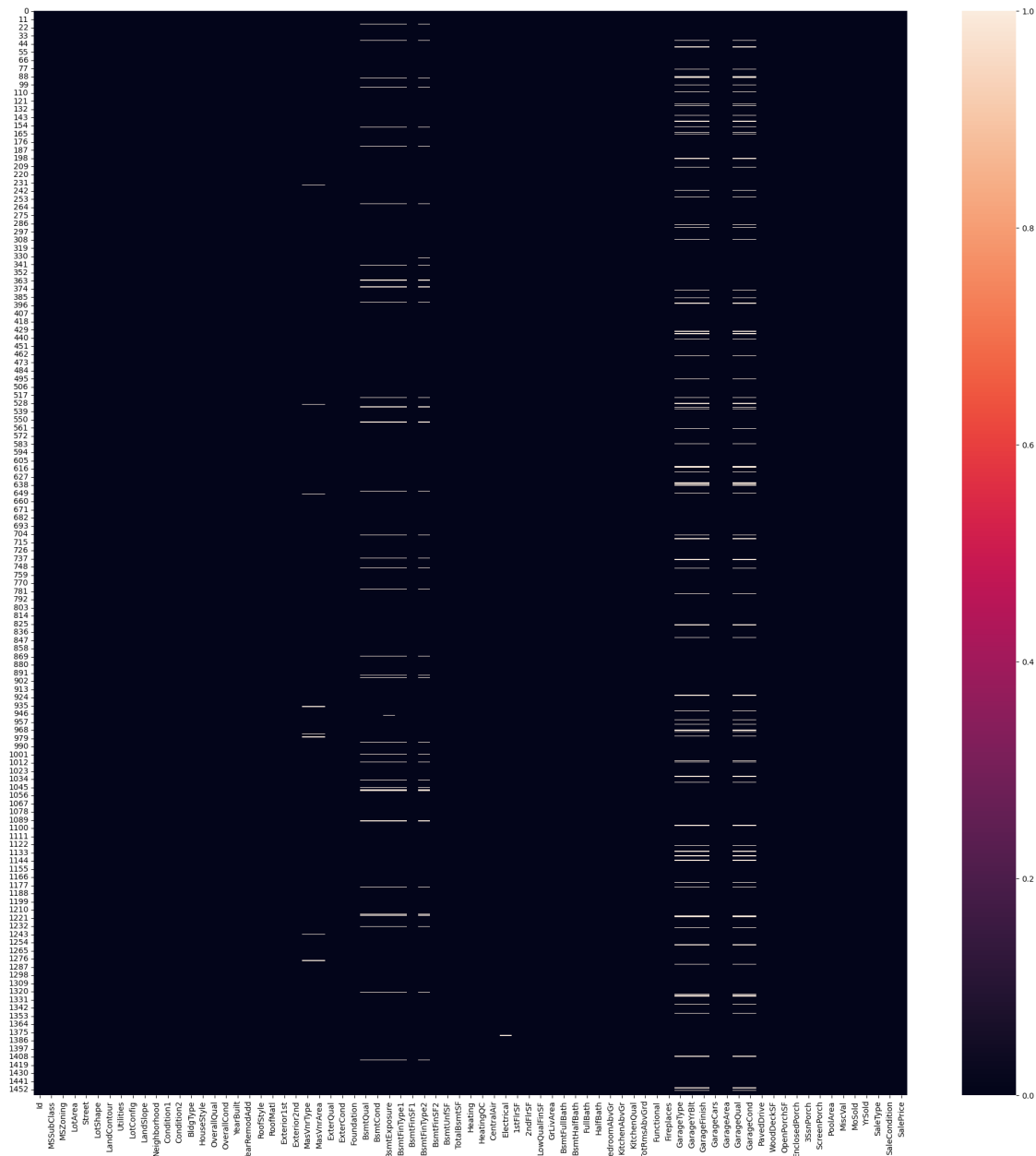
```
[13]: df2.isnull().sum() # All the column that contains majority of the null values
      ↪are dropped
```

```
[13]: Id 0
      MSSubClass 0
      MSZoning 0
      LotArea 0
      Street 0
      LotShape 0
      LandContour 0
      Utilities 0
      LotConfig 0
      LandSlope 0
      Neighborhood 0
      Condition1 0
      Condition2 0
      BldgType 0
      HouseStyle 0
      OverallQual 0
      OverallCond 0
      YearBuilt 0
      YearRemodAdd 0
      RoofStyle 0
      RoofMatl 0
      Exterior1st 0
      Exterior2nd 0
      MasVnrType 8
      MasVnrArea 8
      ExterQual 0
      ExterCond 0
      Foundation 0
      BsmtQual 37
      BsmtCond 37
      BsmtExposure 38
      BsmtFinType1 37
      BsmtFinSF1 0
      BsmtFinType2 38
      BsmtFinSF2 0
      BsmtUnfSF 0
      TotalBsmtSF 0
      Heating 0
      HeatingQC 0
      CentralAir 0
      Electrical 1
      1stFlrSF 0
      2ndFlrSF 0
      LowQualFinSF 0
```

GrLivArea	0
BsmtFullBath	0
BsmtHalfBath	0
FullBath	0
HalfBath	0
BedroomAbvGr	0
KitchenAbvGr	0
KitchenQual	0
TotRmsAbvGrd	0
Functional	0
Fireplaces	0
GarageType	81
GarageYrBlt	81
GarageFinish	81
GarageCars	0
GarageArea	0
GarageQual	81
GarageCond	81
PavedDrive	0
WoodDeckSF	0
OpenPorchSF	0
EnclosedPorch	0
3SsnPorch	0
ScreenPorch	0
PoolArea	0
MiscVal	0
MoSold	0
YrSold	0
SaleType	0
SaleCondition	0
SalePrice	0
dtype:	int64

```
[14]: plt.figure(figsize=(25,25)) # Here we can see that the majority of the columns
      ↪ that contained null values are handled.
      sns.heatmap(df2.isnull())
```

```
[14]: <AxesSubplot:>
```



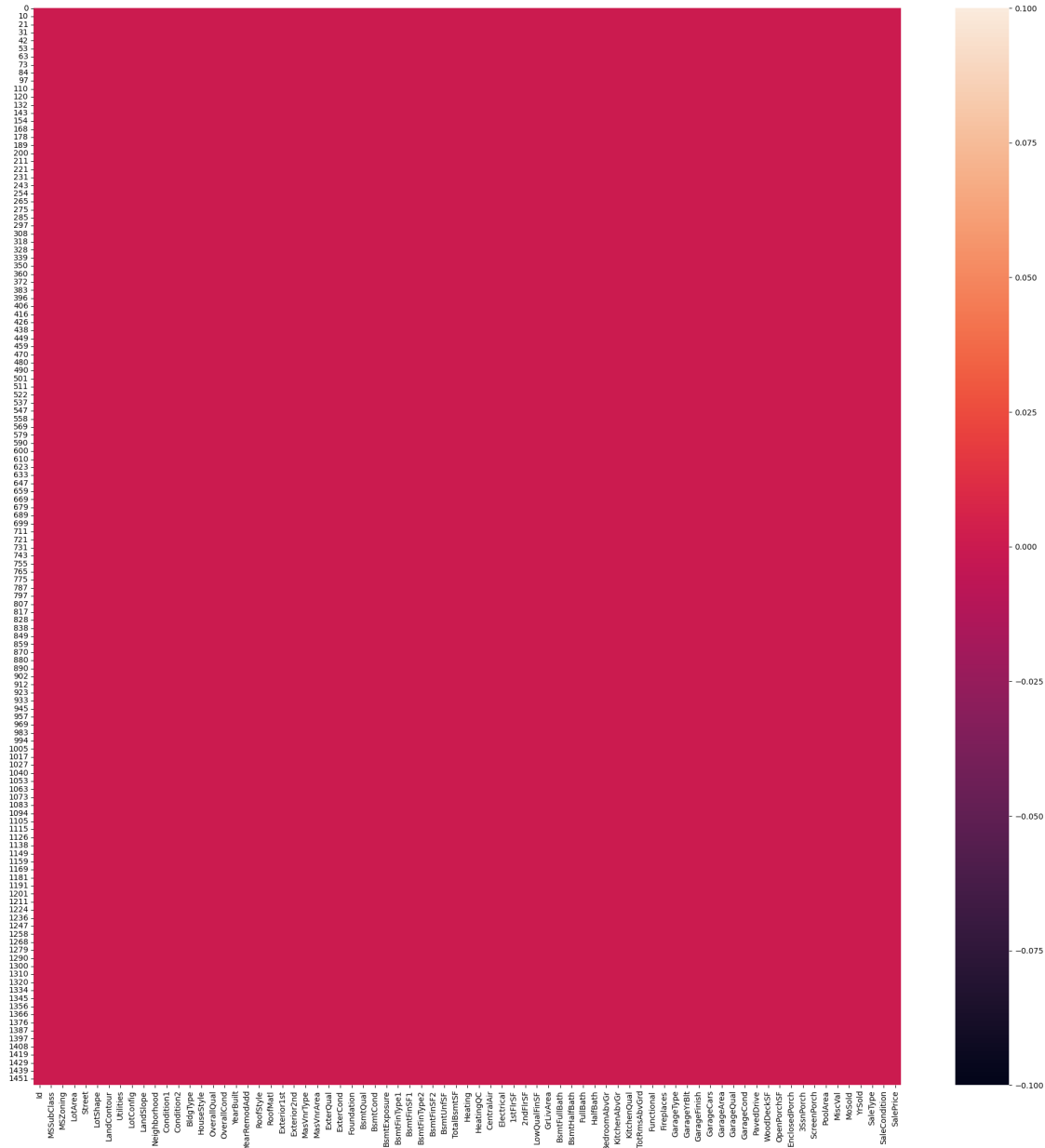
```
[15]: df3= df2.dropna() # we are deleting the rows that contains null values.
```

```
[16]: df3.isnull().sum().sum()
```

[16] : 0

```
[17]: # We can see that null value columns and rows are handled
plt.figure(figsize=(25,25))
sns.heatmap(df3.isnull())
```

```
[17]: <AxesSubplot:>
```



To check data distribution of numerical variables

```
[18]: df3.select_dtypes(include=['int64','float64']).columns.tolist() # Checking for
      ↪ the numerical columns.
```

```
[18]: ['Id',
      'MSSubClass',
      'LotArea',
```



```

'OverallQual',
'OverallCond',
'YearBuilt',
'YearRemodAdd',
'MasVnrArea',
'BsmtFinSF1',
'BsmtFinSF2',
'BsmtUnfSF',
'TotalBsmtSF',
'1stFlrSF',
'2ndFlrSF',
'LowQualFinSF',
'GrLivArea',
'BsmtFullBath',
'BsmtHalfBath',
'FullBath',
'HalfBath',
'BedroomAbvGr',
'KitchenAbvGr',
'TotRmsAbvGrd',
'Fireplaces',
'GarageYrBlt',
'GarageCars',
'GarageArea',
'WoodDeckSF',
'OpenPorchSF',
'EnclosedPorch',
'3SsnPorch',
'ScreenPorch',
'PoolArea',
'MiscVal',
'MoSold',
'YrSold',
'SalePrice']

```

```

[20]: ### Checking data distribution of one column before and after handling the
      ↪missing values

```

```

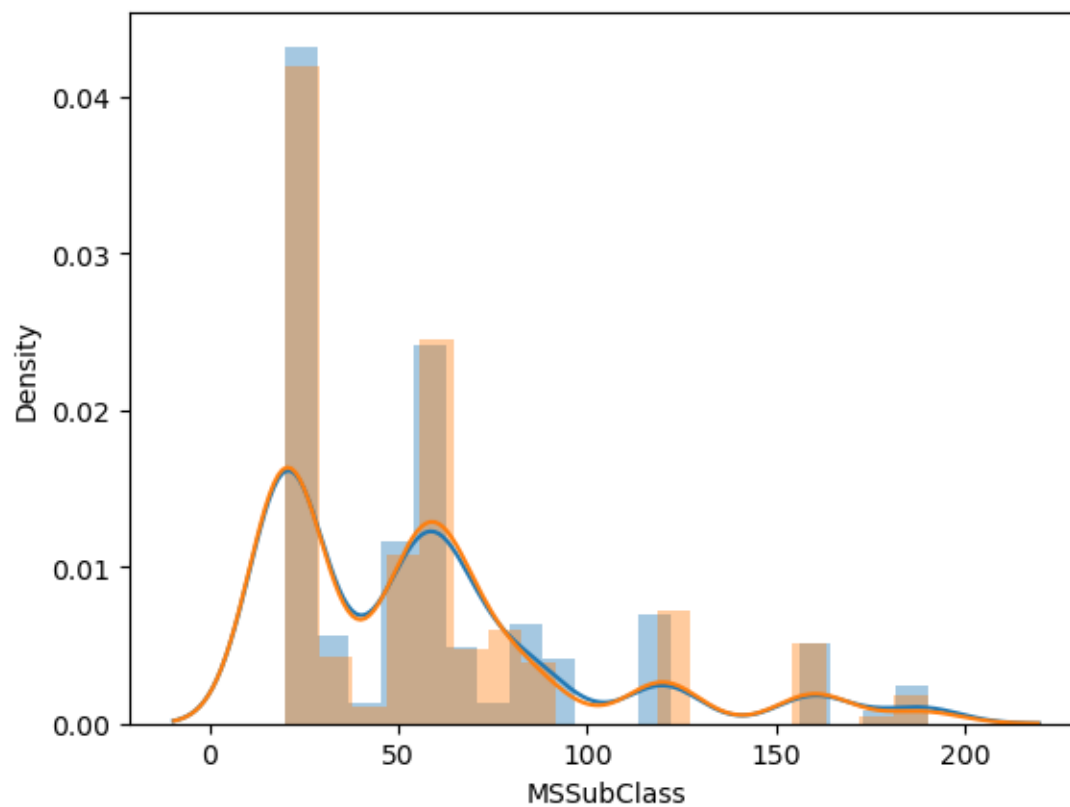
sns.distplot(df['MSSubClass'])
sns.distplot(df3['MSSubClass'])

```

```

[20]: <AxesSubplot:xlabel='MSSubClass', ylabel='Density'>

```



```
[21]: # All the columns with numerical datatype are stored in a varibale.
```

```
num_var= [
    'MSSubClass',
    'LotArea',
    'OverallQual',
    'OverallCond',
    'YearBuilt',
    'YearRemodAdd',
    'MasVnrArea',
    'BsmtFinSF1',
    'BsmtFinSF2',
    'BsmtUnfSF',
    'TotalBsmtSF',
    '1stFlrSF',
    '2ndFlrSF',
    'LowQualFinSF',
    'GrLivArea',
    'BsmtFullBath',
    'BsmtHalfBath',
    'FullBath',
```

```

'HalfBath',
'BedroomAbvGr',
'KitchenAbvGr',
'TotRmsAbvGrd',
'Fireplaces',
'GarageYrBlt',
'GarageCars',
'GarageArea',
'WoodDeckSF',
'OpenPorchSF',
'EnclosedPorch',
'3SsnPorch',
'ScreenPorch',
'PoolArea',
'MiscVal',
'MoSold',
'YrSold',
'SalePrice']

```

```

[23]: # Checking the data distribution of the numerical variable. To know if we
      ↪ cleaned the data in a proper way.
plt.figure(figsize=(25,25))

for i, var in enumerate(num_var):
    plt.subplot(9,4,i+1)
    sns.distplot(df[var], bins=20)
    sns.distplot(df3[var], bins=20)

```

