

```
In [1]: import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
```

```
In [23]: data=pd.read_csv(r"train.csv")
data.head()
```

Out[23]:

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Utilities	LotConfig	LandSlope	Neighborhood
0	1	60	RL	65.0	8450	Pave	NaN	Reg	Lvl	AllPub	Inside	Gtl	Coll
1	2	20	RL	80.0	9600	Pave	NaN	Reg	Lvl	AllPub	FR2	Gtl	Veer
2	3	60	RL	68.0	11250	Pave	NaN	IR1	Lvl	AllPub	Inside	Gtl	Coll
3	4	70	RL	60.0	9550	Pave	NaN	IR1	Lvl	AllPub	Corner	Gtl	Crav
4	5	60	RL	84.0	14260	Pave	NaN	IR1	Lvl	AllPub	FR2	Gtl	NoRi



```
In [4]: data.shape
```

Out[4]: (1460, 81)

```
In [21]: pd.set_option("display.max_columns",None)
pd.set_option("display.max_rows",None)
```

In [25]: `data.head()`

Out[25]:

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Utilities	LotConfig	LandSlope	Neighborhood
0	1	60	RL	65.0	8450	Pave	NaN	Reg	Lvl	AllPub	Inside	Gtl	Collamore
1	2	20	RL	80.0	9600	Pave	NaN	Reg	Lvl	AllPub	FR2	Gtl	Veerht
2	3	60	RL	68.0	11250	Pave	NaN	IR1	Lvl	AllPub	Inside	Gtl	Collamore
3	4	70	RL	60.0	9550	Pave	NaN	IR1	Lvl	AllPub	Corner	Gtl	Cravens
4	5	60	RL	84.0	14260	Pave	NaN	IR1	Lvl	AllPub	FR2	Gtl	NoRidge



In [27]: `data.tail()`

Out[27]:

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Utilities	LotConfig	LandSlope	Neighborhood
1455	1456	60	RL	62.0	7917	Pave	NaN	Reg	Lvl	AllPub	Inside	Gtl	Collamore
1456	1457	20	RL	85.0	13175	Pave	NaN	Reg	Lvl	AllPub	Inside	Gtl	Veerht
1457	1458	70	RL	66.0	9042	Pave	NaN	Reg	Lvl	AllPub	Inside	Gtl	Cravens
1458	1459	20	RL	68.0	9717	Pave	NaN	Reg	Lvl	AllPub	Inside	Gtl	NoRidge
1459	1460	20	RL	75.0	9937	Pave	NaN	Reg	Lvl	AllPub	Inside	Gtl	Collamore



In [30]: data.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: 678.7+ KB

```

In [31]: data.isnull()

1332	False	False	False	False	False	False	True	False	False	False	False	False
1333	False	False	False	False	False	False	True	False	False	False	False	False
1334	False	False	False	False	False	False	True	False	False	False	False	False
1335	False	False	False	False	False	False	True	False	False	False	False	False
1336	False	False	False	False	False	False	True	False	False	False	False	False
1337	False	False	False	False	False	False	False	False	False	False	False	False
1338	False	False	False	False	False	False	True	False	False	False	False	False
1339	False	False	False	False	False	False	True	False	False	False	False	False
1340	False	False	False	False	False	False	True	False	False	False	False	False
1341	False	False	False	False	False	False	True	False	False	False	False	False
1342	False	False	False	True	False	False	True	False	False	False	False	False
1343	False	False	False	False	False	False	True	False	False	False	False	False

```
In [34]: data.isnull().sum()
```

```
Out[34]: 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
```

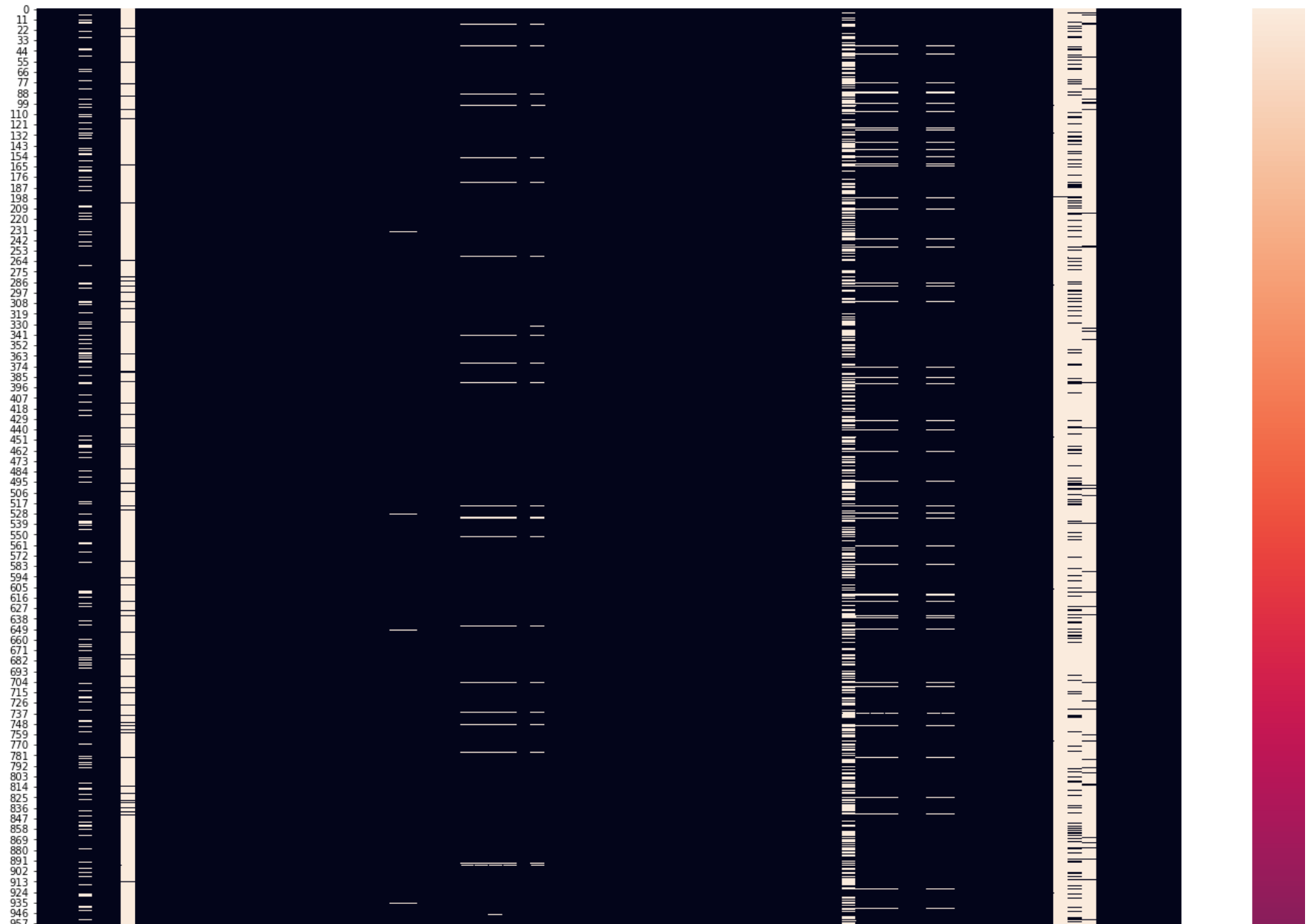
BsmtFinSF2	0
BsmtUnfSF	0
TotalBsmtSF	0
Heating	0
HeatingQC	0
CentralAir	0
Electrical	1
1stFlrSF	0
2ndFlrSF	0
LowQualFinSF	0
GrLivArea	0
BsmtFullBath	0
BsmtHalfBath	0
FullBath	0
HalfBath	0
BedroomAbvGr	0
KitchenAbvGr	0
KitchenQual	0
TotRmsAbvGrd	0
Functional	0
Fireplaces	0
FireplaceQu	690
GarageType	81
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
PoolQC	1453
Fence	1179
MiscFeature	1406

```
MiscVal      0
MoSold       0
YrSold       0
SaleType     0
SaleCondition 0
SalePrice    0
dtype: int64
```



```
In [35]: plt.figure(figsize=(25,25))  
sns.heatmap(data.isnull())
```

Out[35]: <AxesSubplot:>




```
In [41]: null_var=data.isnull().sum()/data.shape[0]*100  
null_var
```

```
Out[41]: 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
GarageYrBlt	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
```

```
In [45]: drop_cloumn = null_var[null_var >17].keys()
drop_cloumn
```

```
Out[45]: Index(['LotFrontage', 'Alley', 'FireplaceQu', 'PoolQC', 'Fence',
               'MiscFeature'],
              dtype='object')
```

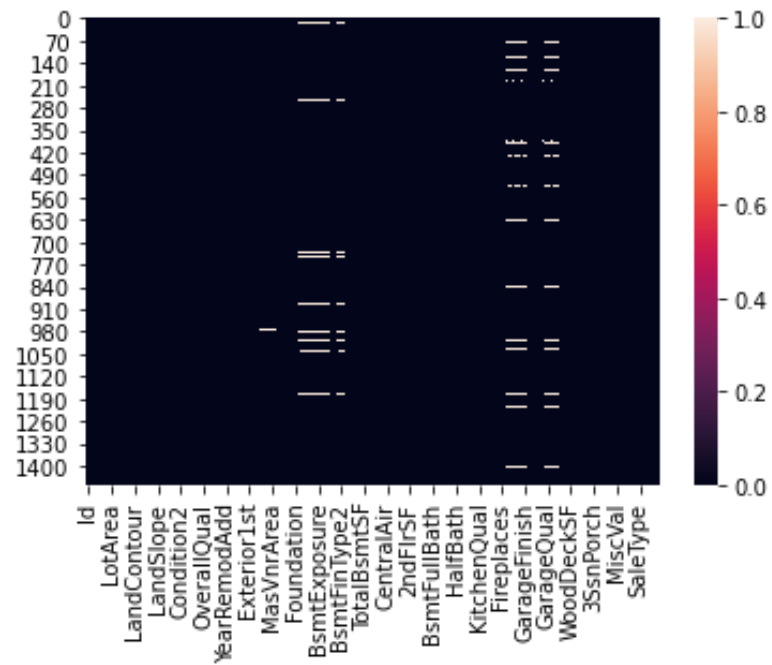
```
In [47]: data2=data.drop(columns=drop_cloumn)
```

```
In [50]: data2.shape
```

```
Out[50]: (1460, 75)
```

```
In [55]: sns.heatmap(data2.isnull())
```

```
Out[55]: <AxesSubplot:>
```



```
In [58]: data3=data2.dropna()
```

```
In [60]: data3.shape
```

```
Out[60]: (1338, 75)
```

```
In [64]: plt.figure(figsize=(25,25))  
sns.heatmap(data3.isnull())
```

Out[64]: <AxesSubplot:>



945	
957	
969	
983	
994	
1005	
1017	
1027	
1040	
1053	
1063	
1073	
1083	
1094	
1105	
1115	
1126	
1136	
1149	
1159	
1169	
1181	
1191	
1201	
1211	
1224	
1236	
1247	
1258	
1268	
1279	
1290	
1300	
1310	
1320	
1334	
1345	
1356	
1366	
1376	
1387	
1397	
1408	
1419	
1429	
1439	
1451	
Id	
MSSubClass	
MSZoning	
LotArea	
Street	
LotShape	
LandContour	
Utilities	
LotConfig	
LandSlope	
Neighborhood	
Condition1	
Condition2	
BldgType	
HouseStyle	
OverallQual	
OverallCond	
YearBuilt	
YearRemodAdd	
RoofStyle	
RoofMatl	
Exterior1st	
Exterior2nd	
MasVnrType	
MasVnrArea	
ExterQual	
ExterCond	
Foundation	
BsmtQual	
BsmtCond	
BsmtExposure	
BsmtFinType1	
BsmtFinSF1	
BsmtFinType2	
BsmtFinSF2	
BsmtUnfSF	
TotalBsmtSF	
Heating	
HeatingQC	
CentralAir	
Electrical	
1stFlrSF	
2ndFlrSF	
LowQualFinSF	
GrLivArea	
BsmtFullBath	
BsmtHalfBath	
FullBath	
HalfBath	
BedroomAbvGr	
KitchenAbvGr	
KitchenQual	
TotRmsAbvGrd	
Functional	
Fireplaces	
GarageType	
GarageYrBlt	
GarageFinish	
GarageCars	
GarageArea	
GarageQual	
GarageCond	
PavedDrive	
WoodDeckSF	
OpenPorchSF	
EnclosedPorch	
3SsnPorch	
ScreenPorch	
PoolArea	
MiscVal	
MoSold	
YrSold	
SaleType	
SaleCondition	
SalePrice	

```
In [66]: data3.isnull().sum()
```

```
Out[66]: 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                0
MasVnrArea                0
ExterQual                 0
ExterCond                 0
Foundation                0
BsmtQual                  0
BsmtCond                  0
BsmtExposure              0
BsmtFinType1              0
BsmtFinSF1                0
BsmtFinType2              0
BsmtFinSF2                0
BsmtUnfSF                 0
```

TotalBsmstSF	0
Heating	0
HeatingQC	0
CentralAir	0
Electrical	0
1stFlrSF	0
2ndFlrSF	0
LowQualFinSF	0
GrLivArea	0
BsmstFullBath	0
BsmstHalfBath	0
FullBath	0
HalfBath	0
BedroomAbvGr	0
KitchenAbvGr	0
KitchenQual	0
TotRmsAbvGrd	0
Functional	0
Fireplaces	0
GarageType	0
GarageYrBlt	0
GarageFinish	0
GarageCars	0
GarageArea	0
GarageQual	0
GarageCond	0
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
```

```
In [68]: data3.isnull().sum().sum()
```

```
Out[68]: 0
```

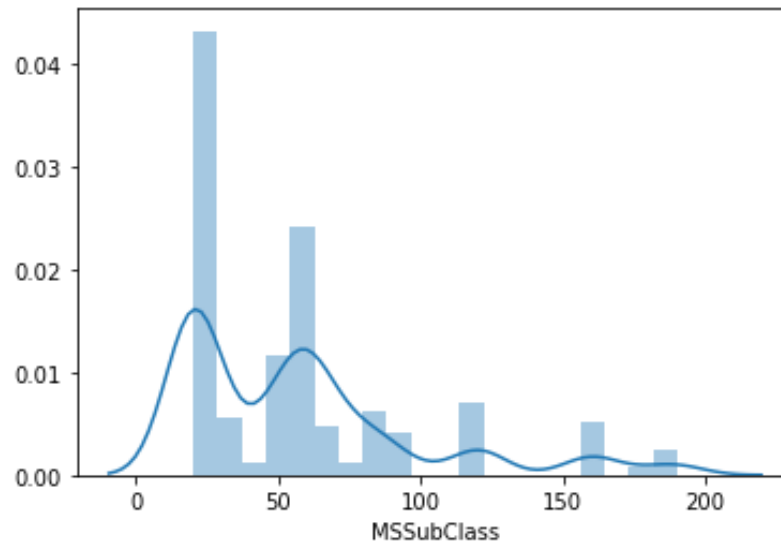
```
In [70]: data3.select_dtypes(include=["int64", "float64"]).columns
```

```
Out[70]: Index(['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'],  
              dtype='object')
```

In [72]: *#Previous data set the are not clens*

```
sns.distplot(data['MSSubClass'])
```

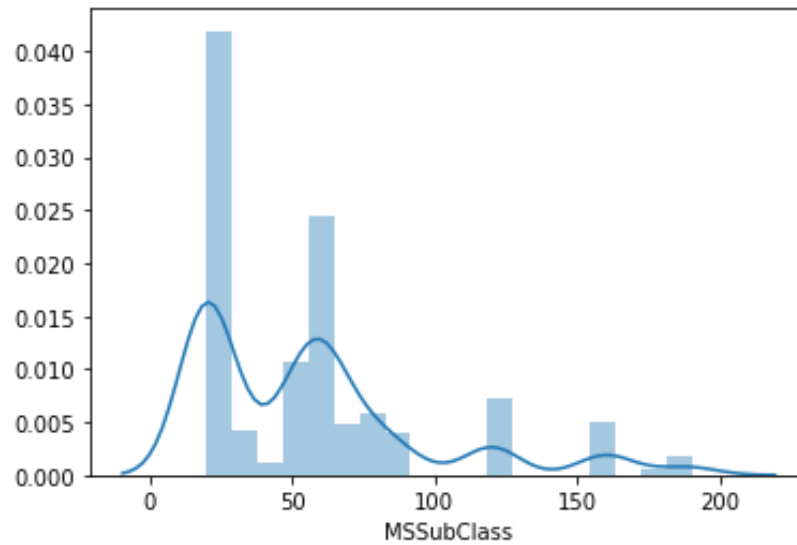
Out[72]: <AxesSubplot:xlabel='MSSubClass'>



In [73]: `# cleans data`

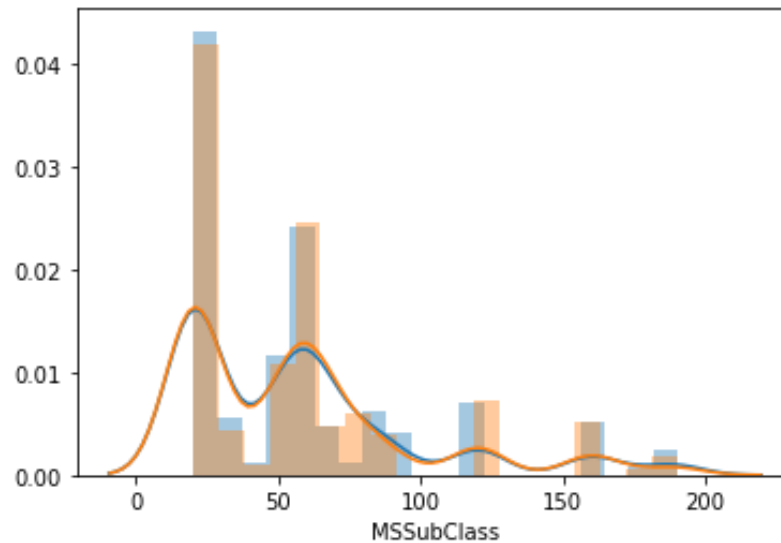
```
sns.distplot(data3['MSSubClass'])
```

Out[73]: `<AxesSubplot:xlabel='MSSubClass'>`



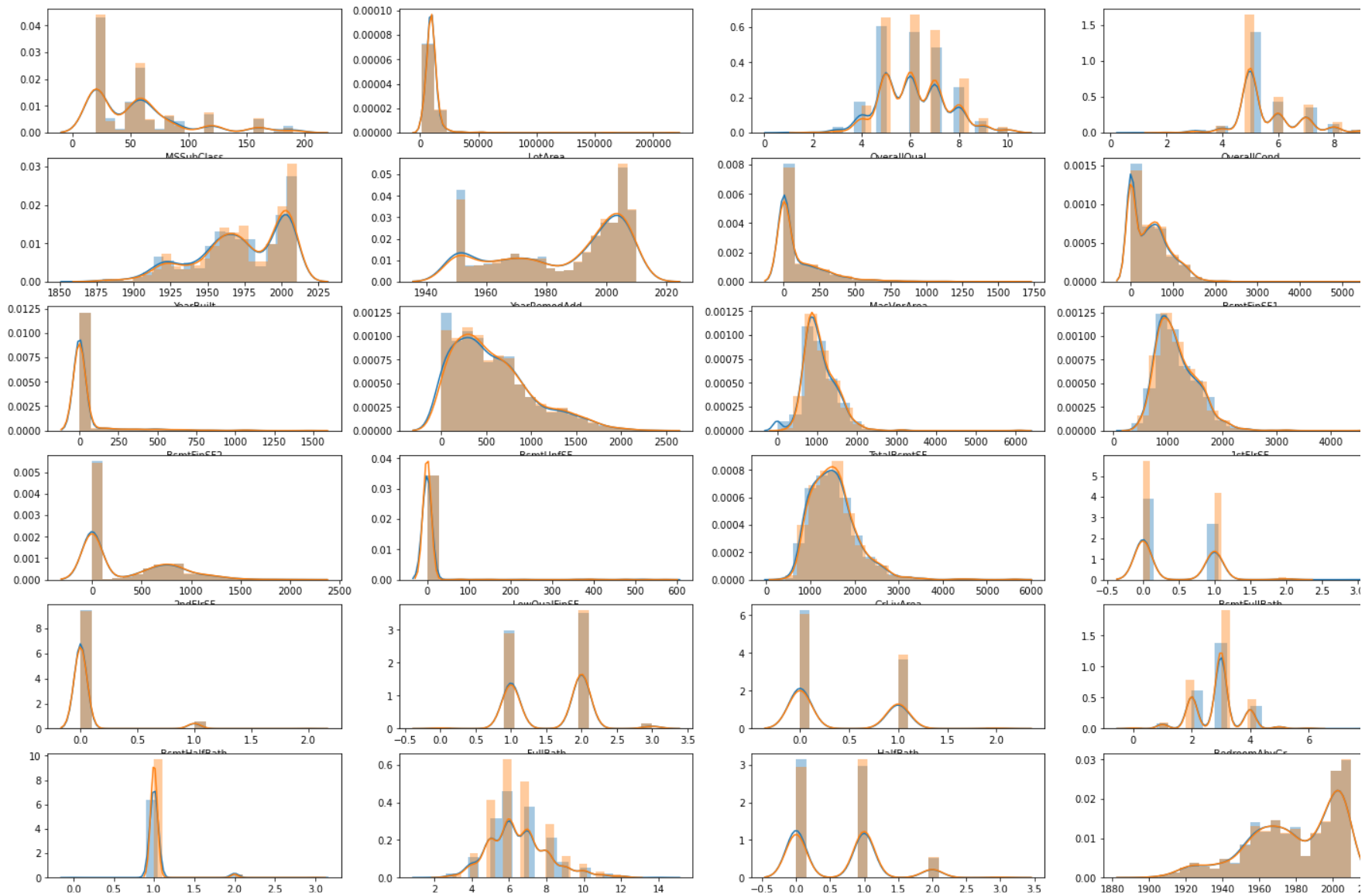
```
In [74]: #overpal 2 data set to comper the data
sns.distplot(data['MSSubClass'])
sns.distplot(data3['MSSubClass'])
```

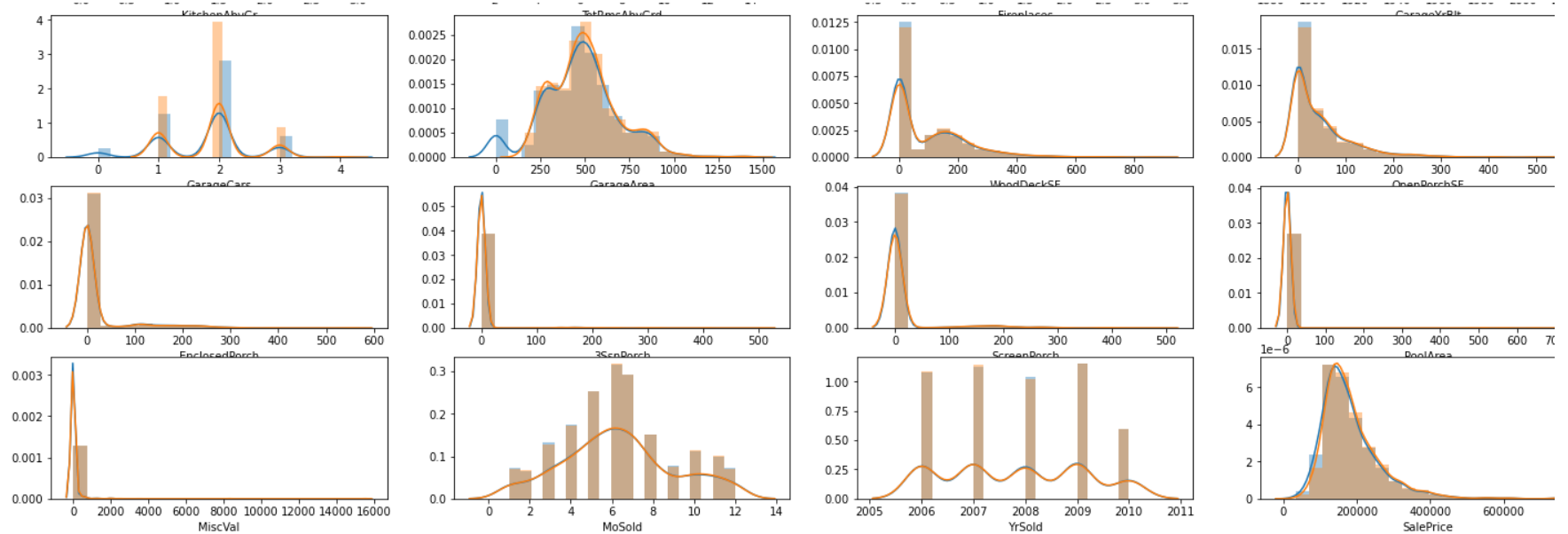
```
Out[74]: <AxesSubplot:xlabel='MSSubClass'>
```



```
In [76]: list_of_column_are_clean=['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']
```

```
In [77]: plt.figure(figsize=(25,25))
for i, var in enumerate(list_of_columns_to_clean):
    plt.subplot(9,4,i+1)
    sns.distplot(data[var],bins=20)
    sns.distplot(data3[var],bins=20)
```





```
In [78]: data3.select_dtypes(include=["object"]).columns
```

```
Out[78]: Index(['MSZoning', 'Street', 'LotShape', 'LandContour', 'Utilities',
               'LotConfig', 'LandSlope', 'Neighborhood', 'Condition1', 'Condition2',
               'BldgType', 'HouseStyle', 'RoofStyle', 'RoofMatl', 'Exterior1st',
               'Exterior2nd', 'MasVnrType', 'ExterQual', 'ExterCond', 'Foundation',
               'BsmtQual', 'BsmtCond', 'BsmtExposure', 'BsmtFinType1', 'BsmtFinType2',
               'Heating', 'HeatingQC', 'CentralAir', 'Electrical', 'KitchenQual',
               'Functional', 'GarageType', 'GarageFinish', 'GarageQual', 'GarageCond',
               'PavedDrive', 'SaleType', 'SaleCondition'],
              dtype='object')
```

```
In [84]: data["MSZoning"].value_counts()/data.shape[0]*100
```

```
Out[84]: RL          78.835616  
RM          14.931507  
FV           4.452055  
RH           1.095890  
C (all)      0.684932  
Name: MSZoning, dtype: float64
```

```
In [85]: data3["MSZoning"].value_counts()/data3.shape[0]*100
```

```
Out[85]: RL          79.671151  
RM          14.275037  
FV           4.633782  
RH           0.822123  
C (all)      0.597907  
Name: MSZoning, dtype: float64
```

```
In [87]: pd.concat([data["MSZoning"].value_counts()/data.shape[0]*100,  
data3["MSZoning"].value_counts()/data3.shape[0]*100],axis=1,keys=["MSZoning_org","MSZoning_clean"])
```

```
Out[87]:
```

	MSZoning_org	MSZoning_clean
RL	78.835616	79.671151
RM	14.931507	14.275037
FV	4.452055	4.633782
RH	1.095890	0.822123
C (all)	0.684932	0.597907

```
In [113]: def comp_clean_data(var):  
  
          return pd.concat([data[var].value_counts()/data.shape[0]*100,  
                           data3[var].value_counts()/data3.shape[0]*100],axis=1,keys=[var+"_org",var+"_clean"])
```

```
In [114]: comp_clean_data('MSZoning')
```

Out[114]:

	MSZoning_org	MSZoning_clean
RL	78.835616	79.671151
RM	14.931507	14.275037
FV	4.452055	4.633782
RH	1.095890	0.822123
C (all)	0.684932	0.597907