

```
1 import pandas as pd
2 import numpy as np
3 import random
4 from random import random
5 import seaborn as sns
6 import matplotlib.pyplot as plt
7 %matplotlib inline
8 import warnings
9 warnings.filterwarnings('ignore')
10 sns.set(style='darkgrid',font_scale=1.2)
11 plt.rcParams['font.family']='SimHei'
12 plt.rcParams['axes.unicode_minus']=False
```

```
1 train = pd.read_csv('train.csv',index_col='Id')
2 train.head()
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

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

5 rows × 80 columns

```
1 test = pd.read_csv('test.csv',index_col='Id')
2 test.head()
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Utilities	LotConfig	...	ScreenPorch	Pool
Id													
1461	20	RH	80.0	11622	Pave	NaN	Reg	Lvl	AllPub	Inside	...	120	0
1462	20	RL	81.0	14267	Pave	NaN	IR1	Lvl	AllPub	Corner	...	0	0
1463	60	RL	74.0	13830	Pave	NaN	IR1	Lvl	AllPub	Inside	...	0	0
1464	60	RL	78.0	9978	Pave	NaN	IR1	Lvl	AllPub	Inside	...	0	0
1465	120	RL	43.0	5005	Pave	NaN	IR1	HLS	AllPub	Inside	...	144	0

5 rows × 79 columns

```
1 | train.shape
```

```
1 | (1460, 80)
```

```
1 | train.describe()
```

```
1 | .dataframe tbody tr th {
2 |     vertical-align: top;
3 | }
4 |
5 | .dataframe thead th {
6 |     text-align: right;
7 | }
```

	MSSubClass	LotFrontage	LotArea	OverallQual	OverallCond	YearBuilt	YearRemodAdd	MasVnrArea	BsmtFinSF1	BsmtFinSF2
count	1460.000000	1201.000000	1460.000000	1460.000000	1460.000000	1460.000000	1460.000000	1452.000000	1460.000000	1460.000000
mean	56.897260	70.049958	10516.828082	6.099315	5.575342	1971.267808	1984.865753	103.685262	443.639726	46.549737
std	42.300571	24.284752	9981.264932	1.382997	1.112799	30.202904	20.645407	181.066207	456.098091	161.314169
min	20.000000	21.000000	1300.000000	1.000000	1.000000	1872.000000	1950.000000	0.000000	0.000000	0.000000
25%	20.000000	59.000000	7553.500000	5.000000	5.000000	1954.000000	1967.000000	0.000000	0.000000	0.000000
50%	50.000000	69.000000	9478.500000	6.000000	5.000000	1973.000000	1994.000000	0.000000	383.500000	0.000000
75%	70.000000	80.000000	11601.500000	7.000000	6.000000	2000.000000	2004.000000	166.000000	712.250000	0.000000
max	190.000000	313.000000	215245.000000	10.000000	9.000000	2010.000000	2010.000000	1600.000000	5644.000000	1474.000000

8 rows × 37 columns

```
1 | train.columns
```

```
1 | Index(['MSSubClass', 'MSZoning', 'LotFrontage', 'LotArea', 'Street', 'Alley',
2 |        'LotShape', 'LandContour', 'Utilities', 'LotConfig', 'LandSlope',
3 |        'Neighborhood', 'Condition1', 'Condition2', 'BldgType', 'HouseStyle',
4 |        'OverallQual', 'OverallCond', 'YearBuilt', 'YearRemodAdd', 'RoofStyle',
5 |        'RoofMatl', 'Exterior1st', 'Exterior2nd', 'MasVnrType', 'MasVnrArea',
6 |        'ExterQual', 'ExterCond', 'Foundation', 'BsmtQual', 'BsmtCond',
7 |        'BsmtExposure', 'BsmtFinType1', 'BsmtFinSF1', 'BsmtFinType2',
8 |        'BsmtFinSF2', 'BsmtUnfsf', 'TotalBsmtSF', 'Heating', 'HeatingQC',
9 |        'CentralAir', 'Electrical', '1stFlrSF', '2ndFlrSF', 'LowQualFinSF',
10 |        'GrLivArea', 'BsmtFullBath', 'BsmtHalfBath', 'FullBath', 'HalfBath',
11 |        'BedroomAbvGr', 'KitchenAbvGr', 'KitchenQual', 'TotRmsAbvGrd',
12 |        'Functional', 'Fireplaces', 'FireplaceQu', 'GarageType', 'GarageYrBlt',
13 |        'GarageFinish', 'GarageCars', 'GarageArea', 'GarageQual', 'GarageCond',
14 |        'PavedDrive', 'WoodDeckSF', 'OpenPorchSF', 'EnclosedPorch', '3SsnPorch',
15 |        'ScreenPorch', 'PoolArea', 'PoolQC', 'Fence', 'MiscFeature', 'MiscVal',
16 |        'MoSold', 'YrSold', 'SaleType', 'SaleCondition', 'SalePrice'],
17 |        dtype='object')
```

Column Info

- MSSubClass: Identifies the type of dwelling involved in the sale.

```
1 | # MSSubClass 用于分类
2 | train['MSSubClass'] = train['MSSubClass'].astype(str)
3 | test['MSSubClass'] = test['MSSubClass'].astype(str)
4 | print(train['MSSubClass'].value_counts(),test['MSSubClass'].value_counts())
```

1	20	536
2	60	299
3	50	144
4	120	87
5	30	69

```

6 160      63
7  70      60
8  80      58
9  90      52
10 190      30
11 85       20
12 75       16
13 45       12
14 180      10
15 40        4
16 Name: MSSubClass, dtype: int64 20      543
17 60       276
18 50       143
19 120      95
20 30       70
21 70       68
22 160      65
23 80       60
24 90       57
25 190      31
26 85       28
27 180       7
28 75       7
29 45       6
30 40       2
31 150      1
32 Name: MSSubClass, dtype: int64

```

```

1 quant = [x for x in train.columns if train[x].dtypes != object]
2 quanli = [x for x in train.columns if train[x].dtypes == object]
3 print('quant: {}, counts: {}'.format(quant,len(quant)))
4 print('-----')
5 print('quanli: {}, counts: {}'.format(quanli,len(quanli)))

```

```

1 quant: ['LotFrontage', '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'], counts: 36
2 -----
3 quanli: ['MSSubClass', 'MSZoning', 'Street', 'Alley', 'LotShape', 'LandContour', 'Utilities', 'LotConfig', 'LandSlope', 'Neighborhood',
'Condition1', 'Condition2', 'BldgType', 'HouseStyle', 'RoofStyle', 'RoofMat1', 'Exterior1st', 'Exterior2nd', 'MasVnrType', 'ExterQual',
'ExterCond', 'Foundation', 'BsmtQual', 'BsmtCond', 'BsmtExposure', 'BsmtFinType1', 'BsmtFinType2', 'Heating', 'HeatingQC', 'CentralAir',
'Electrical', 'KitchenQual', 'Functional', 'FireplaceQu', 'GarageType', 'GarageFinish', 'GarageQual', 'GarageCond', 'PavedDrive', 'PoolQC',
'Fence', 'MiscFeature', 'SaleType', 'SaleCondition'], counts: 44

```

```

1 total = train.isnull().sum().sort_values(ascending=False)
2 total

```

```

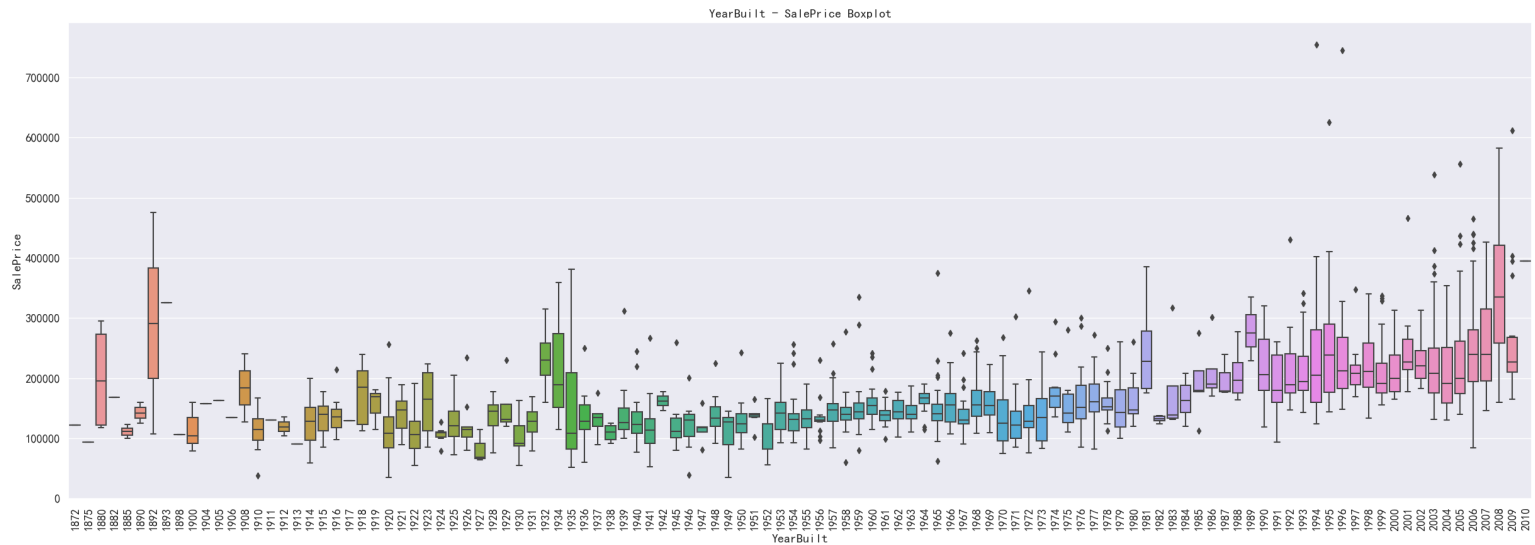
1 PoolQC      1453
2 MiscFeature 1406
3 Alley      1369
4 Fence      1179
5 FireplaceQu    690
6 ...
7 CentralAir      0
8 SaleCondition    0
9 Heating          0
10 TotalBsmtSF      0
11 MSSubClass      0
12 Length: 80, dtype: int64

```

```

1 plt.figure(figsize=(30,10),dpi=100)
2 sns.boxplot(train.YearBuilt,train.SalePrice)
3 plt.title('YearBuilt - SalePrice Boxplot')
4 plt.xticks(rotation=90)
5 plt.savefig('YearBuilt - SalePrice Boxplot.png',dpi=80)
6 plt.show()

```



Insights

- Outlier prices should be trimmed.
- Normal price: [around 100k, around 400k]
- TimeSeries Line plot needed to demonstrate the fluctuation in annual price.

```
1 train['SalePrice'].describe()
```

```
1 count      1460.000000
2 mean       180921.195890
3 std        79442.502883
4 min        34900.000000
5 25%        129975.000000
6 50%        163000.000000
7 75%        214000.000000
8 max        755000.000000
9 Name: SalePrice, dtype: float64
```

```
1 train = train[train.SalePrice >= 40000]
```

```
1 train = train[train.SalePrice <= 500000]
```

```
1 train.shape
```

```
1 (1447, 80)
```

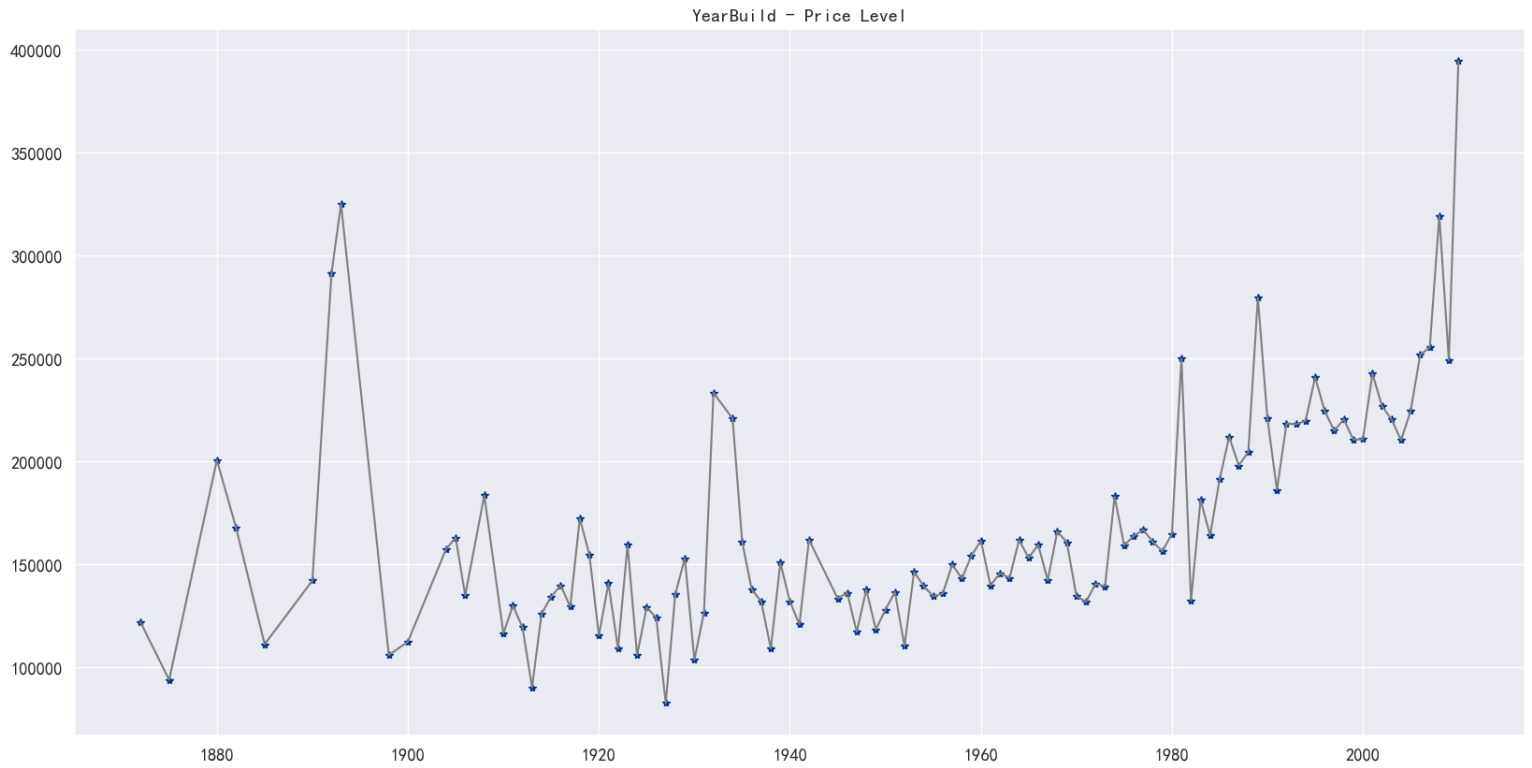
```
1 df1 = train.groupby('YearBuilt').agg({'SalePrice': 'mean'})
2 df1
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	SalePrice
YearBuilt	
1872	122000.000000
1875	94000.000000
1880	200619.750000
1882	168000.000000
1885	111250.000000
...	...
2006	251775.447761
2007	255362.734694
2008	319188.000000
2009	249076.647059
2010	394432.000000

112 rows × 1 columns

```
1 plt.figure(figsize=(20,10),dpi=100)
2 # df1.plot(color='#00338D')
3 plt.plot(df1,"*",color="#00338D")
4 plt.plot(df1,color="gray")
5 plt.title("YearBuild - Price Level")
6 # plt.savefig('YearBuild - Price Level', dpi =100)
7 plt.show()
```



Insights

- General annual trend of *Salesprice*: Up

Conclusion: YearBuilt correlate with SalesPrice.

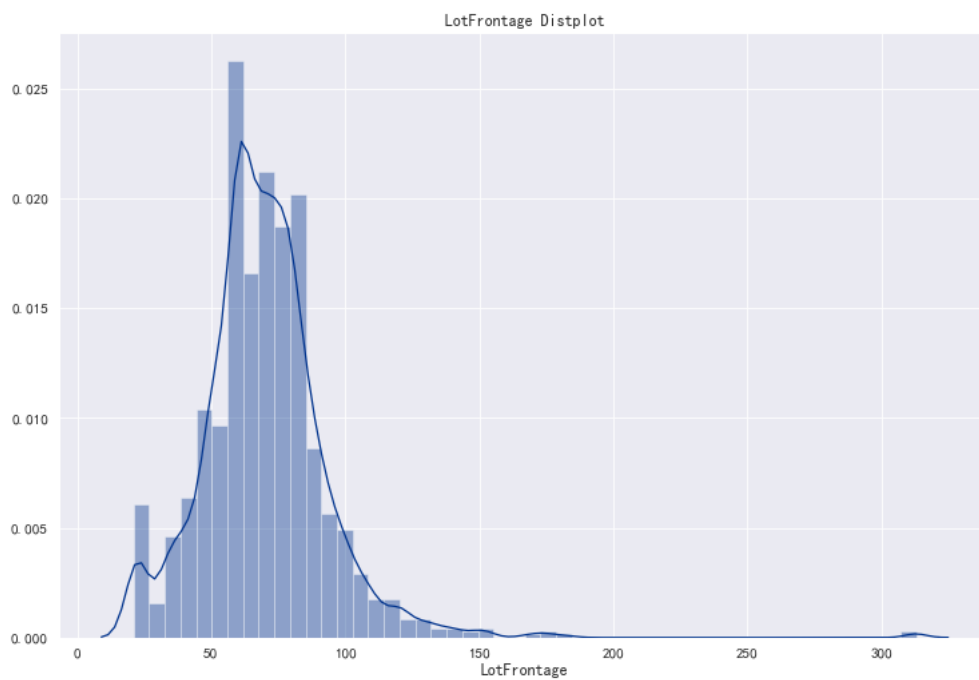
```
1 train.LotFrontage.describe()  #[21,313]
```

```
1 count    1188.000000
2 mean       69.891414
3 std        24.119385
4 min        21.000000
5 25%        59.000000
6 50%        69.000000
7 75%        80.000000
8 max       313.000000
9 Name: LotFrontage, dtype: float64
```

```
1 train.LotFrontage.isnull().sum()
```

```
1 259
```

```
1 plt.figure(figsize=(15,10),dpi=60)
2 sns.distplot(train.LotFrontage,color='#00338D')
3 plt.title('LotFrontage Distplot')
4 plt.savefig('LotFrontage Distplot.png', dpi =100)
5 plt.show()
```



Insights

- Lmt: [0,200]
- Skewed distribution: fillna with median value

```
1 train = train[train['LotFrontage']<=200]
```

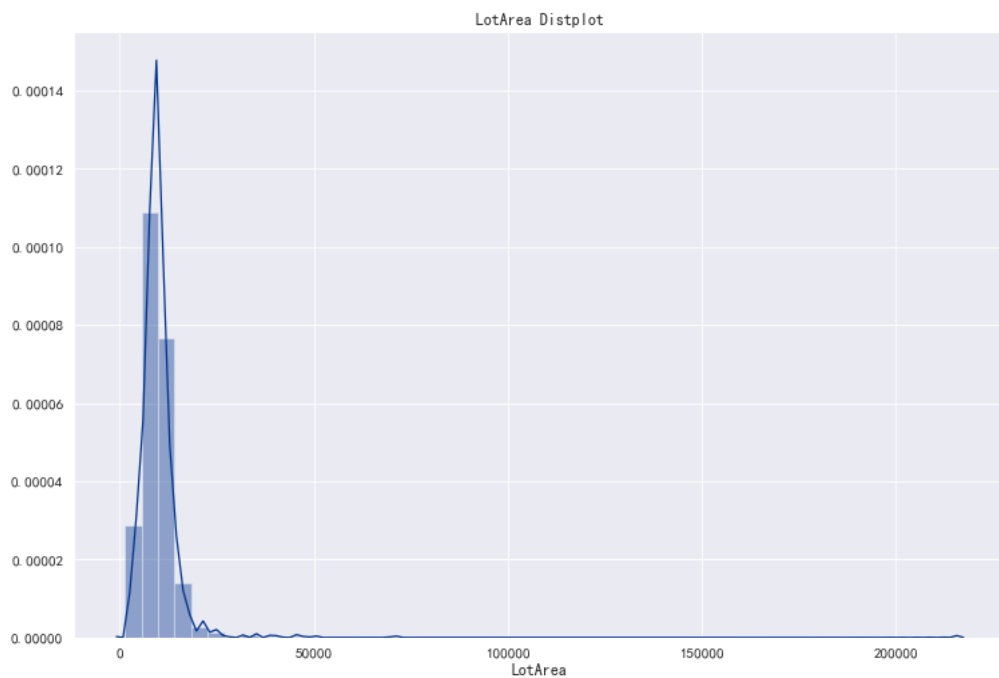
```
1 train.fillna({'LotFrontage':train.LotFrontage.median()},inplace=True)
2 train.LotFrontage.isnull().sum()
```

```
1 0
```

```
1 train.LotArea.describe() # [1.3k,215k+]
```

```
1 count      1186.000000
2 mean       9806.498314
3 std        7640.920262
4 min        1300.000000
5 25%        7409.000000
6 50%        9245.500000
7 75%        11205.250000
8 max        215245.000000
9 Name: LotArea, dtype: float64
```

```
1 plt.figure(figsize=(15,10),dpi=60)
2 sns.distplot(train.LotArea,color='#00338d')
3 plt.title('LotArea Distplot')
4 plt.show()
```



```
1 train = train[train['LotArea'] <= 50000]
```

```
1 train['LotArea'].isnull().sum()
```

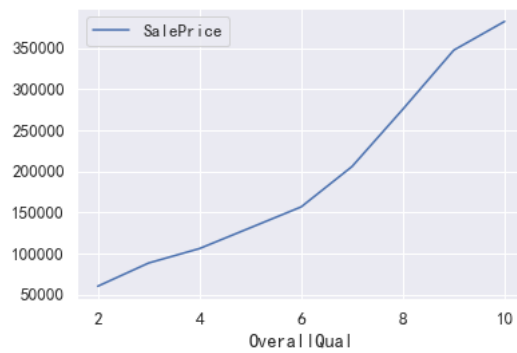
```
1 0
```

```
1 train.OverallQual.isnull().sum()
```

```
1 0
```

```
1 df2 = train.groupby('OverallQual').agg({'SalePrice': 'mean'})
2 df2.plot()
```

```
1 <AxesSubplot: xlabel='OverallQual'>
```



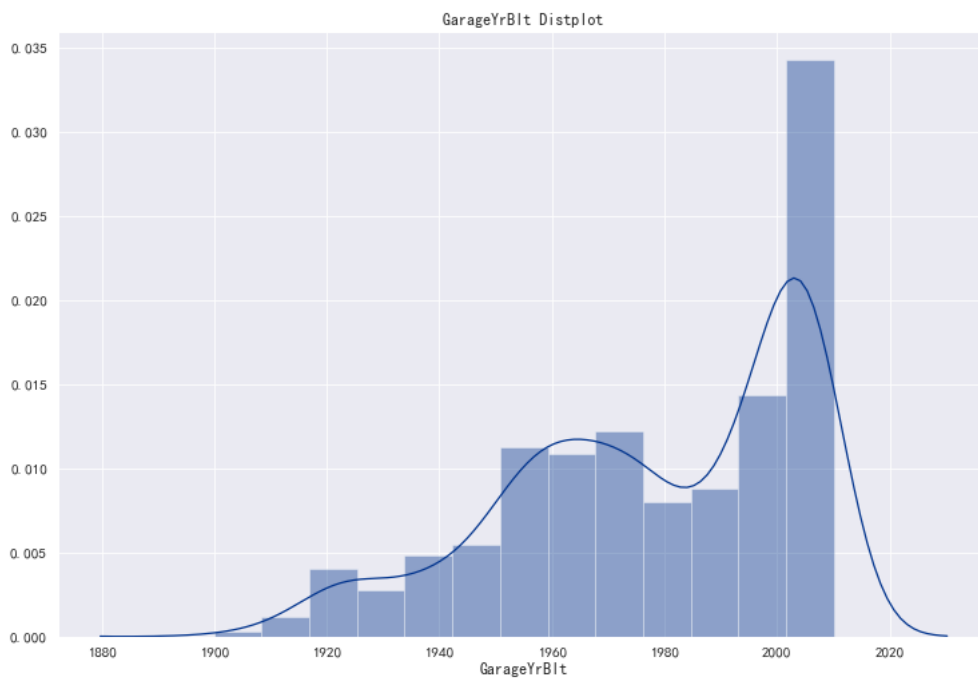
Insights

- SalePrice positively relates to overall rates of material & finish of house (AKA *OverallQual*)

```
1 null_var = [i for i in quant if train[i].isnull().sum() >0]
2 null_var
```

```
1 ['MasVnrArea', 'GarageYrBlt']
```

```
1 plt.figure(figsize=(15,10),dpi=60)
2 sns.distplot(train.GarageYrBlt,color='#00338D')
3 plt.title('GarageYrBlt Distplot')
4 plt.show()
```



```
1 train.fillna({'MasVnrArea':train.MasVnrArea.median(),'GarageYrBlt':train.GarageYrBlt.median()},inplace=True)
```

```
1 null_var = [i for i in quant if train[i].isnull().sum() >0]
2 null_var
```

```
1 []
```

Fillna finished.

```
1 train.Utilities.unique()
```

```
1 array(['AllPub'], dtype=object)
```



```
1 df3= train.groupby('Utilities').agg({'SalePrice':'mean'})
2 df3
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	SalePrice
Utilities	
AllPub	177552.770921

```
1 test.Utilities.unique()
```

```
1 array(['AllPub', nan], dtype=object)
```

- train 和 test Allpub 的 unique值不同，train中都是AllPub，所以对预测没有帮助，需要drop。

```
1 test.drop('Utilities',axis=1,inplace=True)
2 train.drop('Utilities',axis=1,inplace=True)
```

```
1 import sklearn
2 from sklearn import preprocessing
3 from sklearn.preprocessing import LabelEncoder
4 def leo(x):
5     train[x]=LabelEncoder().fit_transform(train[x])
```

```
1 for i in quanli:
2     leo(i)
```

```
1 # 将分类数据转化成整数编码
2 # 获取分类变量的标签值
3 def Labs(x):
4     print(x,"--",train[x].unique())
5 df_obj = train[quanli]
6 print(list(map(Labs,df_obj)))
```

```
1 MSSubClass -- [ 9  4 10  8  3  7 14  0  5 12  1 11  2  6 13]
2 MSZoning -- [3 4 0 1 2]
3 Street -- [1 0]
4 Alley -- [2 0 1]
5 LotShape -- [3 0 1 2]
6 LandContour -- [3 0 1 2]
7 LotConfig -- [4 2 0 1 3]
8 Landslope -- [0 1 2]
9 Neighborhood -- [ 5 24  6 15 11 21 17  3 19 16 20 12  9 10  7 23 22  4  8 14 13  0  2 18
10 1]
11 Condition1 -- [2 1 0 5 8 6 4 3 7]
12 Condition2 -- [2 0 5 1 4 3]
13 BldgType -- [0 1 2 4 3]
14 HouseStyle -- [5 2 0 1 7 4 3 6]
15 RoofStyle -- [1 3 2 4 0]
16 RoofMat1 -- [0 5 1 4 3 2]
17 Exterior1st -- [12  8 13  3  6 14  5  0 11  9  2  1 10  7  4]
18 Exterior2nd -- [13  8 15  6 14 10  5  3 12  0  2  7  1  9 11  4]
19 MasVnrType -- [1 2 3 0 4]
20 ExterQual -- [2 3 0 1]
21 ExterCond -- [4 1 2 3 0]
22 Foundation -- [2 1 0 5 3 4]
23 BsmtQual -- [2 3 0 4 1]
24 BsmtCond -- [3 1 4 0 2]
25 BsmtExposure -- [3 1 2 0 4]
```

[illegible]

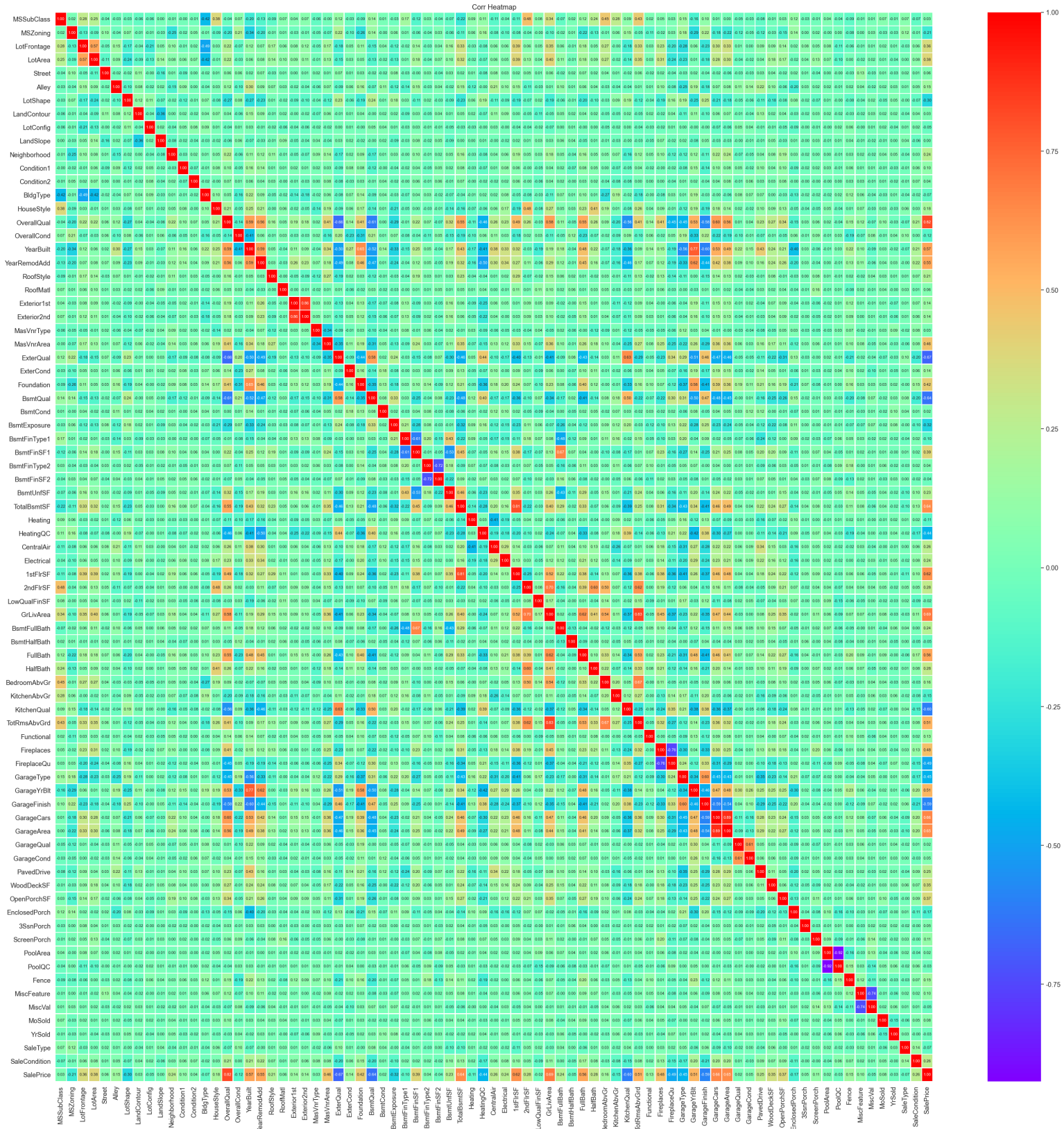
```
1 corr = train.corr()
2 corr
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	LotConfig	LandSlope	...
MSSubClass	1.000000	0.020940	0.279429	0.247561	-0.037199	-0.028659	-0.030395	-0.061403	-0.059057	0.002714	...
MSZoning	0.020940	1.000000	-0.129692	-0.085526	0.103255	-0.035657	0.068120	-0.013183	-0.012228	-0.029553	...
LotFrontage	0.279429	-0.129692	1.000000	0.570965	-0.047611	0.146923	-0.170013	-0.035607	-0.212730	0.049202	...
LotArea	0.247561	-0.085526	0.570965	1.000000	-0.105095	0.087676	-0.235314	-0.091517	-0.131644	0.139614	...
Street	-0.037199	0.103255	-0.047611	-0.105095	1.000000	-0.017300	-0.021099	0.105400	-0.001474	-0.164011	...
...
MoSold	0.067412	-0.030336	0.022697	0.005955	0.013222	0.021470	-0.049747	0.020259	0.016360	-0.008459	...
YrSold	-0.013556	-0.032690	0.005311	-0.036171	-0.029317	0.028094	0.050581	0.023113	-0.035390	0.000823	...
SaleType	0.070457	0.120846	-0.032287	0.003958	0.020851	0.010825	-0.000707	-0.036675	0.014615	0.041055	...
SaleCondition	-0.069217	-0.009377	0.061385	0.081866	0.011092	0.048470	-0.070626	0.041733	0.023852	-0.058860	...
SalePrice	0.030731	-0.208654	0.361636	0.383817	0.059211	0.151276	-0.301186	0.032256	-0.050879	0.002702	...

79 rows x 79 columns

```
1 plt.figure(figsize=(40,40),dpi=100)
2 sns.set(font_scale=1.2)
3 ax=sns.heatmap(corr,
4                 xticklabels=corr.columns,
5                 yticklabels=corr.columns,
6                 linewidths=0.9,annot=True,
7                 cbar=True,cmap="rainbow",fmt='.2f',
8                 annot_kws={'size': 8})
9 plt.title("Corr Heatmap")
10 plt.savefig("Corr Heatmap.png")
11 plt.show()
```



Insights

- 强相关显然保留
- 中相关性和弱相关性，Deep Dive变量之间相关性，选择性保留。

```
1 strong_var = ['OverallQual','YearBuilt','YearRemodAdd','ExterQual','BsmtQual','TotalBsmtSF','1stFlrSF','GrLivArea','KitchenQual',
2              'GarageCars','GarageArea']
```

```
1 other_var = list(train.columns)
```

```
1 for i in strong_var:
2     other_var.remove(i)
```

```
1 print(len(strong_var))
```

```
1 11
```

```
1 k = 15
2 var_ = corr.nlargest(k, 'SalePrice')['SalePrice'].index
```

- strong_var
- corr.nlargest(k,'SalePrice')

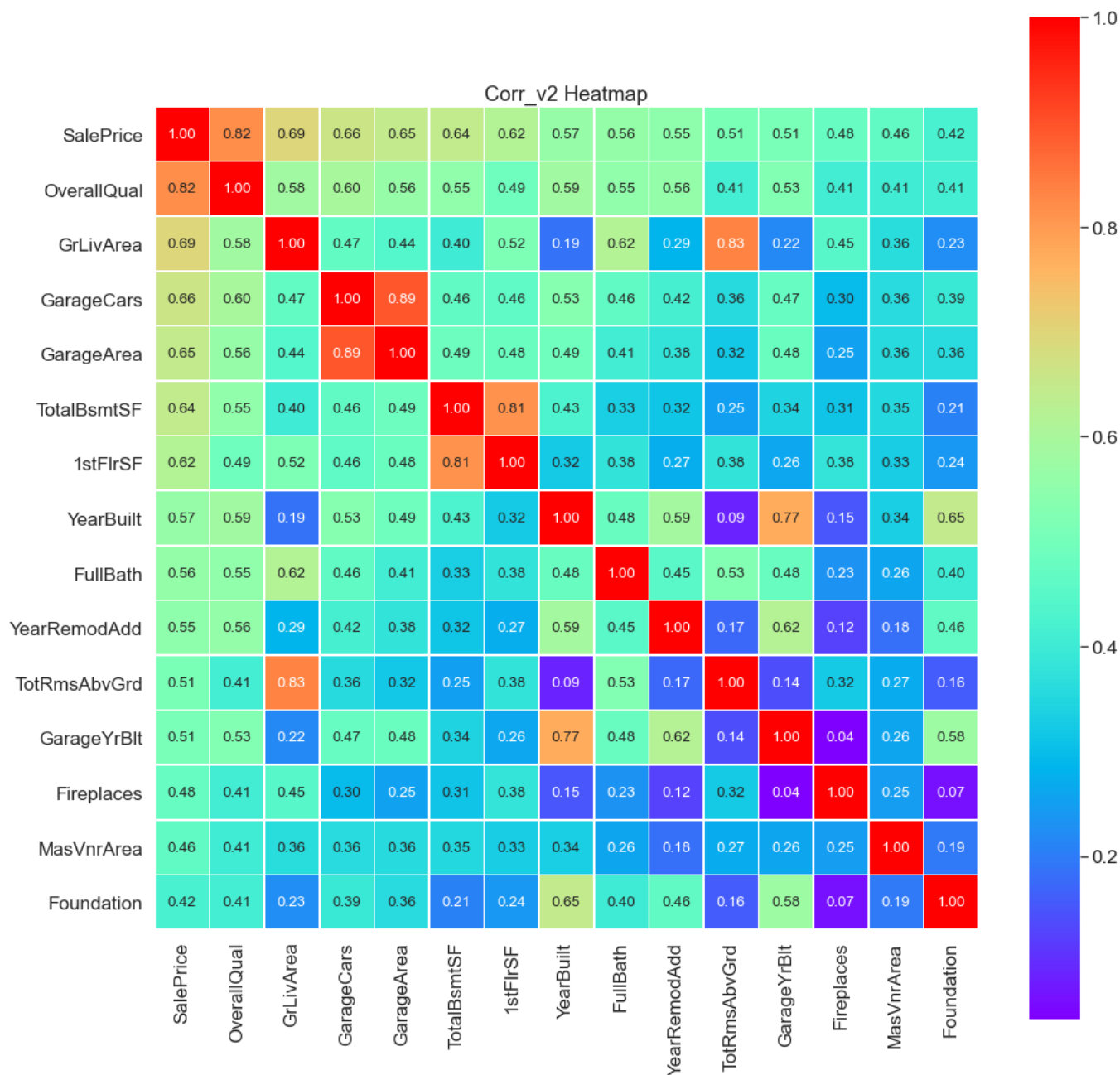
```
1 cm = np.corrcoef(train[var_].values.T)
2 cm
```

```
1 array([[1.          , 0.81836496, 0.69405958, 0.66058917, 0.6462927 ,
2          0.64310743, 0.61768451, 0.56574773, 0.56303443, 0.55088068,
3          0.51319332, 0.51063903, 0.47549164, 0.46349279, 0.4213752 ],
4          [0.81836496, 1.          , 0.58450966, 0.59923389, 0.56356986,
5          0.55246083, 0.489251  , 0.58554145, 0.55424945, 0.55683769,
6          0.41286386, 0.5271104 , 0.40821558, 0.40520085, 0.41226335],
7          [0.69405958, 0.58450966, 1.          , 0.4670706 , 0.4428922 ,
8          0.40127177, 0.51535057, 0.188675  , 0.61736605, 0.28689558,
9          0.83232925, 0.21961226, 0.44688989, 0.36195117, 0.22613848],
10         [0.66058917, 0.59923389, 0.4670706 , 1.          , 0.89365001,
11         0.4630363 , 0.45940744, 0.53472301, 0.46250507, 0.41582948,
12         0.35733104, 0.47382228, 0.29591308, 0.3612021 , 0.39468072],
13         [0.6462927 , 0.56356986, 0.4428922 , 0.89365001, 1.          ,
14         0.48688889, 0.48195849, 0.48984406, 0.40908054, 0.37740308,
15         0.31963221, 0.48175325, 0.2545989 , 0.35714997, 0.35877619],
16         [0.64310743, 0.55246083, 0.40127177, 0.4630363 , 0.48688889,
17         1.          , 0.81032781, 0.42749083, 0.33351637, 0.31553456,
18         0.25034587, 0.34009332, 0.31201215, 0.35269621, 0.20686582],
19         [0.61768451, 0.489251  , 0.51535057, 0.45940744, 0.48195849,
20         0.81032781, 1.          , 0.32287916, 0.37844895, 0.27261388,
21         0.37774798, 0.26200353, 0.38144482, 0.33066204, 0.23862028],
22         [0.56574773, 0.58554145, 0.188675  , 0.53472301, 0.48984406,
23         0.42749083, 0.32287916, 1.          , 0.48007967, 0.59203061,
24         0.08783456, 0.7741855 , 0.15036689, 0.33522184, 0.64560153],
25         [0.56303443, 0.55424945, 0.61736605, 0.46250507, 0.40908054,
26         0.33351637, 0.37844895, 0.48007967, 1.          , 0.45160637,
27         0.52924064, 0.48039193, 0.23192388, 0.25997238, 0.39983469],
28         [0.55088068, 0.55683769, 0.28689558, 0.41582948, 0.37740308,
29         0.31553456, 0.27261388, 0.59203061, 0.45160637, 1.          ,
30         0.17465557, 0.61959647, 0.12399901, 0.1828303 , 0.46318838],
31         [0.51319332, 0.41286386, 0.83232925, 0.35733104, 0.31963221,
32         0.25034587, 0.37774798, 0.08783456, 0.52924064, 0.17465557,
33         1.          , 0.13779946, 0.32415371, 0.26946452, 0.1599862 ],
34         [0.51063903, 0.5271104 , 0.21961226, 0.47382228, 0.48175325,
35         0.34009332, 0.26200353, 0.7741855 , 0.48039193, 0.61959647,
36         0.13779946, 1.          , 0.04441026, 0.26031938, 0.57677663],
37         [0.47549164, 0.40821558, 0.44688989, 0.29591308, 0.2545989 ,
38         0.31201215, 0.38144482, 0.15036689, 0.23192388, 0.12399901,
39         0.32415371, 0.04441026, 1.          , 0.24782868, 0.06650263],
40         [0.46349279, 0.40520085, 0.36195117, 0.3612021 , 0.35714997,
41         0.35269621, 0.33066204, 0.33522184, 0.25997238, 0.1828303 ,
42         0.26946452, 0.26031938, 0.24782868, 1.          , 0.19462898],
43         [0.4213752 , 0.41226335, 0.22613848, 0.39468072, 0.35877619,
44         0.20686582, 0.23862028, 0.64560153, 0.39983469, 0.46318838,
45         0.1599862 , 0.57677663, 0.06650263, 0.19462898, 1.          ]])
```

```

1 plt.figure(figsize=(15,15),dpi=80)
2 sns.set(font_scale=1.4)
3 hm = sns.heatmap(cm, cbar=True, linewidths=0.5,
4                 annot=True, square=True,
5                 fmt='.2f',annot_kws={'size': 12}, cmap="rainbow",
6                 yticklabels=var_.values, xticklabels=var_.values)
7 plt.title("Corr_v2 Heatmap")
8 plt.savefig("Corr_v2 Heatmap.png")
9 plt.show()

```



- 变量之间相关性强的需要drop

```
1 vars_ =list(var_)
```

```
1 vars_.remove('TotRmsAbvGrd')
```

```
1 vars_.remove('GarageArea')
```

```
1 vars_.remove('1stFlrSF')
```

```
1 vars_.remove('GarageYrBlt')
```

```
1 print('Finally var counts: {} {}var: {}'.format(len(vars_),'\n',vars_))
```

```

1 Finally var counts: 11
2 var: ['SalePrice', 'OverallQual', 'GrLivArea', 'GarageCars', 'TotalBsmtSF', 'YearBuilt', 'FullBath', 'YearRemodAdd', 'Fireplaces', 'MasVnrArea',
'Foundation']

```

```
1 vars_.remove('SalePrice')
2 x_train_fn1 = train[vars_]
3 y_train_fn1 = train['SalePrice']
```

```
1 x_test_fn1 = test[vars_]
```

```
1 def leo2(x):
2     x_test_fn1[x]=LabelEncoder().fit_transform(x_test_fn1[x])
```

```
1 for i in vars_:
2     leo2(i)
```

```
1 x_test_fn1.head()
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	OverallQual	GrLivArea	GarageCars	TotalBsmtSF	YearBuilt	FullBath	YearRemodAdd	Fireplaces	MasVnrArea	Foundation
Id										
1461	4	64	1	202	56	1	11	0	0	1
1462	5	313	1	477	53	1	8	0	58	1
1463	4	518	2	234	92	2	48	1	0	2
1464	5	500	2	233	93	2	48	1	6	2
1465	7	282	2	455	87	2	42	0	0	2

```
1 from sklearn.preprocessing import StandardScaler
2 scaler = StandardScaler(copy=False)
3 scaler.fit_transform(x_train_fn1)
4 train2 = scaler.transform(x_train_fn1)
5 train2
```

```
1 array([[ 0.65288347,  0.43354519,  0.31409347, ..., -0.90822251,
2          0.55530153,  0.82762684],
3        [-0.07438691, -0.47798807,  0.31409347, ...,  0.70745753,
4          -0.56020872, -0.52097172],
5        [ 0.65288347,  0.58818029,  0.31409347, ...,  0.70745753,
6          0.36179465,  0.82762684],
7        ...,
8        [ 0.65288347,  1.71538883, -0.98511134, ...,  2.32313758,
9          -0.56020872,  3.52482395],
10       [-0.80165729, -0.8523678 , -0.98511134, ..., -0.90822251,
11        -0.56020872, -0.52097172],
12       [-0.80165729, -0.49019611, -0.98511134, ..., -0.90822251,
13        -0.56020872, -0.52097172]])
```

```
1 from sklearn.linear_model import LinearRegression
```

```
1 scaler2 = StandardScaler(copy=False)
2 scaler2.fit_transform(x_test_fn1)
3 x_test = scaler2.transform(x_test_fn1)
4 lr =LinearRegression()
```

```

1 from sklearn.model_selection import train_test_split #数据集训练集划分
2 # from sklearn import metrics
3 from sklearn.metrics import adjusted_rand_score
4 from sklearn.metrics import accuracy_score
5 from sklearn.metrics import classification_report,precision_score,recall_score,f1_score,r2_score #分类报告
6 x_train,x_test,y_train,y_test=train_test_split(x_train_fn1,y_train_fn1,test_size=.2,random_state=22)
7 lr.fit(x_train,y_train)
8 pred = lr.predict(x_test)

```

```

1 y_test = pd.DataFrame(y_test)
2 pred = pd.DataFrame(pred)
3 print('r2_score:',r2_score(pred,y_test))

```

```

1 r2_score: 0.8504507414086491

```

```

1 # print('f1_score:',precision_score(pred,y_test))
2 xx = np.arange(len(y_test))

```

```

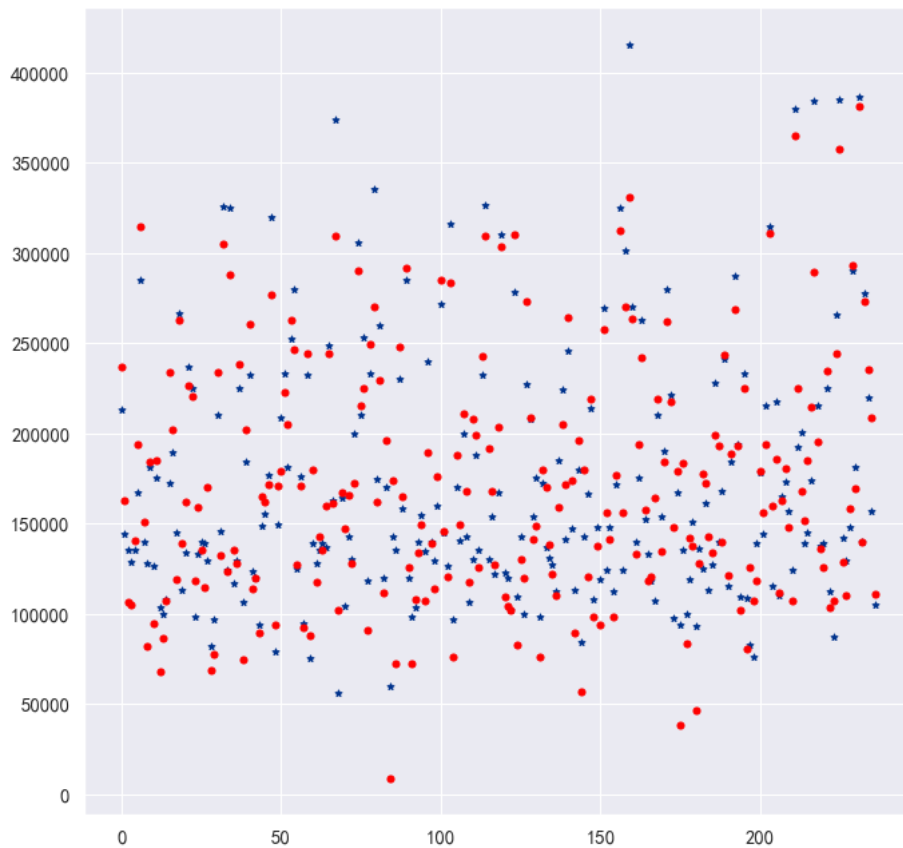
1 plt.figure(figsize = (10,10), dpi=80)
2 plt.scatter(xx, y_test, color="#00338D",s= 20, marker='*')
3 plt.scatter(xx,pred,color="red", s= 20,marker = 'o')

```

```

1 <matplotlib.collections.PathCollection at 0x1129b838850>

```



```

1 lr.fit(train2,y_train_fn1.values)
2 y_train_pred = lr.predict(x_test)
3 y_train_pred

```

```

1 array([ 96036.72599752, 154142.89713714, 181295.75480668, ...,
2        165202.97140824, 112613.91975802, 244607.5935728 ])

```

```
1 ypred = pd.DataFrame(y_train_pred)
2 # ypred.to_csv('Boston ypred v1.csv')
```