

AirBnb listings file modeling

```
In [1]: 1 # Put these at the top of every notebook, to get automatic reloading and inline
2 from IPython.core.display import display, HTML
3 import pandas as pd
4 import warnings
5 import ast
6 warnings.filterwarnings('ignore')
7
8 %reload_ext autoreload
9 %autoreload 1
10 %matplotlib inline
11
12 pd.set_option('display.max_rows', 500)
13 pd.set_option('display.max_columns', 500)
14 pd.set_option('display.width', 1000)
15
16 display(HTML("<style>.container { width:100% !important; }</style>"))
```

```
In [2]: 1 import os
2 import seaborn as sns
3 import pandas as pd
4 import math
5
6 import sklearn.model_selection as cv
7
8 from sklearn.preprocessing import StandardScaler
9 from sklearn.model_selection import train_test_split
10 from sklearn.decomposition import PCA
11 from sklearn.ensemble import RandomForestRegressor, AdaBoostRegressor, GradientBoostingRegressor
12 from sklearn.model_selection import GridSearchCV
13
14 from sklearn.metrics import mean_squared_error as MSE
15
16 from imblearn.over_sampling import SMOTE
17
18 from Utils.UtilsGeoViz import *
19 from Utils.UtilsViz import *
20 from Utils.DataUtils import *
21
22 RANDOM_SEED = 42
```

```
In [3]: 1 # data_path = os.path.join(os.getcwd(), "../data/cleaned_listings.csv")
2 listings = pd.read_csv("cleaned_listings.csv", index_col="id")
3 display(listings.shape)
```

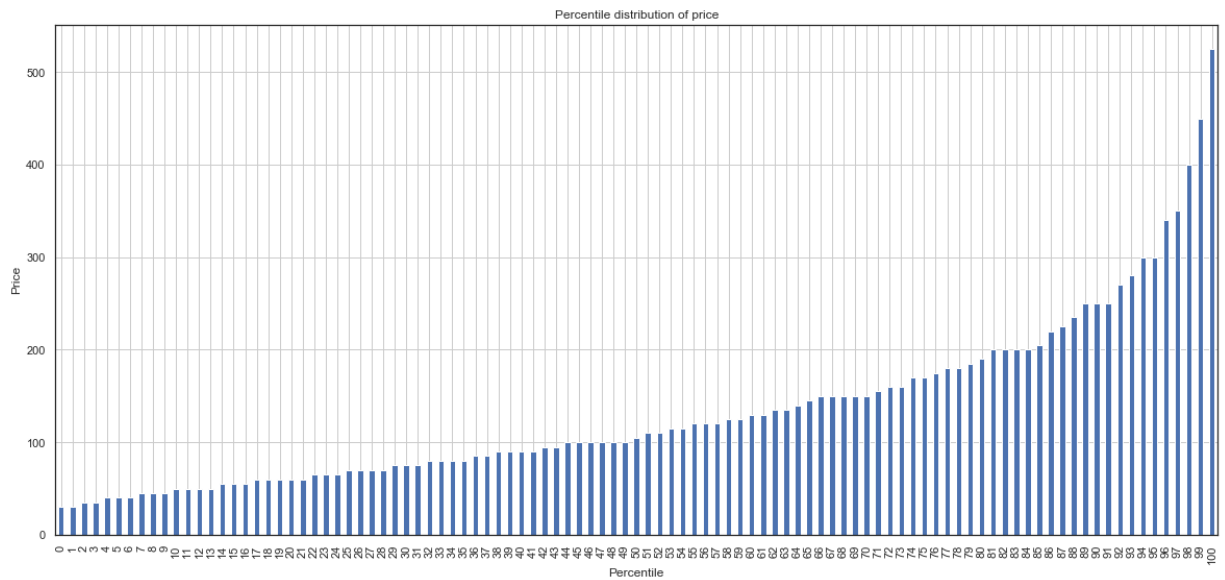
(48855, 65)

Plot the distribution

Let's plot the percentile for price

```
In [4]: 1 percentiles = list(range(0,101, 1))
        2 price_percentile = {}
        3 for p in percentiles:
        4     price_percentile[p] = np.percentile(listings['price'].values, p)
        5
        6 sns.set(style="white")
        7 price_percentile = pd.DataFrame.from_dict(price_percentile, orient='index')
        8 price_percentile.plot(kind='bar', figsize=(20,9), grid=True, legend=False)
        9 plt.title("Percentile distribution of price")
       10 plt.xlabel("Percentile")
       11 plt.ylabel("Price")
```

Out[4]: Text(0, 0.5, 'Price')

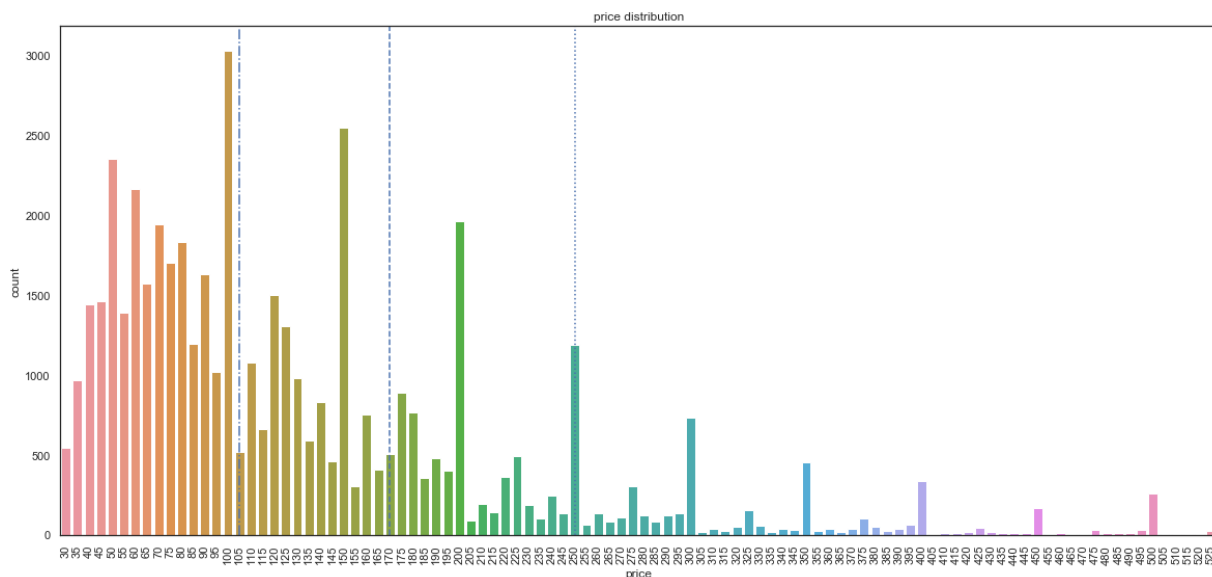


```

In [5]: 1 f, ax = plt.subplots(1,1,figsize=(20,9))
        2 g = sns.countplot(x="price", data=listings, ax=ax)
        3 t = g.set_xticklabels(g.get_xticklabels(), rotation=90)
        4 t = g.set_title("price distribution")
        5 median_idx = np.where(np.sort(listings["price"].unique())==listings["price"]
        6 plt.axvline(x=median_idx, linestyle="-.")
        7 percentile_75_idx = np.where(np.sort(listings["price"].unique())==price_perc
        8 plt.axvline(x=percentile_75_idx, linestyle="--")
        9 percentile_90_idx = np.where(np.sort(listings["price"].unique())==price_perc
        10 plt.axvline(x=percentile_90_idx, linestyle=":")

```

Out[5]: <matplotlib.lines.Line2D at 0x1acb3da5f60>



Quick helper functions

```

In [4]: 1 def roundto(row, base=5):
        2     return int(base * round(float(row) / base))
        3
        4 # Get the index of the price columns
        5 def get_index(vallist, val):
        6     return vallist.index(val)

```

Oversampling using SMOTE

```
In [5]: 1 def check_rep(row):
2         if (row <= 200) | (row==250) | (row==350) | (row==450) | (row==550) :
3             return 0
4         elif (row > 200) & (row < 300) & (row != 250):
5             return 1
6         elif (row > 300) & (row < 400) & (row != 350):
7             return 2
8         else:
9             return 3
10
11 listings["flag_ur"] = listings["price"].apply(check_rep)
```

```
In [6]: 1 vcs = listings["flag_ur"].value_counts()
2         vcs
```

```
Out[6]: 0    43321
1      3093
3      1624
2       817
Name: flag_ur, dtype: int64
```

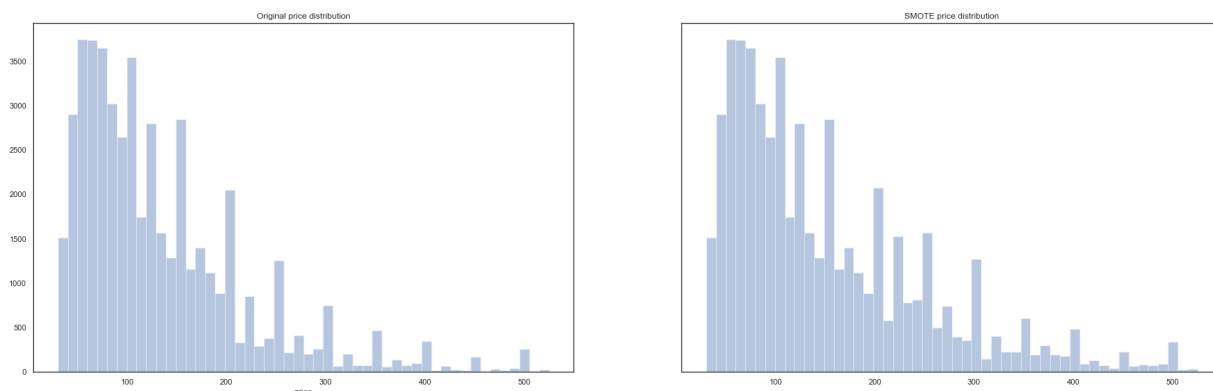
```
In [7]: 1 ycol = ["flag_ur"]
2         xcol = [i for i in listings.columns if i not in ycol]
3
4         x = listings[xcol].values
5         y = listings[ycol].values
6
7         smote_sampling_strategy = {
8             1: int(vcs[1]*2)
9             ,2: int(vcs[2]*2)
10            ,3: int(vcs[3]*2)
11        }
12         sm = SMOTE(random_state=RANDOM_SEED, sampling_strategy=smote_sampling_strate
13         # Fit the smote onto the sample
14         x_new, y_new = sm.fit_sample(x, y)
15
16         # Drop the flag column
17         listings.drop(labels=["flag_ur"], axis=1, inplace=True)
18
19         # -----
20         # Overwrite X and Y
21         price_index = get_index(list(listings.columns), "price")
22
23         y = x_new[:, price_index]
24         x = np.delete(x_new, price_index, axis=1)
25         for i in range(len(y)):
26             y[i] = roundto(y[i])
```

/anaconda3/lib/python3.6/site-packages/sklearn/utils/validation.py:761: DataCon
versionWarning: A column-vector y was passed when a 1d array was expected. Plea
se change the shape of y to (n_samples,), for example using ravel().
y = column_or_1d(y, warn=True)

```
In [8]: 1 print(
2         "Old size :", listings.shape, "\n",
3         "New size :", x.shape
4     )
```

Old size : (48855, 65)
New size : (54389, 64)

```
In [11]: 1 f, ax = plt.subplots(1, 2, figsize=(30, 9), sharey=True)
2 g1 = sns.distplot(listings["price"], ax=ax[0], kde=False)
3 t = g1.set_title("Original price distribution")
4
5 g2 = sns.distplot(y, ax=ax[1], kde=False)
6 t = g2.set_title("SMOTE price distribution")
7
```



Transformation

```
In [12]: 1 x_cols = listings.drop(['price'], axis=1)
2 X = pd.DataFrame(data=x, columns=x_cols.columns)
```

Train test split

```
In [9]: 1 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.30, ra
```

Standardisation

```
In [10]: 1 standard_scaler = StandardScaler()
2 x_train = standard_scaler.fit_transform(x_train)
3 x_test = standard_scaler.transform(x_test)
```

Prediction

a. Ridge Regression

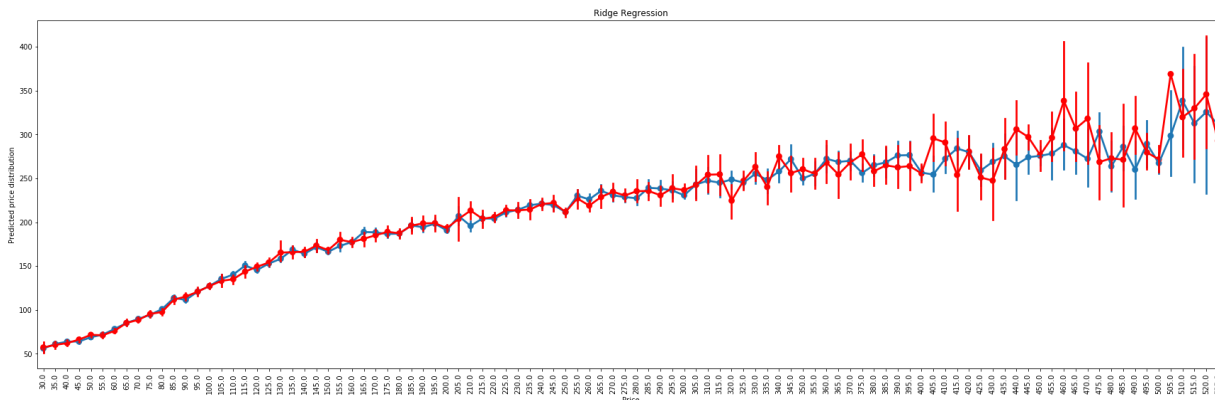
```
In [35]: 1 # Grid Search
2 from sklearn.linear_model import Ridge, Lasso
3 from sklearn.model_selection import GridSearchCV
4
5 rid = Ridge()
6
7 parameters = {'alpha': [1e-15, 1e-12, 1e-10, 1e-08, 0.1, 1, 10, 100]}
8
9 rid_reg = GridSearchCV(rid, parameters, scoring='neg_mean_squared_error', cv=5)
10
11 rid_reg.fit(x_train, y_train)
12
13 print(rid_reg.best_params_)
14
15 rmse = (-rid_reg.best_score_)**(1/2)
16 print(rmse)
```

```
{'alpha': 1e-08}
64.5309608074011
```

```
In [74]: 1 rid = Ridge(alpha=1e-08)
2 rid.fit(X=x_train, y=y_train)
3 y_pred_train = rid.predict(X=x_train)
4 y_pred_test = rid.predict(X=x_test)
5
6 mse_train = MSE(y_train, y_pred_train)
7 mse_test = MSE(y_test, y_pred_test)
8
9 rmse_train = mse_train**(1/2)
10 rmse_test = mse_test**(1/2)
11
12 print("Train set RMSE Scaled: {:.2f}".format(rmse_train))
13 print("Test set RMSE Scaled: {:.2f}".format(rmse_test))
```

```
Train set RMSE Scaled: 64.40
Test set RMSE Scaled: 65.17
```

```
In [75]: 1 f, ax = plt.subplots(1,1, figsize=(30, 9), sharex=False)
2 g = sns.pointplot(x=y_train, y=y_pred_train, ax=ax)
3 g = sns.pointplot(x=y_test, y=y_pred_test, ax=ax, color="red")
4 t = g.set_xlabel("Price")
5 t = g.set_ylabel("Predicted price distribution")
6 t = g.set_xticklabels(g.get_xticklabels(), rotation=90)
7 t = g.set_title("Ridge Regression")
```



```
In [73]: 1 th = 200
2
3 tmpdf = pd.DataFrame({"y_train":y_train, "y_pred_train":y_pred_train})
4 tmpdf = tmpdf[tmpdf["y_train"] <= th]
5
6 tmpdf2 = pd.DataFrame({"y_test":y_test, "y_pred_test":y_pred_test})
7 tmpdf2 = tmpdf2[tmpdf2["y_test"] <= th]
8
9 mse_train = MSE(tmpdf["y_train"].values, tmpdf["y_pred_train"].values)
10 mse_test = MSE(tmpdf2["y_test"].values, tmpdf2["y_pred_test"].values)
11
12 rmse_train = mse_train**(1/2)
13 rmse_test = mse_test**(1/2)
14
15 print("Train set RMSE Scaled: {:.2f}".format(rmse_train))
16 print("Test set RMSE Scaled: {:.2f}".format(rmse_test))
```

Train set RMSE Scaled: 48.04

Test set RMSE Scaled: 49.70

b. Lasso Regression

In [69]:

```
1 # Grid Search
2 from sklearn.linear_model import Ridge, Lasso
3 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
4
5 las = Lasso()
6
7 parameters = {'alpha': [0.15, 0.5, 1]}
8
9 las_reg = GridSearchCV(las, parameters, scoring='neg_mean_squared_error', cv=5)
10
11 las_reg.fit(x_train, y_train)
12
13 print(las_reg.best_params_)
14
15 rmse = (-las_reg.best_score_)**(1/2)
16 print(rmse)
```

```
{'alpha': 0.15}
64.58075921581023
```

In [76]:

```
1 las = Lasso(alpha=0.15)
2 las.fit(X=x_train, y=y_train)
3 y_pred_train = las.predict(X=x_train)
4 y_pred_test = las.predict(X=x_test)
5
6 mse_train = MSE(y_train, y_pred_train)
7 mse_test = MSE(y_test, y_pred_test)
8
9 rmse_train = mse_train**(1/2)
10 rmse_test = mse_test**(1/2)
11
12 print("Train set RMSE Scaled: {:.2f}".format(rmse_train))
13 print("Test set RMSE Scaled: {:.2f}".format(rmse_test))
```

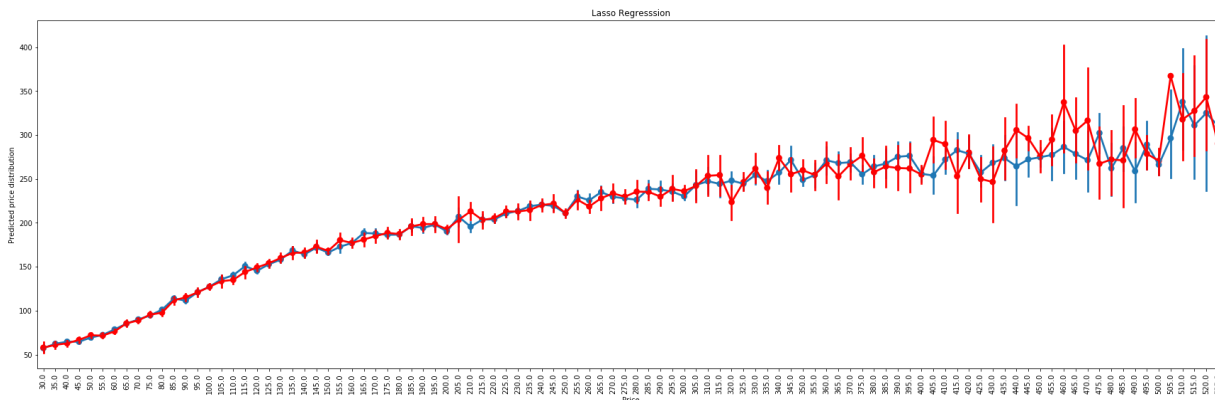
```
Train set RMSE Scaled: 64.46
Test set RMSE Scaled: 64.26
```


In [77]:

```

1 f, ax = plt.subplots(1,1, figsize=(30, 9), sharex=False)
2 g = sns.pointplot(x=y_train, y=y_pred_train, ax=ax)
3 g = sns.pointplot(x=y_test, y=y_pred_test, ax=ax, color="red")
4 t = g.set_xlabel("Price")
5 t = g.set_ylabel("Predicted price distribution")
6 t = g.set_xticklabels(g.get_xticklabels(), rotation=90)
7 t = g.set_title("Lasso Regressionssion")

```



In [71]:

```

1 th = 200
2
3 tmpdf = pd.DataFrame({"y_train":y_train, "y_pred_train":y_pred_train})
4 tmpdf = tmpdf[tmpdf["y_train"] <= th]
5
6 tmpdf2 = pd.DataFrame({"y_test":y_test, "y_pred_test":y_pred_test})
7 tmpdf2 = tmpdf2[tmpdf2["y_test"] <= th]
8
9 mse_train = MSE(tmpdf["y_train"].values, tmpdf["y_pred_train"].values)
10 mse_test = MSE(tmpdf2["y_test"].values, tmpdf2["y_pred_test"].values)
11
12 rmse_train = mse_train**(1/2)
13 rmse_test = mse_test**(1/2)
14
15 print("Train set RMSE Scaled: {:.2f}".format(rmse_train))
16 print("Test set RMSE Scaled: {:.2f}".format(rmse_test))

```

Train set RMSE Scaled: 47.88

Test set RMSE Scaled: 47.86