

Airbnb listings file modeling - Initial probing

This file establishes the baseline metrics for each of the models chosen by running them through the data straight after cleaning.

The effects of additional processing of the data are explored later in subsequent modeling and EDA files.

```
In [151]: 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 [152]: 1 import os
2 import seaborn as sns
3 import pandas as pd
4 import math
5 import sklearn.model_selection as cv
6
7 from sklearn.preprocessing import StandardScaler
8 from sklearn.model_selection import train_test_split
9 from sklearn.decomposition import PCA
10 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
11 from sklearn.tree import DecisionTreeRegressor
12 from sklearn.model_selection import GridSearchCV
13
14 from sklearn.linear_model import Lasso, Ridge
15
16 from sklearn.metrics import mean_squared_error as MSE
17
18 from imblearn.over_sampling import SMOTE
19
20 from Utils.UtilsGeoViz import *
21 from Utils.UtilsViz import *
22 from Utils.DataUtils import *
23
24 RANDOM_SEED = 42
```

```
In [153]: 1 sns.set(font_scale=2, style="whitegrid")
```

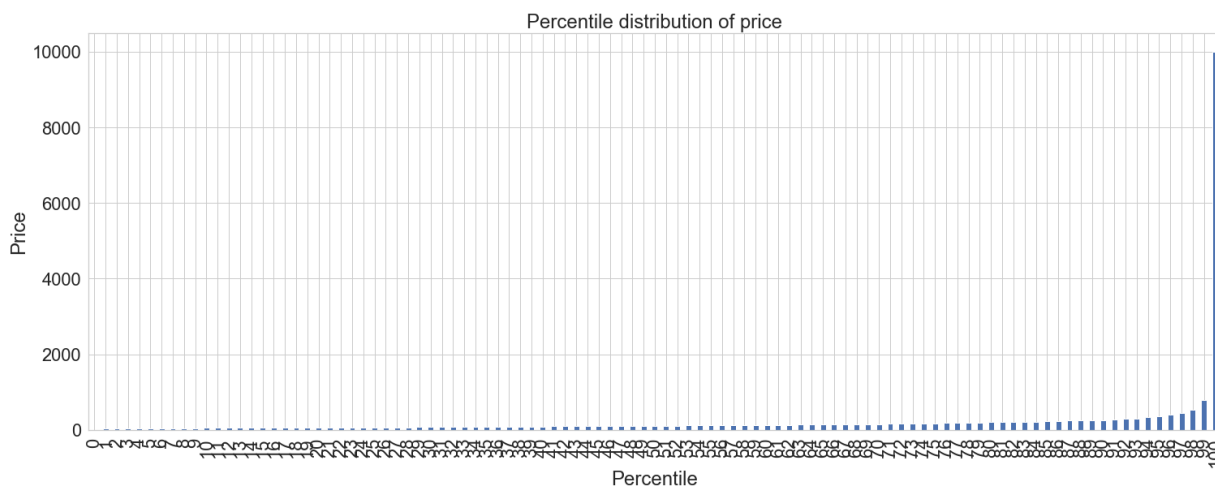
```
In [154]: 1 data_path = os.path.join(os.getcwd(), "../data/cleaned_listings_with_outlier")
2 listings = pd.read_csv(data_path, index_col="id")
3 display(listings.shape)

(50220, 65)
```

Plot the distribution

```
In [155]: 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 price_percentile = pd.DataFrame.from_dict(price_percentile, orient='index')
7 price_percentile.plot(kind='bar', figsize=(25,9), grid=True, legend=False)
8 plt.title("Percentile distribution of price")
9 plt.xlabel("Percentile")
10 plt.ylabel("Price")
```

```
Out[155]: Text(0, 0.5, 'Price')
```



```
In [156]: 1 f, ax = plt.subplots(1,1,figsize=(25,9))
2 g = sns.distplot(a=listings["price"], kde=False)
3 t = g.set_title("price distribution")
```



Round the prices to multiples of 5

```
In [157]: 1 def roundto(row, base=5):
2     return int(base * round(float(row) / base))
3 listings["price"] = listings["price"].apply(roundto)
```

Train Test Split

```
In [158]: 1 ycol = ["price"]
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 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.30, ra
```

Prediction

```
In [159]: 1 # Quick placeholder variable
2 rmse_dict = {}
3 rmse_m_dict = {}
```

a. Ridge Regression

In [160]:

```
1 ridge = Ridge()
2 ridge.fit(X=x_train, y=y_train)
3 y_pred_train = ridge.predict(X=x_train)
4 y_pred_test = ridge.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 rmse_dict["Ridge"] = {"Train":rmse_train,"Test":rmse_test}
12
13 print("Train set RMSE: {:.2f}".format(rmse_train))
14 print("Test set RMSE: {:.2f}".format(rmse_test))
```

Train set RMSE: 189.67

Test set RMSE: 223.51

In [161]:

```
1 th = price_percentile.iloc[80, :].values[0]
2
3 tmpdf = pd.DataFrame({"y_train":y_train.ravel(), "y_pred_train":y_pred_train
4 tmpdf = tmpdf[tmpdf["y_train"] <= th]
5
6 tmpdf2 = pd.DataFrame({"y_test":y_test.ravel(), "y_pred_test":y_pred_test.ra
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 rmse_m_dict["Ridge"] = {"Train":rmse_train,"Test":rmse_test}
16
17 print("Train set RMSE Scaled: {:.2f}".format(rmse_train))
18 print("Test set RMSE Scaled: {:.2f}".format(rmse_test))
```

Train set RMSE Scaled: 62.22

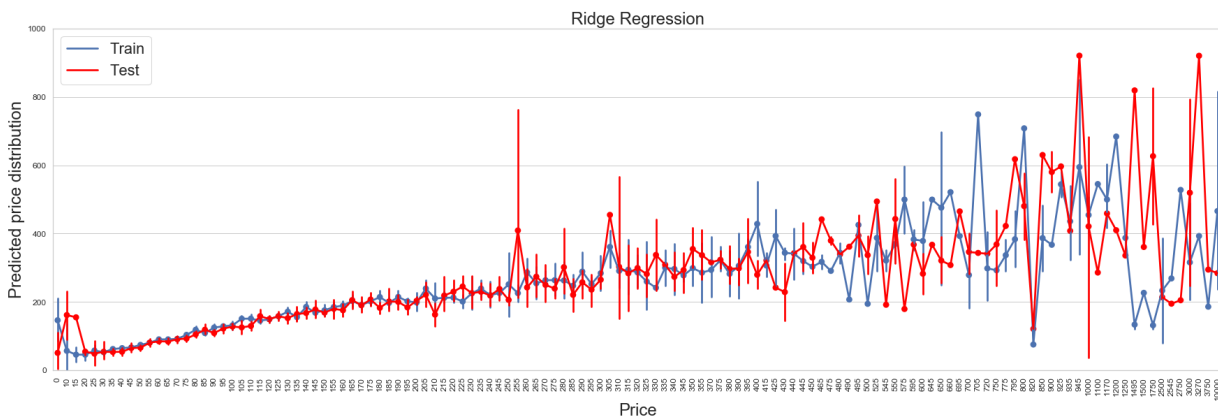
Test set RMSE Scaled: 61.92

In [162]:

```

1 f, ax = plt.subplots(1,1, figsize=(30, 9), sharex=False)
2 # g = sns.pointplot(x=np.unique(y_train), y=np.unique(y_train), color="light
3 g = sns.pointplot(x=y_train.ravel()[::5], y=y_pred_train.ravel()[::5], ax=ax
4 g = sns.pointplot(x=y_test.ravel()[::5], y=y_pred_test.ravel()[::5], ax=ax,
5 t = g.set_xlabel("Price")
6 t = g.set_ylabel("Predicted price distribution")
7 t = g.set_xticklabels(g.get_xticklabels(), rotation=90)
8 t = g.set_title("Ridge Regression")
9 yl = g.set_ylim(0, 1000)
10 l = ax.legend(handles=ax.lines[:,len(np.unique(y_train))+1], labels=["Train"
11 t = g.tick_params(labelsize=12)
12

```



In [164]:

```

1 # Use cross-validation on Train data
2
3 CV_scores = - cv.cross_val_score(ridge, x_train, y_train, scoring='neg_mean_
4
5 # Compute the 10-folds CV
6 CV = CV_scores.mean()*(1/2)
7
8 # Print Train CV accuracy
9 print('Cross Validation RMSE: {:.2f}'.format(CV))

```

Cross Validation RMSE: 190.35

b. Lasso Regression

In [165]:

```
1 lasso = Lasso()
2 lasso.fit(X=x_train, y=y_train)
3 y_pred_train = lasso.predict(X=x_train)
4 y_pred_test = lasso.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 rmse_dict["lasso"] = {"Train":rmse_train,"Test":rmse_test}
12
13 print("Train set RMSE: {:.2f}".format(rmse_train))
14 print("Test set RMSE: {:.2f}".format(rmse_test))
```

Train set RMSE: 190.65

Test set RMSE: 224.78

In [166]:

```
1 th = price_percentile.iloc[80, :].values[0]
2
3 tmpdf = pd.DataFrame({"y_train":y_train.ravel(), "y_pred_train":y_pred_train
4 tmpdf = tmpdf[tmpdf["y_train"] <= th]
5
6 tmpdf2 = pd.DataFrame({"y_test":y_test.ravel(), "y_pred_test":y_pred_test.ra
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 rmse_m_dict["lasso"] = {"Train":rmse_train,"Test":rmse_test}
16
17 print("Train set RMSE Scaled: {:.2f}".format(rmse_train))
18 print("Test set RMSE Scaled: {:.2f}".format(rmse_test))
```

Train set RMSE Scaled: 61.02

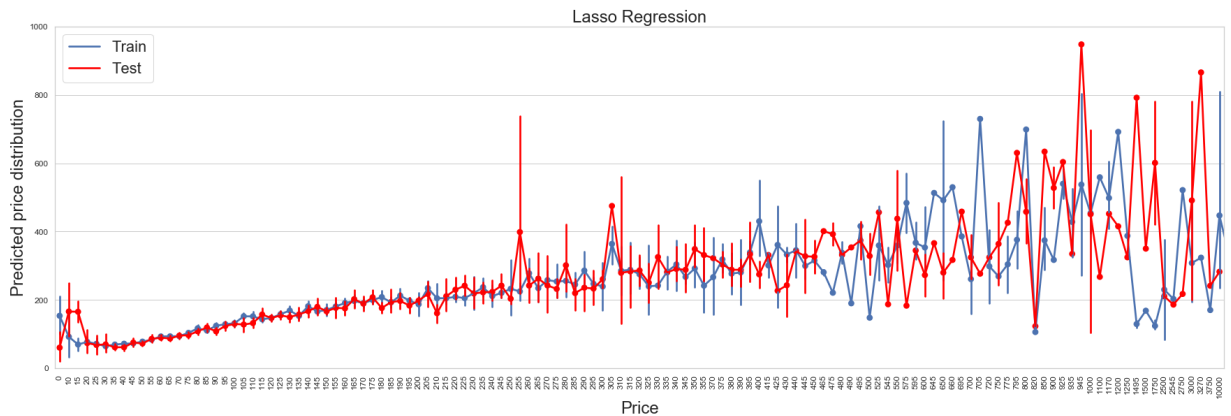
Test set RMSE Scaled: 61.24

In [167]:

```

1 f, ax = plt.subplots(1,1, figsize=(30, 9), sharex=False)
2 # g = sns.pointplot(x=np.unique(y_train), y=np.unique(y_train), color="light
3 g = sns.pointplot(x=y_train.ravel()[::5], y=y_pred_train.ravel()[::5], ax=ax
4 g = sns.pointplot(x=y_test.ravel()[::5], y=y_pred_test.ravel()[::5], ax=ax,
5 t = g.set_xlabel("Price")
6 t = g.set_ylabel("Predicted price distribution")
7 t = g.set_xticklabels(g.get_xticklabels(), rotation=90)
8 t = g.set_title("Lasso Regression")
9 yl = g.set_ylim(0, 1000)
10 l = ax.legend(handles=ax.lines[:,len(np.unique(y_train))+1], labels=["Train"
11 t = g.tick_params(labelsize=12)
12

```



In [169]:

```

1  # Use cross-validation on Train data
2
3  CV_scores = - cv.cross_val_score(lasso, x_train, y_train, scoring='neg_mean_
4
5  # Compute the 10-folds CV
6  CV = CV_scores.mean()**(1/2)
7
8  # Print Train CV accuracy
9  print('Cross Validation RMSE: {:.2f}'.format(CV))

```

```

-----
KeyboardInterrupt                                Traceback (most recent call last)
<ipython-input-169-dabc9a0c5edd> in <module>
      1 # Use cross-validation on Train data
      2
----> 3 CV_scores = - cv.cross_val_score(lasso, x_train, y_train, scoring='neg_
mean_squared_error', cv=10, n_jobs=6)
      4
      5 # Compute the 10-folds CV

c:\users\sriharis\appdata\local\programs\python\python37-32\lib\site-packages\s
klearn\model_selection\_validation.py in cross_val_score(estimator, X, y, group
s, scoring, cv, n_jobs, verbose, fit_params, pre_dispatch, error_score)
    400             fit_params=fit_params,
    401             pre_dispatch=pre_dispatch,
--> 402             error_score=error_score)
    403     return cv_results['test_score']
    404

c:\users\sriharis\appdata\local\programs\python\python37-32\lib\site-packages\s
klearn\model_selection\_validation.py in cross_validate(estimator, X, y, group
s, scoring, cv, n_jobs, verbose, fit_params, pre_dispatch, return_train_score,
return_estimator, error_score)
    238         return_times=True, return_estimator=return_estimator,
    239         error_score=error_score)
--> 240     for train, test in cv.split(X, y, groups))
    241
    242     zipped_scores = list(zip(*scores))

c:\users\sriharis\appdata\local\programs\python\python37-32\lib\site-packages\s
klearn\externals\joblib\parallel.py in __call__(self, iterable)
    928
    929         with self._backend.retrieval_context():
--> 930             self.retrieve()
    931             # Make sure that we get a last message telling us we are do
ne
    932             elapsed_time = time.time() - self._start_time

c:\users\sriharis\appdata\local\programs\python\python37-32\lib\site-packages\s
klearn\externals\joblib\parallel.py in retrieve(self)
    831         try:
    832             if getattr(self._backend, 'supports_timeout', False):
--> 833                 self._output.extend(job.get(timeout=self.timeout))
    834             else:
    835                 self._output.extend(job.get())

c:\users\sriharis\appdata\local\programs\python\python37-32\lib\site-packages\s

```



```

klearn\externals\joblib\_parallel_backends.py in wrap_future_result(future, tim
eout)
    519         AsyncResults.get from multiprocessing."""
    520         try:
--> 521             return future.result(timeout=timeout)
    522         except LokyTimeoutError:
    523             raise TimeoutError()

c:\users\sriharis\appdata\local\programs\python\python37-32\lib\concurrent\futu
res\_base.py in result(self, timeout)
    425         return self.__get_result()
    426
--> 427         self._condition.wait(timeout)
    428
    429         if self._state in [CANCELLED, CANCELLED_AND_NOTIFIED]:

c:\users\sriharis\appdata\local\programs\python\python37-32\lib\threading.py in
wait(self, timeout)
    294         try:      # restore state no matter what (e.g., KeyboardInterrupt)
    295
--> 296             if timeout is None:
    297                 waiter.acquire()
    298                 gotit = True
    299             else:

```

KeyboardInterrupt:

c. Decision Trees

```

In [ ]: 1 dtr = DecisionTreeRegressor()
        2 dtr.fit(X=x_train, y=y_train)
        3 y_pred_train = dtr.predict(X=x_train)
        4 y_pred_test = dtr.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 rmse_dict["dt"] = {"Train":rmse_train,"Test":rmse_test}
       12
       13 print("Train set RMSE: {:.2f}".format(rmse_train))
       14 print("Test set RMSE: {:.2f}".format(rmse_test))

```

```

In [ ]: 1 th = price_percentile.iloc[80, :].values[0]
        2
        3 tmpdf = pd.DataFrame({"y_train":y_train.ravel(), "y_pred_train":y_pred_train
        4 tmpdf = tmpdf[tmpdf["y_train"] <= th]
        5
        6 tmpdf2 = pd.DataFrame({"y_test":y_test.ravel(), "y_pred_test":y_pred_test.ra
        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 rmse_m_dict["dt"] = {"Train":rmse_train,"Test":rmse_test}
       16
       17 print("Train set RMSE Scaled: {:.2f}".format(rmse_train))
       18 print("Test set RMSE Scaled: {:.2f}".format(rmse_test))

```

```

In [ ]: 1 f, ax = plt.subplots(1,1, figsize=(30, 9), sharex=False)
        2 # g = sns.pointplot(x=np.unique(y_train), y=np.unique(y_train), color="light
        3 g = sns.pointplot(x=y_train.ravel()[::5], y=y_pred_train.ravel()[::5], ax=ax
        4 g = sns.pointplot(x=y_test.ravel()[::5], y=y_pred_test.ravel()[::5], ax=ax,
        5 t = g.set_xlabel("Price")
        6 t = g.set_ylabel("Predicted price distribution")
        7 t = g.set_xticklabels(g.get_xticklabels(), rotation=90)
        8 t = g.set_title("Decision Trees")
        9 yl = g.set_ylim(0, 3000)
       10 l = ax.legend(handles=ax.lines[::len(np.unique(y_train))+1], labels=["Train"
       11 t = g.tick_params(labelsize=12)
       12

```

Feature importance

```

In [ ]: 1 x_final = listings.drop(['price'], axis=1)
        2
        3 imp = pd.DataFrame({'Input Variable': x_final.columns, 'Importance (%)': -np
        4 top20 = imp.iloc[:20,:]}
        5
        6 f, ax = plt.subplots(1,1,figsize=(20,9))
        7 g = sns.barplot(x="Importance (%)", y="Input Variable", data=top20, ax=ax)
        8 t = g.set_xlabel("Importance (%)")
        9 t = g.set_ylabel("Input Variable")
       10 t = g.set_title("Feature Importance")

```

```

In [ ]: 1 # Use cross-validation on Train data
        2
        3 CV_scores = - cv.cross_val_score(dtr, x_train, y_train, scoring='neg_mean_sq
        4
        5 # Compute the 10-folds CV
        6 CV = CV_scores.mean()**(1/2)
        7
        8 # Print Train CV accuracy
        9 print('Cross Validation RMSE: {:.2f}'.format(CV))

```

d. Random forests

```

In [ ]: 1 rfr = RandomForestRegressor()
        2 rfr.fit(X=x_train, y=y_train)
        3 y_pred_train = rfr.predict(X=x_train)
        4 y_pred_test = rfr.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 rmse_dict["rf"] = {"Train":rmse_train,"Test":rmse_test}
       12
       13 print("Train set RMSE: {:.2f}".format(rmse_train))
       14 print("Test set RMSE: {:.2f}".format(rmse_test))

```

```

In [ ]: 1 th = price_percentile.iloc[80, :].values[0]
        2
        3 tmpdf = pd.DataFrame({"y_train":y_train.ravel(), "y_pred_train":y_pred_train
        4 tmpdf = tmpdf[tmpdf["y_train"] <= th]
        5
        6 tmpdf2 = pd.DataFrame({"y_test":y_test.ravel(), "y_pred_test":y_pred_test.ra
        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 rmse_m_dict["rf"] = {"Train":rmse_train,"Test":rmse_test}
       16
       17 print("Train set RMSE Scaled: {:.2f}".format(rmse_train))
       18 print("Test set RMSE Scaled: {:.2f}".format(rmse_test))

```

```

In [ ]: 1 f, ax = plt.subplots(1,1, figsize=(30, 9), sharex=False)
        2 # g = sns.pointplot(x=np.unique(y_train), y=np.unique(y_train), color="light
        3 g = sns.pointplot(x=y_train.ravel()[::5], y=y_pred_train.ravel()[::5], ax=ax
        4 g = sns.pointplot(x=y_test.ravel()[::5], y=y_pred_test.ravel()[::5], ax=ax,
        5 t = g.set_xlabel("Price")
        6 t = g.set_ylabel("Predicted price distribution")
        7 t = g.set_xticklabels(g.get_xticklabels(), rotation=90)
        8 t = g.set_title("Random Forest")
        9 yl = g.set_ylim(0, 2000)
       10 l = ax.legend(handles=ax.lines[:,len(np.unique(y_train))+1], labels=["Train"
       11 t = g.tick_params(labelsize=12)
       12

```

Feature importance

```

In [ ]: 1 x_final = listings.drop(['price'], axis=1)
        2
        3 imp = pd.DataFrame({'Input Variable': x_final.columns, 'Importance (%)': -np
        4 top20 = imp.iloc[:20,:])
        5
        6 f, ax = plt.subplots(1,1,figsize=(20,9))
        7 g = sns.barplot(x="Importance (%)", y="Input Variable", data=top20, ax=ax)
        8 t = g.set_xlabel("Importance (%)")
        9 t = g.set_ylabel("Input Variable")
       10 t = g.set_title("Feature Importance")

```

```

In [ ]: 1 # Use cross-validation on Train data
        2
        3 CV_scores = - cv.cross_val_score(rfr, x_train, y_train, scoring='neg_mean_sq
        4
        5 # Compute the 10-folds CV
        6 CV = CV_scores.mean()**(1/2)
        7
        8 # Print Train CV accuracy
        9 print('Cross Validation RMSE: {:.2f}'.format(CV))

```

e. Gradient Boosting regression

```

In [ ]: 1 gbr = GradientBoostingRegressor()
        2 gbr.fit(X=x_train, y=y_train)
        3 y_pred_train = gbr.predict(X=x_train)
        4 y_pred_test = gbr.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 rmse_dict["gbr"] = {"Train":rmse_train,"Test":rmse_test}
       12
       13 print("Train set RMSE: {:.2f}".format(rmse_train))
       14 print("Test set RMSE: {:.2f}".format(rmse_test))

```

```

In [ ]: 1 th = price_percentile.iloc[80, :].values[0]
        2
        3 tmpdf = pd.DataFrame({"y_train":y_train.ravel(), "y_pred_train":y_pred_train
        4 tmpdf = tmpdf[tmpdf["y_train"] <= th]
        5
        6 tmpdf2 = pd.DataFrame({"y_test":y_test.ravel(), "y_pred_test":y_pred_test.ra
        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 rmse_m_dict["gbr"] = {"Train":rmse_train,"Test":rmse_test}
       16
       17 print("Train set RMSE Scaled: {:.2f}".format(rmse_train))
       18 print("Test set RMSE Scaled: {:.2f}".format(rmse_test))

```

```

In [ ]: 1 f, ax = plt.subplots(1,1, figsize=(30, 9), sharex=False)
        2 # g = sns.pointplot(x=np.unique(y_train), y=np.unique(y_train), color="light
        3 g = sns.pointplot(x=y_train.ravel()[::5], y=y_pred_train.ravel()[::5], ax=ax
        4 g = sns.pointplot(x=y_test.ravel()[::5], y=y_pred_test.ravel()[::5], ax=ax,
        5 t = g.set_xlabel("Price")
        6 t = g.set_ylabel("Predicted price distribution")
        7 t = g.set_xticklabels(g.get_xticklabels(), rotation=90)
        8 t = g.set_title("Gradient Boosting Regression")
        9 yl = g.set_ylim(0, 2500)
       10 l = ax.legend(handles=ax.lines[:,len(np.unique(y_train))+1], labels=["Train"
       11 t = g.tick_params(labelsize=12)
       12

```

Feature importance

```

In [ ]: 1 x_final = listings.drop(['price'], axis=1)
        2
        3 imp = pd.DataFrame({'Input Variable': x_final.columns, 'Importance (%)': -np
        4 top20 = imp.iloc[:20,:]}
        5
        6 f, ax = plt.subplots(1,1,figsize=(20,9))
        7 g = sns.barplot(x="Importance (%)", y="Input Variable", data=top20, ax=ax)
        8 t = g.set_xlabel("Importance (%)")
        9 t = g.set_ylabel("Input Variable")
       10 t = g.set_title("Feature Importance")

```

```
In [ ]: 1 # Use cross-validation on Train data
2
3 CV_scores = - cv.cross_val_score(gbr, x_train, y_train, scoring='neg_mean_sq
4
5 # Compute the 10-folds CV
6 CV = CV_scores.mean()**(1/2)
7
8 # Print Train CV accuracy
9 print('Cross Validation RMSE: {:.2f}'.format(CV))
```

Final dataframe to plot the RMSEs

```
In [ ]: 1 rmse_df = pd.DataFrame(rmse_dict)
2 rmse_df = rmse_df.T.reset_index(drop=False)
3 rmse_df.columns = ["Regressor", "Test RMSE", "Train RMSE"]
4 rmse_df
```

```
In [ ]: 1 rmse_m_df = pd.DataFrame(rmse_m_dict)
2 rmse_m_df = rmse_m_df.T.reset_index(drop=False)
3 rmse_m_df.columns = ["Regressor", "Test RMSE", "Train RMSE"]
4 rmse_m_df
```

```

In [ ]: 1 def plot_bar(data, x, y, ax, hue=None, title="", xlabel="", ylabel="",
2         xrot=0, yrot=0, highlight_max_min=True,
3         plot_percentiles=[], plot_mean=True,
4         point_plot=False, annot=True, legend=False):
5     if highlight_max_min:
6         clr = []
7         for v in data[y].values:
8             if v < data[y].max():
9                 if v > data[y].min():
10                    clr.append('lightblue')
11                else:
12                    clr.append('lightgreen')
13            else:
14                clr.append('darksalmon')
15        g = sns.barplot(x=x, y=y, data=data, ax=ax, palette=clr)
16    else:
17        g = sns.barplot(x=x, y=y, data=data, ax=ax, hue=hue)
18
19    if len(plot_percentiles) > 0:
20        for p in plot_percentiles:
21            v = np.percentile(data[y].values, p)
22            plt.axhline(v, 1, 0, color='grey').set_linestyle("--")
23
24    if plot_mean:
25        v = data[y].mean()
26        plt.axhline(v, 1, 0, color='k').set_linestyle("--")
27
28    if point_plot:
29        g1 = sns.pointplot(x=x, y=y, data=data, ax=ax, color="darkblue")
30    if xrot != 0:
31        g.set_xticklabels(rotation=xrot, labels=g.get_xticklabels())
32    if yrot != 0:
33        g.set_yticklabels(rotation=yrot, labels=g.get_yticklabels())
34    if annot:
35        # Add labels to the plot
36        style = dict(size=12, color='darkblue')
37        s1 = data[y].values
38        counter = 0
39        for idx, row in data.iterrows():
40            rx, ry = row[x], row[y]
41            if type('str') == type(idx):
42                ax.text(counter, ry, str(np.round(ry, 2)),
43                        **style, va="bottom", ha='right')
44            else:
45                ax.text(idx*0.99, ry, str(np.round(s1[idx], 2)),
46                        **style, va="bottom", ha='right')
47            counter += 1
48    g.set(xlabel=xlabel, ylabel=ylabel, title=title)
49    ax.set_ylim([0, data[y].max() * 1.2])
50    if legend:
51        ax.legend(handles=ax.lines[:len(data) + 1], labels=[y])
52

```

```
In [ ]: 1 f, ax = plt.subplots(1,2,figsize=(20,8), sharey=False)
2 # g = sns.barplot(x="Regressor", y="Train RMSE", data=rmse_df, ax=ax[0])
3 # t = g.set(title="Train RMSE", ylabel="RMSE", xlabel="")
4 # g = sns.barplot(x="Regressor", y="Test RMSE", data=rmse_df, ax=ax[1])
5 # t = g.set(title="Test RMSE", ylabel="RMSE", xlabel="")
6
7 g = plot_bar(x="Regressor", y="Train RMSE", data=rmse_df, ax=ax[0], plot_mean=True)
8 g = plot_bar(x="Regressor", y="Test RMSE", data=rmse_df, ax=ax[1], plot_mean=True)
9
```

```
In [ ]: 1 f, ax = plt.subplots(1,2,figsize=(20,8), sharey=False)
2 # g = sns.barplot(x="Regressor", y="Train RMSE", data=rmse_m_df, ax=ax[0])
3 # t = g.set(title="Train RMSE", ylabel="RMSE", xlabel="")
4 # g = sns.barplot(x="Regressor", y="Test RMSE", data=rmse_m_df, ax=ax[1])
5 # t = g.set(title="Test RMSE", ylabel="RMSE", xlabel="")
6
7 g = plot_bar(x="Regressor", y="Train RMSE", data=rmse_m_df, ax=ax[0], plot_mean=True)
8 g = plot_bar(x="Regressor", y="Test RMSE", data=rmse_m_df, ax=ax[1], plot_mean=True)
```

```
In [ ]: 1
```