

AirBnb listings file modeling

```
In [6]: 1 # Put these at the top of every notebook, to get automatic reloading and inl
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 [7]: 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 [8]: 1 data_path = os.path.join(os.getcwd(), "../data/cleaned_listings.csv")
2 listings = pd.read_csv(data_path, index_col="id")
3 display(listings.shape)
```

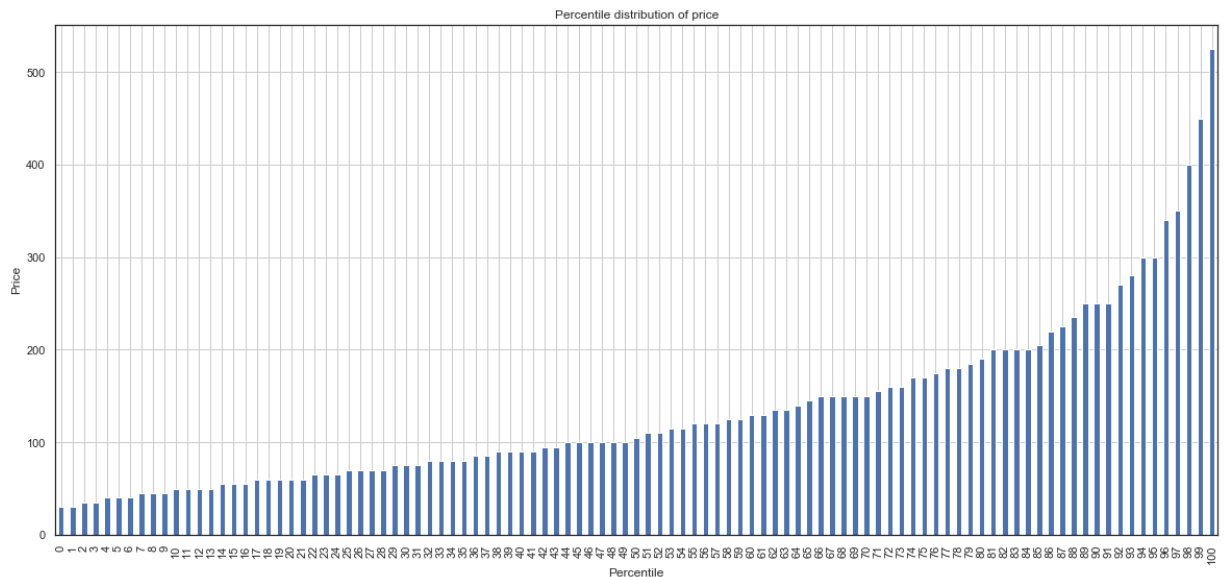
(48855, 65)

Plot the distribution

Let's plot the percentile for price

```
In [9]: 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[9]: Text(0, 0.5, 'Price')

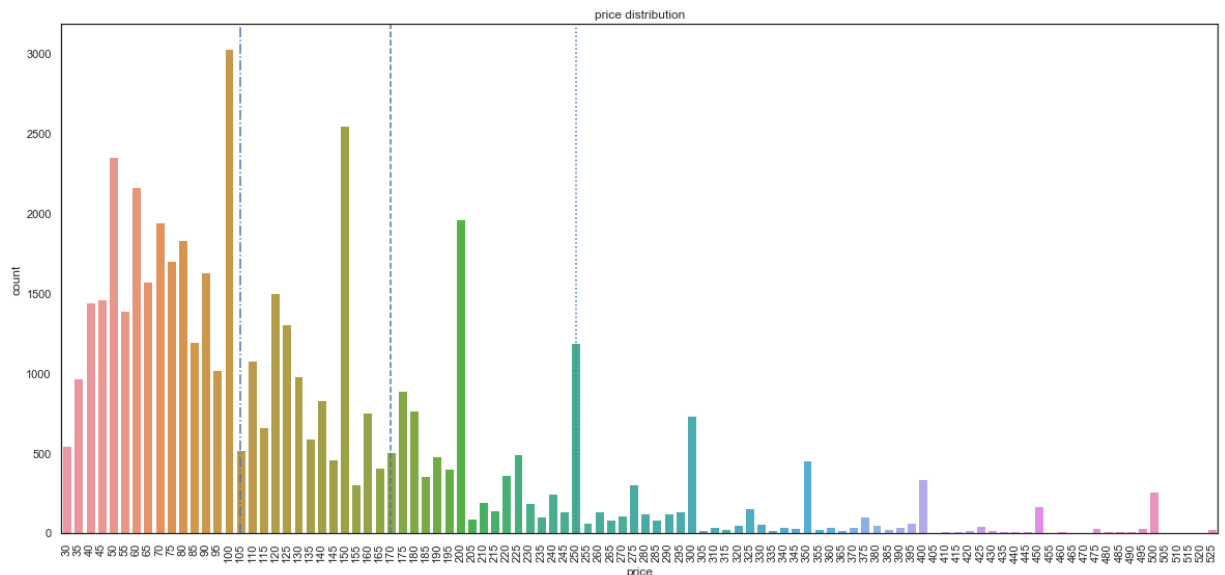


```

In [10]: 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                                .median())[0][0]
7         plt.axvline(x=median_idx, linestyle="-.")
8         percentile_75_idx = np.where(np.sort(listings["price"].unique())==price_perc
9                                       [0.75])[0][0]
10        plt.axvline(x=percentile_75_idx, linestyle="--")
11        percentile_90_idx = np.where(np.sort(listings["price"].unique())==price_perc
12                                       [0.9])[0][0]
13        plt.axvline(x=percentile_90_idx, linestyle=":")

```

Out[10]: <matplotlib.lines.Line2D at 0x1a25e59e80>



Quick helper functions

```

In [11]: 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)

```

1. Oversampling using SMOTE

```

In [12]: 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 [13]: 1 vcs = listings["flag_ur"].value_counts()
2         vcs

```

```

Out[13]: 0    43321
1      3093
3      1624
2       817
Name: flag_ur, dtype: int64

```

```

In [14]: 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])

```

```

In [15]: 1 print(
2         " Old size :", listings.shape, "\n",
3         "New size :", x.shape
4         )

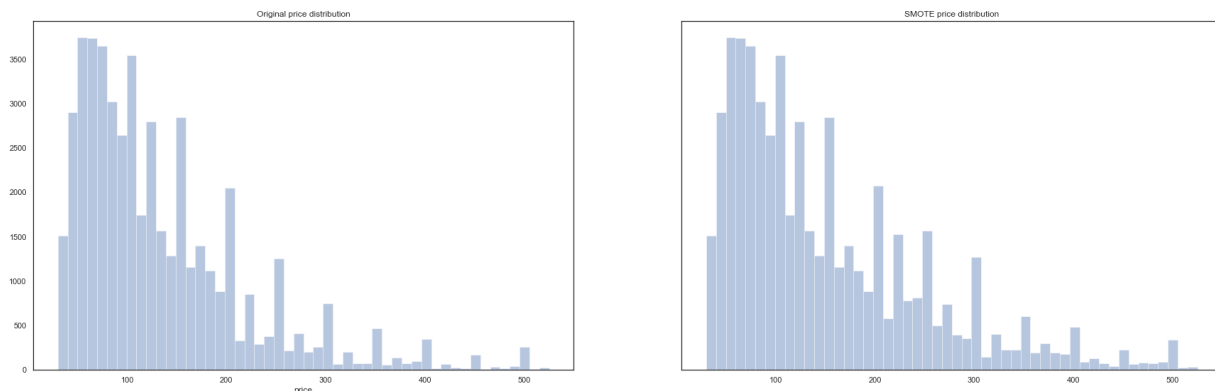
```

```

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

```

```
In [16]: 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 [17]: 1 x_cols = listings.drop(['price'], axis=1)
2         X = pd.DataFrame(data=x, columns=x_cols.columns)
```

```
In [18]: 1 X_log = X.copy()
2         X_sqr = X.copy()
3         X_sqrt = X.copy()
```

```
In [19]: 1 # Taking log, square and square root transformations of x
2         for i in range(X_log.shape[1]):
3             X_log.iloc[:,i] = np.log(X_log.iloc[:,i] + 1)
4             X_sqr.iloc[:,i] = np.square(X_sqr.iloc[:,i])
5             X_sqrt.iloc[:,i] = np.sqrt(X_sqrt.iloc[:,i])
6         X_log.columns = X_log.columns.map(lambda x: x + '_log')
7         X_sqr.columns = X_sqr.columns.map(lambda x: x + '_sqr')
8         X_sqrt.columns = X_sqrt.columns.map(lambda x: x + '_sqrt')
```

```
In [20]: 1 X_sqrt.shape
```

```
Out[20]: (54389, 64)
```

```
In [21]: 1 # Appending X, X_log, X_sqr and X_sqrt
2         X_final = pd.concat([X, X_log, X_sqr, X_sqrt], axis=1)
3         X_final.shape
```

```
Out[21]: (54389, 256)
```

Train test split

```
In [25]: 1 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.30, ra
```

Standardisation

```
In [26]: 1 standard_scaler = StandardScaler()
2 x_train = standard_scaler.fit_transform(x_train)
3 x_test = standard_scaler.transform(x_test)
```

Decision Trees

Initial Baseline

```
In [184]: 1 from sklearn.tree import DecisionTreeRegressor
2
3 dtr = DecisionTreeRegressor(random_state=RANDOM_SEED)
4 dtr.fit(X=x_train, y=y_train)
5 y_pred_train = dtr.predict(X=x_train)
6 y_pred_test = dtr.predict(X=x_test)
7
8 mse_train = MSE(y_train, y_pred_train)
9 mse_test = MSE(y_test, y_pred_test)
10
11 rmse_train = mse_train**(1/2)
12 rmse_test = mse_test**(1/2)
13
14 print("Train set RMSE Scaled: {:.2f}".format(rmse_train))
15 print("Test set RMSE Scaled: {:.2f}".format(rmse_test))
```

Train set RMSE Scaled: 1.75

Test set RMSE Scaled: 78.63

```
In [185]: 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))
```

Cross Validation RMSE: 78.11

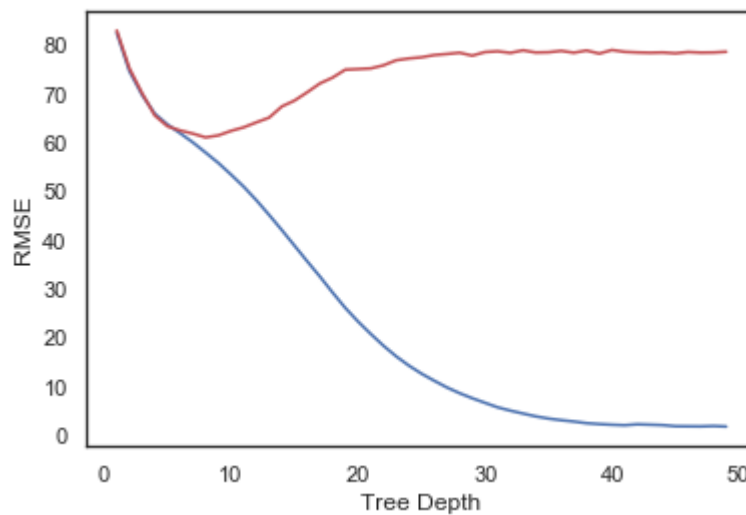
Individual Parameter Tuning

max_depth

```

In [63]: 1 max_depths = np.arange(1, 50, 1)
          2 train_results = []
          3 test_results = []
          4 for test_depth in max_depths:
          5
          6     dt = DecisionTreeRegressor(random_state=RANDOM_SEED,max_depth=test_depth)
          7
          8     dt.fit(X=x_train, y=y_train)
          9
          10    train_pred = dt.predict(X=x_train)
          11    train_mse = MSE(y_train, train_pred)
          12    train_results.append(train_mse**(1/2))
          13
          14    y_pred = dt.predict(x_test)
          15    test_mse = MSE(y_test, y_pred)
          16    test_results.append(test_mse**(1/2))
          17
          18    from matplotlib.legend_handler import HandlerLine2D
          19    line1, = plt.plot(max_depths, train_results, 'b', label='Train RMSE')
          20    line2, = plt.plot(max_depths, test_results, 'r', label='Test RMSE')
          21    plt.ylabel('RMSE')
          22    plt.xlabel('Tree Depth')
          23    plt.show()

```

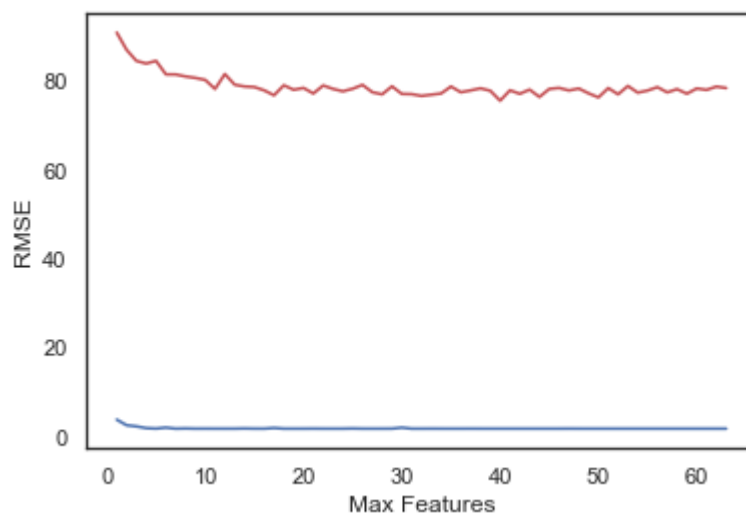


max_features

```

In [195]: 1 max_features = list(range(1,x_train.shape[1]))
          2 train_results = []
          3 test_results = []
          4 for max_feat_test in max_features:
          5
          6     dt = DecisionTreeRegressor(random_state=RANDOM_SEED,max_features=max_fea
          7
          8     dt.fit(X=x_train, y=y_train)
          9
          10    train_pred = dt.predict(X=x_train)
          11    train_mse = MSE(y_train, train_pred)
          12    train_results.append(train_mse**(1/2))
          13
          14    y_pred = dt.predict(x_test)
          15    test_mse = MSE(y_test, y_pred)
          16    test_results.append(test_mse**(1/2))
          17
          18 from matplotlib.legend_handler import HandlerLine2D
          19 line1, = plt.plot(max_features, train_results, 'b', label='Train RMSE')
          20 line2, = plt.plot(max_features, test_results, 'r', label='Test RMSE')
          21 plt.ylabel('RMSE')
          22 plt.xlabel('Max Features')
          23 plt.show()

```

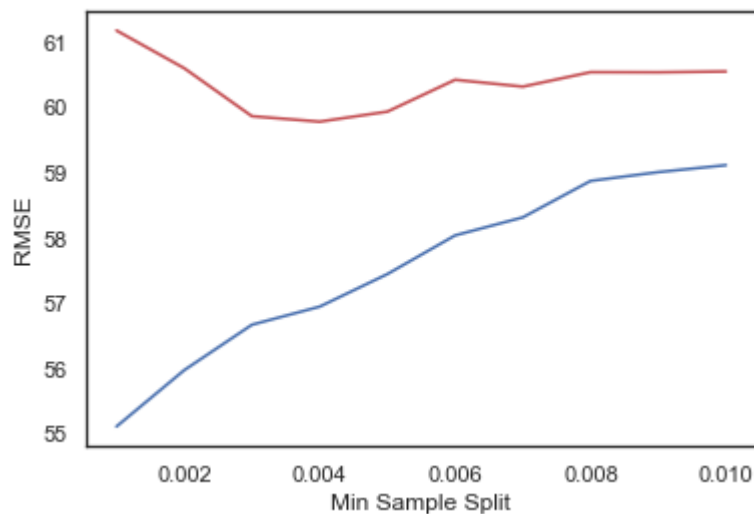


min_samples_splits


```

In [192]: 1 min_samples_splits = np.linspace(0.001, 0.01, 10, endpoint=True)
          2
          3 train_results = []
          4 test_results = []
          5 for min_samp_test in min_samples_splits:
          6
          7     dt = DecisionTreeRegressor(random_state=RANDOM_SEED,max_depth=10,min_sam
          8
          9     dt.fit(X=x_train, y=y_train)
         10
         11     train_pred = dt.predict(X=x_train)
         12     train_mse = MSE(y_train, train_pred)
         13     train_results.append(train_mse**(1/2))
         14
         15     y_pred = dt.predict(x_test)
         16     test_mse = MSE(y_test, y_pred)
         17     test_results.append(test_mse**(1/2))
         18
         19 from matplotlib.legend_handler import HandlerLine2D
         20 line1, = plt.plot(min_samples_splits, train_results, 'b', label='Train RMSE'
         21 line2, = plt.plot(min_samples_splits, test_results, 'r', label='Test RMSE')
         22 plt.ylabel('RMSE')
         23 plt.xlabel('Min Sample Split')
         24 plt.show()

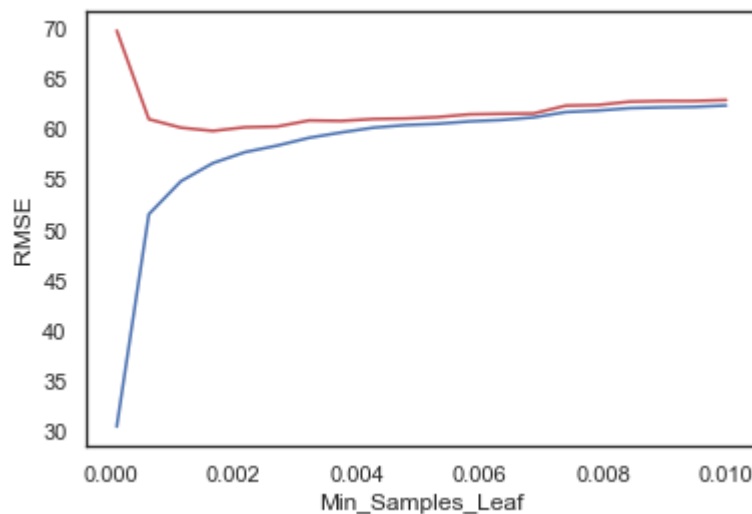
```



```

In [199]: 1 min_samples_leafs = np.linspace(0.0001, 0.01, 20, endpoint=True)
          2
          3 train_results = []
          4 test_results = []
          5 for min_samp_leaf_test in min_samples_leafs:
          6
          7     dt = DecisionTreeRegressor(random_state=RANDOM_SEED,min_samples_leaf=min
          8
          9     dt.fit(X=x_train, y=y_train)
          10
          11     train_pred = dt.predict(X=x_train)
          12     train_mse = MSE(y_train, train_pred)
          13     train_results.append(train_mse**(1/2))
          14
          15     y_pred = dt.predict(x_test)
          16     test_mse = MSE(y_test, y_pred)
          17     test_results.append(test_mse**(1/2))
          18
          19 from matplotlib.legend_handler import HandlerLine2D
          20 line1, = plt.plot(min_samples_leafs, train_results, 'b', label='Train RMSE')
          21 line2, = plt.plot(min_samples_leafs, test_results, 'r', label='Test RMSE')
          22 plt.ylabel('RMSE')
          23 plt.xlabel('Min_Samples_Leaf')
          24 plt.show()

```



GridSearch

In [201]: 1 grid_search.fit(x_test, y_test)

Fitting 3 folds for each of 8000 candidates, totalling 24000 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 37 tasks      | elapsed: 2.4s
[Parallel(n_jobs=-1)]: Done 513 tasks     | elapsed: 7.8s
[Parallel(n_jobs=-1)]: Done 1325 tasks    | elapsed: 20.4s
[Parallel(n_jobs=-1)]: Done 2457 tasks    | elapsed: 46.6s
[Parallel(n_jobs=-1)]: Done 3917 tasks    | elapsed: 1.4min
[Parallel(n_jobs=-1)]: Done 5697 tasks    | elapsed: 2.2min
[Parallel(n_jobs=-1)]: Done 7805 tasks    | elapsed: 2.9min
[Parallel(n_jobs=-1)]: Done 10233 tasks   | elapsed: 3.9min
[Parallel(n_jobs=-1)]: Done 12989 tasks   | elapsed: 4.9min
[Parallel(n_jobs=-1)]: Done 16065 tasks   | elapsed: 6.1min
[Parallel(n_jobs=-1)]: Done 19469 tasks   | elapsed: 7.5min
[Parallel(n_jobs=-1)]: Done 23193 tasks   | elapsed: 9.0min
[Parallel(n_jobs=-1)]: Done 23993 out of 24000 | elapsed: 9.4min remaining:
0.2s
[Parallel(n_jobs=-1)]: Done 24000 out of 24000 | elapsed: 9.4min finished
```

```
Out[201]: GridSearchCV(cv=3, error_score='raise-deprecating',
      estimator=DecisionTreeRegressor(criterion='mse', max_depth=None, max_fea-
      tures=None,
      max_leaf_nodes=None, min_impurity_decrease=0.0,
      min_impurity_split=None, min_samples_leaf=1,
      min_samples_split=2, min_weight_fraction_leaf=0.0,
      presort=False, random_state=None, splitter='best'),
      fit_params=None, iid='warn', n_jobs=-1,
      param_grid={'max_depth': [8, 9, 10, 11, 12, 13, 14, 15], 'max_features':
      [10, 20, 30, 40, 50], 'min_samples_leaf': [0.0001, 0.0006210526315789474, 0.001
      1421052631578948, 0.0016631578947368423, 0.0021842105263157894, 0.0027052631578
      947366, 0.0032263157894736843, 0.0037473684210526316, 0.0042684210526...amples_
      split': [0.001, 0.002, 0.003, 0.004, 0.005, 0.006, 0.007, 0.008, 0.009000000000
      000001, 0.01]}),
      pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
      scoring='neg_mean_squared_error', verbose=2)
```

In [202]: 1 grid_search.best_params_

```
Out[202]: {'max_depth': 14,
      'max_features': 50,
      'min_samples_leaf': 0.0011421052631578948,
      'min_samples_split': 0.01}
```

Final Decision Tree Model

In [203]:

```

1 from sklearn.tree import DecisionTreeRegressor
2
3 dtr = DecisionTreeRegressor(random_state=RANDOM_SEED,max_depth=14,max_featur
4 dtr.fit(X=x_train, y=y_train)
5 y_pred_train = dtr.predict(X=x_train)
6 y_pred_test = dtr.predict(X=x_test)
7
8 mse_train = MSE(y_train, y_pred_train)
9 mse_test = MSE(y_test, y_pred_test)
10
11 rmse_train = mse_train**(1/2)
12 rmse_test = mse_test**(1/2)
13
14 print("Train set RMSE Scaled: {:.2f}".format(rmse_train))
15 print("Test set RMSE Scaled: {:.2f}".format(rmse_test))

```

Train set RMSE Scaled: 58.78

Test set RMSE Scaled: 59.99

In [205]:

```

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))

```

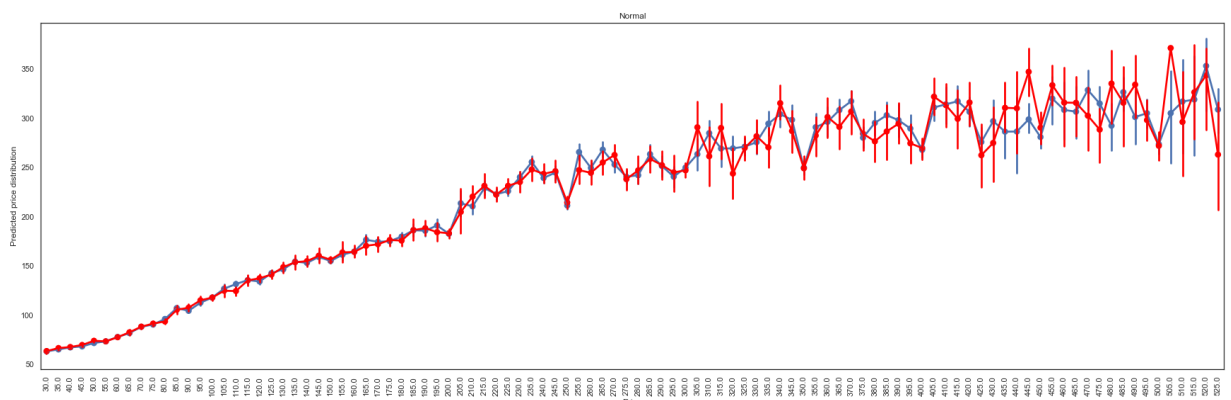
Cross Validation RMSE: 61.04

In [206]:

```

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("Normal")

```



Feature Importance

```
In [106]: 1 x_final = listings.drop(['price'], axis=1)
```

```
In [107]: 1 imp = pd.DataFrame({'Input Variable': x_final.columns, 'Importance (%)': -np
2 top20 = imp.iloc[:20,:]}
3
4 f, ax = plt.subplots(1,1,figsize=(20,9))
5 g = sns.barplot(x="Importance (%)", y="Input Variable", data=top20, ax=ax)
6 t = g.set_xlabel("Importance (%)")
7 t = g.set_ylabel("Input Variable")
8 t = g.set_title("Feature Importance")
```

