# AirBnb listings file modeling

```
In [1]:   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 [1]:   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, Gradi
         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(data_path, index_col="id")
          3  display(listings.shape)
```
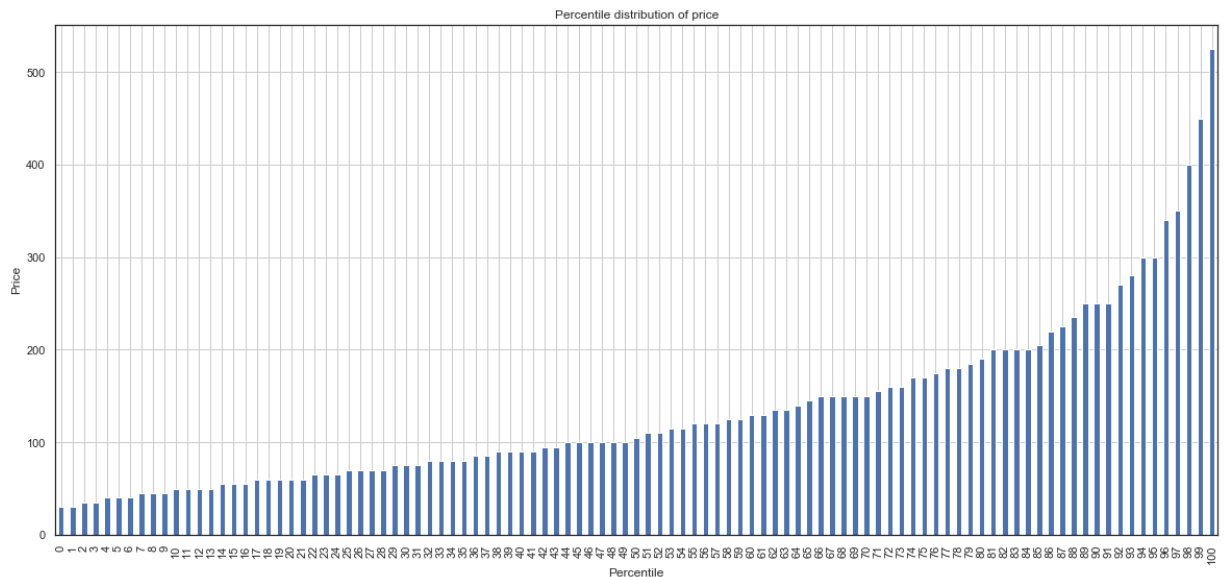
(48855, 65)

## Plot the dstribution

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")
```

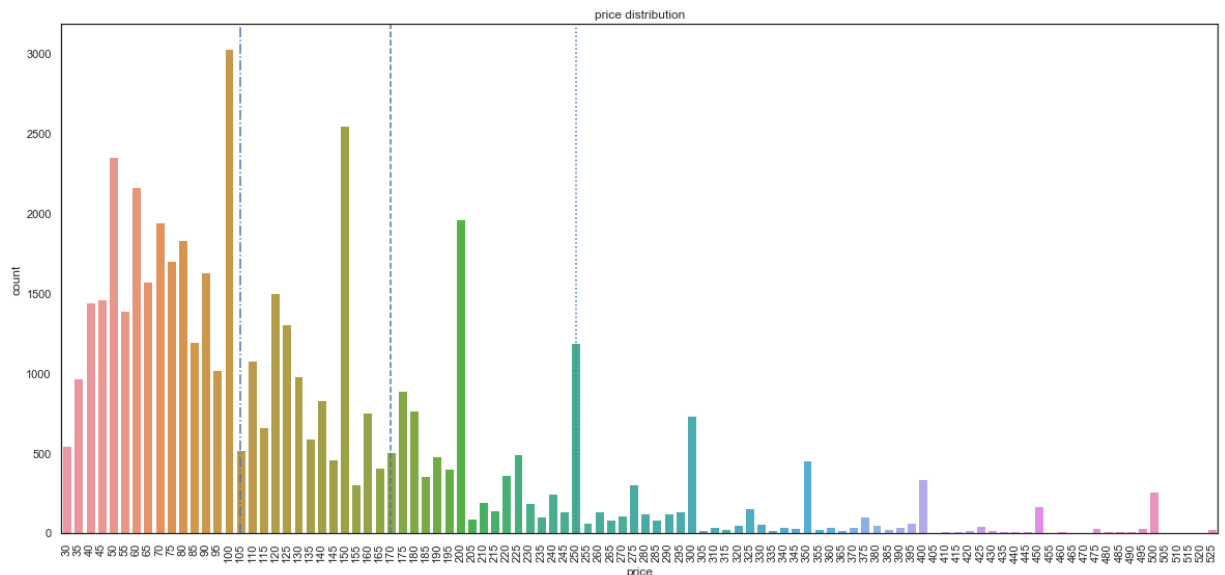Out[4]:  Text(0, 0.5, 'Price')

```
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 0x12b4d870>



Quick helper functions

```
In [6]:    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 [7]:
```python
def check_rep(row):
    if (row <= 200) | (row==250) | (row==350) | (row==450) | (row==550) :
        return 0
    elif (row > 200) & (row < 300) & (row != 250):
        return 1
    elif (row > 300) & (row < 400) & (row != 350):
        return 2
    else:
        return 3

listings["flag_ur"] = listings["price"].apply(check_rep)
```

In [8]:
```python
vcs = listings["flag_ur"].value_counts()
vcs
```

Out[8]:
```
0    43321
1     3093
3     1624
2      817
Name: flag_ur, dtype: int64
```

In [9]:
```python
ycol = ["flag_ur"]
xcol = [i for i in listings.columns if i not in ycol]

x = listings[xcol].values
y = listings[ycol].values

smote_sampling_strategy = {
    1: int(vcs[1]*2)
    ,2: int(vcs[2]*2)
    ,3: int(vcs[3]*2)
}
sm = SMOTE(random_state=RANDOM_SEED, sampling_strategy=smote_sampling_strate
# Fit the smote onto the sample
x_new, y_new = sm.fit_sample(x, y)

# Drop the flag column
listings.drop(labels=["flag_ur"], axis=1, inplace=True)

# --------------------------------------------------------------------------
# Overwrite X and Y
price_index = get_index(list(listings.columns), "price")

y = x_new[:, price_index]
x = np.delete(x_new, price_index, axis=1)
for i in range(len(y)):
    y[i] = roundto(y[i])
```
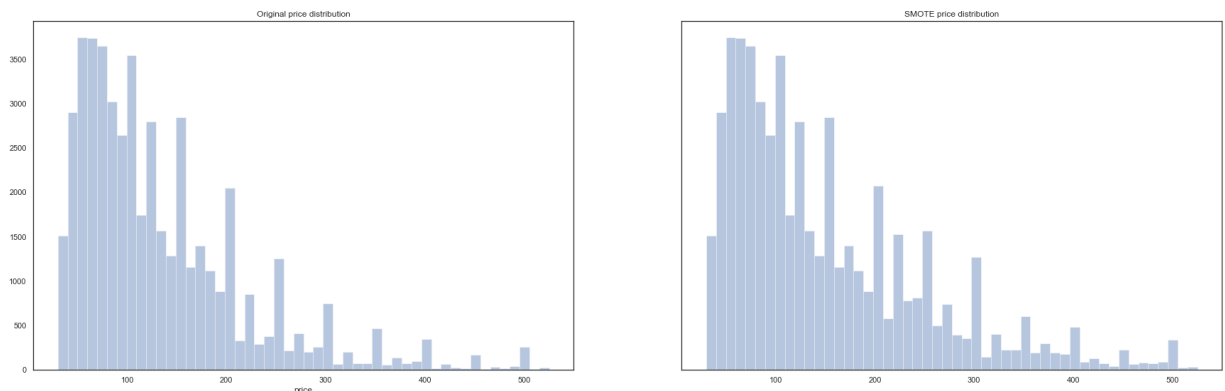
```
c:\users\sriharis\appdata\local\programs\python\python37-32\lib\site-packages\s
klearn\utils\validation.py:761: DataConversionWarning: A column-vector y was pa
ssed when a 1d array was expected. Please change the shape of y to (n_samples,
), for example using ravel().
  y = column_or_1d(y, warn=True)
```

In [10]:
```python
print(
    " Old size :", listings.shape, "\n",
    "New size :", x.shape
)
```

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

In [11]:
```python
f, ax = plt.subplots(1, 2, figsize=(30, 9), sharey=True)
g1 = sns.distplot(listings["price"], ax=ax[0], kde=False)
t = g1.set_title("Original price distribution")

g2 = sns.distplot(y, ax=ax[1], kde=False)
t = g2.set_title("SMOTE price distribution")
```



## 2. Transformation

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

In [13]:
```python
X_log = X.copy()
X_sqr = X.copy()
X_sqrt = X.copy()
```

In [14]:
```python
# Taking log, square and square root transformations of x
for i in range(X_log.shape[1]):
    X_log.iloc[:,i] = np.log(X_log.iloc[:,i] + 1)
    X_sqr.iloc[:,i] = np.square(X_sqr.iloc[:,i])
    X_sqrt.iloc[:,i] = np.sqrt(X_sqrt.iloc[:,i])
X_log.columns = X_log.columns.map(lambda x: x + '_log')
X_sqr.columns = X_sqr.columns.map(lambda x: x + '_sqr')
X_sqrt.columns = X_sqrt.columns.map(lambda x: x + '_sqrt')
```

In [15]:
```python
X_sqrt.shape
```

Out[15]: (54389, 64)

```
In [16]:  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[16]: (54389, 256)

## Prediction

```
In [17]:  1  x_train, x_test, y_train, y_test = train_test_split(X, y, test_size=0.30, ra
```

```
In [18]:  1  rfr = RandomForestRegressor(random_state=RANDOM_SEED)
          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
         12  print("Train set RMSE with Transformation: {:.2f}".format(rmse_train))
         13  print("Test set RMSE with Transformation: {:.2f}".format(rmse_test))
```

```
c:\users\sriharis\appdata\local\programs\python\python37-32\lib\site-packages\s
klearn\ensemble\forest.py:246: FutureWarning: The default value of n_estimators
will change from 10 in version 0.20 to 100 in 0.22.
  "10 in version 0.20 to 100 in 0.22.", FutureWarning)

Train set RMSE with Transformation: 24.25
Test set RMSE with Transformation: 56.60
```

```
In [19]:  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 with Transformation: {:.2f}'.format(CV))
```

```
Cross Validation RMSE with Transformation: 57.38
```

## 3. Train test split

```
In [20]:  1  x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.30, ra
```
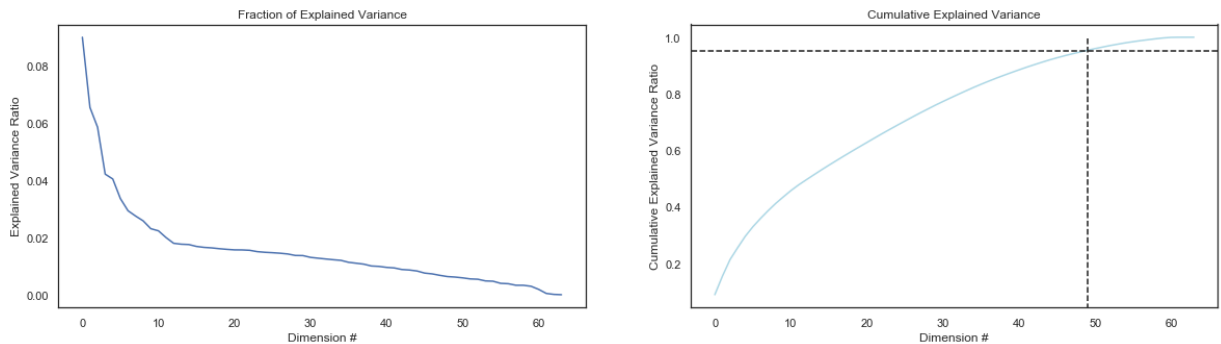
## 4. Standardisation

In [21]:
```python
1  standard_scaler = StandardScaler()
2  x_train = standard_scaler.fit_transform(x_train)
3  x_test = standard_scaler.transform(x_test)
```

## 5. PCA

In [22]:
```python
1  pca_naive = PCA(n_components=None)
2  pca_naive.fit(x_train)
3  f, ax = plt.subplots(1,2,figsize=(20,5))
4  plot_naive_variance(pca_naive, ax[0])
5  num_pc = plot_pca_var_cum(pca_naive, ax=ax[1], cutoff=0.95)
6  print("Number of components : ", num_pc)
```

```
c:\users\sriharis\appdata\local\programs\python\python37-32\lib\site-packages\m
atplotlib\cbook\__init__.py:1725: UserWarning: Saw kwargs ['c', 'color'] which
are all aliases for 'color'.  Kept value from 'color'
  seen=seen, canon=canonical, used=seen[-1]))
```

```
Number of components :   49
```



In [23]:
```python
1  # pca = PCA(n_components=num_pc)
2  pca = PCA(n_components=None)
3  pca.fit(x_train)
4  x_train_pca = pca.transform(x_train)
5  x_test_pca = pca.transform(x_test)
```

## Prediction

## Random Forest Regressor

**With PCA**

In [24]:
```python
rfr = RandomForestRegressor(random_state=RANDOM_SEED)
rfr.fit(X=x_train_pca, y=y_train)
y_pred_train = rfr.predict(X=x_train_pca)
y_pred_test = rfr.predict(X=x_test_pca)

mse_train = MSE(y_train, y_pred_train)
mse_test = MSE(y_test, y_pred_test)

rmse_train = mse_train**(1/2)
rmse_test = mse_test**(1/2)

print("Train set RMSE with PCA: {:.2f}".format(rmse_train))
print("Test set RMSE with PCA: {:.2f}".format(rmse_test))
```

```
c:\users\sriharis\appdata\local\programs\python\python37-32\lib\site-packages\s
klearn\ensemble\forest.py:246: FutureWarning: The default value of n_estimators
will change from 10 in version 0.20 to 100 in 0.22.
  "10 in version 0.20 to 100 in 0.22.", FutureWarning)

Train set RMSE with PCA: 26.72
Test set RMSE with PCA: 62.94
```

In [25]:
```python
# Use cross-validation on Train data

CV_scores = - cv.cross_val_score(rfr, x_train_pca, y_train, scoring='neg_mea

# Compute the 10-folds CV
CV = CV_scores.mean()**(1/2)

# Print Train CV accuracy
print('Cross Validation RMSE with PCA: {:.2f}'.format(CV))
```

```
Cross Validation RMSE with PCA: 63.11
```

In [26]:
```python
f, ax = plt.subplots(1,1, figsize=(30, 9))
g = sns.pointplot(x=y_train, y=y_pred_train, ax=ax)
g = sns.pointplot(x=y_test, y=y_pred_test, ax=ax, color="red")

t = g.set_xlabel("Price")
t = g.set_ylabel("Predicted price distribution")
t = g.set_xticklabels(g.get_xticklabels(), rotation=90)
t = g.set_title("PCA")
```