# ▾ Imports

```
!pip install fake_useragent

    Collecting fake_useragent
      Downloading https://files.pythonhosted.org/packages/d1/79/af647635d6968e2deb57a208d30g
    Building wheels for collected packages: fake-useragent
      Building wheel for fake-useragent (setup.py) ... done
      Created wheel for fake-useragent: filename=fake_useragent-0.1.11-cp37-none-any.whl siz
      Stored in directory: /root/.cache/pip/wheels/5e/63/09/d1dc15179f175357d3f5c00cbffbac37
    Successfully built fake-useragent
    Installing collected packages: fake-useragent
    Successfully installed fake-useragent-0.1.11
```

```python
from bs4 import BeautifulSoup as bs4
from requests import get
import json
import pandas as pd
import requests
import matplotlib.pyplot as plt
import seaborn as sns
import mpl_toolkits
import numpy as np
%matplotlib inline
#from fake_useragent import UserAgent
```

# ▾ Data preparation (Web scraping)

Preparing the data by extracting information from the first three pages of the website.

```python
url_1 = 'https://www.pararius.com/apartments/eindhoven'
url_2 = 'https://www.pararius.com/apartments/eindhoven/page-2'
url_3 = 'https://www.pararius.com/apartments/eindhoven/page-3'
urls= ['https://www.pararius.com/apartments/eindhoven','https://www.pararius.com/apartments/e
#var_urls = url_1, url_2, url_3
#urls.append(var_urls)
#result=get(my_url)

page_1 = requests.get(url_1)
page_2 = requests.get(url_2)
page_3 = requests.get(url_3)
#page_2.content
```

```
# html parsing
page1_soup= bs4(page_1.text, "html.parser")
page2_soup= bs4(page_2.text, "html.parser")
page3_soup= bs4(page_3.text, "html.parser")
```

## Page 1

```
# grab each product
allHouses1 = page1_soup.findAll("li", {"class": "search-list__item search-list__item--listing
houses1 = page1_soup.findAll("ul", {"class": "search-list"})[0].text
#data_rows = table.findAll('')[2:]
print(len(allHouses1))
print(len(houses1))
```

```
    32
    13585
```

## Page 2

```
allHouses2 = page2_soup.findAll("li", {"class": "search-list__item search-list__item--listing
houses2 = page2_soup.findAll("ul", {"class": "search-list"})[0].text
#data_rows = table.findAll('')[2:]
print(len(allHouses2))
print(len(houses2))
```

```
    32
    13782
```

## Page 3

```
allHouses3 = page3_soup.findAll("li", {"class": "search-list__item search-list__item--listing
houses3 = page3_soup.findAll("ul", {"class": "search-list"})[0].text
#data_rows = table.findAll('')[2:]
print(len(allHouses3))
print(len(houses3))
```

```
    32
    13702
```

### Data description

- house_type - house/apartment/room
- street - the name of the street where the property is placed

- postcode - the postcode of the property
- price
- year - the year of construction
- living_area
- rooms - how many rooms are in the property

Page 1 - saving the extracted data

```
catalog=[]
for h in allHouses1:
    #data['houses'].append({
        name = h.findAll('a',class_='listing-search-item__link listing-search-item__link--tit
        _name = name.split()
        house_type = _name[0]
        street = _name[1]
        _address= h.findAll('div', class_='listing-search-item__location')[0].text
        #String manipulation to remove the unwanted signs from the address
        __address = _address.replace("\nnew\n ", "")
        address = __address.replace("\n ", "")   #actual address after string manipulation -
        new_address = address.split()
        postcode = new_address[0]
        price_ = h.findAll('span', class_='listing-search-item__price')[0].text
        #splitting the string to find the price
        p=price_.split()
        _price = p[0] #actual price before string manipulation
        __price = _price.replace("€", "")  #actual price before full string manipulation
        price = __price.replace(",", "")   #actual price after string manipulation - ready to

        #finding the whole element from the web page
        ylr= h.findAll('section', class_= 'illustrated-features illustrated-features--vertica

        #splitting the string to find the living are, rooms and year
        lry= ylr.split()
#may use another for
        #living_area after taking the indexes that define it
        year = h.findAll('span', class_= 'illustrated-features__description')[0].text

        #living_area after taking the indexes that define it
        living_area = lry[0]

        #rooms after taking the index that defines the variable
        rooms = lry[4]

        vars = house_type, street, postcode,price,year,living_area,rooms
        catalog.append(vars)

print(year)
```

```
    34 m²
```

## Page 2 - saving the extracted data

```
catalog2=[]
for h in allHouses2:
    #data['houses'].append({
        name = h.findAll('a',class_='listing-search-item__link listing-search-item__link--tit
        _name = name.split()
        house_type = _name[0]
        street = _name[1]
        _address= h.findAll('div', class_='listing-search-item__location')[0].text
        #String manipulation to remove the unwanted signs from the address
        __address = _address.replace("\nnew\n ", "")
        address = __address.replace("\n ", "")   #actual address after string manipulation -
        new_address = address.split()
        postcode = new_address[0]
        price_ = h.findAll('span', class_='listing-search-item__price')[0].text
        #splitting the string to find the price
        p=price_.split()
        _price = p[0] #actual price before string manipulation
        __price = _price.replace("€", "")  #actual price before full string manipulation
        price = __price.replace(",", "")   #actual price after string manipulation - ready to

        #finding the whole element from the web page
        ylr= h.findAll('section', class_= 'illustrated-features illustrated-features--vertica

        #splitting the string to find the living are, rooms and year
        lry= ylr.split()
#may use another for
        #living_area after taking the indexes that define it
        year = h.findAll('span', class_= 'illustrated-features__description')[0].text

        #living_area after taking the indexes that define it
        living_area = lry[0]

        #rooms after taking the index that defines the variable
        rooms = lry[4]

        vars = house_type, street, postcode,price,year,living_area,rooms
        catalog2.append(vars)

print(year)

    100 m²
```

## Page 3 - saving the extracted data

```
catalog3=[]
for h in allHouses3:
    #data['houses'].append({
        name = h.findAll('a',class_='listing-search-item__link listing-search-item__link--tit
        _name = name.split()
        house_type = _name[0]
        street = _name[1]
        _address= h.findAll('div', class_='listing-search-item__location')[0].text
        #String manipulation to remove the unwanted signs from the address
        __address = _address.replace("\nnew\n ", "")
        address = __address.replace("\n ", "")   #actual address after string manipulation -
        new_address = address.split()
        postcode = new_address[0]
        price_ = h.findAll('span', class_='listing-search-item__price')[0].text
        #splitting the string to find the price
        p=price_.split()
        _price = p[0] #actual price before string manipulation
        __price = _price.replace("€", "")  #actual price before full string manipulation
        price = __price.replace(",", "")   #actual price after string manipulation - ready to

        #finding the whole element from the web page
        ylr= h.findAll('section', class_= 'illustrated-features illustrated-features--vertica

        #splitting the string to find the living are, rooms and year
        lry= ylr.split()
#may use another for
        #living_area after taking the indexes that define it
        year = h.findAll('span', class_= 'illustrated-features__description')[0].text

        #living_area after taking the indexes that define it
        living_area = lry[0]

        #rooms after taking the index that defines the variable
        rooms = lry[4]

        vars = house_type, street, postcode,price,year,living_area,rooms
        catalog3.append(vars)

print(street)

    Edenstraat
```

> Saving the scraped data to pandas dataframe (creating the table and giving names to the cokumns)

```
# Create DataFrame
df1 = pd.DataFrame(catalog)
df2 = pd.DataFrame(catalog2)
df3 = pd.DataFrame(catalog3)
```

```
df1.columns=['TYPE', 'STREET NAME', 'POSTCODE', 'PRICE', 'YEAR', 'LIVING_AREA', 'ROOMS']
df2.columns=['TYPE', 'STREET NAME', 'POSTCODE', 'PRICE', 'YEAR', 'LIVING_AREA', 'ROOMS']
df3.columns=['TYPE', 'STREET NAME', 'POSTCODE', 'PRICE', 'YEAR', 'LIVING_AREA', 'ROOMS']
```

**Data integration**

Using `Union` to integrate the scraped data from the three web pages.

```
frames = [df1, df2, df3]
```

```
df = pd.concat(frames)
df
```

|  | TYPE | STREET NAME | POSTCODE | PRICE | YEAR | LIVING_AREA | ROOMS |
|---|---|---|---|---|---|---|---|
| 0 | Apartment | Kruisstraat | 5612 | 1000 | 42 m² | 42 | 2 |
| 1 | House | Nieuwe | 5612 | 1150 | 65 m² | 65 | 3 |
| 2 | Apartment | 1e | 5614 | 1195 | 93 m² | 93 | 3 |
| 3 | House | Henkenshage | 5653 | 1425 | 138 m² | 138 | 5 |
| 4 | House | Count | 5629 | 3495 | 279 m² | 279 | 5 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 27 | Apartment | Hertog | 5611 | 1100 | 65 m² | 65 | 3 |
| 28 | Apartment | Aalsterweg | 5615 | 1400 | 60 m² | 60 | 3 |
| 29 | Apartment | St | 5616 | 895 | 25 m² | 25 | 2 |
| 30 | Apartment | De | 5611 | 1150 | 80 m² | 80 | 2 |
| 31 | Apartment | Edenstraat | 5611 | 995 | 68 m² | 68 | 2 |

96 rows × 7 columns

# Data analysis

Here we can see the shape of our data with the .shape. Here we see that we have 31 rows and 7 columns. However, they are always changing because the data is alway up to date by using the web scraping technique.

Checking the dimension of the dataset

```
df.shape
```

```
(96, 7)
```

```
df.head()
```

|   | TYPE | STREET NAME | POSTCODE | PRICE | YEAR | LIVING_AREA | ROOMS |
|---|------|-------------|----------|-------|------|-------------|-------|
| 0 | Apartment | Kruisstraat | 5612 | 1000 | 42 m² | 42 | 2 |
| 1 | House | Nieuwe | 5612 | 1150 | 65 m² | 65 | 3 |
| 2 | Apartment | 1e | 5614 | 1195 | 93 m² | 93 | 3 |
| 3 | House | Henkenshage | 5653 | 1425 | 138 m² | 138 | 5 |
| 4 | House | Count | 5629 | 3495 | 279 m² | 279 | 5 |

```
df.describe()
```

|       | TYPE | STREET NAME | POSTCODE | PRICE | YEAR | LIVING_AREA | ROOMS |
|-------|------|-------------|----------|-------|------|-------------|-------|
| count | 96 | 96 | 96 | 96 | 96 | 96 | 96 |
| unique | 3 | 66 | 18 | 54 | 55 | 55 | 6 |
| top | Apartment | Kruisstraat | 5611 | 1500 | 45 m² | 45 | 2 |
| freq | 75 | 7 | 24 | 8 | 7 | 7 | 38 |

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 96 entries, 0 to 31
Data columns (total 7 columns):
 #   Column        Non-Null Count   Dtype
---  ------        --------------   -----
 0   TYPE          96 non-null      object
 1   STREET NAME   96 non-null      object
 2   POSTCODE      96 non-null      object
 3   PRICE         96 non-null      object
 4   YEAR          96 non-null      object
 5   LIVING_AREA   96 non-null      object
 6   ROOMS         96 non-null      object
dtypes: object(7)
memory usage: 6.0+ KB
```

```
df.tail()
```

|  | TYPE | STREET NAME | POSTCODE | PRICE | YEAR | LIVING_AREA | ROOMS |
|---|---|---|---|---|---|---|---|
| 27 | Apartment | Hertog | 5611 | 1100 | 65 m² | 65 | 3 |
| 28 | Apartment | Aalsterweg | 5615 | 1400 | 60 m² | 60 | 3 |
| 29 | Apartment | St | 5616 | 895 | 25 m² | 25 | 2 |

`df.iloc[0]`

```
TYPE            Apartment
STREET NAME    Kruisstraat
POSTCODE             5612
PRICE                1000
YEAR               42 m²
LIVING_AREA            42
ROOMS                   2
Name: 0, dtype: object
```

`df.sort_values('TYPE', ascending = True)`

|  | TYPE | STREET NAME | POSTCODE | PRICE | YEAR | LIVING_AREA | ROOMS |
|---|---|---|---|---|---|---|---|
| 0 | Apartment | Kruisstraat | 5612 | 1000 | 42 m² | 42 | 2 |
| 6 | Apartment | Hoefkestraat | 5611 | 950 | 70 m² | 70 | 3 |
| 5 | Apartment | Cornelis | 5642 | 650 | 25 m² | 25 | 2 |
| 3 | Apartment | Petrus | 5613 | 750 | 22 m² | 22 | 1 |
| 2 | Apartment | Alpenroosstraat | 5644 | 1095 | 47 m² | 47 | 2 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 8 | Room | Boschdijk | 5612 | 795 | 24 m² | 24 | 1 |
| 4 | Room | Petrus | 5613 | 650 | 18 m² | 18 | 1 |
| 7 | Room | St | 5616 | 495 | 16 m² | 16 | 1 |
| 18 | Room | Tongelresestraat | 5642 | 595 | 19 m² | 19 | 1 |
| 21 | Room | Heydaalweg | 5616 | 410 | 14 m² | 14 | 1 |

96 rows × 7 columns

There are no missing values in the dataset.

`df.isnull().all()`

```
TYPE           False
STREET NAME    False
POSTCODE       False
PRICE          False
```

```
    YEAR            False
    LIVING_AREA     False
    ROOMS           False
    dtype: bool
```

```
# Find columns with missing values and their percent missing
df.isnull().sum()
miss_val = df.isnull().sum().sort_values(ascending=False)
miss_val = pd.DataFrame(data=df.isnull().sum().sort_values(ascending=False), columns=['Missva

# Add a new column to the dataframe and fill it with the percentage of missing values
miss_val['Percent'] = miss_val.MissvalCount.apply(lambda x : '{:.2f}'.format(float(x)/df.shap
miss_val = miss_val[miss_val.MissvalCount > 0]
miss_val
```

**MissvalCount Percent**

**Pre Processing**

Handling Outlier

> An **outlier** is a data point in a data set that is distant from all other observations (a
> data point that lies outside the overall distribution of the dataset.)

```
plt.figure(figsize=(10, 10), dpi=80)
plt.scatter(df.LIVING_AREA, df.PRICE, c= 'red')
plt.title("Outliers")
plt.xlabel("LivArea")
plt.ylabel("Price")
plt.show()
```

Outliers

```python
df.boxplot(column=['PRICE'])
plt.show
```

```
<function matplotlib.pyplot.show>
```



```python
sorted(df)
```

```
['LIVING_AREA', 'POSTCODE', 'PRICE', 'ROOMS', 'STREET NAME', 'TYPE', 'YEAR']
```

Using scatter plots to visualize the relationship between the variables and the targeted variable - PRICE.

```python
plt.figure(figsize=(20, 5), dpi=80)
plt.scatter(df['PRICE'],df['ROOMS'])
plt.title("Price vs Rooms")
```

```
plt.xlabel("Price")
plt.ylabel("Rooms")
plt.show()
sns.despine
```



```
<function seaborn.utils.despine>
```

> It can be noticed that there is a positive correlation between the price and the living area, which means that the variables move in tandem—that is, in the same direction. This means that whenever one variable increases, the other decreases. For instance, the price increases with the more rooms the housing has.
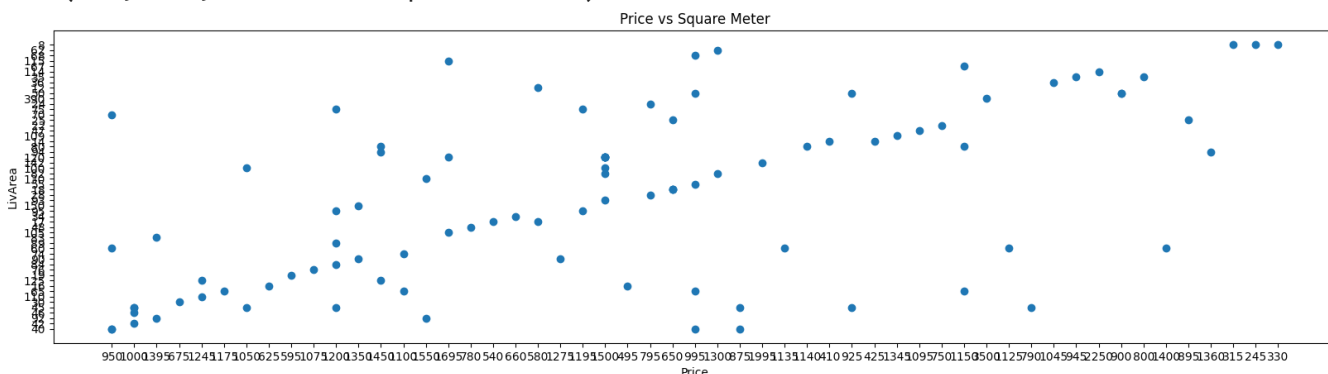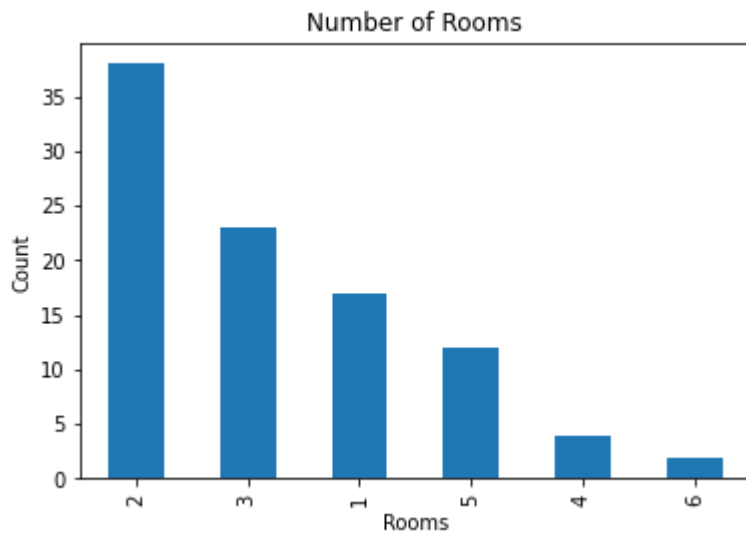
```
plt.figure(figsize=(20, 5), dpi=80)
plt.scatter(df['PRICE'],df['POSTCODE'])
plt.xlabel("Price")
plt.ylabel("Postcode")
plt.title("Price vs Postcode")
```

```
Text(0.5, 1.0, 'Price vs Postcode')
```

Price vs Postcode



```
plt.figure(figsize=(20, 5), dpi=100)
plt.scatter(df['PRICE'],df['LIVING_AREA'])
plt.xlabel("Price")
plt.ylabel("LivArea")
plt.title("Price vs Square Meter")
```

```
Text(0.5, 1.0, 'Price vs Square Meter')
```

Price vs Square Meter



> It can be noticed that there is a positive correlation between the price and the living
> area, which means that the variables move in tandem—that is, in the same direction.
> This means that whenever one variable increases, the other decreases. For instance,
> the price increases with the increase in the living area.

```
df['PRICE'] =df['PRICE'].astype(float)
df['ROOMS'] =df['ROOMS'].astype(int)
df['LIVING_AREA'] =df['LIVING_AREA'].astype(int)
```

```
df['ROOMS'].value_counts().plot(kind='bar')
plt.title('Number of Rooms')
plt.xlabel('Rooms')
plt.ylabel('Count')
sns.despine
```

```
<function seaborn.utils.despine>
```



```
print(df['PRICE'])
```

```
0        1000.0
1        1150.0
2        1195.0
3        1425.0
4        3495.0
          ...
27       1100.0
28       1400.0
29        895.0
30       1150.0
31        995.0
Name: PRICE, Length: 96, dtype: float64
```

Changing the type of the variable `price` in order to plot it in the next diagram.

```
df['PRICE'] =df['PRICE'].astype(float)
df['POSTCODE'] =df['POSTCODE'].astype(int)
df['LIVING_AREA'] =df['LIVING_AREA'].astype(int)
df['ROOMS'] =df['ROOMS'].astype(int)
code_numeric = {'Apartment': 1, 'Room': 2, 'House': 3}
df ['TYPE'] = df['TYPE'].map(code_numeric)
```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 96 entries, 0 to 31
Data columns (total 7 columns):
 #   Column       Non-Null Count  Dtype
---  ------       --------------  -----
 0   TYPE         96 non-null     int64
 1   STREET NAME  96 non-null     object
 2   POSTCODE     96 non-null     int64
 3   PRICE        96 non-null     float64
```

```
 4   YEAR          96 non-null     object
 5   LIVING_AREA   96 non-null     int64
 6   ROOMS         96 non-null     int64
dtypes: float64(1), int64(4), object(2)
memory usage: 6.0+ KB
```
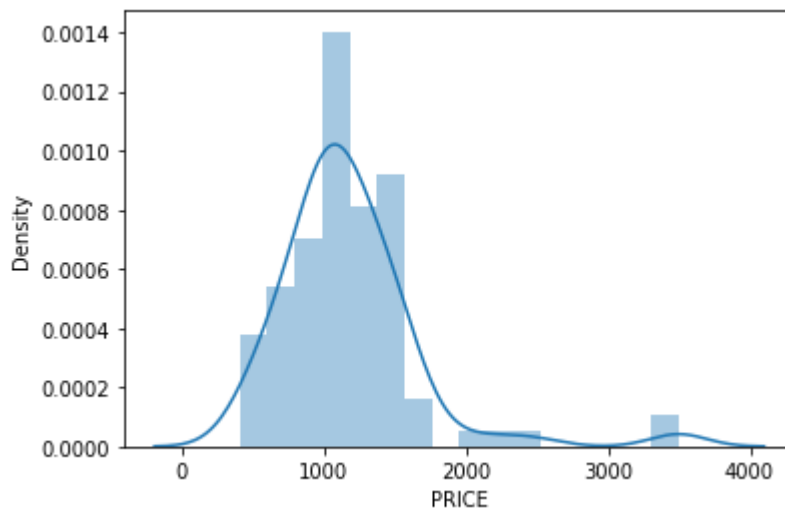
Examining the data distributions of the features. We will start with the target variable, `PRICE`, to make sure it's normally distributed.

```
sns.distplot(df['PRICE'])
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2557: FutureWarning: `di
  warnings.warn(msg, FutureWarning)
<matplotlib.axes._subplots.AxesSubplot at 0x7effa05c7910>
```



We can see that the `PRICE` distribution is not skewed, but normally distributed.

**Normally distributed** means that the data is symmetric about the mean, showing that data near the mean are more frequent in occurrence than data far from the mean.

Reviewing the skewness of each feature

```
df.skew().sort_values(ascending=False)
```

```
LIVING_AREA    2.899261
PRICE          2.333455
TYPE           1.619129
POSTCODE       1.559712
ROOMS          0.863245
dtype: float64
```

> Values closer to zero are less skewed. The results show some features having a
> positive (right-tailed) or negative (left-tailed) skew.

Factor plot is informative when we have multiple groups to compare.

```
sns.factorplot('ROOMS', 'PRICE', data=df,kind='bar',size=3,aspect=3)
fig, (axis1) = plt.subplots(1,1,figsize=(10,3))
sns.countplot('ROOMS', data=df)
df['PRICE'].value_counts()
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass th
  FutureWarning
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass th
  FutureWarning
1500.0     8
995.0      5
1200.0     5
950.0      4
1000.0     4
650.0      4
1150.0     3
1695.0     3
1450.0     3
875.0      2
925.0      2
795.0      2
580.0      2
900.0      2
1550.0     2
1100.0     2
1350.0     2
1395.0     2
1245.0     2
1195.0     2
1050.0     2
1300.0     2
780.0      1
540.0      1
660.0      1
595.0      1
495.0      1
625.0      1
1175.0     1
675.0      1
2250.0     1
1140.0     1
410.0      1
1125.0     1
1095.0     1
1345.0     1
1135.0     1
1995.0     1
1275.0     1
1075.0     1
330.0      1
245.0      1
315.0      1
1360.0     1
895.0      1
```

> Real estate with 5 rooms has the highest `Price` while the sales of others with rooms
> of 2 is the most sold ones.

```
------   -
750.0     1
```

```
#g = sns.factorplot(x='POSTCODE', y='Skewed_SP', col='PRICE', data=df, kind='bar', col_wrap=4
sns.factorplot('POSTCODE', 'PRICE', data=df,kind='bar',size=3,aspect=6)
```
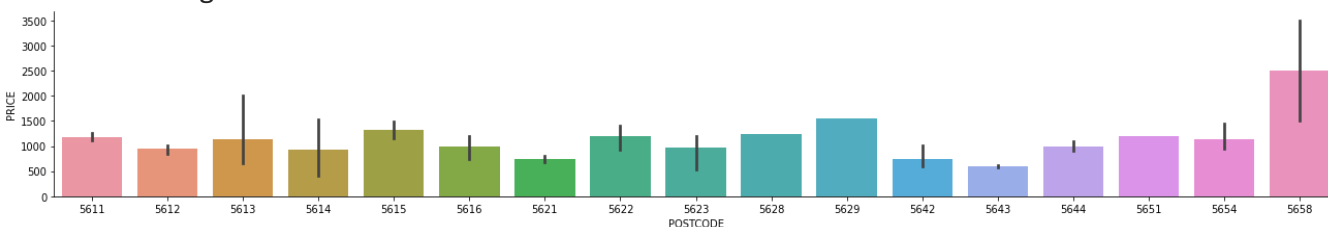
```
/usr/local/lib/python3.7/dist-packages/seaborn/categorical.py:3714: UserWarning: The `fa
  warnings.warn(msg)
/usr/local/lib/python3.7/dist-packages/seaborn/categorical.py:3720: UserWarning: The `si
  warnings.warn(msg, UserWarning)
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass th
  FutureWarning
<seaborn.axisgrid.FacetGrid at 0x7f0fabaa5f10>
```

The diagram represents the `price` of a rpoperty, depending on its `postcode`.

# ▾ Train-Test Split dataset

Necessary imports

```
from sklearn.linear_model import LinearRegression, Ridge, Lasso
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import cross_val_score, train_test_split, GridSearchCV
```

Analyzing the numeric features.

```
numeric_features = df.select_dtypes(include=[np.number])
```

```
numeric_features.columns
```

```
       Index(['TYPE', 'POSTCODE', 'PRICE', 'LIVING_AREA', 'ROOMS'], dtype='object')
```

```
# set the target and predictors
y = df.PRICE  # target
```

```
# use only those input features with numeric data type
df_temp = df_select_dtypes(include=["int64", "float64"])
```

```
df_temp = df.select_dtypes(include=[ int64 , float64 ])

X = df_temp.drop(["PRICE"],axis=1)  # predictors


X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20, random_state=0)
```

# Modelling

### Linear Regression

```
lr = LinearRegression()
# fit optimal linear regression line on training data
lr.fit((X_train),y_train)
```

```
    LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

```
#predict y_values using X_test set
yr_hat = lr.predict(X_test)
```

```
lr_score =lr.score((X_test),y_test)
print("Accuracy: ", lr_score)
```

```
    Accuracy:  0.8860749530938599
```

Using cross-validation to see whether the model is over-fitting the data.

```
# cross validation to find 'validate' score across multiple samples, automatically does Kfold
lr_cv = cross_val_score(lr, X, y, cv = 5, scoring= 'r2')
print("Cross-validation results: ", lr_cv)
print("R2: ", lr_cv.mean())
```

```
    Cross-validation results:  [0.7633578  0.76243208 0.79135206 0.74946125 0.68485638]
    R2:  0.7502919129374968
```

### Random Forest

```
rfr = RandomForestRegressor()
rfr.fit(X_train, y_train) # gets the parameters for the rfr model
rfr_cv = cross_val_score(rfr,X, y, cv = 5, scoring = 'r2')
print("R2: ", rfr_cv.mean())
```

```
    R2:  0.6020438259058933
```

```
rfr.feature_importances_
```

```
    array([0.00537268, 0.05812923, 0.88598102, 0.05051707])
```

# ▾ Plotting the Feature Importance

```
importance = rfr.feature_importances_

# map feature importance values to the features
feature_importances = zip(importance, X.columns)

#list(feature_importances)
sorted_feature_importances = sorted(feature_importances, reverse = True)

#print(sorted_feature_importances)
top_15_predictors = sorted_feature_importances[0:15]
values = [value for value, predictors in top_15_predictors]
predictors = [predictors for value, predictors in top_15_predictors]
print(predictors)
```
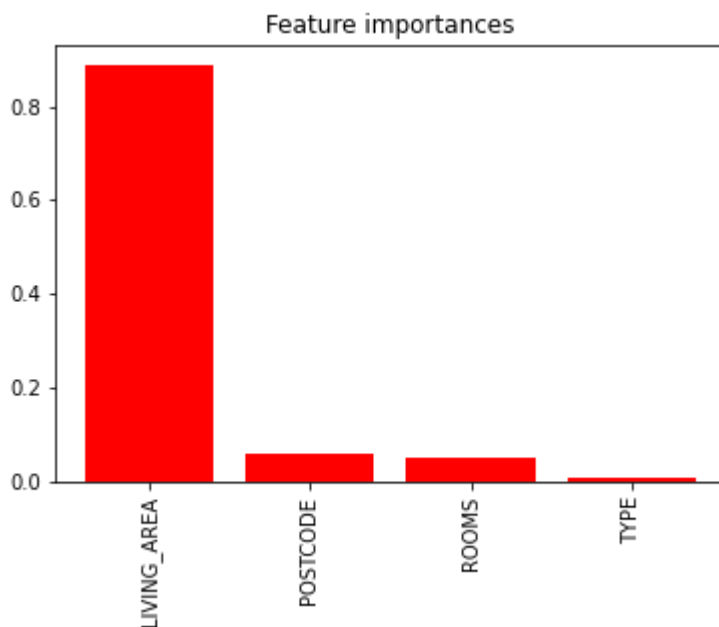
```
[→    ['LIVING_AREA', 'POSTCODE', 'ROOMS', 'TYPE']
```

**Plotting the feauture importance of the Random forest.**

```
plt.figure()
plt.title("Feature importances")
plt.bar(range(len(predictors)), values,color="r", align="center");
plt.xticks(range(len(predictors)), predictors, rotation=90);
```



Feature importances

# Conclusion

**Data collection:**

For the data collection part, I decided to use `web scraping` as e technique because it gives the opportunity to work with a data set that is up to date and therefore, make more accurate summaries.

**Data preprocessing:**

I tried different types of data transfoms to expose the data structure better, so we may be able to improve model accuracy later.

- `Standardizing` was made to the data set so as to reduce the effects of differing distributions.
- `The skewness` of the feautures was checked in order to see how distorted a data sample is from the normal distribution.
- `Rescaling (normalizing)` the dataset was also included to reduce the effects of differing scales

**Modelling:**

I used two models to determine the accuracy - Linear Regression and Random Forest.

Linear Regression turns out to be the more accurate model for predicting the house price. It scored an estimated accuracy of 75%, out performing the Random Forest. Random Forest determined that overall the living area of a home is by far the most important predictor. Following are the size of above rooms and postcode.