

Data Analysis with Python

House Sales in King County, USA

This dataset contains house sale prices for King County, which includes Seattle. It includes homes sold between May 2014 and May 2015.

id : A notation for a house

date: Date house was sold

price: Price is prediction target

bedrooms: Number of bedrooms

bathrooms: Number of bathrooms

sqft_living: Square footage of the home

sqft_lot: Square footage of the lot

floors :Total floors (levels) in house

waterfront: House which has a view to a waterfront

view: Has been viewed

condition :How good the condition is overall

grade: overall grade given to the housing unit, based on King County grading system

sqft_above : Square footage of house apart from basement

sqft_basement: Square footage of the basement

yr_built : Built Year

yr_renovated : Year when house was renovated

zipcode: Zip code

lat: Latitude coordinate

long: Longitude coordinate

sqft_living15: Living room area in 2015(implies-- some renovations) This might or might not have affected the lotsize area

sqft_lot15 : LotSize area in 2015(implies-- some renovations)

You will require the following libraries:

```
In [12]: import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler, PolynomialFeatures
from sklearn.linear_model import LinearRegression
matplotlib inline
```

Module 1: Importing Data Sets

Load the csv:

```
In [13]: file_name='https://s3-apl.us-gco.objectstorage.softlayer.net/cf-courses-data/CognitiveClass/DA0101EN/co
ursera/project/kc_house_data_NaN.csv'
df=pd.read_csv(file_name)
```

We use the method `head` to display the first 5 columns of the dataframe.

```
In [14]: df.head()
Out[14]:
```

	Unnamed: 0	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	...	grade	sqft_abv
0	0	7129300520	20141013T000000	221900.0	3.0	1.00	1180	5650	1.0	0	...	7	1'
1	1	6414100192	20141029T000000	538000.0	3.0	2.25	2570	7242	2.0	0	...	7	2'
2	2	5631500400	20150225T000000	180000.0	2.0	1.00	770	10000	1.0	0	...	6	;
3	3	2487200875	20141209T000000	604000.0	4.0	3.00	1960	5000	1.0	0	...	7	1f
4	4	1954400510	20150216T000000	510000.0	3.0	2.00	1680	8080	1.0	0	...	8	1f

5 rows × 22 columns

Question 1

Display the data types of each column using the attribute `dtype`, then take a screenshot and submit it, include your code in the image.

```
In [15]: df.dtypes
Out[15]:
```

	Unnamed: 0	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	...	grade	sqft_abv
count	21613.00000	2.161300e+04	2.161300e+04	21600.000000	21603.000000	21603.000000	21613.000000	21613.000000	21613.000000	21613.000000	21613.000000	21613.000000	2
mean	10806.00000	4.580302e+09	5.400881e+05	3.372870	3.372870	2.115736	2079.899736	1.510697e+04	1.494309	0.007542	0.234303	3.409430	
std	6239.28002	2.876566e+09	3.671272e+05	0.926657	0.768996	0.768996	918.440897	4.142051e+04	0.539989	0.086517	0.766318	0.650743	
min	0.00000	1.000102e+06	7.500000e+04	1.000000	0.500000	0.500000	290.000000	5.200000e+02	1.000000	0.000000	0.000000	0.000000	
25%	5403.00000	2.123049e+09	3.219500e+05	3.000000	1.750000	1.750000	1427.000000	5.040000e+03	1.000000	0.000000	0.000000	3.000000	
50%	10806.00000	3.904930e+09	4.500000e+05	3.000000	2.250000	1910.000000	7.618000e+03	1.500000	0.000000	0.000000	0.000000	3.000000	
75%	16209.00000	7.308900e+09	6.450000e+05	4.000000	2.500000	2550.000000	1.068800e+04	2.000000	0.000000	0.000000	4.000000	5.000000	
max	21612.00000	9.900000e+09	7.700000e+06	33.000000	8.000000	8.000000	13540.000000	1.651359e+06	3.500000	1.000000	4.000000	5.000000	

8 rows × 21 columns

Module 2: Data Wrangling

Question 2

Drop the columns `"id"` and `"Unnamed: 0"` from axis 1 using the method `drop()`, then use the method `describe()` to obtain a statistical summary of the data. Take a screenshot and submit it, make sure the `inplace` parameter is set to `True`

```
In [19]: df=df.drop('id',1)
df=df.drop('Unnamed: 0',1)
df.describe()
Out[19]:
```

	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	condition
count	21613.00000	21600.000000	21603.000000	21613.000000	21613.000e+04	21613.000000	21613.000000	21613.000000	2
mean	5.400881e+05	3.372870	2.115736	2079.899736	1.510697e+04	1.494309	0.007542	0.234303	3.409430
std	3.671272e+05	0.926657	0.768996	918.440897	4.142051e+04	0.539989	0.086517	0.766318	0.650743
min	7.500000e+04	1.000000	0.500000	290.000000	5.200000e+02	1.000000	0.000000	0.000000	1.000000
25%	3.219500e+05	3.000000	1.750000	1427.000000	5.040000e+03	1.000000	0.000000	0.000000	3.000000
50%	4.500000e+05	3.000000	2.250000	1910.000000	7.618000e+03	1.500000	0.000000	0.000000	3.000000
75%	6.450000e+05	4.000000	2.500000	2550.000000	1.068800e+04	2.000000	0.000000	0.000000	4.000000
max	7.700000e+06	33.000000	8.000000	13540.000000	1.651359e+06	3.500000	1.000000	4.000000	5.000000

We can see we have missing values for the columns `bedrooms` and `bathrooms`

```
In [20]: print("number of NaN values for the column bedrooms :", df['bedrooms'].isnull().sum())
print("number of NaN values for the column bathrooms :", df['bathrooms'].isnull().sum())

number of NaN values for the column bedrooms : 13
number of NaN values for the column bathrooms : 10
```

We can replace the missing values of the column `'bedrooms'` with the mean of the column `'bedrooms'` using the method `replace()`. Don't forget to set the `inplace` parameter to `True`

```
In [21]: mean=df['bedrooms'].mean()
df['bedrooms'].replace(np.nan,mean, inplace=True)
```

We also replace the missing values of the column `'bathrooms'` with the mean of the column `'bathrooms'` using the method `replace()`. Don't forget to set the `inplace` parameter to `True`

```
In [22]: mean=df['bathrooms'].mean()
df['bathrooms'].replace(np.nan,mean, inplace=True)
```

```
In [23]: print("number of NaN values for the column bedrooms :", df['bedrooms'].isnull().sum())
print("number of NaN values for the column bathrooms :", df['bathrooms'].isnull().sum())

number of NaN values for the column bedrooms : 0
number of NaN values for the column bathrooms : 0
```

Module 3: Exploratory Data Analysis

Question 3

Use the method `value_counts` to count the number of houses with unique floor values, use the method `.to_frame()` to convert it to a dataframe.

```
In [27]: df['floors'].value_counts()
Out[27]:
```

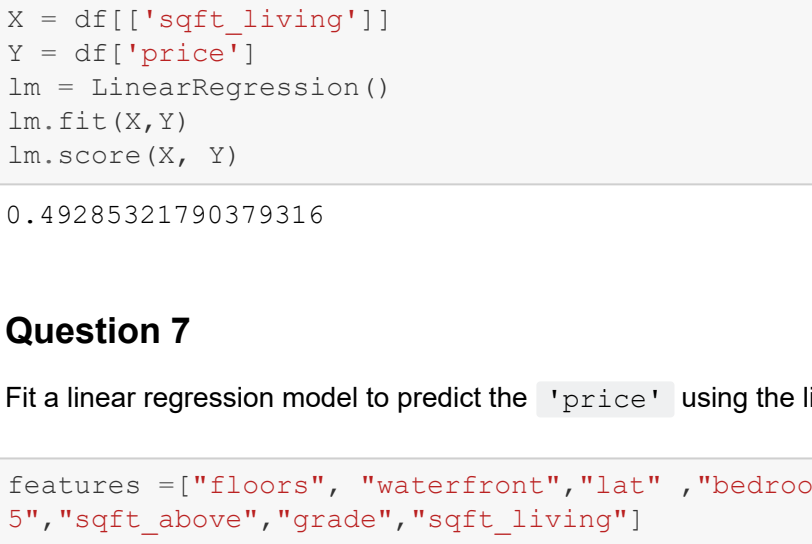
	price
1.0	10680
2.0	8241
1.5	1910
3.0	613
2.5	161
3.5	8

Name: floors, dtype: int64

Question 4

Use the function `boxplot` in the seaborn library to determine whether houses with a waterfront view or without a waterfront view have more price outliers.

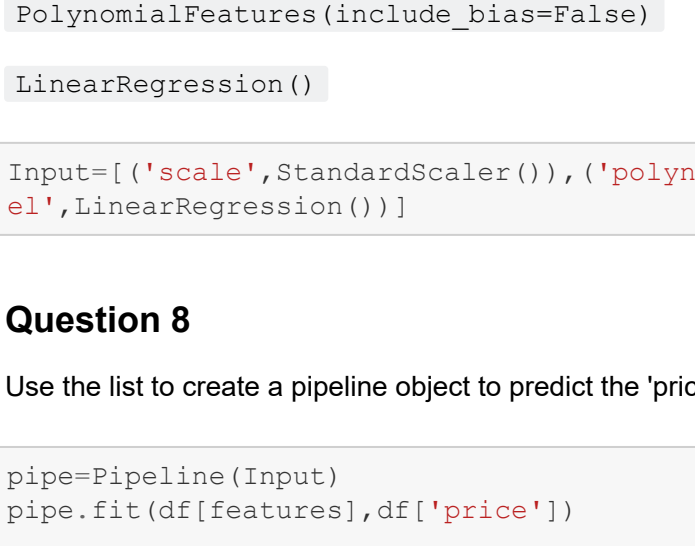
```
In [34]: df.boxplot(by = 'waterfront', column = ['price'])
Out[34]: <matplotlib.axes._subplots.AxesSubplot at 0x2084dfc5b88>
```



Question 5

Use the function `regplot` in the seaborn library to determine if the feature `sqft_above` is negatively or positively correlated with price.

```
In [36]: sns.set_style('whitegrid')
sns.lmplot(x = 'sqft_above', y = 'price', data = df)
Out[36]: <seaborn.axisgrid.FacetGrid at 0x2084e1340c8>
```



We can use the Pandas method `corr()` to find the feature other than price that is most correlated with price.

```
In [37]: df.corr()['price'].sort_values()
Out[37]:
```

	price
zipcode	-0.053203
long	0.021626
condition	0.036362
yr_built	0.054012
sqft_lot15	0.062447
sqft_lot	0.089661
yr_renovated	0.126434
floors	0.256794
waterfront	0.266369
lat	0.307003
bedrooms	0.308797
sqft_basement	0.323816
view	0.397293
bathrooms	0.525378
sqft_living15	0.585379
sqft_above	0.605567
grade	0.667434
sqft_living	0.702035
price	1.000000

Name: price, dtype: float64

Module 4: Model Development

We can Fit a linear regression model using the longitude feature `'long'` and calculate the R^2 .

```
In [38]: X = df[['long']]
Y = df['price']
lm = LinearRegression()
lm.fit(X,Y)
lm.score(X, Y)
Out[38]: 0.00046769430149029567
```

Question 6

Fit a linear regression model to predict the `'price'` using the feature `'sqft_living'` then calculate the R^2 . Take a screenshot of your code and the value of the R^2 .

```
In [39]: X = df[['sqft_living']]
Y = df['price']
lm = LinearRegression()
lm.fit(X,Y)
lm.score(X, Y)
Out[39]: 0.49285321790379316
```

Question 7

Fit a linear regression model to predict the `'price'` using the list of features:

```
In [44]: features = ["floors", "waterfront","lat", "bedrooms", "sqft_basement", "view", "bathrooms","sqft_living15", "sqft_above","grade","sqft_living"]

Then calculate the  $R^2$ . Take a screenshot of your code.
```

```
In [46]: X = df[features]
Y = df['price']
lm = LinearRegression()
lm.fit(X,Y)
lm.score(X, Y)
Out[46]: 0.6576853050765703
```

This will help with Question 8

Create a list of tuples, the first element in the tuple contains the name of the estimator:

```
'scale'
'polynomial'
'model'
```

The second element in the tuple contains the model constructor

```
StandardScaler()
PolynomialFeatures(include_bias=False)
LinearRegression()
```

```
In [69]: Input=[('scale',StandardScaler()),('polynomial', PolynomialFeatures(include_bias=False,degree=2)),('model',LinearRegression())]
```

Question 8

Use the list to create a pipeline object to predict the `'price'`, fit the object using the features in the list `features`, and calculate the R^2 .

```
In [77]: pipe=Pipeline(Input)
pipe.fit(df[features],df['price'])

pipe.score(df[features],df['price'])
Out[77]: 0.7513409690477972
```

Module 5: Model Evaluation and Refinement

Import the necessary modules:

```
In [71]: from sklearn.model_selection import cross_val_score
from sklearn.model_selection import train_test_split
print("done")

done
```

We will split the data into training and testing sets:

```
In [72]: features = ["floors", "waterfront","lat", "bedrooms", "sqft_basement", "view", "bathrooms","sqft_living15", "sqft_above","grade","sqft_living"]
X = df[features]
Y = df['price']

X_train, x_test, y_train, y_test = train_test_split(X, Y, test_size=0.15, random_state=1)

print("number of test samples:", x_test.shape[0])
print("number of training samples:",x_train.shape[0])

number of test samples: 3242
number of training samples: 18371
```

Question 9

Create and fit a Ridge regression object using the training data, set the regularization parameter to 0.1, and calculate the R^2 using the test data.

```
In [73]: from sklearn.linear_model import Ridge

In [76]: ridge=Ridge()
ridge.fit(x_train,y_train,0.1)
y_pred=ridge.predict(x_test)
ridge.score(x_test,y_test)
Out[76]: 0.6470831724882585
```

Question 10

Perform a second order polynomial transform on both the training data and testing data. Create and fit a Ridge regression object using the training data, set the regularisation parameter to 0.1, and calculate the R^2 utilising the test data provided. Take a screenshot of your code and the R^2 .

```
In [78]: pr = PolynomialFeatures(degree = 2)
X_train_pr = pr.fit_transform(x_train)
X_test_pr = pr.fit_transform(x_test)

rr = Ridge(alpha = 0.1)
rr.fit(X_train_pr, y_train)
rr.score(X_test_pr, y_test)
Out[78]: 0.7002744250224031
```

Once you complete your work you will have to share it. Select the icon on the top right a marked in red in the image below, a dialogue box should open, and select the option all content excluding sensitive code cells.

You can then share the notebook via a URL by scrolling down as shown in the following image:

About the Authors:

[Joseph Santacangelo](#) has a PhD in Electrical Engineering, his research focused on using machine learning, signal processing, and computer vision to determine how videos impact human cognition. Joseph has been working for IBM since he completed his PhD.

Other contributors: [Michelle Carey](#), [Mavis Zhou](#)

```
In [ ] :
```