



## 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/House

**bathrooms**: Number of bathrooms/bedrooms

**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 [2]: 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
%matplotlib inline
```

## 1.0 Importing the Data

Load the csv:

```
In [3]: file_name='https://s3-api.us-geo.objectstorage.softlayer.net/cf-courses-data/C
df=pd.read_csv(file_name)
```

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

```
In [4]: df.head()
```

Out [4]:

	Unnamed: 0	id	date	price	bedrooms	bathrooms	sqft_living	sqft
0	0	7129300520	20141013T000000	221900.0	3.0	1.00	1180	
1	1	6414100192	20141209T000000	538000.0	3.0	2.25	2570	
2	2	5631500400	20150225T000000	180000.0	2.0	1.00	770	1
3	3	2487200875	20141209T000000	604000.0	4.0	3.00	1960	
4	4	1954400510	20150218T000000	510000.0	3.0	2.00	1680	

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 [5]: `print(df.dtypes)`

```

Unnamed: 0      int64
id              int64
date           object
price          float64
bedrooms       float64
bathrooms      float64
sqft_living    int64
sqft_lot       int64
floors         float64
waterfront     int64
view           int64
condition      int64
grade          int64
sqft_above     int64
sqft_basement  int64
yr_built       int64
yr_renovated   int64
zipcode        int64
lat            float64
long           float64
sqft_living15  int64
sqft_lot15     int64
dtype: object

```

We use the method `describe` to obtain a statistical summary of the dataframe.

In [6]: `df.describe()`

Out [6]:

	Unnamed: 0	id	price	bedrooms	bathrooms	sqft_living
<b>count</b>	21613.00000	2.161300e+04	2.161300e+04	21600.000000	21603.000000	21613.000000
<b>mean</b>	10806.00000	4.580302e+09	5.400881e+05	3.372870	2.115736	2079.899736
<b>std</b>	6239.28002	2.876566e+09	3.671272e+05	0.926657	0.768996	918.440897
<b>min</b>	0.00000	1.000102e+06	7.500000e+04	1.000000	0.500000	290.000000
<b>25%</b>	5403.00000	2.123049e+09	3.219500e+05	3.000000	1.750000	1427.000000
<b>50%</b>	10806.00000	3.904930e+09	4.500000e+05	3.000000	2.250000	1910.000000
<b>75%</b>	16209.00000	7.308900e+09	6.450000e+05	4.000000	2.500000	2550.000000
<b>max</b>	21612.00000	9.900000e+09	7.700000e+06	33.000000	8.000000	13540.000000

8 rows x 7 columns

## 2.0 Data Wrangling

In [22]:

```
#### Question 2
Drop the columns <code>"id"</code> and <code>"Unnamed: 0"</code> from axis 1 u
```

```
File "<ipython-input-22-2139ceb94885>", line 2
    Drop the columns <code>"id"</code> and <code>"Unnamed: 0"</code> from axis
s 1 using the method <code>drop()</code>, then use the method <code>describe()
</code> to obtain a statistical summary of the data. Take a screenshot and sub
mit it, make sure the inplace parameter is set to <code>True</code>
      ^
SyntaxError: invalid syntax
```

In [7]:

```
df=pd.read_csv(file_name)

df.drop(["id", "Unnamed: 0"], axis=1, inplace = True)

df.describe()
```

Out [7]:

	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors
<b>count</b>	2.161300e+04	21600.000000	21603.000000	21613.000000	2.161300e+04	21613.000000
<b>mean</b>	5.400881e+05	3.372870	2.115736	2079.899736	1.510697e+04	1.494309
<b>std</b>	3.671272e+05	0.926657	0.768996	918.440897	4.142051e+04	0.539989
<b>min</b>	7.500000e+04	1.000000	0.500000	290.000000	5.200000e+02	1.000000
<b>25%</b>	3.219500e+05	3.000000	1.750000	1427.000000	5.040000e+03	1.000000
<b>50%</b>	4.500000e+05	3.000000	2.250000	1910.000000	7.618000e+03	1.500000
<b>75%</b>	6.450000e+05	4.000000	2.500000	2550.000000	1.068800e+04	2.000000
<b>max</b>	7.700000e+06	33.000000	8.000000	13540.000000	1.651359e+06	3.500000

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

```
In [8]: 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 [9]: 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 [10]: mean=df['bathrooms'].mean()
df['bathrooms'].replace(np.nan,mean, inplace=True)
```

```
In [11]: 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

## 3.0 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 [12]: df['floors'].value_counts()
```

```
Out[12]: 1.0    10680
2.0     8241
1.5     1910
3.0      613
2.5      161
3.5         8
Name: floors, dtype: int64
```

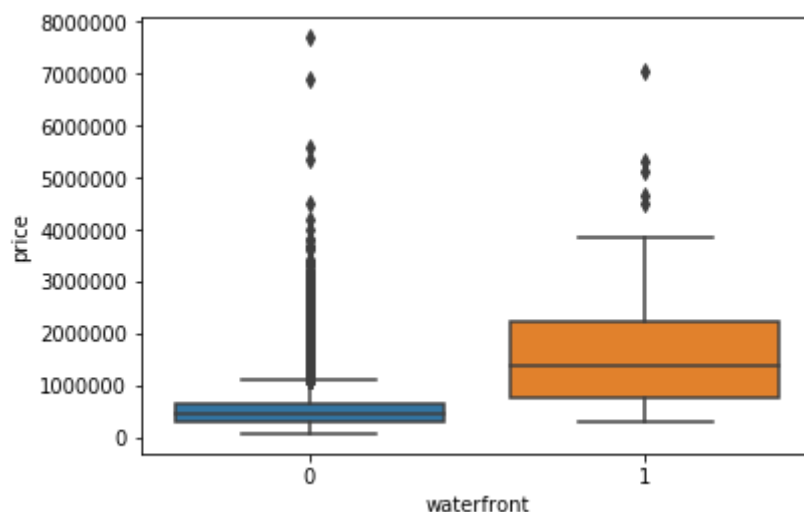
```
In [13]: df['floors'].value_counts().to_frame()
```

Out[13]:

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

## Question 4

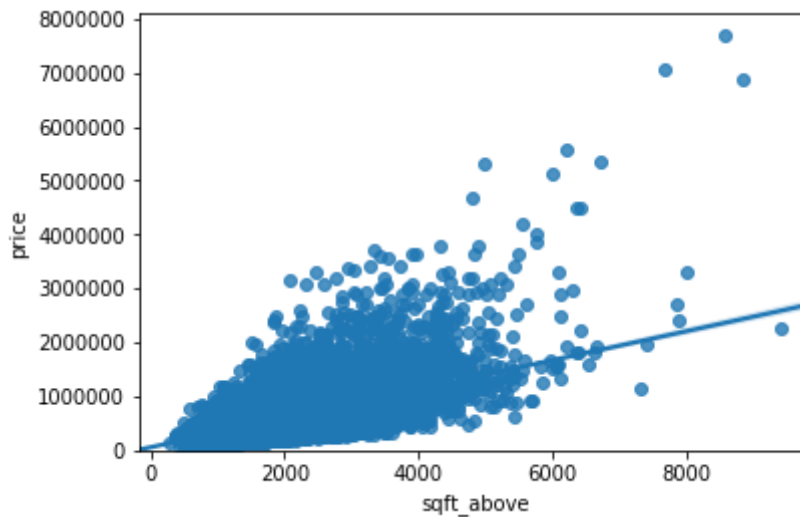
Use the function `boxplot` in the seaborn library to determine whether houses with a waterfront view or without a waterfront view have more price outliers .

In [14]: `sns.boxplot(x="waterfront", y="price", data=df)`Out[14]: `<matplotlib.axes._subplots.AxesSubplot at 0x7fb1c0a32da0>`

## 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 [15]: `sns.regplot(x="sqft_above", y="price", data=df)  
plt.ylim(0,)`Out[15]: `(0, 8086161.400594347)`



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

```
In [16]: df.corr()['price'].sort_values()
```

```
Out[16]: zipcode      -0.053203
long              0.021626
condition         0.036362
yr_built         0.054012
sqft_lot15       0.082447
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        0.525738
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

Import libraries

```
In [17]: import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression
```

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

```
In [18]: X = df[['long']]
Y = df['price']
lm = LinearRegression()
lm
lm.fit(X,Y)
lm.score(X, Y)
```

Out[18]: 0.00046769430149007363

## 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 [19]: lm = LinearRegression()
lm

X = df[['sqft_living']]
Y = df['price']

lm.fit(X,Y)

lm.score(X,Y)
```

Out[19]: 0.49285321790379316

```
In [20]: y_data = df['price']

x_data=df.drop('price',axis=1)

from sklearn.model_selection import train_test_split

x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_size=0.1)

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

lre=LinearRegression()

lre.fit(x_train[['sqft_living']], y_train)
lre.score(x_test[['sqft_living']], y_test)
```

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

Out[20]: 0.4910058627910614

## Question 7

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

```
In [21]: features =["floors", "waterfront","lat" ,"bedrooms" ,"sqft_basement" ,"view" ,'
```



```
In [22]: lm = LinearRegression()  
lm  
  
X = df[['floors']]  
Y = df['price']  
  
lm.fit(X,Y)  
lm.score(X,Y)
```

Out[22]: 0.06594310068341092

```
In [23]: lm = LinearRegression()  
lm  
  
X = df[['waterfront']]  
Y = df['price']  
  
lm.fit(X,Y)  
lm.score(X,Y)
```

Out[23]: 0.07095267538578309

```
In [24]: lm = LinearRegression()  
lm  
  
X = df[['lat']]  
Y = df['price']  
  
lm.fit(X,Y)  
lm.score(X,Y)
```

Out[24]: 0.09425113672917484

```
In [25]: lm = LinearRegression()  
lm  
  
X = df[['bedrooms']]  
Y = df['price']  
  
lm.fit(X,Y)  
lm.score(X,Y)
```

Out[25]: 0.09535546506131365

```
In [26]: lm = LinearRegression()  
lm  
  
X = df[['sqft_basement']]  
Y = df['price']  
  
lm.fit(X,Y)  
lm.score(X,Y)
```

Out[26]: 0.104856815269744

```
In [27]: lm = LinearRegression()
lm

X = df[['view']]
Y = df['price']

lm.fit(X,Y)

lm.score(X,Y)
```

Out[27]: 0.15784211584121544

```
In [28]: lm = LinearRegression()
lm

X = df[['bathrooms']]
Y = df['price']

lm.fit(X,Y)

lm.score(X,Y)
```

Out[28]: 0.2763999306031437

```
In [29]: lm = LinearRegression()
lm

X = df[['sqft_living15']]
Y = df['price']

lm.fit(X,Y)

lm.score(X,Y)
```

Out[29]: 0.3426684607560172

```
In [30]: lm = LinearRegression()
lm

X = df[['sqft_above']]
Y = df['price']

lm.fit(X,Y)

lm.score(X,Y)
```

Out[30]: 0.3667117528382793

```
In [31]: lm = LinearRegression()
lm

X = df[['grade']]
Y = df['price']

lm.fit(X,Y)
```

```
lm.score(X,Y)
```

Out[31]: 0.4454684861092873

```
In [32]: lm = LinearRegression()
lm

X = df[['sqft_living']]
Y = df['price']

lm.fit(X,Y)
lm.score(X,Y)
```

Out[32]: 0.49285321790379316

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

In [ ]:

**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 [33]: Input=[('scale',StandardScaler()),('polynomial', PolynomialFeatures(include_bias=False))]
```

## Question 8

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

```
In [34]: pipe=Pipeline(Input)
pipe
```

```
Out[34]: Pipeline(memory=None,
                 steps=[('scale', StandardScaler(copy=True, with_mean=True, with_std=True)), ('polynomial', PolynomialFeatures(degree=2, include_bias=False, interaction_only=False)), ('model', LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False))])
```

```
In [36]: pipe.fit(X,Y)
```

```
/opt/conda/envs/Python36/lib/python3.6/site-packages/sklearn/preprocessing/data.py:645: DataConversionWarning: Data with input dtype int64 were all converted to float64 by StandardScaler.
    return self.partial_fit(X, y)
/opt/conda/envs/Python36/lib/python3.6/site-packages/sklearn/base.py:467: DataConversionWarning: Data with input dtype int64 were all converted to float64 by StandardScaler.
    return self.fit(X, y, **fit_params).transform(X)
```

```
Out[36]: Pipeline(memory=None,
                 steps=[('scale', StandardScaler(copy=True, with_mean=True, with_std=True)), ('polynomial', PolynomialFeatures(degree=2, include_bias=False, interaction_only=False)), ('model', LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False))])
```

```
In [37]: pipe.score(X,Y)
```

```
/opt/conda/envs/Python36/lib/python3.6/site-packages/sklearn/pipeline.py:511: DataConversionWarning: Data with input dtype int64 were all converted to float64 by StandardScaler.
    Xt = transform.transform(Xt)
```

```
Out[37]: 0.5327430940591443
```

## Module 5: MODEL EVALUATION AND REFINEMENT

import the necessary modules

```
In [38]: 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 set

```
In [39]: features = ["floors", "waterfront", "lat", "bedrooms", "sqft_basement", "view", "
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=42)
```

```
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, setting the regularization parameter to 0.1 and calculate the  $R^2$  using the test data.

```
In [40]: from sklearn.linear_model import Ridge
```

```
In [41]: pr=PolynomialFeatures(degree=2)  
x_train_pr=pr.fit_transform(x_train[['floors', 'waterfront','lat' ,'bedrooms' ,  
x_test_pr=pr.fit_transform(x_test[['floors', 'waterfront','lat' ,'bedrooms' , 's
```

```
In [43]: RidgeModel=Ridge(alpha=0.1)  
  
RidgeModel.fit(x_train_pr, y_train)
```

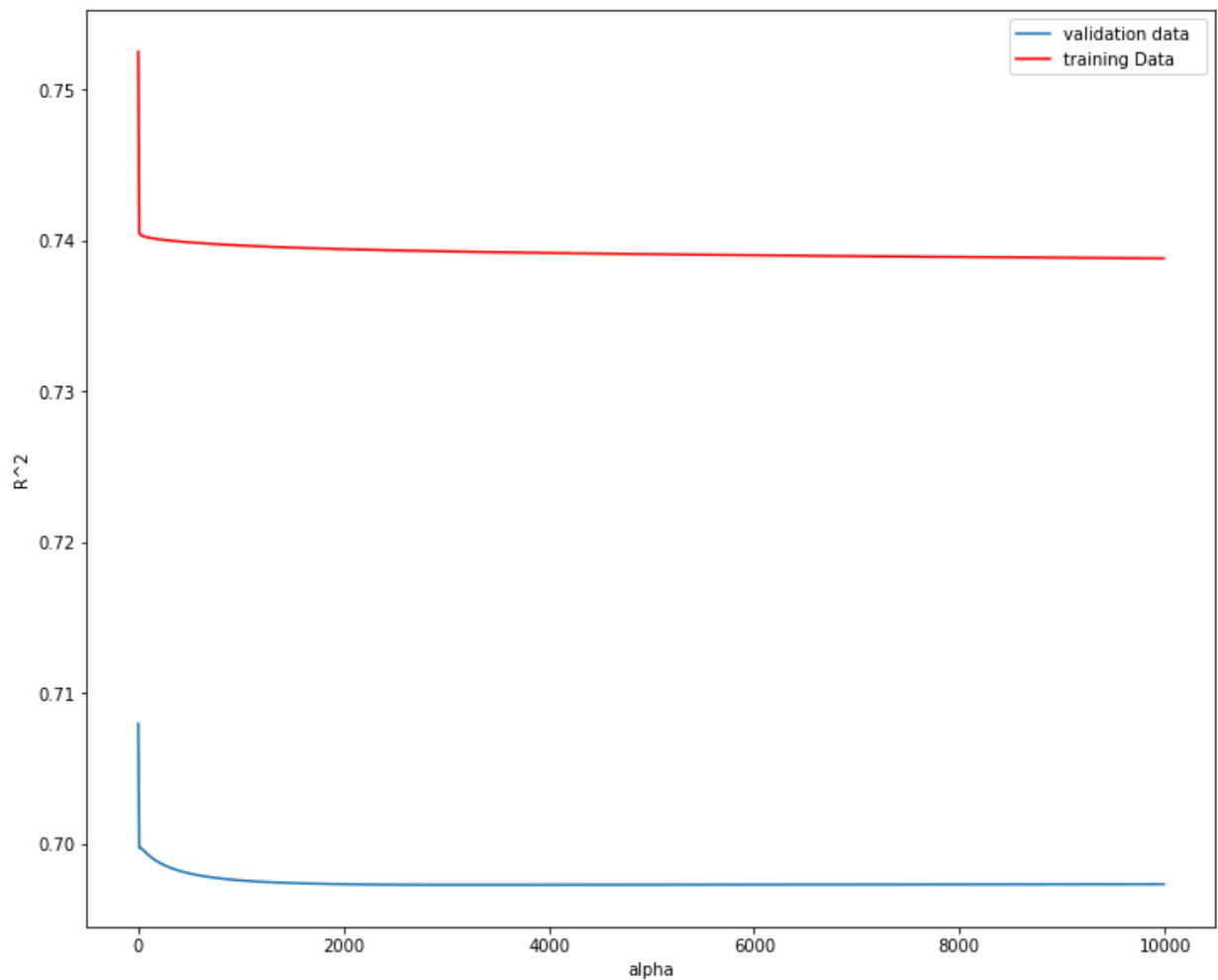
```
Out[43]: Ridge(alpha=0.1, copy_X=True, fit_intercept=True, max_iter=None,  
              normalize=False, random_state=None, solver='auto', tol=0.001)
```

```
In [50]: RidgeModel.score(x_train_pr, y_train)
```

```
Out[50]: 0.741816743868634
```

```
In [51]: width = 12  
height = 10  
plt.figure(figsize=(width, height))  
  
plt.plot(ALFA, Rsqu_test, label='validation data ' )  
plt.plot(ALFA, Rsqu_train, 'r', label='training Data ' )  
plt.xlabel('alpha')  
plt.ylabel('R^2')  
plt.legend()
```

```
Out[51]: <matplotlib.legend.Legend at 0x7f607850c898>
```



## 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, setting the regularisation parameter to 0.1. Calculate the  $R^2$  utilising the test data provided. Take a screenshot of your code and the  $R^2$ .

```
In [80]: from sklearn.preprocessing import PolynomialFeatures
```

```
In [81]: pr=PolynomialFeatures(degree=2)
pr
```

```
Out[81]: PolynomialFeatures(degree=2, include_bias=True, interaction_only=False)
```

```
In [82]: x_train_pr=pr.fit_transform(x_train[['floors', 'waterfront','lat' ,'bedrooms' ,
```

```
In [83]: x_polly=pr.fit_transform(x_train[['floors', 'waterfront','lat' ,'bedrooms' ,'sc
```

```
In [88]: RidgeModel=Ridge(alpha=0.1)

RidgeModel.fit(x_train_pr, y_train)

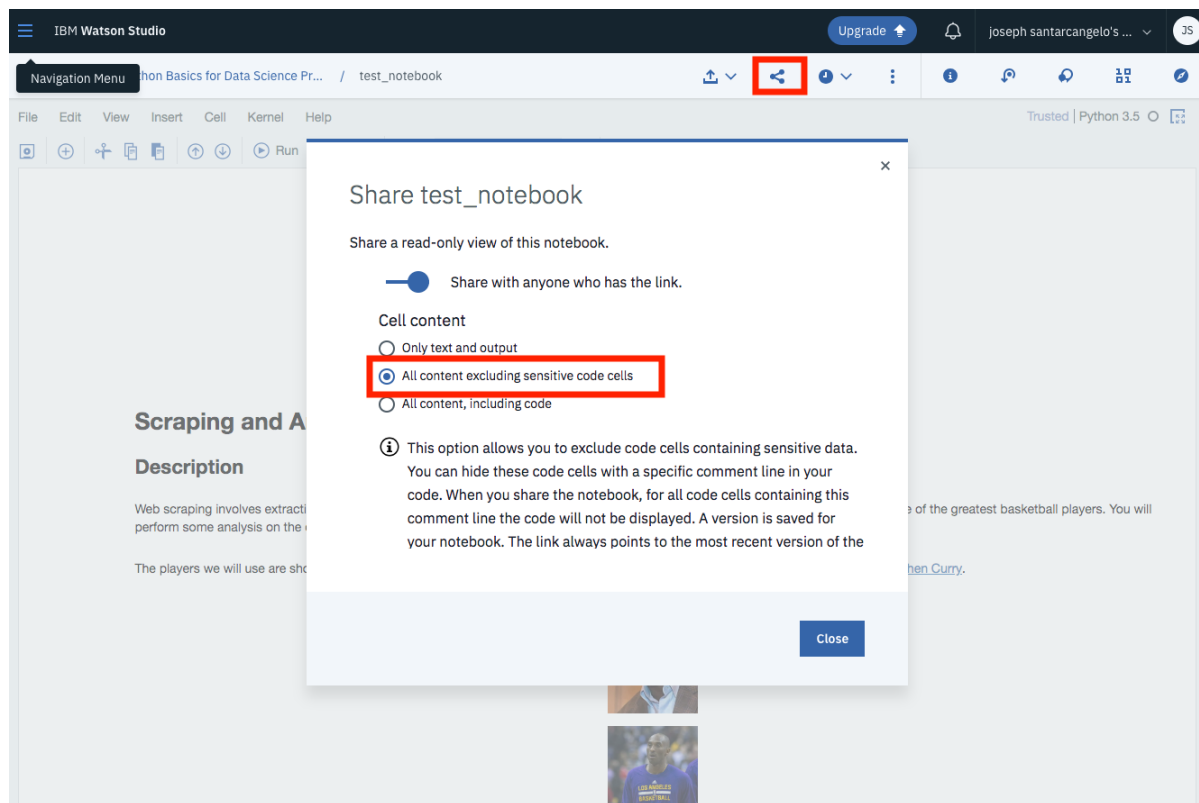
RidgeModel.score(x_train_pr, y_train)
```

Out [88]: 0.741816743868634

```
In [89]: x_test_pr=pr.fit_transform(x_test[['floors', 'waterfront','lat' ,'bedrooms' ,'sqft  
x_polly=pr.fit_transform(x_test[['floors', 'waterfront','lat' ,'bedrooms' ,'sqft  
RidgeModel=Ridge(alpha=0.1)  
RidgeModel.fit(x_test_pr, y_test)  
RidgeModel.score(x_test_pr, y_test)
```

Out [89]: 0.7666545737165752


Once you complete your notebook 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, 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:


×

## Share test\_notebook



 This option allows you to exclude code cells containing sensitive data. You can hide these code cells with a specific comment line in your code. When you share the notebook, for all code cells containing this comment line the code will not be displayed. A version is saved for your notebook. The link always points to the most recent version of the notebook.

Permalink to view notebook

<https://dataplatform.cloud.ibm.com/analytics/notebooks/v2/106a6db4->



Share on social media



Close

## About the Authors:

[Joseph Santarcangelo](#) 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 [ ]: