

**Skills Network Labs**

File Edit View Run Kernel Diagram Skills Network Tabs Settings Help

Launcher X House\_Sales\_in\_King\_County\_X launcher X kc\_house\_data\_HaN.csv X + Python (Pydide) ○

**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.

Variable	Description
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	Lots size area in 2015 (implies-- some renovations)

If you run the lab locally using Anaconda, you can load the correct library and versions by uncommenting the following:

```
[35]: # All libraries required for this lab are listed below. The libraries pre-installed on Skills Network Labs are commented.
# !nmaeo install -y pandas==1.1.4 numpy==1.21.4 seaborn==0.9.0 matplotlib==3.5.0 scikit-learn==0.20.2
# Note: If your environment doesn't support "nmaeo install", use "pip install"
```

```
[36]: # Supress warnings:
def warn(*args, **kwargs):
    pass
import warnings
warnings.warn = warn
```

You will require the following libraries:

```
[37]: import piplite
await piplite.install(['pandas','matplotlib','scikit-learn','seaborn', 'numpy'])
```

```
[38]: import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import warnings
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler, PolynomialFeatures
from sklearn.linear_model import LinearRegression
%matplotlib inline
```

## Module 1: Importing Data Sets

The functions below will download the dataset into your browser:

```
[39]: from pyodide.http import pyfetch
async def download(url, filename):
    response = await pyfetch(url)
    if response.status == 200:
        with open(filename, "wb") as f:
            f.write(await response.bytes())
```

```
[40]: file_name='https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/100DeveloperSkillsNetwork-DA0101EN-SkillNetwork/labs/FinalModule_Courses/data/kc_house_data_HaN.csv'
```

You will need to download the dataset; if you are running locally, please comment out the following code:

```
[41]: await download(file_name, "kc_house_data_HaN.csv")
file_name="kc_house_data_HaN.csv"
```

Use the Pandas method `read_csv()` to load the data from the web address.

```
[42]: df = pd.read_csv(file_name)
```

We use the method `head()` to display the first 5 columns of the data frame.

```
[43]: df.head()
```

Unnamed: 0	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	... grade	sqft_above	sqft_basement	yr_built	yr_renovated	zipcode	lat	long	sqft_living15	sqft_lot15	
0	0	712930050	20140101T000000	221900.0	3.0	1.00	1180	5550	1.0	0 ...	7	1180	0	1955	0	98178	47.5112	-122.2537	1340	5650
1	1	641410192	20141209T000000	538000.0	3.0	2.25	2570	7242	2.0	0 ...	7	2170	400	1951	1991	98125	47.7212	-122.319	1690	7639
2	2	5631500400	20150225T000000	1800000.0	2.0	1.00	770	10000	1.0	0 ...	6	770	0	1933	0	98028	47.7379	-122.233	2720	8062
3	3	2487200875	20141209T000000	604000.0	4.0	3.00	1960	5000	1.0	0 ...	7	1050	910	1965	0	98136	47.5208	-122.393	1380	5000
4	4	1954400510	20150218T000000	510000.0	3.0	2.00	1680	8080	1.0	0 ...	8	1680	0	1987	0	98074	47.6168	-122.045	1800	7503

5 rows × 22 columns

### Question 1

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

```
[44]: df.dtypes
```

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

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

```
[45]: df.describe()
```

Unnamed: 0	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	... grade	sqft_above	sqft_basement	yr_built	yr_renovated	zipcode	lat	long	sqft_living15	sqft_lot15		
count	21613.000000	2.161300e+04	2.161300e+04	21600.000000	21613.000000	2.161300e+04	21613.000000	21613.000000	21613.000000	... grade	21613.000000	21613.000000	21613.000000	21613.000000	21613.000000	21613.000000	21613.000000	21613.000000	21613.000000		
mean	10800.000000	4.50302e+09	5.40088e+05	3.727870	2.15736	2079.899736	1.510697e+04	1.494209	0.07542	... grade	7.656873	1788.390691	291.509045	1971.005136	84.402258	98077.939805	47.560033	-122.213896	1986.552492	1276.455652	
std	6239.2802	2.87656e+09	3.67127e+05	0.592657	0.768995	918.40937	4.14205e+04	0.539969	0.06517	... grade	0.766318	1.175459	828.099778	442.575043	29.373141	401.679240	53.505026	0.138544	0.140823	685.391304	273041.79631
min	0.000000	1.000102e+06	7.500000e+04	1.000000	0.500000	290.000000	5.200000e+02	1.000000	0.000000	... grade	0.000000	0.000000	190.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
25%	5403.000000	2.123049e+09	3.219500e+05	3.000000	1.750000	1427.000000	5.040000e+03	1.500000	0.000000	... grade	0.000000	7.000000	1150.000000	0.000000	0.000000	1951.000000	0.000000	0.000000	0.000000	0.000000	
50%	10806.000000	3.940492e+09	4.500000e+05	3.000000	2.250000	1910.000000	7.618000e+03	1.500000	0.000000	... grade	0.000000	7.000000	1560.000000	0.000000	0.000000	1975.000000	0.000000	0.000000	0.000000	0.000000	
75%	16209.000000	7.308950e+09	6.450000e+05	4.000000	2.500000	2500.000000	1.068000e+04	2.000000	0.000000	... grade	0.000000	8.000000	2210.000000	560.000000	1997.000000	0.000000	0.000000	9811.000000	47.678000		
max	21612.000000	9.800000e+09	7.700000e+05	33.000000	8.000000	13540.000000	1.651359e+05	3.500000	1.000000	... grade	4.000000	... grade	9410.000000	4820.000000	2015.000000	9819.000000	47.777600	0.000000	6210.000000	871200.000000	

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`.

```
[53]: df.drop(['id', 'Unnamed: 0'], axis=1, inplace=True)
df.describe()

      price    bedrooms   bathrooms    sqft_living    sqft_lot     floors    waterfront       view      condition        grade    sqft_above    sqft_basement    yr_built    yr_renovated      zipcode      lat        long    sqft_living15    sqft_lots15
count  2.161300e+04  2.1600.000000  2.1603.000000  2.1613.000000  2.1613.000000  2.1613.000000  2.1613.000000  2.1613.000000  2.1613.000000  2.1613.000000  2.1613.000000  2.1613.000000  2.1613.000000  2.1613.000000  2.1613.000000
mean  5.40081e+05  3.372870  2.115736  2.079.899736  1.510697e+04  1.494309  0.007542  0.23403  3.409430  7.658673  2.91.50945  1971.005136  48.402258  98077.939805  47.560053  -122.13896  1986.532496  12768.456562
std   3.67172e+05  0.926657  0.768996  7.91.400997  4.142051e+04  0.539989  0.065517  0.766318  0.650743  1.175459  828.09978  442.575043  29.373411  401.679240  53.505026  0.138564  0.140828  665.391309  27304.179631
min   2.500000e+04  1.000000  0.500000  2.90.000000  5.200000e+02  1.000000  0.000000  1.000000  0.000000  2.90.00000  0.000000  1900.00000  0.000000  98001.00000  47.155900  -122.519000  399.000000  651.000000
25%  3.219500e+05  3.000000  1.750000  1.427.000000  5.040000e+03  1.000000  0.000000  0.000000  3.000000  7.000000  1190.00000  0.000000  1931.00000  0.000000  98033.00000  47.471000  -122.328000  1490.000000  5100.00000
50%  4.300000e+05  3.000000  2.250000  1.910.000000  7.618000e+04  1.500000  0.000000  3.000000  7.000000  1560.00000  0.000000  1975.00000  0.000000  98065.00000  47.571800  -122.230000  1840.000000  7620.00000
75%  6.450000e+05  4.000000  2.550000  2.5000.000000  1.068800e+04  2.000000  0.000000  4.000000  8.000000  2210.00000  560.00000  1997.00000  0.000000  98118.00000  47.678000  -122.125000  2360.000000  10083.00000
max  7.700000e+06  33.000000  8.000000  13540.000000  1.651359e+06  3.500000  1.000000  4.000000  5.000000  13.00000  9410.00000  4820.00000  2015.00000  0.000000  98199.00000  47.777600  -121.315000  6210.00000  67120.00000
```

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

```
[54]: 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

[55]: 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

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

[57]: 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

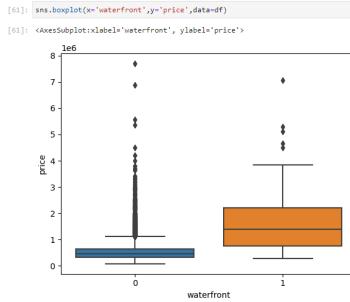
Use the method `value_counts` to count the number of houses with unique floor values, use the method `.to_frame()` to convert it to a dataframe.

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

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

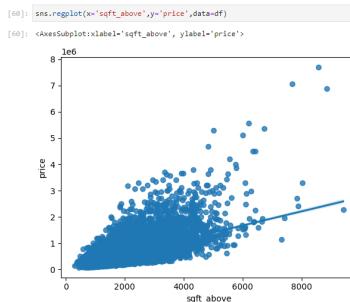
### Question 4

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.



### Question 5

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



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

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

      zipCode  -0.952509
long   0.021626
condition  0.036362
yr_built  0.054012
sqft_lots15  0.044747
sqft_above  0.080661
yr_renovated  0.126434
floors   0.256794
waterfront  0.300000
lat    0.307003
bedrooms  0.308797
sqft_basement  0.311141
view   0.397293
bathrooms  0.525738
sqft_living15  0.585379
sqft_above15  0.607260
grade   0.667424
sqft_living  0.702835
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.

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

[65]: 0.00046769493149007363
```

### 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.

```
[66]: X=df[['sqft_living']]
Y=df['price']
lm=LinearRegression()
lm.fit(X,Y)
lm.score(X,Y)
```

```
[65]: 0.4928552179037931
```

#### Question 7

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

```
[66]: 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.

```
[67]: X=df[features]
y=df['price']
lm=LinearRegression()
lm.fit(X,Y)
lm.score(X,Y)
```

```
[67]: 0.6576950029060081
```

#### 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()
```

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

```
[71]: pipe=Pipeline(Input)
pipe
pipe.fit(X,Y)
pipe.score(X,Y)
```

```
[71]: 0.7512786321941719
```

## Module 5: Model Evaluation and Refinement

Import the necessary modules:

```
[72]: 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:

```
[73]: 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: 16871
```

#### 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.

```
[74]: from sklearn.linear_model import Ridge
[75]: RidgeModel=Ridge(alpha=0.1)
RidgeModel.fit(X_train,Y_train)
RidgeModel.score(X_test,Y_test)
```

```
[75]: 0.647875916393905
```

#### 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.

```
[76]: from sklearn.preprocessing import PolynomialFeatures
from sklearn.linear_model import Ridge
pr = PolynomialFeatures(degree=2)
X_train_pr = pr.fit_transform(X_train)
X_test_pr = pr.fit_transform(X_test)
poly = Ridge(alpha=0.1)
poly.fit(X_train_pr, Y_train)
poly.score(X_test_pr, Y_test)
```

```
[76]: 0.7002744263398642
```

Once you complete your notebook, you can download the notebook. To download the notebook, navigate to **File** and click **Download**.

#### 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

#### Change Log

Date (YYYY-MM-DD)	Version	Changed By	Change Description
2020-12-01	2.2	Aje Eswarikide	Converted Data description from text to table
2020-10-06	2.1	Lakshmi Holla	Changed markdown instruction of Question1
2020-08-27	2.0	Malika Singla	Added lab to GitLab

© IBM Corporation 2020. All rights reserved.