

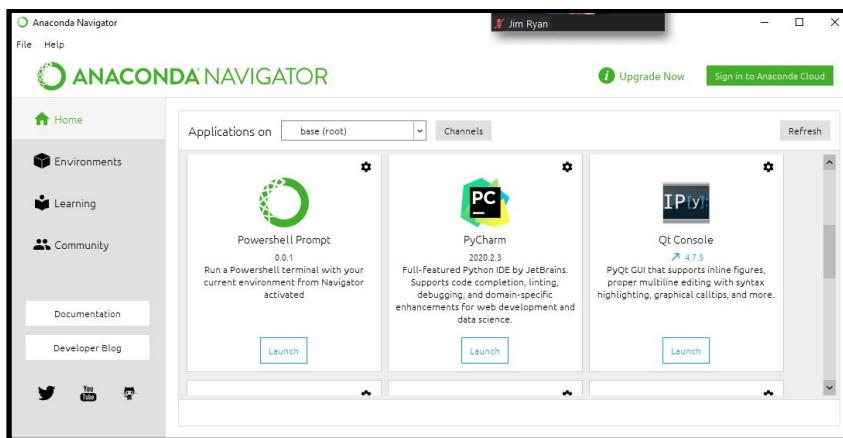
Data Mining Fundamentals

What is data mining?

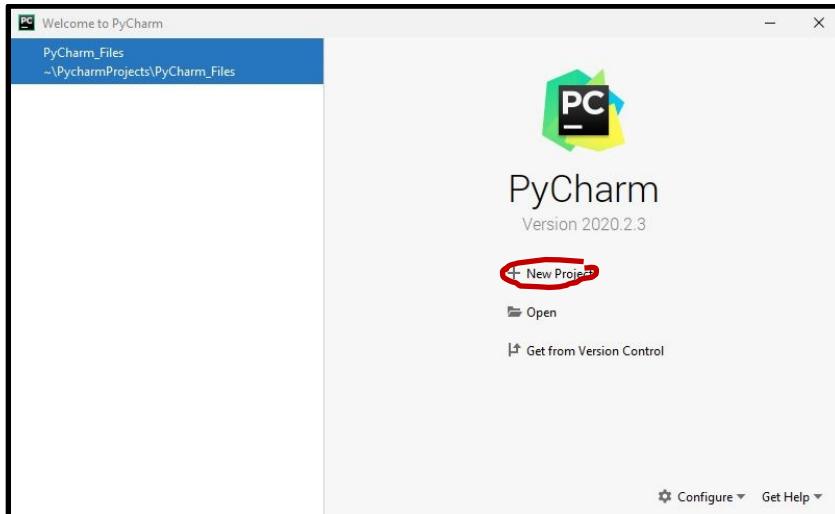
Data mining is a combination of several different processes, tools, or methodologies used to extract relevant data from a larger set of raw data. Data mining implies analyzing data patterns in large batches using one or more software applications. With respect to this tutorial, we will narrow our focus on Python data mining processes of cluster analysis, regression and classification, association and correlation, as well as dimensionality reduction. Python has built-in features to help in data mining methodologies for various data sets. For more on data mining not in this tutorial, refer to the URL: https://www.tutorialspoint.com/data_mining/index.htm

Create a Data Mining Project in PyCharm

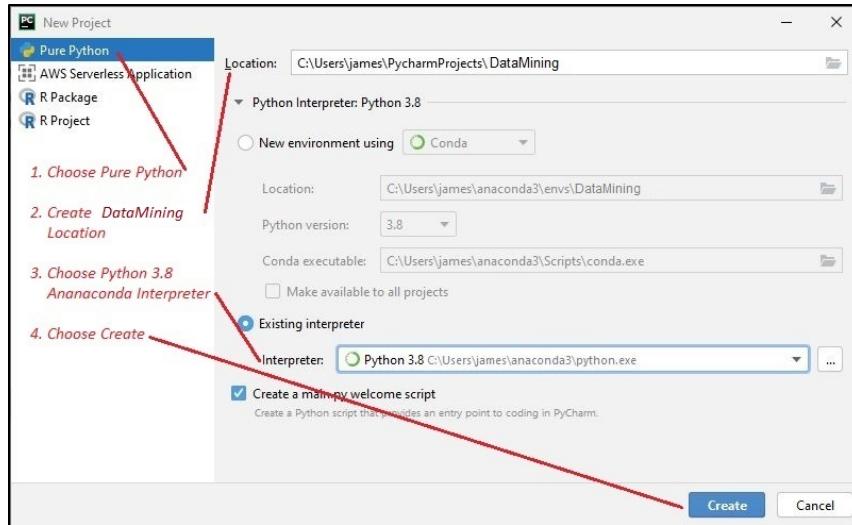
- 0) To ensure Python package libraries are available for import, use Anaconda Navigator to launch PyCharm



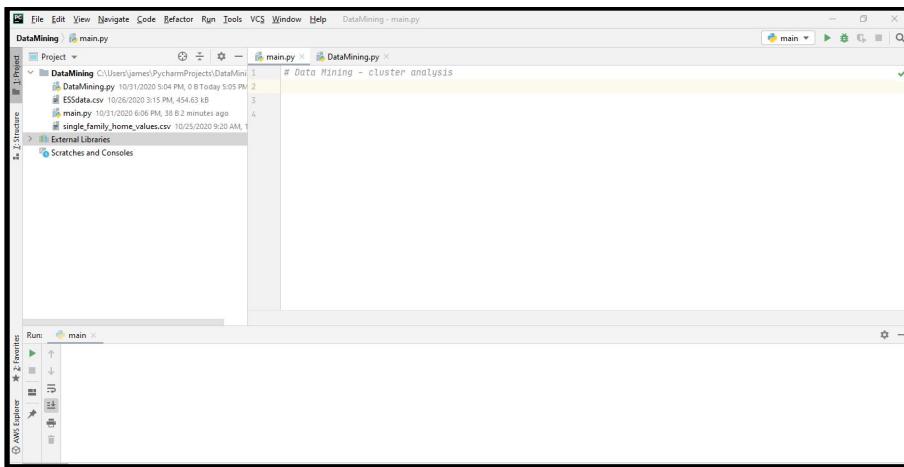
- 1) Open PyCharm and choose a New Project.



- 1) Choose a Pure Python project,
- 2) Change your project location to create a DataMining folder
- 3) Be sure to choose the Python interpreter as Python 3.8 associated with Anaconda.
- 4) Choose Create to open PyCharm with these settings.



- 5) Download the single_family_home_values.csv file from the Week 10 Canvas Module to the DataMining directory created in the last step. We used this CSV file during Week 09 Data Wrangling.



All visualizations may be viewed at 400% magnification for clarity.

Reading a CSV file into a Pandas DataFrame:

We learned about Pandas in Tutorial Seven. Again this tutorial, we will look at DataFrame or df. For more specific details about Pandas not covered in this tutorial, please refer to the following URL: https://www.tutorialspoint.com/python_data_science/python_pandas.htm

- 1) Anaconda will load the libraries in PyCharm for you. To access the four libraries we are using for Tutorial Eight (e.g., pandas, seaborn, matplotlib.pyplot, and KMeans), enter the following code into Main.py in the DataMining project directory of PyCharm:

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
```

- 2) To format the run window to handle the size of the csv data file we will be data mining in Tutorial Eight, enter the following Python code:

```
desired_width = 400
pd.set_option('display.width', desired_width)
pd.set_option('display.max_columns', 20)
```

- 3) To read the single_family_home_values.csv file into a pandas dataframe, type the following Python code into main.py via PyCharm in the line below the previous code above.

```
df = pd.read_csv(r'single_family_home_values.csv')      # reads Zillow file
```

- 4) Once the csv file is loaded into a dataframe (df), we can begin data mining. Any of the following Python functions may be used with a df, once it has been assigned.

```
print(df.head(2))      # returns top two rows
print(df.tail(5))      # returns bottom five rows
print(type(df))        # returns type of df
print(df.shape)         # returns row & column count
print(df.describe())    # returns data profiling metrics
print(df.info())        # returns column attributes
```



Cluster Analysis:

Cluster analysis is a type of unsupervised learning method and a common data mining technique for statistical data analysis. Cluster analysis is a task of dividing the data set of observations into subsets, called clusters, in such a way that observations in the same cluster are similar in one sense and dissimilar to observations in other clusters. The main goal of clustering is to group the data on the basis of similarity and dissimilarity. For cluster analysis in Tutorial Eight, we will use a well-known algorithm for data clustering the Scikit-learn package library called KMeans. KMeans takes the number of K-groups assigned, so it assumes that a known K of clusters exist. This is also called flat clustering. KMeans is an iterative algorithm, where data points are randomly assigned to each K-group and the K-groups' centroids in each iteration find optimized separation until the K-groups' centroids reach their optimal location separation. For more information and details on Cluster Analysis, refer to the following URLs:

https://www.tutorialspoint.com/artificial_intelligence_with_python/artificial_intelligence_with_python_unsupervised_learning_clustering.htm

https://www.tutorialspoint.com/data_mining/dm_cluster_analysis.htm

- 1) We imported pandas, numpy, seaborn, matplotlib.pyplot, and KMeans from sklearn.cluster. Looking at the PyCharm screenshot below, you will see how to (1) create a new data frame df2 that dropped a column (estimated_value) from a data frame, (2) slice the new df2 into a new data frame df3 that contains seven columns, (3) replace the df3 fields in columns with missing data with 0, (4) uses KMeans to create a new data frame k_groups that separates all rows of df3 data observations into 5 similar cluster centroids, (5) displays the k_groups' labels [0 to 4] for each of the top 3 rows and bottom 3 rows, and (6) compare the # rows in k_groups to the # of rows and columns in the reduced data frame df3.

```

1 # Data Mining - cluster analysis
2
3 import pandas as pd
4 import numpy as np
5 import seaborn as sns
6 import matplotlib.pyplot as plt
7 from sklearn.cluster import KMeans
8
9 desired_width = 400
10 pd.set_option('display.width', desired_width) # sets run screen width to 400
11 pd.set_option('display.max_columns', 20) # sets run screen column display to 20
12 df = pd.read_csv(r'single_family_home_values.csv') # reads Zillow file
13 (1) df2 = df.drop(['estimated_value', 'axis = 1']) # any data frame column can be dropped
14 (2) df3 = df2[['bedrooms', 'bathrooms', 'rooms', 'squareFootage', 'lotSize', 'yearBuilt', 'priorSaleAmount']] # reduced df
15 (3) df3.fillna(0, inplace=True) # replaces the NaN in priorSaleAmount with 0 -- may get a warning, but better than NaN
16 (4) k_groups = KMeans(n_clusters=5, random_state=0).fit(df3) # separates data set into 5 distinguishable groups
17 (5) print(k_groups.labels_) # displays k_groups' label (0 to 4) for each row
18 (6) print(len(k_groups.labels_), df3.shape) # displays rows in k_groups as well as rows, columns in df3
19
20
21
22
23

```

Process finished with exit code 0

- 2) In the following PyCharm screenshot, we (7) print out the cluster centroid coordinates for each of the seven k_groups columns. Each cluster centroid show coordinates respectively for bedroom, bathroom, rooms, squareFootage, lotSize, yearBuilt', and priorSaleAmount.

```

1 # Data Mining - cluster analysis
2
3 import pandas as pd
4 import numpy as np
5 import seaborn as sns
6 import matplotlib.pyplot as plt
7 from sklearn.cluster import KMeans
8
9 desired_width = 400
10 pd.set_option('display.width', desired_width) # sets run screen width to 400
11 pd.set_option('display.max_columns', 20) # sets run screen column display to 20
12 df = pd.read_csv(r'single_family_home_values.csv') # reads Zillow file
13 (1) df2 = df.drop(['estimated_value', 'axis = 1']) # any data frame column can be dropped
14 (2) df3 = df2[['bedrooms', 'bathrooms', 'rooms', 'squareFootage', 'lotSize', 'yearBuilt', 'priorSaleAmount']] # reduced df
15 (3) df3.fillna(0, inplace=True) # replaces the NaN in priorSaleAmount with 0 -- may get a warning, but better than NaN
16 (4) print(df3.head(2))
17 (5) k_groups = KMeans(n_clusters=5, random_state=0).fit(df3) # separates data set into 5 distinguishable groups
18 (6) print(k_groups.labels_)
19 (7) print(k_groups.cluster_centers_)
20 (8) print(k_groups.cluster_centers_[0])
21
22
23

```

Process finished with exit code 0

- 3) The PyCharm screenshot below identifies how to show the k_groups cluster centroid coordinates for k_groups cluster centroid [0].

```

1 # Data Mining - cluster analysis
2
3 import pandas as pd
4 import numpy as np
5 import seaborn as sns
6 import matplotlib.pyplot as plt
7 from sklearn.cluster import KMeans
8
9 desired_width = 400
10 pd.set_option('display.width', desired_width) # sets run screen width to 400
11 pd.set_option('display.max_columns', 20) # sets run screen column display to 20
12 df = pd.read_csv(r'single_family_home_values.csv') # reads Zillow file
13 (1) df2 = df.drop(['estimated_value', 'axis = 1']) # any data frame column can be dropped
14 (2) df3 = df2[['bedrooms', 'bathrooms', 'rooms', 'squareFootage', 'lotSize', 'yearBuilt', 'priorSaleAmount']] # reduced df
15 (3) df3.fillna(0, inplace=True) # replaces the NaN in priorSaleAmount with 0 -- may get a warning, but better than NaN
16 (4) print(df3.head(2))
17 (5) k_groups = KMeans(n_clusters=5, random_state=0).fit(df3) # separates data set into 5 distinguishable groups
18 (6) print(k_groups.labels_)
19 (7) print(len(k_groups.labels_), df3.shape)
20 (8) print(k_groups.cluster_centers_)
21 (9) print(k_groups.cluster_centers_[0])
22
23

```

Process finished with exit code 0

- 4) The PyCharm screenshot below identifies how to add the k_groups cluster centroid number as a column of the reduced data frame df3.

The screenshot shows a PyCharm interface with a code editor and a terminal window. The code in main.py performs the following steps:

- Imports pandas, numpy, seaborn, and matplotlib.pyplot.
- Reads a CSV file 'single_family_home_values.csv' into df.
- Drops the 'estimated_value' column from df.
- Reduces df to columns: 'bedrooms', 'bathrooms', 'rooms', 'squareFootage', 'lotSize', 'yearBuilt', and 'priorSaleAmount'.
- Fills NaN values in 'priorSaleAmount' with 0.
- Prints the first two rows of df3.
- Creates a KMeans object with 5 clusters and fits it to df3.
- Prints the cluster labels for each row.
- Prints the shape of df3.
- Adds a new column 'cluster' to df3 based on the cluster labels.
- Prints the top three rows of df3.

Annotations in red highlight specific parts of the code:

- 'print(df3.head(3))' is annotated with 'Prints out the top three rows of df3'.
- 'df3['cluster'] = k_groups.labels_' is annotated with 'Adds cluster as a column to df3'.

```

1 # Data Mining - cluster analysis
2 import pandas as pd
3 import numpy as np
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 from sklearn.cluster import KMeans
7 desired_width = 480
8 pd.set_option('display.width', desired_width) # sets run screen width to 480
9 pd.set_option('display.max_columns', 28) # sets run screen column display to 28
10 df = pd.read_csv(r'single_family_home_values.csv') # reads Zillow file
11 df2 = df.drop(['estimated_value'], axis = 1) # any data frame column can be dropped
12 df3 = df2[['bedrooms', 'bathrooms', 'rooms', 'squareFootage', 'lotSize', 'yearBuilt', 'priorSaleAmount']] # reduced df
13 df3.fillna(0, inplace=True) # replaces the NaN in priorSaleAmount with 0 -- may get a warning, but better than NaN
14 print(df3.head(2)) # prints top two rows of df3
15 k_groups = KMeans(n_clusters=5, random_state=0).fit(df3) # separates data set into 5 distinguishable groups
16 print(k_groups.labels_) # displays K-groups' label (0 to 4) for each row
17 print(len(k_groups.labels_),df3.shape) # displays rows in K-groups as well as rows, columns in df3
18 df3['cluster'] = k_groups.labels_ # add a new column to df3 called 'cluster', the k-group #
19 print(df3.head(3)) # display the top three rows of data frame df3

```

- 5) The PyCharm screenshot below depicts how to get k_groups means (i.e., df3.groupby('cluster').mean()) from data frame df3. You can substitute means with minimum (i.e., df3.groupby('cluster').min()), maximum (i.e., df3.groupby('cluster').max()), or median (i.e., df3.groupby('cluster').median()).

The screenshot shows a PyCharm interface with a code editor and a terminal window. The code in main.py performs the following steps:

- Imports pandas and KMeans.
- Reads a CSV file 'single_family_home_values.csv' into df.
- Drops the 'estimated_value' column from df.
- Reduces df to columns: 'bedrooms', 'bathrooms', 'rooms', 'squareFootage', 'lotSize', 'yearBuilt', and 'priorSaleAmount'.
- Fills NaN values in 'priorSaleAmount' with 0.
- Prints the first three rows of df3.
- Creates a KMeans object with 5 clusters and fits it to df3.
- Prints the cluster labels for each row.
- Prints the shape of df3.
- Prints the cluster centers.
- Prints the average values for each cluster center.
- Prints the cluster centers again.
- Adds a new column 'cluster' to df3 based on the cluster labels.
- Prints the top three rows of df3.
- Prints the mean values for each cluster.

Annotations in red highlight specific parts of the code:

- 'print(df3.groupby('cluster').mean())' is annotated with 'Mean(s)'.

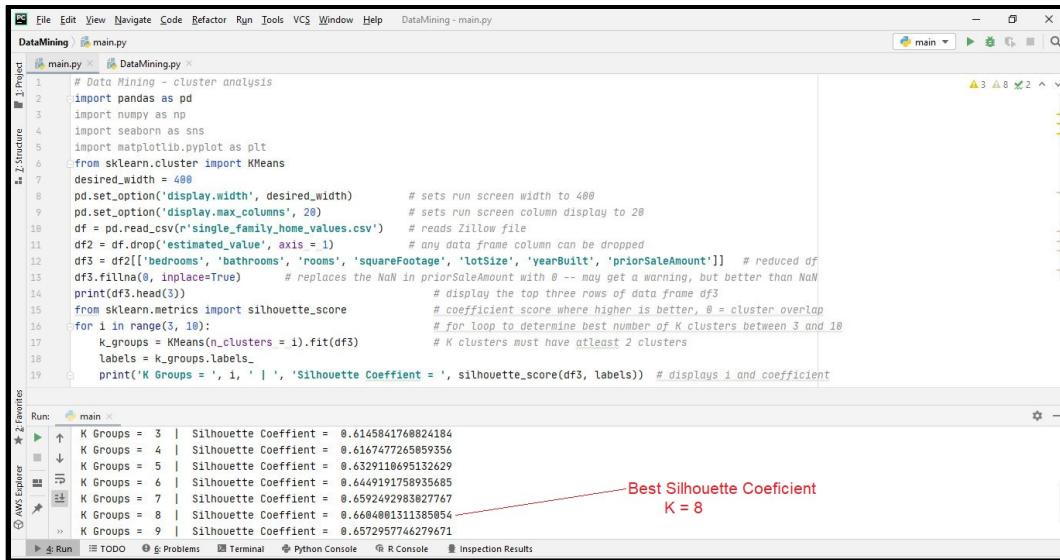
```

1 # Data Mining - cluster analysis
2 import pandas as pd
3 from sklearn.cluster import KMeans
4
5 desired_width = 480
6 pd.set_option('display.width', desired_width) # sets run screen width to 480
7 pd.set_option('display.max_columns', 28) # sets run screen column display to 28
8 df = pd.read_csv(r'single_family_home_values.csv') # reads Zillow file
9 df2 = df.drop(['estimated_value'], axis = 1) # any data frame column can be dropped
10 df3 = df2[['bedrooms', 'bathrooms', 'rooms', 'squareFootage', 'lotSize', 'yearBuilt', 'priorSaleAmount']] # reduced df
11 df3.fillna(0, inplace=True) # replaces the NaN in priorSaleAmount with 0 -- may get a warning, but better than NaN
12 print(df3.head(3)) # display the top three rows of data frame df3
13 k_groups = KMeans(n_clusters=5, random_state=0).fit(df3) # separates data set into 5 distinguishable groups
14 print(k_groups.labels_) # displays K-groups' label (0 to 4) for each row
15 print(len(k_groups.labels_),df3.shape) # displays rows in K-groups as well as rows, columns in df3
16 print(k_groups.cluster_centers_) # displays averages of the seven columns for each cluster centroid [0, 1, 2, 3, 4]
17 print(k_groups.cluster_centers_[0]) # displays averages for each of the seven columns in the cluster centroid [0]
18 df3['cluster'] = k_groups.labels_ # add a new column to df3 called 'cluster', the k-group #
19 print(df3.head(3)) # display the top three rows of data frame df3
20 print(df3.groupby('cluster').mean())

```

- 6) The PyCharm screenshot below depicts how to use a 'for loop' to calculate a silhouette_score library package from sklearn.metrics. The silhouette_score returns a silhouette coefficient (i.e., $-1 < \text{silhouette coefficient} < 1$) over all samples in the clusters. A silhouette coefficient near 1 is preferred and near 0 reflects cluster overlap. More information on silhouette_score from sklearn.metrics is available via the URL: https://scikit-learn.org/stable/modules/generated/sklearn.metrics.silhouette_score.html.

Note: In the screenshot below, df3 has dropped the 'cluster' column.



```

DataMining > main.py
File Edit View Navigate Code Refactor Run Tools VCS Window Help DataMining - main.py
main.py DataMining.py
Project Z-Structure
1 Project 1 main.py
1 main.py
# Data Mining - cluster analysis
2 import pandas as pd
3 import numpy as np
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6
7 from sklearn.cluster import KMeans
8 desired_width = 400
9 pd.set_option('display.width', desired_width) # sets run screen width to 400
10 pd.set_option('display.max_columns', 20) # sets run screen column display to 20
11 df = pd.read_csv('single_family_home_values.csv') # reads Zillow file
12 df2 = df.drop(['estimated_value', 'axis = 1']) # any data frame column can be dropped
13 df3 = df2[['bedrooms', 'bathrooms', 'rooms', 'squareFootage', 'lotSize', 'yearBuilt', 'priorSaleAmount']] # reduced df
14 df3.fillna(0, inplace=True) # replaces the NaN in priorSaleAmount with 0 -- may get a warning, but better than NaN
15 print(df3.head(3)) # display the top three rows of data frame df3
16 from sklearn.metrics import silhouette_score
17 for i in range(3, 10):
18     k_groups = KMeans(n_clusters = i).fit(df3) # coefficient score where higher is better, 0 = cluster overlap
19     labels = k_groups.labels_
20     print('K Groups = ', i, ' | ', silhouette_score(df3, labels)) # displays i and coefficient
21
22
Run: main
K Groups = 3 | Silhouette Coefficient =  0.6145841768824184
K Groups = 4 | Silhouette Coefficient =  0.6167477265895356
K Groups = 5 | Silhouette Coefficient =  0.6329118695132629
K Groups = 6 | Silhouette Coefficient =  0.6449191758935685
K Groups = 7 | Silhouette Coefficient =  0.6592492983827767
K Groups = 8 | Silhouette Coefficient =  0.6604001311385054
K Groups = 9 | Silhouette Coefficient =  0.6572957746279671

```

Best Silhouette Coefficient
K = 8

The DataMining.py file has the Python code for all screenshots in this Tutorial.

- 5) **Exercise deliverable:** Construct an example of a cluster analysis in Python code. Use either the single_family_home_values.csv or your data project csv. Determine the correct number of K clusters to use (i.e., see previous PyCharm screenshot in step 6) and then let Python calculate the means of all the columns separated into K clusters (i.e., steps 1 to 5). Capture two PyCharm screenshots for use later: (1) one of preferred K clusters due to silhouette coefficient and (2) a listing of dataframe.groupby('clusters').mean() for the preferred K clusters.

Classification & Regression:

Data mining to discover patterns in data that lead to actionable insights can be classified into two groups based on how the algorithms “learn” about data to make predictions: supervised and unsupervised learning. Practical machine learning uses supervised learning. Supervised learning is where you have input variables (x) and an output variable (Y) and you use an algorithm to learn the mapping function from the input to the output $Y = f(X)$. The goal is to approximate the mapping function so well that when you have new input data (x) that you can predict the output variables (Y) for that data. Supervised learning requires that the data used to train the algorithm is already labeled with correct answers. For example, a classification algorithm will learn to identify animals after being trained on a dataset of images that are properly labeled with the species of the animal and some identifying characteristics. Supervised learning problems can be further grouped into Regression and Classification problems. Both problems have as goal the construction of a succinct model that can predict the value of the dependent attribute from the attribute variables. The difference between the two tasks is the fact that the dependent attribute is numerical for regression (i.e., linear regression) and categorical for classification (i.e., logistic regression). For more discussion on the differentiation between classification and regression, refer to the URL: <https://machinelearningmastery.com/classification-versus-regression-in-machine-learning/>

[learning/#:~:text=Fundamentally%2C%20classification%20is%20about%20predicting,is%20about%20predicting%20a%20quantity.&text=That%20classification%20is%20the%20problem,quantity%20output%20for%20an%20example.](#)

- 1) The following Python commands establish the PyCharm libraries and variables we will use in understanding classification and regression:

```
# Data Mining - classification & regression
import pandas as pd
from sklearn.linear_model import LinearRegression, LogisticRegression
from sklearn.svm import SVC, SVR
from sklearn.neighbors import KNeighborsClassifier, KNeighborsRegressor
from sklearn.model_selection import train_test_split
desired_width = 400
pd.set_option('display.width', desired_width)      # sets run screen width to 400
pd.set_option('display.max_columns', 20)           # sets run screen column display to 20
df = pd.read_csv(r'single_family_home_values.csv')    # reads Zillow file
df.fillna(0, inplace = True)                      # replaces the NaN with 0
X = df[['bedrooms', 'bathrooms', 'rooms', 'squareFootage', 'lotSize', 'yearBuilt',
'priorSaleAmount']] # reduced df as X
y = df.estimated_value
```

Several of the library packages from sklearn have classifiers (i.e., categorical variables) and regressors (i.e., numeric variables) listed. Logistic regression is a classifier that has regression in the name. sklearn expects X to be a feature matrix (Pandas Dataframe) and y to be a response vector (Pandas Series). The train_test_split function randomly divides a data set into an 80% train and 20% test data sets. The following PyCharm screenshots demonstrate the Python code above using LinearRegression() as a regressor model and LogisticRegression() as a classifier model. Note a categorical variable was created for estimated_value in the classifier model. The DataMining.py file has the Python code for all screenshots in this Tutorial.

Python Package Imports / X,y

```
File Edit View Navigate Code Refactor Run Tools VCS Window Help DataMining - main.py
DataMining  main.py DataMining.py
1 # Data Mining - classification & regression
2 import pandas as pd
3 from sklearn.linear_model import LinearRegression, LogisticRegression
4 from sklearn.svm import SVC, SVR
5 from sklearn.neighbors import KNeighborsClassifier, KNeighborsRegressor
6 desired_width = 400
7 pd.set_option('display.width', desired_width)      # sets run screen width to 400
8 pd.set_option('display.max_columns', 20)           # sets run screen column display to 20
9 df = pd.read_csv(r'single_family_home_values.csv')    # reads Zillow file
10 df.fillna(0, inplace = True)                      # replaces the NaN with 0
11 X = df[['bedrooms', 'bathrooms', 'rooms', 'squareFootage', 'lotSize', 'yearBuilt',
12 'priorSaleAmount']] # reduced df as X
13 y = df.estimated_value
14 print(X.info())
15
```

LinearRegression() X,y



```
# Data Mining - classification & regression
import pandas as pd
from sklearn.linear_model import LinearRegression, LogisticRegression
from sklearn.svm import SVC, SVR
from sklearn.neighbors import KNeighborsClassifier, KNeighborsRegressor
from sklearn.model_selection import train_test_split
df = pd.read_csv('family_home_values.csv')
df.fillna(0, inplace = True)
X = df[['bedrooms', 'bathrooms', 'rooms', 'squareFootage', 'lotSize', 'yearBuilt', 'priorSaleAmount']] # reduced df as upper case X matrix
y = df.estimated_value # predictor variable lower case y as array
lg = LinearRegression()
lg.fit(X,y)
print(lg.score(X,y))
print(lg.score(X,y))
X_train, X_test, y_train, y_test = train_test_split(X,y)
print(X_train.shape, y_train.shape, X_test.shape, y_test.shape) # we used default settings 80% is train to 20% test
print(lg.fit(X_train, y_train))
print(lg.score(X_test, y_test))
print(lg.score(X_test, y_test))
```

Run: main
C:\Users\james\Anaconda3\python.exe C:/Users/james/PycharmProjects/DataMining/main.py

Coefficient of Determination R2: -0.1 <= R2 <= 1.0
R2 of all 11,250 data points
X_train, y_train, X_test, y_test data sets from train_test_split()
R2 of 3,500 test data points using regression model of 11,250 train data points

LogisticRegression() X,y2

```
File Edit View Navigate Code Refactor Run Tools VCS Window Help DataMining - main.py
DataMining main.py
main.py DataMining x
1 from sklearn.neighbors import KNeighborsClassifier, KNeighborsRegressor
2 from sklearn.model_selection import train_test_split
3 desired_width = 400
4 pd.set_option('display.width', desired_width) # sets run screen width to 400
5 pd.set_option('display.max_columns', 20) # sets max columns displayed to 20
6 df = pd.read_csv('small_homes.csv') # loads Zillow file
7 df.fillna(0, inplace = True) # replaces the NaN with 0 to have even 15,880 in all 7 variables
8 X = df[['bedrooms', 'bathrooms', 'rooms', 'squareFootage', 'lotSize', 'yearBuilt', 'priorSaleAmount']] # reduced of as upper case X matrix
9 y = df.estimated_value # predictor variable lower case y as array
10 df['estimated_value_class'] = df.estimated_value.apply(lambda x: 'low' if x < 58800 else 'high')
11 print(df.estimated_value_class.value_counts()) # displays distribution of estimated_value_class
12 y2 = df.estimated_value_class # assigns y2 as cat variable estimated_value_class
13 log = LogisticRegression() # assigning alias lg to LogisticRegression() function
14 print(log.fit(X,y))
15 print(log.score(X,y))
16 X_train, y_train, X_test, y_test = train_test_split(X,y, test_size=0.2) # test size 0.2 means 20% of data
17 print(X_train.shape, y_train.shape, X_test.shape, y2_test.shape) # we used default settings 80% is train to 20% test
18 print(log.fit(X_train, y2_train)) # Using 11,250 data points to train model.
19 print(log.score(X_test, y2_test)) # Using 1,750 data points to test/evaluate R2 of model
20 print(log.coef_)

Run main
Name: estimated_value_class, dtype: int64
Coefficient of Determination R2
R2 of all 15,000 data observations
R2 of 3,570 data observations in the model trained by 11,250 data observations
Categorical variable y2 = df.estimated_value_class count of high values vs. low values
```

Enzyme expression on classification

- 2) **Prediction Check:** When developing regression or classification models, besides using R², you can use a Python confusion matrix on classifications to check what predictions would be and the associated Type I and Type II errors. For regressions, you can evaluate metrics (e.g., mean absolute error, mean squared error, and root mean squared error). The following two PyCharm screenshots show comparison metrics and a confusion matrix for regression and classification, respectively.

Regression – Linear Regression()

The screenshot displays two PyCharm code editors side-by-side.

Top Editor (Linear Regression):

```

1 Project DataMining main.py
2 Structure 1 main.py
3 main.py
4 import numpy as np
5 from sklearn.linear_model import LinearRegression, LogisticRegression
6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import mean_squared_error
8 from sklearn.metrics import mean_absolute_error
9 desired_width = 400
10 pd.set_option('display.width', desired_width) # sets run screen width to 400
11 pd.set_option('display.max_columns', 20) # sets run screen column display to 20
12 df = pd.read_csv('single_family_home_values.csv') # reads Zillow file
13 df.fillna(0, inplace = True) # replaces the NaN with 0 to have even 15,000 in all 7 variables
14 X = df[['bedrooms', 'bathrooms', 'rooms', 'squareFootage', 'lotSize', 'yearBuilt', 'priorSaleAmount']] # reduced df as upper case X matrix
15 y = df.estimated_value # predictor variable lower case y as array
16 lg = LinearRegression() # assigning alias lg to LinearRegression() function
17 X_train, X_test, y_train, y_test = train_test_split(X,y) # randomly split X,y data to 2 X,y (train,test) sets
18 lg.fit(X_train, y_train) # use 11,250 data points to train model
19 print(lg.score(X_test, y_test)) # displays R2 measure from 3,750 data points used in 11,250 trained model
20 y_prime = lg.predict(X_test) # assigns y_pred to prediction of X_test data
21 print(y_prime, np.array(y_test)) # displays y_pred versus y of test data set
22 print("MAE =", mean_absolute_error(y_test, y_prime)) # displays mean absolute error of actual vs. predicted
23 print("MSE =", mean_squared_error(y_test, y_prime)) # displays mean squared error of actual vs. predicted
24 print("RMSE =", np.sqrt(mean_squared_error(y_test, y_prime))) # displays root mean squared error of actual vs. predicted

```

Bottom Editor (Logistic Regression):

```

1 Project DataMining main.py
2 Structure 1 main.py
3 main.py
4 # Data Mining - Classification & regression - Prediction (Logistic)
5 import pandas as pd
6 import numpy as np
7 from sklearn.linear_model import LinearRegression, LogisticRegression
8 from sklearn.model_selection import train_test_split
9 from sklearn.metrics import confusion_matrix
10 desired_width = 400
11 pd.set_option('display.width', desired_width) # sets run screen width to 400
12 pd.set_option('display.max_columns', 20) # sets run screen column display to 20
13 df = pd.read_csv('single_family_home_values.csv') # reads Zillow file
14 df.fillna(0, inplace = True) # replaces the NaN with 0 to have even 15,000 in all 7 variables
15 X = df[['bedrooms', 'bathrooms', 'rooms', 'squareFootage', 'lotSize', 'yearBuilt', 'priorSaleAmount']] # reduced df as upper case X matrix
16 y = df.estimated_value # predictor variable lower case y as array
17 df['estimated_value_class'] = df.estimated_value.apply(lambda x: 'low' if x < 500000 else 'high')
18 y2 = df.estimated_value_class # assigns y2 as cat variable estimated_value_class
19 log = LogisticRegression() # assigning alias lg to LogisticRegression() function
20 X_train, X_test, y2_train, y2_test = train_test_split(X,y2) # randomly split X,y2 data to 2 X,y2 (train,test) sets
21 # we used default settings 80% is train to 20% test
22 log.fit(X_train, y2_train) # use 11,250 data points to train model
23 log.score(X_test, y2_test) # measures R2 from 3,750 data points used in 11,250 trained model
24 y2_pred = log.predict(X_test) # assigns y2_pred to prediction of X_test data
25 print(y2_pred, np.array(y2_test)) # displays y2_pred versus y2 of test data set
26 print(confusion_matrix(y2_test, y2_pred)) # displays correct on diagonal, typeI right, typeII left, all 3,750

```

The bottom editor's Run tab shows the output of the confusion matrix:

```

[[ 'low' 'high' 'high' ... 'low' 'high' 'low'] [ 'low' 'high' 'high' ... 'low' 'high' 'high']
 [[1543 394]--[high predicted vs. high test data , high predicted vs. low test data - typeI error]
 [ 287 1526]]--[low predicted vs. high test data - type II error , low predicted vs. low test data ]]

```

The bottom editor also includes a note: "Confusion Matrix - all numbers add-up to X_test n = 3,750".

- 3) **Support Vector Machine - Regression (SVR) / Classification (SVC):** SVR.score() or SVC.score() return the coefficient of determination R2 for the regression or classification model. The same prediction check on classification can be run using the confusion matrix. For more detail on SVR, please refer to the following URL: <https://scikit-learn.org/stable/modules/svm.html#svm-regression> The following PyCharm screenshot depict Python code using SVR, SVC, and a predict check on classification.

The screenshot shows the PyCharm IDE interface with the following details:

- Project:** DataMining
- File:** main.py
- Code Content:** The code performs various data processing and machine learning tasks. It includes:
 - Importing `train_test_split` from `sklearn.model_selection`.
 - Importing `confusion_matrix` from `sklearn.metrics`.
 - Setting screen width to 400.
 - Setting column display to 28.
 - Reading a CSV file named 'single_family_home_values.csv'.
 - Filling NaN values with 0.
 - Replacing 'NaN' with '0' in all 7 variables.
 - Creating X and y variables.
 - Creating X1 and X2.
 - Predicting y based on X1.
 - Assigning 'estimated_value_class' based on X2.
 - Splitting data into X_train, X_test, y_train, y_test.
 - Splitting X into X1_train, X1_test, X2_train, X2_test.
 - Splitting X1 into X1_train, X1_test.
 - Creating a SVR model (svr_reg).
 - Creating an SVC model (svc_class).
 - Fitting SVR to X_train, y_train.
 - Fitting SVC to X1_train, y1_train.
 - Scoring SVR on X_test, y_test.
 - Scoring SVC on X1_test, y2_test.
 - Predicting y2 using SVC.
 - Printing confusion matrix for y2.
- AWS Explorer Results:** A red box highlights the following results:
 - Poor R² on the SVR() regression model
 - Good R² on the SVC() classification model
 - Type I and Type II errors too high in the confusion matrix
- Run:** Run main
- Output:** svr_score = -0.055934722658349534
svr_score = 0.7168
[[654 1332]
 [565 1199]]
- Event Log:** Packages installed successfully: Installed packages: scikit-learn (yesterday, 149 PM)

- 4) **K Nearest Neighbors (KNN) – Regressor & Classifier:** Provides functionality for unsupervised and supervised neighbors-based learning methods. Supervised neighbors-based learning comes in two flavors: classification for data with discrete labels, and regression for data with continuous labels. The principle behind nearest neighbor methods is to find a predefined number of training samples closest in distance to the new point and predict the label from these. Being a non-parametric method, it is often successful in classification situations where the decision boundary is very irregular. The classes in `sklearn.neighbors` can handle either NumPy arrays or `scipy.sparse` matrices as input. A large number of possible distance metrics are supported for dense matrices, but fewer for sparse matrices. For more information on KNN, please refer to the following URL: <https://scikit-learn.org/stable/modules/neighbors.html#unsupervised-neighbors>

The following PyCharm screenshot demonstrates a KNN regression, classification, and classification prediction check example using Python code.

The screenshot shows the PyCharm IDE interface with the following details:

- Project:** DataMining
- File:** DataMining.py
- Code Content:** The script imports necessary libraries from sklearn, performs data loading and splitting, and trains two KNN models (regressor and classifier) on a dataset of house prices. It prints the R-squared values for both models.

```
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix
desired_width = 400
pd.set_option('display.width', desired_width) # sets run screen width to 400
pd.set_option('display.max_columns', 20) # sets run screen column display to 20
df = pd.read_csv('single_family_home_values.csv') # reads Zillow file
df.fillna(0, inplace=True) # replaces the Nahn with 0 to have even 15,000 in all 7 variables
X = df[['bedrooms', 'bathrooms', 'rooms', 'squareFootage', 'lotSize', 'yearBuilt', 'priorSaleAmount']] # reduced df as upper case X matrix
X1 = X # duplicate of X
y = df.estimated_value # predictor variable lower case y as array
df['estimated_value_class'] = df.estimated_value.apply(lambda x: 'low' if x < 500000 else 'high')
y2 = df.estimated_value_class # assigns y2 as cat variable estimated_value_class
X_train, X_test, y_train, y_test = train_test_split(X,y) # randomly split X,y data to 2 X,y (train,test) sets
X1_train, X1_test, y2_train, y2_test = train_test_split(X1,y2) # randomly split X,y data to 2 X,y (train,test) sets
knn_reg = KNeighborsRegressor() # assign knn_reg to the KNeighborsRegressor() function
knn_class = KNeighborsClassifier() # assign knn_class to the KNeighborsClassifier function
knn_reg.fit(X_train, y_train) # fit a knn_reg model using 11,250 data points
print("knn_reg score = ", knn_reg.score(X_test, y_test)) # score the knn_reg model based on 11,250 data points using the 3,750 data points
knn_class.fit(X1_train, y2_train) # fit a knn_class model using 11,250 data points
print("knn_class score = ", knn_class.score(X1_test, y2_test)) # score the knn_class model base on 11,250 data points using the 3,750 data points
y2_pred = knn_class.predict(X_test) # assigns y2_pred to knn_class prediction of X_test data
print(confusion_matrix(y2_test, y2_pred)) # displays correct on diagonal, typeI right, typeII left, all 3,750
```

- Run Tab:** Shows the main module is selected.
- R2 Scores:** The output shows the R-squared values for both models:
 - R2 of knn_reg model: 0.6865855366517732
 - R2 of knn_class model: 0.832
- Console Output:** Displays the confusion matrix:

| Actual | Predicted | Low | High |
|--------|-----------|-----|------|
| Low | 994 | 966 | |
| High | 917 | 873 | |
- Bottom Status Bar:** Shows the event log, Python version, and other system information.

- 5) **Exercise deliverable:** Construct a regression, classification, and prediction check example using Python code. Use either the single_family_home_values.csv or your data project csv. If you are using the csv above, then mix-up the columns used for X. (1) Choose from Linear/Logistic, SVR/SVC, or K Nearest Neighbor methods of regression and classification. (2) Be sure to use the train_test_split() function in both regression and classification. (3) Choose either of the regression or classification example for the prediction check. (4) Capture PyCharm screenshots for each of the three examples (e.g., regressor, classifier, & prediction check) to use later in this tutorial.

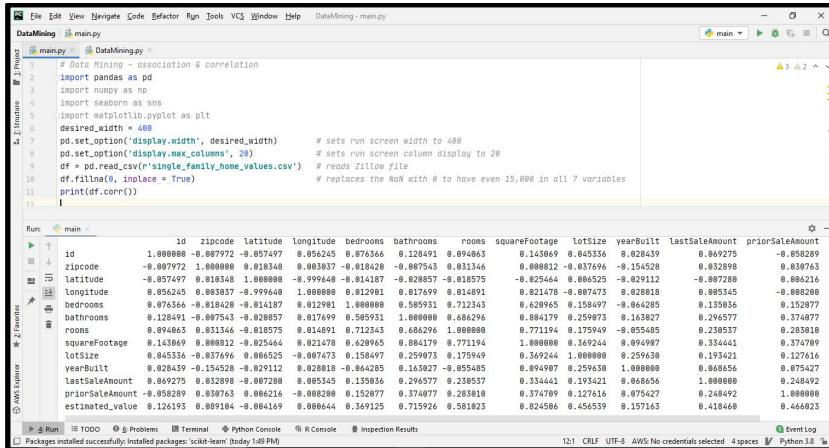
Association & Correlation:

Association is measured by coefficients of correlation or determination. As a measure of association, the coefficient of correlation (i.e., r) is part of a covariance matrix, where correlation measures if a relationship between two variables is direct (i.e., positive) or inverse (i.e., negative). A coefficient of correlation falls between $-1 \leq r \leq 1$. A measure close to zero reflects no association, a measure near -1 reflects an inverse relationship, and a measure near 1 reflects a direct relationship. Coefficient of determination (i.e., r^2) reflects the strength of association between variables that explain variance, where near 0 (i.e., none) reflects little explanation while near to 1 (i.e., strong) reflects stronger explanation. Correlation is part of the covariance matrix, dispersion matrix, variance matrix, or variance-covariance matrix. It is a matrix in which i-j position defines the correlation between the ith and jth parameter of the given data-set.

When the data points follow a roughly straight-line trend, the variables are said to have an approximately linear or direct relationship. In some cases, the data points fall close to a straight line, but more often there is quite a bit of variability of the points around the straight-line trend. Correlation describes the strength of the linear association. Correlation summarizes the strength and direction of the linear (straight-line) association between two quantitative variables. Denoted by r, it takes values between -1 and +1. A positive value for r indicates a positive association, and a negative value for r indicates a negative association.

The closer r is to 1 the closer the data points fall to a straight line, thus, the linear association is stronger. The closer r is to 0, making the linear association weaker.

- 1) The Python code using the df.corr() function in the PyCharm screenshot below, provides a simple correlation matrix between all non-string column attributes:

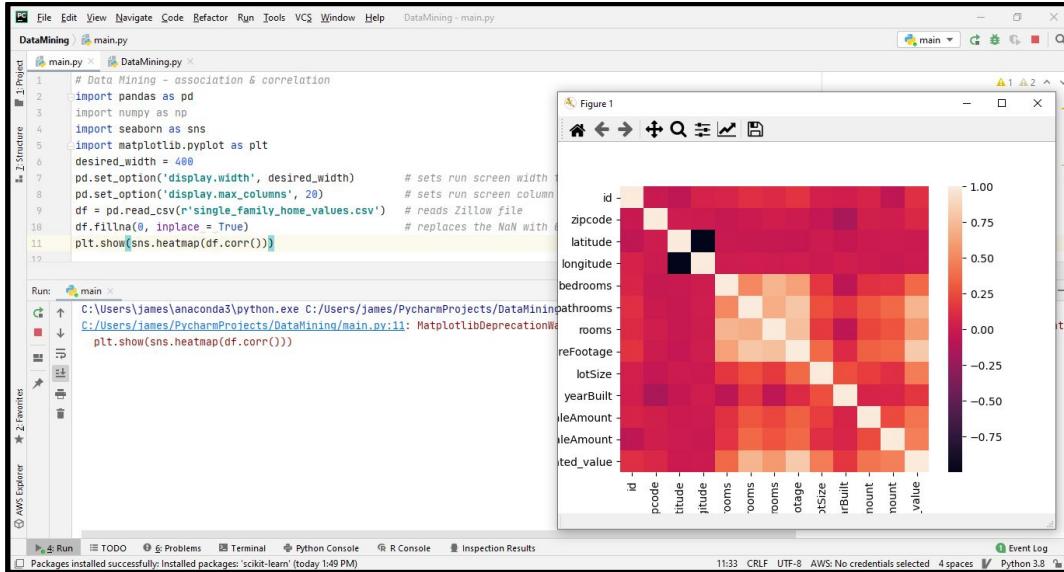


```

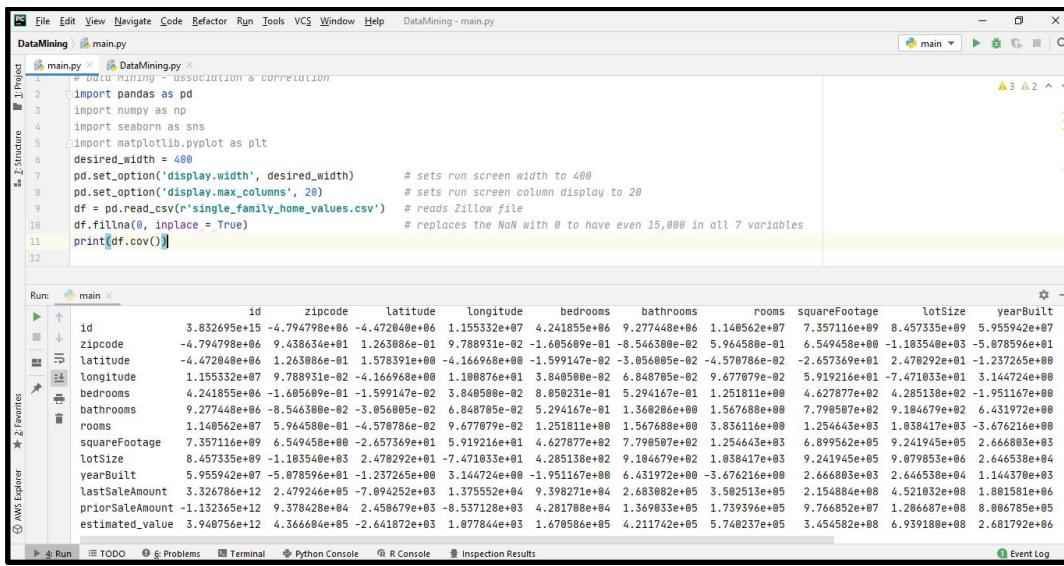
File Edit View Navigate Code Refactor Run Tools VCS Window Help DataMining - main.py
DataMining main.py DataMining
1 # Data Mining - association & correlation
2 import pandas as pd
3 import numpy as np
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 desired_width = 480
7 pd.set_option('display.width', desired_width) # sets run screen width to 480
8 pd.set_option('display.max_columns', 20) # sets run screen column display to 20
9 df = pd.read_csv("single_family_home_values.csv") # reads Zillow file
10 df.fillna(0, inplace = True) # replaces the NaN with 0 to have even 15,000 in all 15 variables
11 print(df.corr())
12
Run main
      id      zipcode      latitude      longitude      bedrooms      bathrooms      rooms      squareFootage      lotSize      yearBuilt      lastSaleAmount      priorSaleAmount
id  1.000000  -0.089792  -0.057497  0.056245  0.076366  0.128491  0.094863  0.143669  0.045356  0.028439  0.069275  -0.058289
zipcode  -0.089792  1.000000  0.018340  0.003837  -0.018426  -0.007543  0.013146  0.008812  -0.037694  -0.154528  0.032898  0.038763
latitude  -0.057497  0.018340  1.000000  -0.099648  -0.014187  -0.028857  -0.018575  -0.025464  0.006525  -0.029112  -0.007288  0.006216
longitude  0.056245  -0.099648  -0.014187  1.000000  0.018340  0.013146  0.008812  0.008812  0.006525  0.006525  0.006525  -0.006525
bedrooms  0.076366  0.018340  0.013146  0.018340  1.000000  0.059593  0.712343  0.429465  0.158497  0.162295  0.059836  0.152977
bathrooms  0.003837  -0.018426  -0.007543  0.018340  0.059593  1.000000  0.464094  0.884179  0.259873  0.163807  0.294577  0.374977
rooms  0.018340  0.003837  -0.018426  0.018340  0.059593  0.059593  1.000000  0.484294  0.168294  0.175949  0.238537  0.283018
squareFootage  0.013146  0.006525  0.007543  0.018340  0.059593  0.059593  0.059593  1.000000  0.369244  0.894987  0.334461  0.374789
lotSize  0.008812  0.006525  0.006525  0.018340  0.059593  0.059593  0.059593  0.059593  1.000000  0.259638  0.193421  0.127651
yearBuilt  0.008812  0.006525  0.006525  0.018340  0.059593  0.059593  0.059593  0.059593  0.059593  1.000000  0.068656  0.075427
lastSaleAmount  0.008812  0.006525  0.006525  0.018340  0.059593  0.059593  0.059593  0.059593  0.059593  0.059593  1.000000  0.248492
priorSaleAmount  0.008812  0.006525  0.006525  0.018340  0.059593  0.059593  0.059593  0.059593  0.059593  0.059593  0.059593  1.000000
estimated_value  0.008812  0.006525  0.006525  0.018340  0.059593  0.059593  0.059593  0.059593  0.059593  0.059593  0.059593  0.059593

```

- 2) For more visualization, the PyCharm screenshot below depicts a heat map of correlation between variables and the darker squares indicate higher correlations between variables. Note the diagonal of the matrix is white because it is the correlation of a variable to itself.

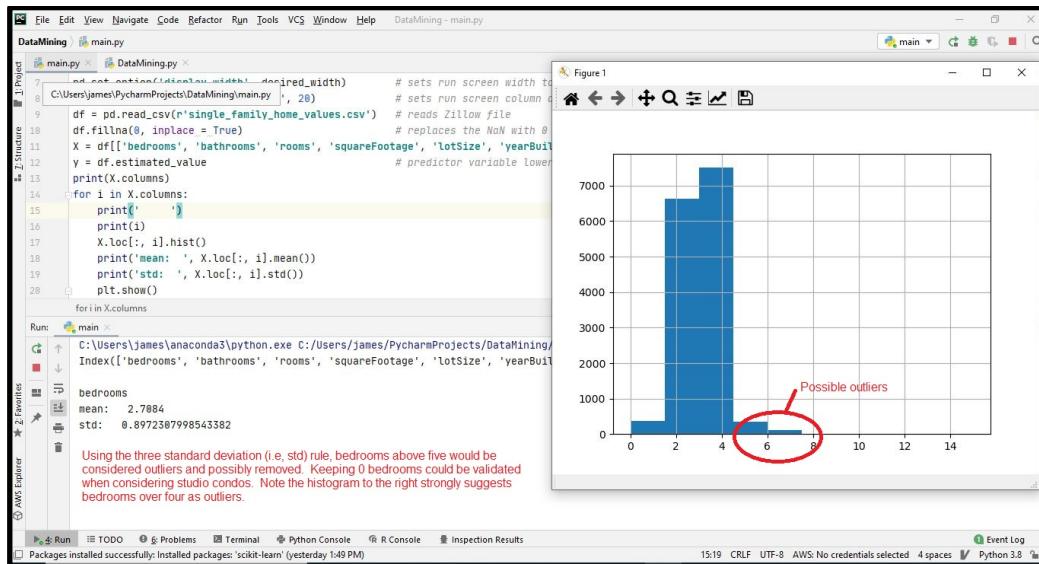


- 3) Covariance is the unstandardized correlation and can be identified in a matrix by the Python code `dataframe.cov()` as illustrated in the PyCharm screenshot below.

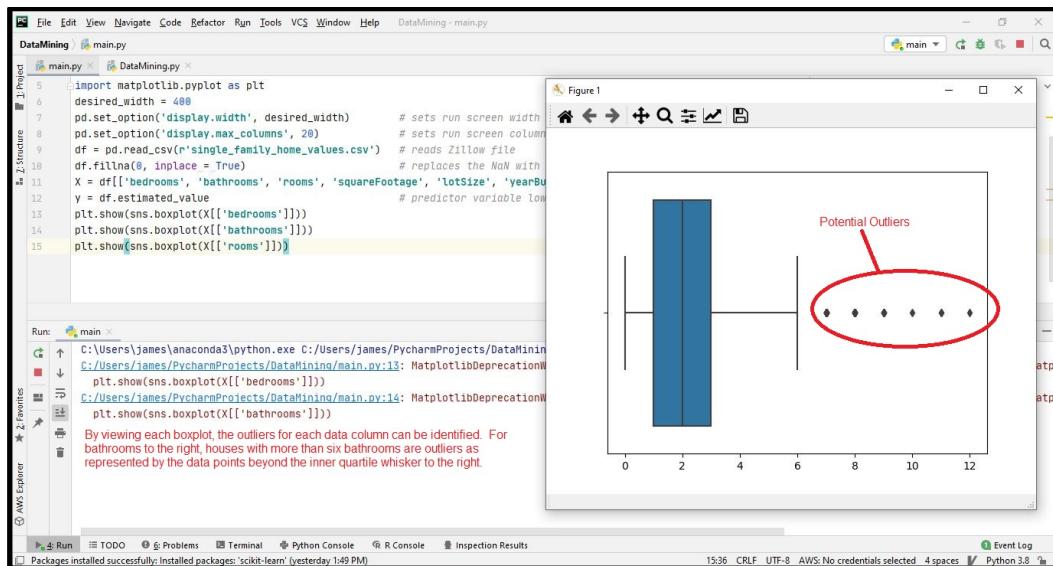


- 4) Measures of association help data scientists identify which variables explain other variables in terms of variance. Coefficient of correlation, covariance, and a coefficient of determination are all calculated in terms of variance. Outliers generate noise by artificially amplifying variance. Beyond correlation and covariance, two measures to identify outliers are the three standard deviation rule and boxplots identifying IQRs. Using the three standard deviation rule, we would exclude all outliers that are three standard deviations from the mean. The PyCharm

screenshot below uses a ‘for loop’ to display histograms for each of our seven columns (e.g., bedrooms, bathrooms, rooms, squareFootage, lotSize, yearBuilt, and priorSaleAmount). For bedrooms, the three standard deviation rule suggests $2.784 \pm (3)(0.892307998543382)$. Per the justification in the PyCharm screenshot below, a case can be made to keep data points of houses having between 0 to 4 and 0 to 5 bedrooms.



Another tool to use in outlier identification is a boxplot that uses inter-quartile ranges IRQs. The following PyCharm screenshot shows the Python code to produce boxplots for the X column data points.



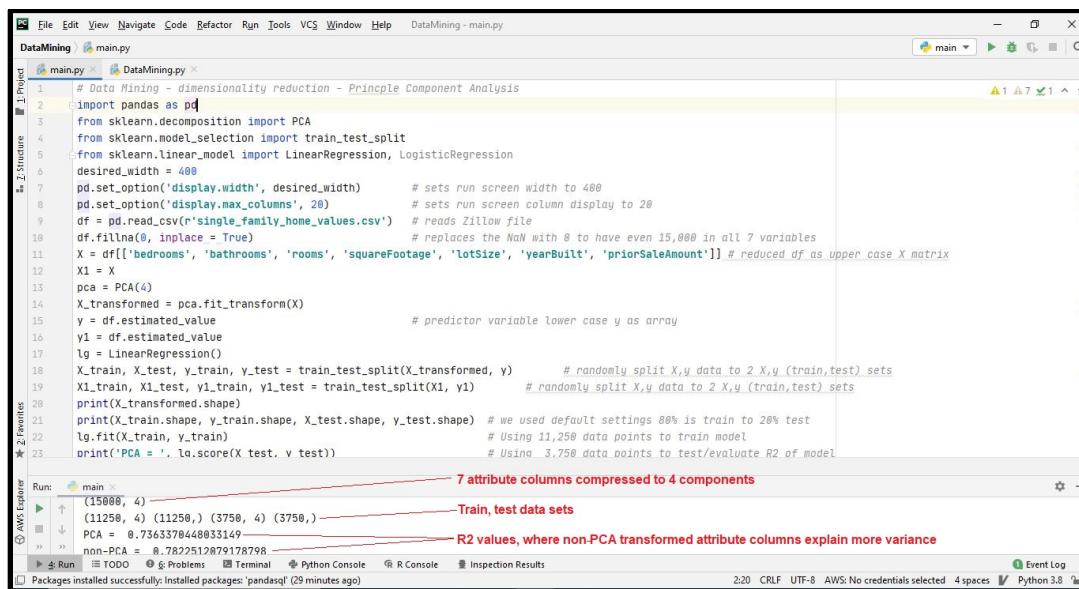
- 5) **Exercise deliverable:** (1)Construct an example of `df.corr()` function or `df.cov()` in Python code using the matrix or a heat map of the matrix to display the correlation or covariance between column attributes. (2) As a second PyCharm screenshot, use the three standard deviation rule

or a boxplot to identify outliers. Use either the single_family_home_values.csv or your data project csv. In using the CSV file above, choose different column attributes than those shown in the previous two screenshots. Save the two PyCharm screenshots for use later in this tutorial.

Dimensionality Reduction:

Principle Component Analysis (PCA): PCA is used to decompose or compress a multivariate dataset into a set of successive orthogonal components that explain a maximum amount of the variance. In sklearn, PCA is implemented as a transformer object that learns components in its fit method, and can be used on new data to project it on these components. The focus of PCA is to compress the number of explanatory variables by taking portions of starting variables and compressing them into a given number of components. For more details on PCA, refer to the following URL: <https://scikit-learn.org/stable/modules/decomposition.html#pca>

- 1) The following PyCharm screenshot demonstrates how PCA can condense seven attribute columns into four components. The R^2 of the PCA compressed components show how explained variance can be lost when compressing explanatory variables into fewer components, especially when the reduction is small (i.e., 7 to 4).



The screenshot shows a PyCharm IDE window with the following details:

- Code Editor:** The main.py file contains Python code for dimensionality reduction using PCA. It imports pandas, sklearn.decomposition, and sklearn.model_selection. It reads a CSV file (single_family_home_values.csv) and performs PCA on the features (bedrooms, bathrooms, rooms, squareFootage, lotSize, yearBuilt, priorSaleAmount). It splits the data into training and testing sets and fits a LinearRegression model. The output shows that 7 attribute columns were compressed into 4 components, with an R^2 value of 0.7363378448033149.
- Run Tab:** Shows the run configuration for 'main' with arguments: (15000, 4). It also shows the output: Train, test data sets and R2 values, where non-PCA transformed attribute columns explain more variance.
- Bottom Status Bar:** Shows the current time (2:20), file encoding (CRLF, UTF-8), AWS status (No credentials selected), and Python version (Python 3.8).

- 2) **Exercise deliverable:** Use PCA to reduce the number of explanatory variables into fewer components. Use either the single_family_home_values.csv or your data project csv. If you choose the csv file above, do not use the exact same attributes. Save the PyCharm screenshot for use later in this tutorial.

Deliverables from Exercises Above:

Take the noted deliverables (i.e., PyCharm screenshots) from each exercise above (i.e., total of 8 screenshots), insert the screenshots into a MS Word document as you label each screenshot deliverable from the corresponding exercise, save the document as

one PDF, and submit the PDF to the Tutorial Eight-A Canvas Assignment uplink having the following qualities:

- a. Your screenshots for each of the four exercise deliverables above is included.
- b. Each screenshot is labeled appropriately for each exercise.
 - i. Cluster Analysis 2 screenshots
 - ii. Regression and Classification 3 screenshots
 - iii. Association and Correlation 2 screenshots
 - iv. Dimensionality Reduction 1 screenshotTotal PyCharm screenshots to submit 8 screenshots
- c. All screenshots are legible output for each PyCharm exercise deliverable.