

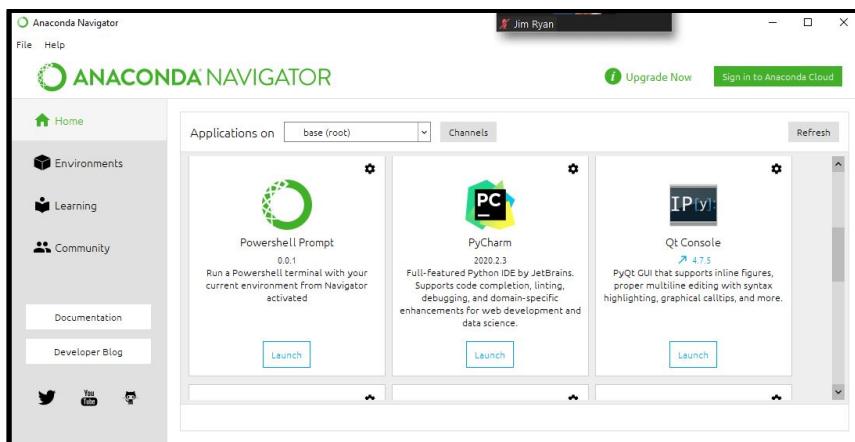
## Mining & Data Storing

### Data mining continued...

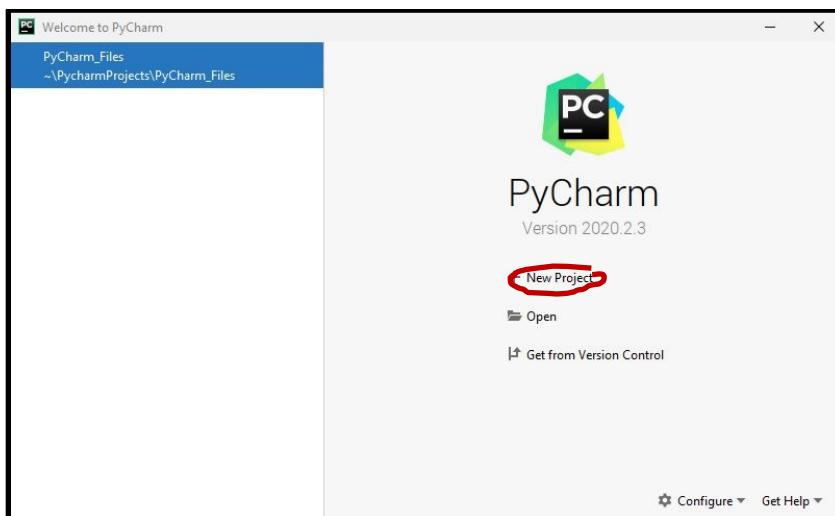
Data mining is a combination of several different processes, tools, or methodologies used to extract relevant data from a larger set of raw data. As we saw in Tutorial Eight-A, data mining implies analyzing data patterns in large batches using one or more software applications. With respect to this tutorial, we will continue our focus on Python data mining processes for text mining and network mining, as well as introducing array and matrix math along with data input and output (reading, writing, and appending) via SQL scripts, NumPy, and Pandas. Tutorial Ten will explore more text mining via Natural Language Processing (NLP).

### Create a MiningStoring Project in PyCharm

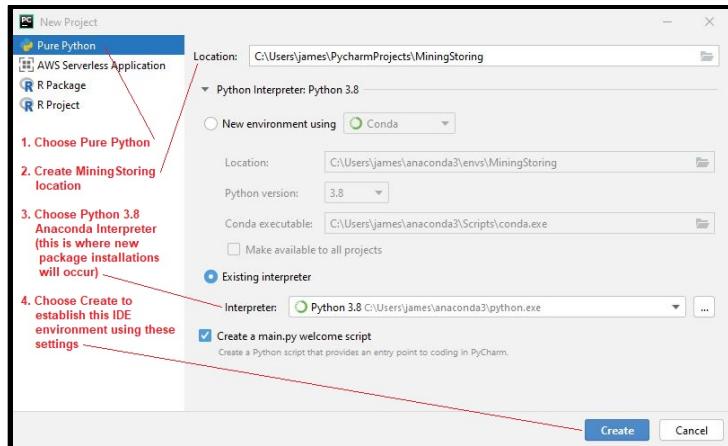
- 0) To ensure your various Python package libraries are available to import, be sure to use Anaconda Navigator to launch PyCharm.



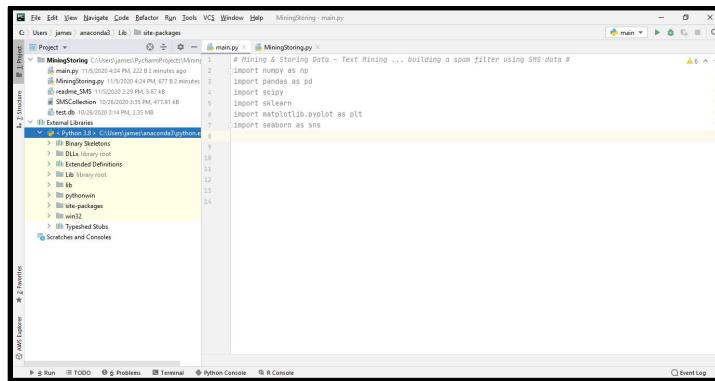
- 1) Launch PyCharm and choose a New Project.



- 1) Choose a Pure Python project,
- 2) Change your project location to create a MiningStoring 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.

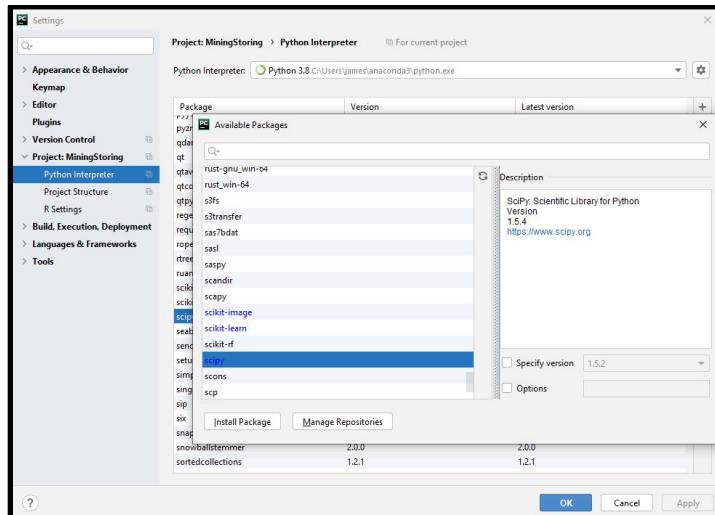


- 5) Download the SMSCollection, readme\_SMS, and test.db files from the Week 11 Canvas Module to the MiningStoring directory created in the last step. The readme\_SMS file explains the data in the SMSCollection file. The test.db file is an SQLite database.



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

- 6) When importing a new package library, check the Python Interpreter to ensure the package is installed via Files->Settings->Python Interpreter or left-click Python3.8 in the lower right screen corner choosing Interpreter Settings... Once in the Python Interpreter, pan down to find the particular Python package library and double click on package version. In the resulting window, all installed packages are highlighted in blue.



 **Reading a text file into a Pandas DataFrame:**

We learned about Pandas in Tutorials Seven and Eight-A. Again in this tutorial, we will look at DataFrame. For more specific details about Pandas not covered in this tutorial, please refer to the following URL:

[https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.read\\_table.html](https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.read_table.html)

- 1) Anaconda will load the libraries in PyCharm. To access the libraries we are using for text mining in Tutorial Eight-B (e.g., numpy, pandas, scipy, seaborn, sklearn, and matplotlib.pyplot), enter the following code into Main.py in the MiningStoring project directory of PyCharmProjects:

```
import numpy as np
import pandas as pd
import scipy
import sklearn
import matplotlib.pyplot as plt
import seaborn as sns
```

- 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 SMSCollection file into a pandas dataframe, type the following Python code into main.py via PyCharm in the line below the previous code above.

```
data_file = 'SMSCollection'      # assigns file name variable data_file
sms_raw = pd.read_table(data_file, header=None)      # reads SMS file
sms_raw.columns = ('spam', 'message')      # assigns column names
```

- 4) Once the SMSCollection file is loaded into a dataframe table, we can begin text mining. Any of the following Python functions may be used with the dataframe sms\_raw, once it has been assigned.

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

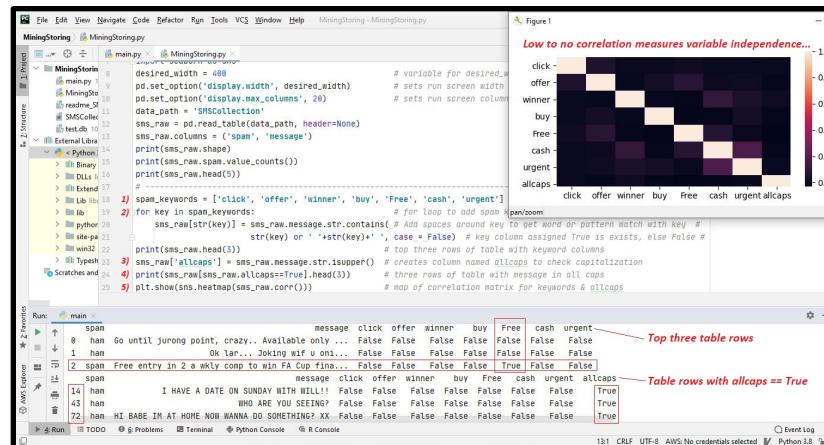
- 5) The MiningStoring.py file has the Python code for all screenshots in Tutorial Eight-B.

 **Text Mining:**

Text mining, also referred to as text analytics, transforms the free (unstructured) text in documents and databases into normalized, structured data suitable for analysis or to drive machine learning (ML) algorithms. Widely used in knowledge-driven organizations, text mining is the process of examining large collections of documents to discover new information or help answer specific research questions. Text mining identifies facts, relationships and assertions that would otherwise remain buried in the mass of textual big data. Tutorial Ten explains Natural Language Processing (NLP), which is similar to text mining. Text mining and NLP have difference in the end goal and difference in methods. Text mining techniques are usually shallow and do not consider the text structure, while using bag-of-words, n-grams, and possibly stemming. NLP methods usually involve

the text structure, sentence splitting, part-of-speech tagging, and parse tree construction. For this exercise in Tutorial Eight-B, we will demonstrate general text mining methods via Python code towards developing a spam filter.

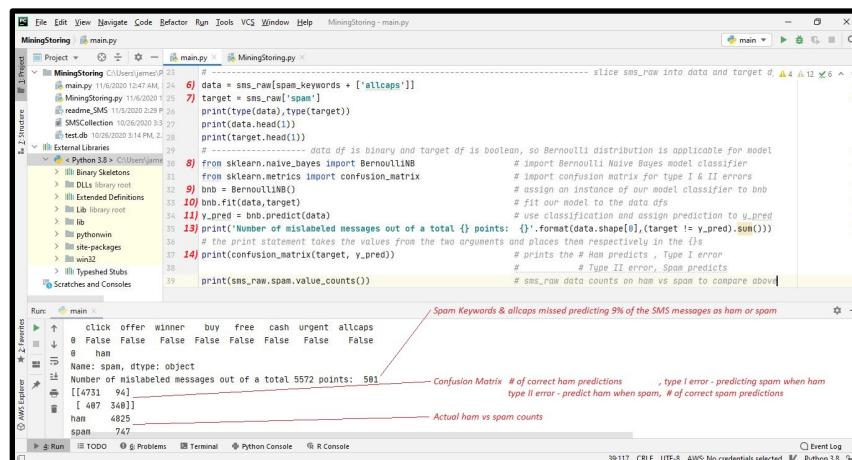
- 1) We imported the five libraries and added Python code from the prior reading df section. In the PyCharm screenshot below, you will see how to (1) assign spam keywords to a list, (2) use a 'for loop' to check each message column row to see if a spam keyword is in it, (3) add an allcaps column in sms\_raw to check if message is all upper case, (4) print sms\_raw rows with allcaps == True, and (5) plot a heat map to identify correlation between spam keywords and allcaps. The resulting heatmap plot shows little to no correlation and suggests strong independence between spam keyword and 'allcaps' variables.



*All screenshots may be enlarged 400x for clarity*

- 2) In the following PyCharm screenshot, we (6) slice the sms\_raw df columns spam\_keywords and 'allcaps' into data df, (7) slice the sms\_raw df column 'spam' into target df, (8) import a Bernoulli Naïve Bayesian classifier model and a confusion matrix prediction checker, (9) assign a variable to the classifier model, (10) fit the classifier model to data and target, (11) assign the classifier model prediction to y\_pred, (12) count the number of total data points vs. mislabeled predictions (type I & II errors), (13) calculate the prediction check confusion matrix, and (14) print the sms\_raw df column 'spam' value counts. More information on the Bernoulli Naïve Bayes classifier is available via the URL:

[https://scikit-learn.org/stable/modules/naive\\_bayes.html#bernoulli-naive-bayes](https://scikit-learn.org/stable/modules/naive_bayes.html#bernoulli-naive-bayes)



- 3) In the two prior PyCharm screenshots, the spam keywords and allcaps column attributes are used as rules to detect spam SMS messages from regular (ham) SMS messages. The mispredictions of type I and type II errors are calculated. In the following PyCharm screenshot, we (15) create a data frame sms\_raw2 which holds all the mispredictions of type I and type II errors, (16) write the sms\_raw2 data frame to a .CSV file for review in Notepad or Excel, (17) use a ‘for loop’ to list all the type I and type II errors where a spam keyword was found in the SMS message, (18) listing of type I and type II errors where allcaps were true, and (19) a count of rows where type I and type II errors showed allcaps were true.

The screenshot shows a PyCharm interface with the following details:

- Code Editor:** Displays main.py with code related to SMS classification and analysis. Lines 15-19 are highlighted in red, corresponding to the numbered steps in the list below.
- Run Tab:** Shows the output of the script. It includes:
  - Number of mislabeled messages out of a total 5572 points: 581
  - Type I errors: 4731 (spam)
  - Type II errors: 3461 (ham)
  - spam keyword is click: A table showing counts for various columns (message, click, offer, winner, buy, Free, cash, urgent, allcaps) across different categories (spam, ham).
  - spam keyword is offer: Similar table for the 'offer' keyword.
- Bottom Status Bar:** Shows the Python version (Python 3.8), encoding (UTF-8), and AWS credentials status.

- 4) **Exercise deliverable:** Review the sms\_mispredic.csv file using Notepad, Excel, or another text editor. Using the Python code in this text mining example (steps 1 to 14), add or change spam\_keywords to improve the model fit and minimize type II errors. If you prefer to duplicate this text mining example from your project data set, then please do. Capture a PyCharm screenshot showing the number of mislabeled predictions as well as the confusion matrix for use later in this tutorial.

## Network Mining:

Graph and network mining has leaped to the forefront of data mining research, spurred by an avalanche of structured data from applications such as bioinformatics, cheminformatics, online social networks, and sensor networks. This structured data is often best represented either as a set of independent graphs or as a large network of interconnected nodes. The proliferation of graph and network data is both an opportunity and a challenge. Graph mining is playing an increasingly important role in the analysis of highly structured data such as chemical compounds, proteins, very large scale integration (VLSI) designs, and program execution traces. Examples of graph mining applications include predicting protein function (graph classification), searching for compounds with certain substructures (graph similarity search), and detecting software bugs in programs

(graph anomaly detection). Network mining offers the opportunity for analyzing large collections of inter-related objects such as Web graphs, social networks, and biological networks. Common applications are assessing the spread of epidemics (influence maximization), predicting future collaboration between authors (link prediction), and finding authoritative web pages (link-based node ranking). Traditional algorithms, which typically assume that objects are independent and identically distributed, are not appropriate for mining such data. Furthermore, some of the data mining tasks of interest (particularly in network data) have no counterparts in record-based data (e.g., influence maximization and link prediction) .For more information on network mining, refer to the URL: <https://networkx.org/> .

- 1) **Network Graphs:** The following Python screenshot uses Python code and the networkx package library to show relationships within a family unit {Fiancée, Child, Mom, Dad, Uncle}. Several different graph layouts are used to plot a network graph, specifically circular, spiral, and spring layouts. For more information on the numerous networkz graph layouts, refer to the URL:

[https://networkx.org/documentation/stable//\\_modules/networkx/drawing/layout.html](https://networkx.org/documentation/stable//_modules/networkx/drawing/layout.html)

The figure consists of three sub-diagrams arranged horizontally, each showing a network of six nodes: Child, Dad, Mom, Uncle, Fiancee, and Uncle. The edges represent relationships between these nodes.

- Figure 1 – Circular Layout:** The nodes are arranged in a circular pattern. Edges connect Child to Mom, Child to Uncle, Uncle to Fiancee, Uncle to Uncle, and Fiancee to Dad.
- Figure 2 – Spiral Layout:** The nodes are arranged in a spiral pattern. Edges connect Child to Uncle, Uncle to Uncle, Uncle to Fiancee, Fiancee to Dad, Dad to Mom, and Child to Mom.
- Figure 3 – Spring Layout:** The nodes are arranged in a triangular pattern. Edges connect Child to Uncle, Uncle to Uncle, Uncle to Fiancee, Fiancee to Dad, Dad to Mom, and Child to Mom.

- 2) ***Network Graph Information:*** Once the network graph is plotted, Python code allows you to summarize the nodes, edges (i.e., network paths or relationships), betweenness centrality to identify the shortest network path, and number of degrees (i.e., relationships) for each node.

The following PyCharm screenshot illustrates the Python code to retrieve network graph information.

A screenshot of the PyCharm IDE. The top menu bar includes File, Edit, View, Navigate, Code, Refactor, Run, Tools, VCS, Window, Help. The title bar says "MiningStorage - main.py". The left sidebar shows a project structure with files like fig1.png, fig2.png, fig3.png, main.py, MiningStorage.py, and test.db. The main editor window contains Python code for printing network graph statistics. The run output window at the bottom shows the execution results, including node degrees and betweenness centrality scores. Red annotations highlight the 'Betweenness Centrality score' as 'shortest path' and 'higher score denotes influence', and the 'Number of relationships' as 'the number of edges (node relationships)'.

```

#----- # pulling information from the network graph
print('This family network graph has {} nodes and {} edges.'.format(G.number_of_nodes(), G.number_of_edges())) # placing numbers in {}
print('The nodes are: {}'.format(G.nodes()))
print('The edges are: {}'.format(G.edges()))
print('-----')
print('The Child node has in-degree of {} and an out-degree of {}'.format(G.in_degree('Child'), G.out_degree('Child')))
print('The Fiancee node has in-degree of {} and an out-degree of {}'.format(G.in_degree('Fiancee'), G.out_degree('Fiancee')))
print('The Mom node has in-degree of {} and an out-degree of {}'.format(G.in_degree('Mom'), G.out_degree('Mom')))
print('The Dad node has in-degree of {} and an out-degree of {}'.format(G.in_degree('Dad'), G.out_degree('Dad')))
print('The Uncle node has in-degree of {} and an out-degree of {}'.format(G.in_degree('Uncle'), G.out_degree('Uncle')))
print('-----') # line separator
print('The betweenness centrality scores are: {}'.format(nx.betweenness_centrality(G))) # shortest path - most influence
print('The node degrees are: {}'.format(G.degree())) # Node degree is the number of edges (node relationships)

```

- 3) **Exercise deliverable:** Construct a network graph from your project data set, if applicable. If a network graph is not applicable, then choose four friends from Facebook or four colleagues from your LinkedIn network. (0) Code the names of your friends with their initials and the number of friends or colleagues in their networks. (1) Identify if any of the four have relationships with each other on the social platform. (2) Build a network graph and try a few different graph layouts and choose the one you like. (3) Save the network graph as a png file. (4) Gather the same network graph information as detailed in the PyCharm screenshot above. (5) Capture a PyCharm screenshot displaying the png file above the run window that displays the information. An example of this PyCharm screenshot configuration is below.

A screenshot of the PyCharm IDE. The top menu bar includes File, Edit, View, Navigate, Code, Refactor, Run, Tools, VCS, Window, Help. The title bar says "MiningStorage - fig1.png". The left sidebar shows a project structure with files like fig1.png, fig2.png, fig3.png, main.py, MiningStorage.py, and test.db. The main editor window contains Python code for printing network graph statistics. The run output window at the bottom shows the execution results, including node degrees and betweenness centrality scores. Above the run window, a network graph visualization is displayed, showing nodes for Child, Fiancee, Mom, and Dad, with edges connecting them. The graph has 5 nodes and 10 edges.

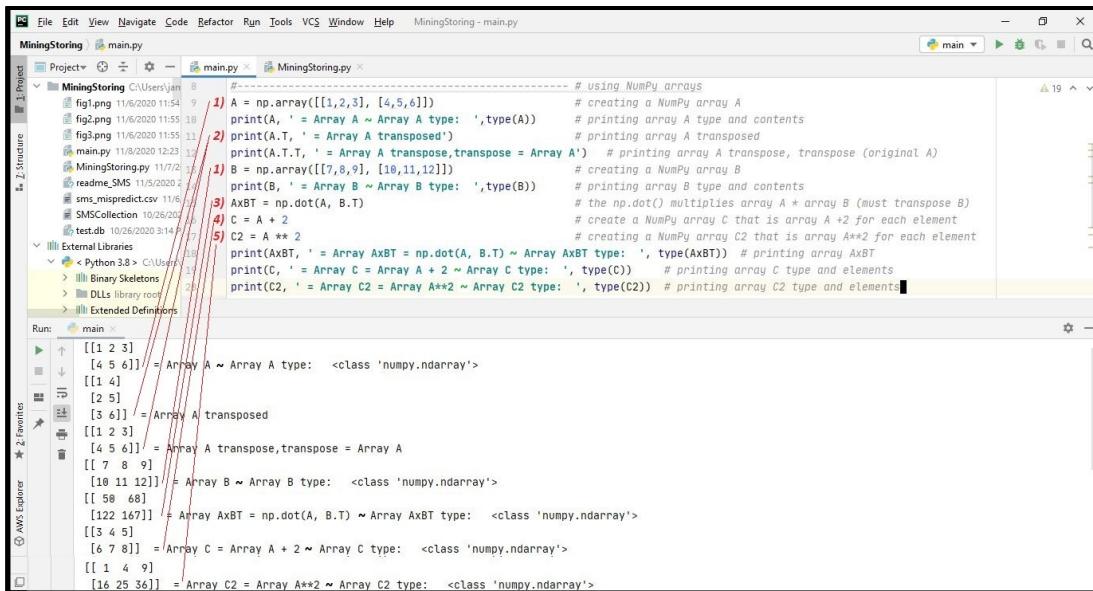
## Matrix:

*Should you take the red pill or blue pill?* In Python, we can use the Numpy or Pandas libraries to help with performing matrix math. NumPy is a powerful python library that expands Python's functionality by allowing users to create multi-dimensional array objects (ndarray). In addition to the creation of ndarray objects, NumPy provides a large set of mathematical functions that can operate quickly on the entries of the ndarray without the need of for loops. The Pandas (i.e., PAnel DAta) Python library allows for easy and fast data analysis and manipulation tools by providing numerical tables and time series data structures called DataFrame and Series, respectively. Pandas was created to (1) provide data structures that can handle both time and non-time series data, (2) allow mathematical operations on the data structures while ignoring the metadata of the data structures, (3) use relational operations like those found in SQL (join, group by, etc.), and (4) handle missing data. Pandas is built on top of NumPy, relying on ndarray and its mathematical functions. Likewise, Pandas relies heavily on the NumPy array for the implementation of pandas data objects and shares many of its features. In addition, pandas builds upon functionality provided by NumPy. For more information on matrix math under NumPy and Pandas in Python, refer to the following URLs:

[https://www.python-course.eu/matrix\\_arithmetic.php](https://www.python-course.eu/matrix_arithmetic.php) (NumPy)

<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.dot.html> (Pandas)

- 1) The Python code in the following PyCharm screenshot reviews some basic matrix math using NumPy two dimensional arrays. Specifically, we execute (1) np.array to create a NumPy array A and array B, (2) A.T to transpose array A (i.e., array A.T.T is the same as array A), (3) np.dot() that multiplies two arrays, (4) adding 2 to each element in array C, and (5) raising each element in array C2 to the power of 2 (i.e., squaring each item).



The screenshot shows a PyCharm IDE window with the following code in the main.py file:

```

#----- # using NumPy arrays
A = np.array([[1,2,3], [4,5,6]]) # creating a NumPy array A
print(A, ' = Array A ~ Array A type: ', type(A)) # printing array A type and contents
print(A.T, ' = Array A transpose') # printing array A transposed
print(A.T.T, ' = Array A transpose,transpose = Array A') # printing array A transpose, transpose (original A)
B = np.array([[7,8,9], [10,11,12]]) # creating a NumPy array B
print(B, ' = Array B ~ Array B type: ', type(B)) # printing array B type and contents
AXBT = np.dot(A, B.T) # np.dot() multiplies array A * array B (must transpose B)
print(AXBT, ' = Array AXBT = np.dot(A, B.T) ~ Array AXBT type: ', type(AXBT)) # printing array AXBT
C = A + 2 # create a NumPy array C that is array A +2 for each element
print(C, ' = Array C = Array A + 2 ~ Array C type: ', type(C)) # printing array C type and elements
C2 = A ** 2 # creating NumPy array C2 that is array A**2 for each element
print(C2, ' = Array C2 = Array A**2 ~ Array C2 type: ', type(C2)) # printing array C2 type and elements

```

The code defines two 3x3 arrays, A and B, and performs various operations on them, including dot product, addition, and exponentiation.

- 2) For more examples of arrays, matrices, and data frames, the Python code in the following PyCharm screenshot reviews more functionality using NumPy arrays, a NumPy matrix, and a NumPy matrix converted to a Pandas data frame. For precision in math calculations, a NumPy

array or matrix is recommended. The following PyCharm screenshot specifically demonstrates the execution of (6) np.array to create a 3x3 NumPy two dimensional array C3; (7) use of the array.diagonal() function to display NumPy array C4 as the diagonal of C3; (8) use of np.matrix() to convert array A into matrix A; (9) use of the matrix.sum(), matrix.mean(), matrix.std(), and matrix.var() functionality; (10) use of the matrix.nonzero() function; (11) subtracting array B from matrix A; (12) adding array B to matrix A; (13) use of the matrix.dot() function to multiply matrix A by array B transposed; and (14) converting a NumPy matrix to a Pandas data frame.

```

# more NumPy arrays, NumPy Matrix, & NumPy Matrix as Pandas df
C3 = np.array([[1,2,3], [4,5,6], [7,8,9]]) # creating a NumPy two dimensional array C3 that is a 3x3
C4 = C3.diagonal() # np.array.diagonal() returns the diagonal elements of an array
print(C3, ' ~ Array C3 is a 3x3 ~ Array C3 type: ', type(C3)) # print NumPy two dimensional array C3
print(C4, ' ~ Array C4 is the C3 diagonal ~ Array C4 type: ', type(C4)) # print NumPy two dimensional array C4
A = np.matrix(A)
print(A, ' = matrix A = array A ~ matrix A type: ', type(A)) # print NumPy matrix A
print('..... matrix A sum = ', A.sum(), ' | mean = ', A.mean(), ' | std = ', A.std(), ' | variance = ', A.var())
print('..... matrix nonzero() returns indices of nonzero elements ~ ',A.nonzero()) # print indices of nonzero elements
print(A-B, ' = A-B results when array B is subtracted from matrix A') # print element results of matrix A - array B
print(A+B, ' = A+B results when array B is added to matrix A') # print element results of matrix A + array B
print(A.dot(B.T), ' = A.dot(B.T) results when matrix A is multiplied by array B transposed') # print element results
A_df = pd.DataFrame(A) # print matrix A as a Pandas Dataframe ~ type(A_df) = , type(A_df)
print(A_df, ' = matrix A as a Pandas Dataframe ~ type(A_df) = <class \'pandas.core.frame.DataFrame\'>')

[[1 2 3]
 [4 5 6]] = matrix A = array A ~ matrix A type: <class 'numpy.matrix'>
..... matrix A sum = 21 | mean = 3.5 | std = 1.707825127659933 | variance = 2.9166666666666665
..... matrix nonzero() returns indices of nonzero elements ~ (array([0, 0, 0, 1, 1, 1], dtype=int64), array([0, 1, 2, 0, 1, 2], dtype=int64))
[[6 -6 -6]
 [-6 -6 -6]] = A-B results when array B is subtracted from matrix A
[[1 0 12]
 [14 16 18]] = A+B results when array B is added to matrix A
[[ 50  68]
 [122 167]] = A.dot(B.T) results when matrix A is multiplied by array B transposed
[[ 0  1  2]
 [ 1  2  3]
 [ 1  4  5  6]] = matrix A as a Pandas Dataframe ~ type(A_df) = <class 'pandas.core.frame.DataFrame'>

```

- 3) Exercise deliverable:** Use your project data set or a data set you derive to illustrate array, matrix, and Pandas data frame functionality. In a PyCharm screenshot, show the array/matrix and data frame elements, mean(), sum(), std(), var(), and if applicable diagonal(). Note that an array or matrix applies aggregate functions to all elements. The data frame will apply the aggregate functions to a column (i.e., in the screenshot above print(A\_df.mean()) would print 2.5 for column 0, 3.5 for column 1, and 4.5 for column 2). The diagonal() function is not applicable to data frames as the matrix is segregated by rows and columns. Save the PyCharm screenshot for use later in this tutorial.

## Storing and Reading Data:

In prior Tutorials, we have looked at reading CSV or flat text into Panda data frames. This section of the tutorial will look at reading as well as storing data using Python code, specifically in SQLite databases and CSV files. For reading and writing data to SQLite databases, we will import and use the sqlite3 package library. For more information on the sqlite3 package library not covered in this tutorial, please refer to the URL: <https://docs.python.org/3/library/sqlite3.html>

- 1) SQL Script Execution:** The following PyCharm screenshot demonstrates how to (1) open a SQLite database, (2) execute SQL scripts, and (3) close the SQLite database connection. In this particular example, two SQL scripts are executed from the SQLite database test.db located in the ...PyCharmProject/MiningStoring. The SQLite database should always be in the subdirectory where the Python program is executed or the path to the database should be specified in the .connect function (i.e., sqlite.connect(path/database.db)).

```

# Mining & Storing Data - reading and storing data - using sqlite3 package to access a SQL database#
import pandas as pd
import sqlite3
desired_width = 400
pd.set_option('display.width', desired_width)
pd.set_option('display.max_columns', 20)
pd.set_option('display.max_rows', 100)
#-----#
# using sqlite3 package library to read a SQLite database#
# establish connection to SQLite test.db#
con = sqlite3.connect('test.db')#
cur = con.cursor()#
for row in cur.execute("SELECT * FROM physician LIMIT 5;"):
    print(row)
    print()
for row in cur.execute("SELECT * FROM procedure LIMIT 5;"):
    print(row)
    print()
con.close()
    
```

Run: main

```

('id', 'specialty')
('0', 'General Surgery')
('1', 'Unknown')
('2', 'Family Practice')
('3', 'Emergency Medicine')

('physician_id', 'procedure_code', 'procedure', 'number_of_patients')
('0', '99202', 'new_patient_office_or_other_outpatient_visit,_typically_20_minutes', '14')
('0', '99203', 'new_patient_office_or_other_outpatient_visit,_typically_30_minutes', '15')
('0', '99205', 'new_patient_office_or_other_outpatient_visit,_typically_60_minutes', '12')
('0', '99212', 'established_patient_office_or_other_outpatient_visit,_typically_10_minutes', '27')
    
```

The following PyCharm screenshot displays Python code using a user-defined function to execute the same two SQL scripts from the previous example. It should be noted SQL scripts may be data manipulation language to project or insert data, data definition language to create database objects (i.e., tables, views, indexes, etc.), or data control language to assign rights and security.

```

# using a user defined function to execute a SQL script#
def exec_sql():
    print('')
    print(SQL_script)
    cur = con.cursor()
    for row in cur.execute(SQL_script):
        print(row)
        print()
    con = sqlite3.connect('test.db')
    SQL_script = "SELECT * FROM physician LIMIT 5;"
    exec_sql()
    SQL_script = "SELECT * FROM procedure LIMIT 5;"
    exec_sql()
    con.close()
    
```

Run: main

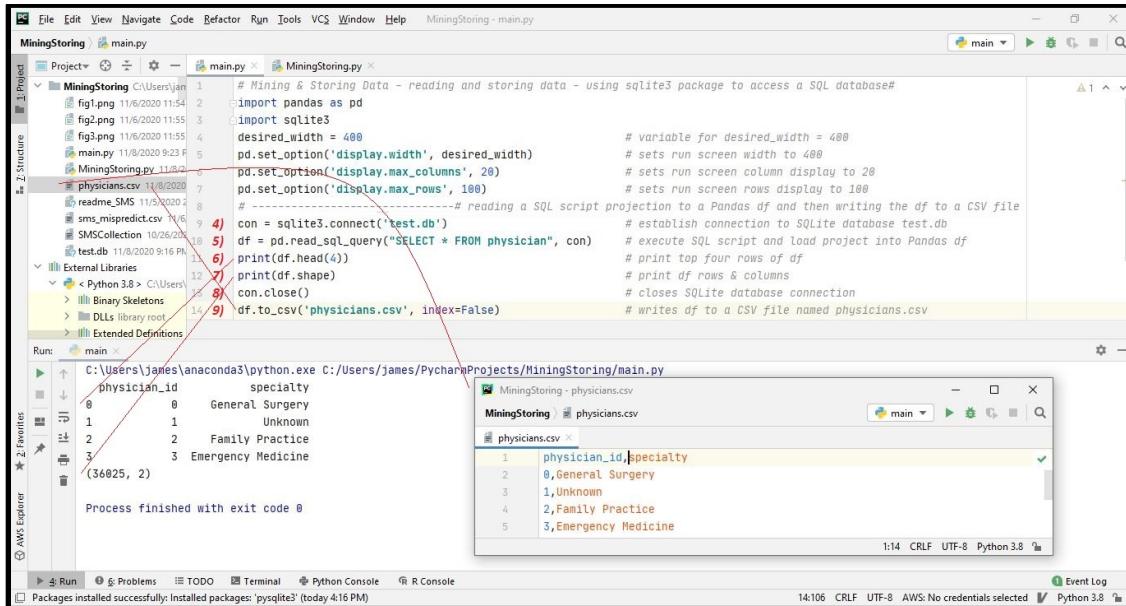
```

SELECT * FROM physician LIMIT 5;
('id', 'specialty')
('0', 'General Surgery')
('1', 'Unknown')
('2', 'Family Practice')
('3', 'Emergency Medicine')

SELECT * FROM procedure LIMIT 5;
('physician_id', 'procedure_code', 'procedure', 'number_of_patients')
('0', '99202', 'new_patient_office_or_other_outpatient_visit,_typically_20_minutes', '14')
('0', '99203', 'new_patient_office_or_other_outpatient_visit,_typically_30_minutes', '15')
('0', '99205', 'new_patient_office_or_other_outpatient_visit,_typically_60_minutes', '12')
('0', '99212', 'established_patient_office_or_other_outpatient_visit,_typically_10_minutes', '27')
    
```

- 2) **Reading a SQL Script Projection into a Pandas Data File and Saving the df as a CSV:** The following PyCharm screenshot demonstrates how to (4) open a SQLite database, (5) read an SQL script SELECT statement into a Pandas df, (6) check the first four rows of the df, (7) list the df's number of columns and rows, (8) close the SQLite database connection, and (9) write the df to a CSV file. In this particular example, the SQL script projection is read directly into a

Pandas df and the resulting df is written to a CSV file in the ...PyCharmProjects/MiningStoring directory.



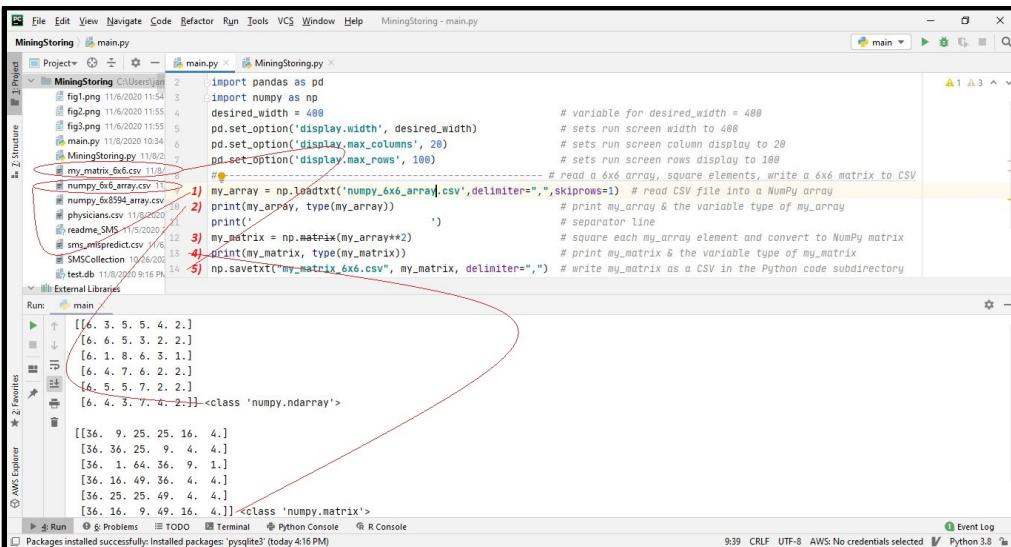
```

# Mining & Storing Data - reading and storing data - using sqlite3 package to access a SQL database#
import pandas as pd
import sqlite3
desired_width = 400
pd.set_option('display.width', desired_width)
pd.set_option('display.max_columns', 20)
pd.set_option('display.max_rows', 100)
# ---#
con = sqlite3.connect('test.db')
df = pd.read_sql_query("SELECT * FROM physician", con)
print(df.head(4))
print(df.shape)
con.close()
df.to_csv('physicians.csv', index=False)

```

The screenshot shows the PyCharm interface with the code editor containing the provided Python script. The code connects to a SQLite database, reads the 'physician' table into a Pandas DataFrame, prints the first four rows, prints the shape of the DataFrame, closes the database connection, and then writes the DataFrame to a CSV file named 'physicians.csv'. The Run tab shows the output of the script, which includes the DataFrame's head and its shape, followed by the message 'Process finished with exit code 0'. To the right, there are two tabs showing the contents of 'physicians.csv': one tab shows the raw text of the CSV, and another shows the CSV data as a table with columns 'physician\_id' and 'specialty' and rows 0 through 5.

- 3) **Reading a CSV file into NumPy array; Writing manipulated matrix to a CSV file.** In this example, we read a CSV called numpy\_6x6\_array.csv that is stored in the Python code subdirectory ...PycharmProjects/MiningStoring. As a general rule, using the Pandas import method is a more 'forgiving', so if you have trouble reading directly into a NumPy array, try loading in a Pandas data frame and then converting it to a NumPy array. NumPy's loadtxt method reads delimited text. We specify the separator as a comma. The data we are loading also has a text header, so we use skiprows=1 to skip the header row, which would cause problems for NumPy. We use the savetxt method to save the resulting matrix to a CSV file. In the following PyCharm screenshot, we specifically (1) load a CSV file into a NumPy array, (2) print the array and variable type, (3) square the elements of the array and convert the array to a matrix, (4) print the matrix and the variable type, and (4) save the matrix to a CSV file.



```

# variable for desired_width = 400
# sets run screen width to 400
# sets run screen column display to 20
# sets run screen rows display to 100
# ---#
my_array = np.loadtxt('numpy_6x6_array.csv', delimiter=',', skiprows=1) # read CSV file into a NumPy array
print(my_array, type(my_array)) # print my_array & the variable type of my_array
print('') # separator line
my_matrix = np.matrix(my_array**2) # square each my_array element and convert to NumPy matrix
print(my_matrix, type(my_matrix)) # print my_matrix & the variable type of my_matrix
np.savetxt("my_matrix_6x6.csv", my_matrix, delimiter=",") # write my_matrix as a CSV in the Python code subdirectory

```

The screenshot shows the PyCharm interface with the code editor containing the provided Python script. The code reads a CSV file 'numpy\_6x6\_array.csv' into a NumPy array 'my\_array', prints the array and its type, creates a NumPy matrix 'my\_matrix' by squaring each element of 'my\_array', prints the matrix and its type, and finally saves the matrix to a CSV file 'my\_matrix\_6x6.csv' in the same directory. The Run tab shows the output of the script, which includes the original 6x6 array, the squared matrix, and the saved CSV file. A red circle highlights the 'my\_array' assignment in line 11, and a red oval highlights the entire block of code from line 11 to line 14.

- 4) **Reading a CSV file into Pandas data frame in chunks; Writing/appending to a CSV file.** In this example, we read a CSV file numpy\_6x8594\_array.csv that is stored in the Python code subdirectory ...PycharmProjects/MiningStoring. The CSV file has 8,594 rows, yet let us imagine the file is too large to load all the array elements. We can load a data frame in chunks and then convert the chunk data frame into a NumPy array. We use the savetxt method to save the resulting matrix to a CSV file, but after the first chunk write we need to append the latter chunk writes. In the following PyCharm screenshot, we specifically (0) establish chunks size, i counter, and eof\_size variables to use in an ‘if statement’ (1) load a CSV file into a chunk data frame using a ‘for loop’, (2) pass the chunk data frame to a NumPy array, (3) square the NumPy array and save as a NumPy matrix, (4) print the array first five rows and variable type, (5) print the matrix first five rows and variable type, (6) write the NumPy matrix to a CSV file for the first pass chunk, (7) create a file handle to open the CSV file for appending, (8) append successive chunk matrices to the CSV file, and (9) close the file handle.

```

# Mining & Storing Data - reading and storing data - reading chunks of a df from CSV; writing/appending to a CSV file #
import pandas as pd
desired_width = 400
pd.set_option('display.width', desired_width)
pd.set_option('display.max_columns', 20)
pd.set_option('display.max_rows', 100)
# --- ----- # read a 6x8594 array in chunks; write/ append to a CSV file
# declare chunks variable and assign value of 1000
# end of CSV file is 8594
chunks = 2000
eof = 8594
i = 0
if i <= eof:
    for chunk in pd.read_csv('numpy_6x8594_array.csv', skiprows=i, header=0, chunksize=chunks):
        my_array = chunk.to_numpy(dtype=int)
        my_matrix = np.matrix(my_array**2)
        print(my_array[0:5,:], type(my_array))
        print('')
        print(my_matrix[0:5,:], type(my_matrix))
        if i < 1:
            np.savetxt('my_matrix_6x8594.csv', my_matrix, delimiter=',') # write my_matrix_6x8594.csv as a CSV file
        else:
            csv_file = open('my_matrix_6x8594.csv', 'a')
            np.savetxt(csv_file, my_matrix)
            csv_file.close()
        i = i + chunks

```

**All Python code in PyCharm screenshots Is available in MiningStoring.py**

- 5) **Exercise deliverable:** (1) Practice reading and writing CSV files from your project data set. Choose a SQLite database file from Tutorials One or Four and (2) execute a SQL script as well as read an SQL SELECT statement into a Pandas dataframe. Save PyCharm screenshots from both deliverables and save for use later in the 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-B 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. Text Mining ..... 1 screenshot
  - ii. Network Mining ..... 1 screenshot
  - iii. Matrix ..... 1 screenshot
  - iv. SQL, NumPy, & Pandas Read/Write CSV ..... 2 screenshots  
Total PyCharm screenshots to submit ..... 5 screenshots
- c. All screenshots are legible output for each PyCharm exercise deliverable.