# **Purpose**

This homework will give you practice using two very important Python modules for data manipulation: Numpy and Pandas. Numpy is a powerful library for performing mathematical and logical operations on Arrays. It provides an abundance of useful features for operations on n-arrays and matrices in Python. Pandas builds upon Numpy adding even more features and the ability to handle tabular data from nonhomogeneous types. Both of these modules are commonly used by data scientists and analysts.

#### Skills

After learning Numpy and Pandas, you'll be able to input, clean, and aggregate large quantities of data, and then use that data with other Python modules such as SciPy that is used for statistical analysis or Matplotlib that is used for visualizing the data.

# Knowledge

Knowledge of these two very powerful modules will add to your overall Python and data manipulation skills. When given a new data set in your future career, you will have some powerful options on how to find answers to questions about the data.

# **Important**

- 1. Due Date: **04/01/2021** at **11:59** pm
- 2. This homework is graded out of **100** points.
- 3. This is an individual assignment. You may collaborate with other students in this class. Collaboration means talking through problems, assisting with debugging, explaining a concept, etc. Students may only collaborate with fellow students currently taking CS 2316, the TA's and the lecturer. You should not exchange code or write code for others. For individual assignments, each student must turn in a unique program. Your submission must not be substantially similar to another student's submission. Collaboration at a reasonable level will not result in substantially similar code.
- 4. For Help:
  - TA Helpdesk (Schedule posted on class website)
  - o Email TA's or use Ed Discussion Notes
  - o Numpy Reference Guide PDF, Numpy User Guide PDF
  - o [Pandas Documentation]
  - Handouts
- 5. Comment out or delete all your function calls. Only global variables, and comments are okay to be outside the scope of a function. When your code is run, all it should do is run without any errors.
- 6. HAVING FUNCTION CALLS OR EXTRANEOUS CODE OUTSIDE THE SCOPE OF FUNCTIONS WILL RESULT IN AN AUTOMATIC 0.

#### 7. IF YOUR CODE CANNOT RUN BECAUSE OF AN ERROR, IT IS A 0%

# Introduction

Please read through the entire document before starting this homework. The goal of this homework is to demonstrate your understanding of Numpy and Pandas by using Numpy arrays and Pandas dataframes. You will be defining a few functions that create and manipulate arrays and dataframes. You should test them out first on your own computer, then when you have one or more of them working, upload the entire file (which must be named HW03.py) to GradeScope to see how well your code performs on the test cases we have created. The score shown on Gradescope may not reflect your final grade on this assignment. You can submit the homework file as many times as you'd like before the deadline.

Allowed imports: numpy, pandas

# **One-line Formatting**

#### **Correct:**

```
def function_a(param1):
    return [i for i in param1]
```

or

```
def function_a(param1):
    return type(param1)
```

#### **Incorrect:**

```
def function_a(param1):
    val = [i for i in param1]
    return val
```

or

```
def function_a(param1):
    return helper(param1)

def helper(param1):
```

```
return [i for i in param1]
```

# **Functions**

Note: The grading rubric at the end of the document for point distribution.

Function name: spring weather

Parameter(s): temperatures ( np.array ), minTemp(int), maxTemp(int)

Return Type: int

**Description:** It's the start of spring, and you're curious what the average weather is like this month. Write a function that takes in temperatures (a np.array representing the temperature each week) and finds the average temperature for all values that are greater than or equal to minTemp and less than maxTemp.

#### There is a one-line maximum requirement.

#### **Python Test Cases**

```
Function name: random sale
```

Parameter(s): items (np.array), cost (np.array), tax(np.array), discount(

np.array) budget (int)
Return Type: np.array

#### **Description:**

Write a function that takes in items (a np.array of items), cost (a np.array of the cost of each items in items), tax (a np.array of the tax rate corresponding to each item in items), discount (a np.array of the discount corresponding to each item in items) and the budget (an int). It is only possible to purchase items if their total cost (cost including tax minus discount) is under the budget. Return a np.array of all the items that can be purchased. Each item corresponds to the same index for all arrays.

## There is a one-line maximum requirement.

#### **Python Test Cases**

Function name: no\_null
Parameter(s): data (np.array)

**Return Type:** np.array

#### **Description:**

Write a function that takes in data (a np.array of the scores). The function should return a np.array where the np.nan values **are replaced with the mean of all values that are NOT np.nan**.

Hint: Use np.isnan()

# There is a one-line maximum requirement.

#### **Python Test Case**

```
82., 108., 100., 108., 18.])
```

Function name: wordle

Parameter(s): player1 ( np.array ), player2 (np.array)

Return Type: bool

**Description:** You and your friend never miss the daily wordle. In order to see who is the superior player, you decide to compare your lowest weekly total guesses. Given player1 and player2 (np.arrays with each row representing a different week and each column representing the number of guesses it took per day), return True if your (player1's) lowest total weekly guesses is less than your friend's (player2's) and False otherwise.

# There is a one-line maximum requirement.

# **Python Test Ca**

Function name: csv\_parser
Parameters: filename (str)
Return Type: pd.DataFrame

**Description**: Write a function that takes in the filename as a string and returns the

Pandas DataFrame representing the corresponding csv file.

# There is a one-line maximum requirement.

**Do not hardcode the filename.** The csv file that will be entered as a parameter is titled cars.csv and contains the following information about each car:

Column name	Data Type	Description
vehicle_id	string	The vehicle identification number is a collection of 17 characters (digits and capital letters)
price	int	Sale price of vehicle in ad
brand	string	Brand of car

model	string	Model of car	
year	int	Car registration year	
condition	string	Either 'clean vehicle' or 'salvage insurance'	
mileage	float	Miles the car has traveled	
color	string	Color of the vehicle	
lot	int	Identification number assigned to a particular quantity or lot of material from a single manufacturer.	
state	string	The state in which the car is being available for purchase	
country	string	The country in which the car is being available for purchase. Either 'usa' or 'canada'.	
bidding_time	string	Time remaining to bid on the vehicle	

## **Python Test Case:**

```
>>> df = csv parser("cars.csv")
               vehicle id price
                                       brand
                                                                country
                                                                          bidding time
                                                        state
0
        jtezu11f88k007763
                                                                          10 days left
                             6300
                                      toyota
                                                   new jersey
                                                                    usa
                                              . . .
1
        2fmdk3gc4bbb02217
                            2899
                                        ford
                                                    tennessee
                                                                    usa
                                                                           6 days left
2
        3c4pdcgg5jt346413
                                                                           2 days left
                            5350
                                       dodge
                                              . . .
                                                      georgia
                                                                    usa
3
        1ftfw1et4efc23745 25000
                                                     virginia
                                                                    usa 22 hours left
                                        ford
4
        3gcpcrec2jg473991 27700
                                  chevrolet
                                                      florida
                                                                    usa 22 hours left
                                              . . .
                              . . .
        3n1cn7ap9k1880319
                            7800
                                              ... california
                                                                           1 days left
2494
                                      nissan
                                                                    นรล
        3n1cn7ap5j1884088
                            9200
                                                                    usa 21 hours left
2495
                                                      florida
                                      nissan ...
2496
        3n1cn7ap9j1884191
                            9200
                                      nissan
                                                      florida
                                                                    usa 21 hours left
                                              . . .
2497
        3n1cn7ap3j1883263
                            9200
                                                      florida
                                                                           2 days left
                                      nissan ...
                                                                    usa
2498
        3n1cn7ap4j1884311
                             9200
                                      nissan ...
                                                      florida
                                                                    usa 21 hours left
>>> df.shape
(2499, 12)
```

Function name: add\_new\_col\_1
Parameters: df (pd.DataFrame)

**Return Type:** NoneType

**Description**: Write a function that takes in the Pandas DataFrame corresponding to the output from  $csv_parser$  and returns the DataFrame with an added column titled score, which is equal to (-1/30000) \* mileage + (year/2022) rounded to **three decimal** places for each car.

#### **Python Test Case:**

```
>>> df = csv parser("cars.csv")
>>> df2 = add new col 1(df)
>>> df2
                                                              bidding time score
                vehicle id price
                                        brand
                                                ... country
0
        jtezu11f88k007763
                              6300
                                                               10 days left -8.144
                                       toyota
                                                        usa
1
        2fmdk3gc4bbb02217
                              2899
                                                                6 days left -5.357
                                         ford
                                                . . .
                                                        usa
2
        3c4pdcqq5jt346413
                                                                2 days left -0.322
                              5350
                                        dodge
                                                        usa
3
        1ftfw1et4efc23745
                            25000
                                         ford
                                                . . .
                                                        usa
                                                              22 hours left -1.142
        3qcpcrec2jq473991
                            27700 chevrolet
                                                              22 hours left, 0.776
4
                                                        usa
                                                . . .
                              . . .
                                                        . . .
2494
        3n1cn7ap9k1880319
                              7800
                                                                1 days left
                                                                             0.212
                                       nissan
                                                . . .
                                                        usa
2495
                                                             21 hours left -0.154
        3n1cn7ap5j1884088
                              9200
                                       nissan
                                                        usa
2496
        3n1cn7ap9j1884191
                              9200
                                       nissan
                                                              21 hours left -0.055
                                                . . .
                                                        usa
2497
        3n1cn7ap3j1883263
                              9200
                                                                2 days left -0.087
                                       nissan
                                               . . .
                                                        usa
2498
        3n1cn7ap4j1884311
                              9200
                                                             21 hours left -0.048
                                       nissan
                                                        บรล
                                                . . .
```

Function name: add\_new\_col\_2
Parameters: df (pd.DataFrame)
Return Type: pd.DataFrame

**Description**: You've developed a very specific taste for cars. You are only interested in a car if its brand is "chevrolet", "dodge", or "ford" and if its color is either "blue" or "silver". Additionally, you would like to have at least one day for bidding.

Write a function that takes in the Pandas DataFrame corresponding to the output from <code>csv\_parser</code> and returns the Pandas DataFrame with an added column titled relative\_price. If the requirements above are not satisfied for a specific car, this column should have the value "Not Interested". Otherwise, you should calculate the signed difference <code>price</code> – <code>avg.price</code> truncated as type integer. (Where <code>avg.price</code> represents the mean price of ALL cars in the list).

# **Python Test Case:**

```
>>> df = csv parser("cars.csv")
>>> df3 = add new col 2(df)
>>> df3
        vehicle id
                                        bidding time relative price
                            price ...
0
        jtezu11f88k007763
                             6300
                                         10 days left Not Interested
                                   . . .
1
        2fmdk3gc4bbb02217
                             2899
                                          6 days left
                                                                -15868
                                   . . .
2
        3c4pdcqq5jt346413
                           5350
                                          2 days left
                                                                -13417
3
        1ftfw1et4efc23745
                            25000
                                        22 hours left Not Interested
                                   . . .
4
        3qcpcrec2jq473991
                            27700
                                        22 hours left Not Interested
                                   . . .
                              . . .
2494
        3n1cn7ap9k1880319
                                        1 days left Not Interested
                             7800
                                   . . .
2495
        3n1cn7ap5j1884088
                             9200
                                        21 hours left Not Interested
                                   . . .
2496
        3n1cn7ap9j1884191
                             9200
                                        21 hours left Not Interested
2497
        3n1cn7ap3j1883263
                             9200
                                          2 days left Not Interested
2498
        3n1cn7ap4j1884311
                             9200
                                        21 hours left Not Interested
                                   . . .
>>> len(df3[df3["relative price"] != "Not Interested"])
189
```

Function name: <a href="mailto:year\_data">year\_data</a>
Parameters: <a href="mailto:df">df</a> ( pd.DataFrame )
Return Type: <a href="mailto:pd.DataFrame">pd.DataFrame</a>

**Description**: You want to see how many different years are available for each make and model of car. For example, if there are only 2008 and 2009 Acura Doors available, the number of years available for that make and model is 2.

**Hint:** Use .agg() or .aggregate()

Note: You should return a pd.DataFrame, not pd.Series.

#### There is a one-line maximum requirement.

#### **Python Test Case**

```
>>> df = csv parser("cars.csv")
>>> df4 = year data(df)
>>> pd.set option('display.max rows', df4.shape[0]+1)
>>> df4
                    year
brand
          model
            door
                      2
acura
           mdx
                        1
           5
audi
                         1
           door
                        2
            q5
                        1
                        1
bmw
            coupe
            door
                        1
            m
            series
                        8
           x3
peterbilt truck 4
ram
           door
                        1
toyota cruiser
>>> df4.shape
(180,1)
```

Function name: cannot\_a\_ford
Parameters: df (pd.DataFrame)
Return Type: pd.DataFrame

**Description**: Retrieve the average price of the ford cars per year of production. Convert the

average to integers.

Hint: Use .agg() or .aggregate()

Note: You should return a pd.DataFrame, not pd.Series.

There is a one-line maximum requirement.

## **Python Test Case:**

```
>>> df = csv parser("cars.csv")
>>> ford_df = cannot_a_ford(df)
>>> print(ford_df)
    price
year
1984
        25
1994
        12
1996
        0
1997
        0
1998
2017 20286
2018 25315
2019 28045
2020 27088
```

Function name: export excel

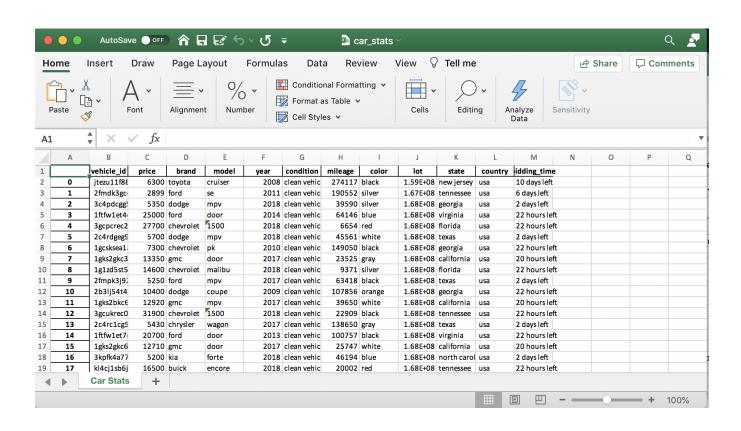
Parameters: df ( pd.DataFrame ), filename ( str ), sheetname ( str )

**Return Type:** NoneType

**Description**: Write a function that takes in the given DataFrame called df and writes to a new excel file with the given filename. You should then name the corresponding sheet name to the name given by sheetname.

## **Python Test Case:**

```
>>> df = csv_parser("cars.csv")
>>> export_excel(df, "car_stats.xlsx", "Car Stats")
```



# **Gradescope Requirements**

- o NO PRINT STATEMENTS this will break the autograder
- NO FUNCTION CALLS outside of function definitions this will also break the autograder
- Make sure the file submitted is titled **HW03.py**
- Do not import any modules, packages, or libraries other than numpy, pandas, pprint
- o Only submit the HW03.py file to Gradescope

# **Grading Rubric**

spring_weather	10	pts
random_sale	10	pts
no_null	10	pts
wordle	10	pts
csv_parser	5	pts
add_new_col_1	10	pts
add_new_col_2	15	pts
year_data	10	pts
cannot_a_ford	15	pts
export_excel	5	pts
Total	100	pts