

# Udacity: Data Science

## Contents

<b>1</b>	<b>Introduction to Data Science</b>	<b>1</b>
1.1	Data Wrangling . . . . .	1
1.2	Data Analysis . . . . .	3

# 1 Introduction to Data Science

## 1.1 Data Wrangling

### SQL Queries in Python:

[Click here](#) to download csv for data.

```
1 import pandas
2 import pandasql
3
4 def num_rainy_days(filename):
5     # The SQL query should return the number of days it rained.
6     weather_data = pandas.read_csv(filename)
7
8     q = """
9     SELECT COUNT(rain)
10    FROM weather_data
11    WHERE CAST(rain AS INTEGER) = 1
12    """
13    rainy_days = pandasql.sqldf(q.lower(), locals())
14    return rainy_days
15
16 def max_temp_aggregate_by_fog(filename):
17     # The SQL query should return the maximum max temp for both foggy and non-foggy
18     # days
19     weather_data = pandas.read_csv(filename)
20
21     q = """
22     SELECT fog, MAX(maxtempi)
23    FROM weather_data
24    GROUP BY 1
25    """
26    foggy_days = pandasql.sqldf(q.lower(), locals())
27    return foggy_days
28
29 def avg_weekend_temperature(filename):
30     # The SQL query should return the average mean temperature on weekends.
31     weather_data = pandas.read_csv(filename)
32
33     q = """
34     SELECT AVG(CAST(meantempi AS INTEGER))
35    FROM weather_data
36    WHERE strftime('%w', date) IN ('0', '6')
37    """
38    mean_temp_weekends = pandasql.sqldf(q.lower(), locals())
39    return mean_temp_weekends
40
41 def avg_min_temperature(filename):
42     # The SQL query should find the average minimum temperature on
43     # rainy days where the minimum temperature is greater than 55 degrees.
44     weather_data = pandas.read_csv(filename)
45
46     q = """
47     SELECT AVG(mintempi)
48    FROM weather_data
49    WHERE (rain = 1) AND (mintempi > 55)
50    """
51    avg_min_temp_rainy = pandasql.sqldf(q.lower(), locals())
52    return avg_min_temp_rainy
```

## Working with Turnstile Data:

[Click here](#) to see an example of an input turnstile file.

[Click here](#) to see an example of the function output turnstile file.

```
1 import csv
2
3 def fix_turnstile_data(filenamees):
4     """
5     Filenames is a list of MTA Subway turnstile text files. As you can see, there
6     are numerous data points included in each row. You want to write a function
7     that will update each row in the text file so there is only one entry per row.
8     """
9     for name in filenamees: # go through each file
10         with open(name, 'r') as f_in, open(''.join(['updated_'+name]), 'w') as f_out:
11             reader = csv.reader(f_in) # csv object to read file
12             writer = csv.writer(f_out) # csv object to write file
13             for row in reader: # go through each row in file
14                 for i in range(3, len(row), 5): # 3 elements repeated, 5 elements unique
15                     writer.writerow(row[0:3] + row[i:i+5])
16     return None
17
18 def create_master_turnstile_file(filenamees, output_file):
19     """
20     Write a function that takes the files in the list filenamees, which all have the
21     columns 'C/A, UNIT, SCP, DATEn, TIMEn, DESCn, ENTRIESn, EXITSn', and consolidates
22     them into one file located at output_file. There should be ONE row with the column
23     headers, located at the top of the file. The input files do not have column header
24     rows of their own.
25     """
26     with open(output_file, 'w') as master_file: # open file to write to
27         master_file.write('C/A,UNIT,SCP,DATEn,TIMEn,DESCn,ENTRIESn,EXITSn\n')
28         for filename in filenamees: # go through each file
29             with open(filename, 'r') as f_in: # open the file to be read
30                 master_file.write(f_in.read()) # write all rows in file to output
31     return None
32
```

## Filtering, Creating, and Re-Formatting Data:

```
1 import pandas
2 import datetime
3
4 def filter_by_regular(filename):
5     """
6     Read the csv file located at filename into a pandas dataframe, and filter
7     the dataframe to only rows where the 'DESCn' column has the value 'REGULAR'.
8     """
9     turnstile_data = pandas.read_csv(filename)
10    turnstile_data = turnstile_data[turnstile_data['DESCn'] == 'REGULAR']
11    return turnstile_data
12
13 def get_hourly_entries(df):
14     # This function should count the number of entries since the last reading
15     df['ENTRIESn_hourly'] = df['ENTRIESn'].diff().fillna(1)
16     return df
17
18 def get_hourly_exits(df):
19     # This function should count the number of exits since the last reading
20     df['EXITSn_hourly'] = df['EXITSn'].diff().fillna(0)
21     return df
22
```

```

1 def time_to_hour(time):
2     # Given a time in the format of "00:00:00" (H:M:S) return hour as an integer.
3     hour = int(time[0:2])
4     return hour
5
6 def reformat_subway_dates(date):
7     '''
8     The dates in our subway data are formatted in the format month-day-year.
9     The dates in our weather underground data are formatted year-month-day.
10    Write a function that takes as its input a date in the MTA Subway
11    data format, and returns a date in the weather underground format.
12    '''
13    date_formatted = datetime.datetime.strptime(date, '%M-%d-%y').strftime('%Y-%M-%d')
14    return date_formatted
15

```

## 1.2 Data Analysis

[Click here](#) to download the data we will be using.

```

1 import numpy
2 import scipy.stats as s
3 import pandas
4 import matplotlib.pyplot as plt
5
6 def entries_histogram(turnstile_weather):
7     '''
8     Let's examine the hourly entries in our NYC subway data and determine what
9     distribution the data follows. Plot two histograms on the same axes to show
10    hourly entries when raining vs. when not raining.
11    '''
12    plt.figure()
13    turnstile_weather['ENTRIESn_hourly'].loc[turnstile_weather['rain'] == 1].hist(bins=30)
14    turnstile_weather['ENTRIESn_hourly'].loc[turnstile_weather['rain'] == 0].hist(bins=30)
15    plt.title('Histogram of ENTRIESn_hourly')
16    plt.legend(['Rain', 'No Rain'])
17    plt.xlabel('ENTRIESn_hourly')
18    plt.ylabel('Frequency')
19    return plt
20
21 def mann_whitney_plus_means(turnstile_weather):
22     '''
23     Take the means and run the Mann Whitney U-test on the ENTRIESn_hourly column.
24     We use this test instead of the Welch t-Test since our data does not follow
25     the normal distribution (seen from graph outputted above).
26     '''
27     with_rain = turnstile_weather['ENTRIESn_hourly'].loc[
28         turnstile_weather['rain'] == 1].values
29     without_rain = turnstile_weather['ENTRIESn_hourly'].loc[
30         turnstile_weather['rain'] == 0].values
31
32     with_rain_mean = np.mean(with_rain)
33     without_rain_mean = np.mean(without_rain)
34
35     U, p = scipy.stats.mannwhitneyu(with_rain, without_rain)
36     return with_rain_mean, without_rain_mean, U, p

```

Note that there is a statistical difference since our  $p < 0.05$ , so we reject  $H_0 : \mu_1 = \mu_2$

**Linear Regression:**

1	
2	
3	
4	