```
1.   Problem Set 3.1
2.   import numpy as np
3.   import pandas
4.   import matplotlib.pyplot as plt
5.
6.   def entries_histogram(turnstile_weather):
7.
8.
9.      plt.figure()
10.     turnstile_weather[turnstile_weather['rain'] == 1]['ENTRIESn_hourly'].hist() # your
     code here to plot a historgram for hourly entries when it is raining
11.     turnstile_weather[turnstile_weather['rain'] == 0]['ENTRIESn_hourly'].hist() # your
     code here to plot a historgram for hourly entries when it is not raining
12.     return plt
13.
14.
15.  Problem Set 3.2
```

Does entries data from the previous exercise seem normally distributed?

◯ Yes

◉ No

Can we run Welch's T test on entries data? Why or why not?

Not until the data is normalized. Suggest a log() transformation. Otherwise, the data is too skewed. A Mann--Whitney U test could be implemented in it's place.

```
16.
17.
18.
19.
20.
21.
22.
23.
24.
25.
26.
27.
28.
29.
30.
31.
32.
33.
34.
```

```
35. Problem Set 3.3
36.
37. import numpy as np
38. import scipy
39. import scipy.stats
40. import pandas
41.
42. def mann_whitney_plus_means(turnstile_weather):
43.
44.     ### YOUR CODE HERE ###
45.     with_rain_mean = np.mean(turnstile_weather[turnstile_weather['rain'] == 1]['ENTRIES
    n_hourly'])
46.     without_rain_mean = np.mean(turnstile_weather[turnstile_weather['rain'] == 0]['ENTR
    IESn_hourly'])
47.     U, p = scipy.stats.mannwhitneyu(turnstile_weather[turnstile_weather['rain'] == 1]['
    ENTRIESn_hourly'],
48.                                     turnstile_weather[turnstile_weather['rain'] == 0]['
    ENTRIESn_hourly'])
49.
50.
51.     return with_rain_mean, without_rain_mean, U, p # leave this line for the grader
52.
53.
54.
55.
56.
57.
58.
59.
60.
61.
62.
63.
64.
65. Problem Set 3.4
```

## Is the distribution of the number of entries statistically different between rainy & non rainy days?

◉ Yes

○ No

## Describe your results and the methods used

I performed a Kolmogorov–Smirnov test which rejected the null hypothesis that the two samples were from the same distribution.

```
66.
67.
```

```
68. Problem Set 3.5
69.
70. import numpy as np
71. import pandas
72. from ggplot import *
73.
74.
75. def normalize_features(df):
76.     mu = df.mean()
77.     sigma = df.std()
78.
79.     if (sigma == 0).any():
80.         raise Exception("One or more features had the same value for all samples, and t
    hus could " + \
81.                         "not be normalized. Please do not include features with only a
     single value " + \
82.                         "in your model.")
83.     df_normalized = (df - df.mean()) / df.std()
84.
85.     return df_normalized, mu, sigma
86.
87. def compute_cost(features, values, theta):
88.
89.     # your code here
90.     cost = np.square(np.dot(features, theta) - values).sum()
91.
92.     return cost
93.
94. def gradient_descent(features, values, theta, alpha, num_iterations):
95.     """
96.     Perform gradient descent given a data set with an arbitrary number of features.
97.
98.     This can be the same gradient descent code as in the lesson #3 exercises,
99.     but feel free to implement your own.
100.         """
101.
102.         m = len(values) * 1.0
103.         cost_history = []
104.
105.         for i in range(num_iterations):
106.             # your code here
107.             update = compute_cost(features, values, theta)
108.             cost_history.append(update)
109.             theta = theta + (1/m) * alpha * np.dot((values - np.dot(features, theta)
    ), features)
110.             return theta, pandas.Series(cost_history)
111.
112.     def predictions(dataframe):
113.         # Select Features (try different features!)
114.         features = dataframe[['Hour', 'maxdewpti', 'maxtempi', 'mindewpti']]
115.
116.         # Add UNIT to features using dummy variables
117.         dummy_units = pandas.get_dummies(dataframe['UNIT'], prefix='unit')
118.         features = features.join(dummy_units)
119.
120.         # Values
121.         values = dataframe['ENTRIESn_hourly']
122.         m = len(values)
123.
124.         features, mu, sigma = normalize_features(features)
125.         features['ones'] = np.ones(m) # Add a column of 1s (y intercept)
```

```
126.
127.            # Convert features and values to numpy arrays
128.            features_array = np.array(features)
129.            values_array = np.array(values)
130.
131.            # Set values for alpha, number of iterations.
132.            alpha = 0.3 # please feel free to change this value
133.            num_iterations = 100 # please feel free to change this value
134.
135.            # Initialize theta, perform gradient descent
136.            theta_gradient_descent = np.zeros(len(features.columns))
137.            theta_gradient_descent, cost_history = gradient_descent(features_array,
138.                                                      values_array,
139.                                                      theta_gradient_desce
     nt,
140.                                                      alpha,
141.                                                      num_iterations)
142.
143.        plot = None
144.        # ------------------------------------------------
145.        # Uncomment the next line to see your cost history
146.        # ------------------------------------------------
147.        plot = plot_cost_history(alpha, cost_history)
148.        #
149.        # Please note, there is a possibility that plotting
150.        # this in addition to your calculation will exceed
151.        # the 30 second limit on the compute servers.
152.
153.        predictions = np.dot(features_array, theta_gradient_descent)
154.        return predictions, plot
155.
156.
157.    def plot_cost_history(alpha, cost_history):
158.        """This function is for viewing the plot of your cost history.
159.        You can run it by uncommenting this
160.
161.            plot_cost_history(alpha, cost_history)
162.
163.        call in predictions.
164.
165.        If you want to run this locally, you should print the return value
166.        from this function.
167.        """
168.        cost_df = pandas.DataFrame({
169.            'Cost_History': cost_history,
170.            'Iteration': range(len(cost_history))
171.        })
172.        return ggplot(cost_df, aes('Iteration', 'Cost_History')) + \
173.            geom_point() + ggtitle('Cost History for alpha = %.3f' % alpha )
174.
175.
176.
177.
178.
179.
180.
181.
182.
183.
184.
185.
```

```
186.        Problem Set 3.6
187.
188.        import numpy as np
189.        import scipy
190.        import matplotlib.pyplot as plt
191.
192.        def plot_residuals(turnstile_weather, predictions):
193.
194.            plt.figure()
195.            (turnstile_weather['ENTRIESn_hourly'] - predictions).hist()
196.            return plt
197.
198.
199.        Problem Set 3.7
200.
201.        import numpy as np
202.        import scipy
203.        import matplotlib.pyplot as plt
204.        import sys
205.        def compute_r_squared(data, predictions):
206.            r_squared = 1 - ((np.square(data - predictions).sum())/(np.square(data - np.m
      ean(data)).sum()))
207.
208.            return r_squared
209.
210.
211.        Problem Set 3.8
212.
213.        import numpy as np
214.        import pandas
215.        import scipy
216.        import statsmodels.api as sm
217.        import datetime
218.
219.        def predictions(weather_turnstile):
220.            #
221.            # Your implementation goes here. Feel free to write additional
222.            # helper functions
223.            #
224.            df = weather_turnstile
225.            dummyunit = pandas.get_dummies(df['UNIT'])
226.            dummyhour = pandas.get_dummies(df['Hour'])
227.            features = dummyunit.join([dummyhour])
228.
229.            entries = df['ENTRIESn_hourly']
230.            prediction = sm.OLS(entries, features).fit().predict(features)
231.
232.            return prediction
```