

Analyzing the NYC Subway Dataset

Questions

Overview

This project consists of two parts. In Part 1 of the project, you should have completed the questions in Problem Sets 2, 3, and 4 in the Introduction to Data Science course. This document addresses part 2 of the project. Please use this document as a template and answer the following questions to explain your reasoning and conclusion behind your work in the problem sets. You will attach a document with your answers to these questions as part of your final project submission.

Section 0. References

https://en.wikipedia.org/wiki/Nonparametric_statistics (https://en.wikipedia.org/wiki/Nonparametric_statistics) https://en.wikipedia.org/wiki/Mann-Whitney_U_test (https://en.wikipedia.org/wiki/Mann-Whitney_U_test) http://dss.princeton.edu/online_help/analysis/interpreting_regression.htm (http://dss.princeton.edu/online_help/analysis/interpreting_regression.htm) https://en.wikipedia.org/wiki/Polynomial_regression (https://en.wikipedia.org/wiki/Polynomial_regression)

Section 1. Statistical Test

1.1 Which statistical test did you use to analyze the NYC subway data? Did you use a one-tail or a two-tail P value? What is the null hypothesis? What is your p-critical value?

ANS :

We choose Mann-Whitney U test two-tail P value.

H_0 is $\mu(\text{rain}) = \mu(\text{notrain})$

We define $\alpha = 0.05$

1.2 Why is this statistical test applicable to the dataset? In particular, consider the assumptions that the test is making about the distribution of ridership in the two samples.

```
In [1]: import pandas as pd
import os
import scipy
import scipy.stats
import numpy as np
import statsmodels.api as sm
alpha = 0.05
filename = r'~/WorkSpace/Udemy-Datascience/IntroductionToDataScience/P2-Analyzing-NYC/turnstile_weather_v2.csv'
turnstile_weather = pd.read_csv(filename)
```


tempi	ENTRIES_hourly
50	1500
60	1850
70	2020
80	2300
90	1700

```
We reject the null hypothesis ( $\mu(\text{with low temp}) \neq \mu(\text{with high temp})$ )  
p-values is 0.01127793160554062598666114070056210039183497428894042968750000000000000000000000  
with low temp mean is 2021.6490019960  
with high temp mean is 1695.4392523364
```

A scatter plot showing the relationship between hourly wind speed (wspeed) on the x-axis and hourly wind energy (ENTRESn_hourly) on the y-axis. The x-axis ranges from -5 to 30, and the y-axis ranges from 1000 to 2400. There are 7 data points plotted, showing a positive correlation. The points are approximately at (0, 1550), (5, 1700), (10, 2250), (15, 2250), (20, 1550), and (25, 1100).

```
We reject the null hypothesis ( $\mu(\text{with low wspdi}) \neq \mu(\text{with high wspdi})$ )
p-values is 0.000000000000530516934301853578265511129698338247816409574131313320322078652679920196533203
with low wspdi mean is 1549.0320479863
with high wspdi mean is 1715.1301188904
```

ANS: We use rain,fog as features in our model and add UNIT, hour, tempi, wspdi and day_week to features using dummy variable.

```
In [10]: features = turnstile_weather[['rain','fog']]
dummy_units = pd.get_dummies(turnstile_weather['UNIT'], prefix='unit')
features = features.join(dummy_units)
dummy_hour = pd.get_dummies(turnstile_weather['hour'], prefix='hour')
features = features.join(dummy_hour)
dummy_day_week = pd.get_dummies(turnstile_weather['day_week'], prefix='day_week')
features = features.join(dummy_day_week)
dummy_tempi = pd.get_dummies(turnstile_weather['tempi_group'], prefix='tempi_group')
features = features.join(dummy_tempi)
dummy_wspdi = pd.get_dummies(turnstile_weather['wspdi_group'], prefix='wspdi_group')
features = features.join(dummy_wspdi)
print len(features.columns )
```

266

2.3 Why did you select these features in your model? We are looking for specific reasons that lead you to believe that

the selected features will contribute to the predictive power of your model.

Your reasons might be based on intuition. For example, response for fog might be: "I decided to use fog because I thought that when it is very foggy outside people might decide to use the subway more often."

Your reasons might also be based on data exploration and experimentation, for example: "I used feature X because as soon as I included it in my model, it drastically improved my R2 value."

ANS:

UNIT,Hour,day of week has obviously effect the ENTRIESn_hourly.

We decide to do not use station because it corresponding to the remote unit.

"rain,fog,tempi and precipi" has an effect by the t-test conclusion.

2.4 What are the parameters (also known as "coefficients" or "weights") of the non-dummy features in your linear regression model?

```
In [11]: X = features
y = np.array(turnstile_weather['ENTRIESn_hourly'])
X = sm.add_constant(X)
olsmod = sm.OLS(y, X)
olsres = olsmod.fit()
ynewpred = olsres.predict(X)
intercept , params = olsres.params[0],olsres.params[1:]
print(olsres.summary())
```

```
OLS Regression Results
=====
Dep. Variable:          y      R-squared:          0.546
Model:                  OLS    Adj. R-squared:      0.544
Method:                 Least Squares      F-statistic:      195.5
Date:                   Sat, 31 Oct 2015    Prob (F-statistic):      0.00
Time:                   17:35:24           Log-Likelihood:      -3.8444e+05
No. Observations:      42649             AIC:              7.694e+05
Df Residuals:          42387             BIC:              7.717e+05
Df Model:              261
Covariance Type:       nonrobust
=====
               coef      std err          t      P>|t|      [95.0% Conf. Int.]
-----
const          1.756e+14   2.74e+14      0.640    0.522   -3.62e+14  7.13e+14
rain           -54.5781     26.934    -2.026    0.043   -107.369   -1.787
fog           -215.4450    101.790    -2.117    0.034   -414.955  -15.935
unit_R003       1.11e+14   4.21e+14     0.264    0.792   -7.14e+14  9.36e+14
unit_R004       1.11e+14   4.21e+14     0.264    0.792   -7.14e+14  9.36e+14
unit_R005       1.11e+14   4.21e+14     0.264    0.792   -7.14e+14  9.36e+14
unit_R006       1.11e+14   4.21e+14     0.264    0.792   -7.14e+14  9.36e+14
unit_R007       1.11e+14   4.21e+14     0.264    0.792   -7.14e+14  9.36e+14
unit_R008       1.11e+14   4.21e+14     0.264    0.792   -7.14e+14  9.36e+14
unit_R009       1.11e+14   4.21e+14     0.264    0.792   -7.14e+14  9.36e+14
unit_R011       1.11e+14   4.21e+14     0.264    0.792   -7.14e+14  9.36e+14
unit_R012       1.11e+14   4.21e+14     0.264    0.792   -7.14e+14  9.36e+14
```

Page 6 of 14

Page 7 of 14

Page 8 of 14


```

day_week_3      -4.988e+14  1.75e+15  -0.285  0.776  -3.93e+15  2.93e+15
day_week_4      -4.988e+14  1.75e+15  -0.285  0.776  -3.93e+15  2.93e+15
day_week_5      -4.988e+14  1.75e+15  -0.285  0.776  -3.93e+15  2.93e+15
day_week_6      -4.988e+14  1.75e+15  -0.285  0.776  -3.93e+15  2.93e+15
tempi_group_50.0  7.354e+14  2.36e+15  0.312  0.755  -3.88e+15  5.35e+15
tempi_group_60.0  7.354e+14  2.36e+15  0.312  0.755  -3.88e+15  5.35e+15
tempi_group_70.0  7.354e+14  2.36e+15  0.312  0.755  -3.88e+15  5.35e+15
tempi_group_80.0  7.354e+14  2.36e+15  0.312  0.755  -3.88e+15  5.35e+15
tempi_group_90.0  7.354e+14  2.36e+15  0.312  0.755  -3.88e+15  5.35e+15
wspdi_group_0.0  -8.626e+14  1.89e+15  -0.456  0.648  -4.57e+15  2.84e+15
wspdi_group_5.0  -8.626e+14  1.89e+15  -0.456  0.648  -4.57e+15  2.84e+15
wspdi_group_10.0 -8.626e+14  1.89e+15  -0.456  0.648  -4.57e+15  2.84e+15
wspdi_group_15.0 -8.626e+14  1.89e+15  -0.456  0.648  -4.57e+15  2.84e+15
wspdi_group_20.0 -8.626e+14  1.89e+15  -0.456  0.648  -4.57e+15  2.84e+15
wspdi_group_25.0 -8.626e+14  1.89e+15  -0.456  0.648  -4.57e+15  2.84e+15
=====
Omnibus:          30229.221  Durbin-Watson:          1.530
Prob(Omnibus):    0.000  Jarque-Bera (JB):      1106543.906
Skew:             2.975  Prob(JB):              0.00
Kurtosis:         27.234  Cond. No.              9.87e+14
=====

```

Warnings:

```

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The smallest eigenvalue is 9.07e-26. This might indicate that there are
strong multicollinearity problems or that the design matrix is singular.

```

ANS:

```

In [12]: print "rain coefficients is %.10f"%params[0]
         print "fog coefficients is %.10f"%params[1]

rain coefficients is -54.5780521669
fog coefficients is -215.4449661619

```

2.5 What is your model's R2 (coefficients of determination) value?

ANS:

```

In [13]: print olsres.rsquared

0.54629540704

```

2.6 What does this R2 value mean for the goodness of fit for your regression model? Do you think this linear model to predict ridership is appropriate for this dataset, given this R2 value?

ANS:

R2 value higher mean that the model is better (the prediction is closer).

As you can see this R2(0.546) is greater than Problem Set 3-5 minimum requirement (0.4) so this linear model is appropriate.

Section 3. Visualization

3.1 One visualization should contain two histograms: one of *ENTRIESn_hourly* for rainy days and one of *ENTRIESn_hourly* for non-rainy days.

You can combine the two histograms in a single plot or you can use two separate plots.

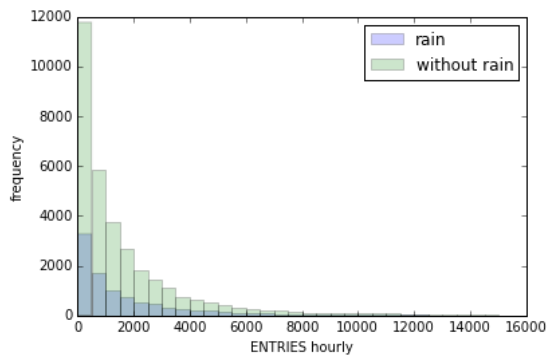
If you decide to use two separate plots for the two histograms, please ensure that the x-axis limits for both of the plots are identical. It is much easier to compare the two in that case.

For the histograms, you should have intervals representing the volume of ridership (value of *ENTRIESn_hourly*) on the x-axis and the frequency of occurrence on the y-axis. For example, each interval (along the x-axis), the height of the bar for this interval will represent the number of records (rows in our data) that have *ENTRIESn_hourly* that falls in this interval.

Remember to increase the number of bins in the histogram (by having larger number of bars). The default bin width is not sufficient to capture the variability in the two samples.

```
In [14]: rain_df = turnstile_weather[turnstile_weather.rain == 1]['ENTRIESn_hourly']
norain_df = turnstile_weather[turnstile_weather.rain == 0]['ENTRIESn_hourly']
df_compare = pd.DataFrame({'rain': rain_df,
                           'without rain': norain_df})
df_compare.plot(kind='hist', alpha=0.2, bins = range(0, 15000 + 500, 500) ) ##,bins = range(0, 37 + 3, 3)

plt.xlabel("ENTRIES hourly")
plt.ylabel("frequency")
plt.show()
```



3.2 One visualization can be more freeform. You should feel free to implement something that we discussed in class (e.g., scatter plots, line plots) or attempt to implement something more advanced if you'd like. Some suggestions are:

Ridership by time-of-day

Ridership by day-of-week

```
In [15]: agg = turnstile_weather.groupby(['day_week'], as_index=False).mean()

x_data = agg['day_week']
y_data = agg['ENTRIESn_hourly']
plt.plot(x_data, y_data, label="average all day")
plt.scatter(x=agg['day_week'], y=agg['ENTRIESn_hourly'])

for hour in range(0,25):
    agg = turnstile_weather[turnstile_weather.hour == hour].groupby(['day_week'], as_index=False).mean()
    x_data = agg['day_week']
    y_data = agg['ENTRIESn_hourly']
    if np.mean(y_data) > 0 :
        plt.plot(x_data, y_data, label="average hour %s"%hour)
        plt.scatter(x=agg['day_week'], y=agg['ENTRIESn_hourly'])

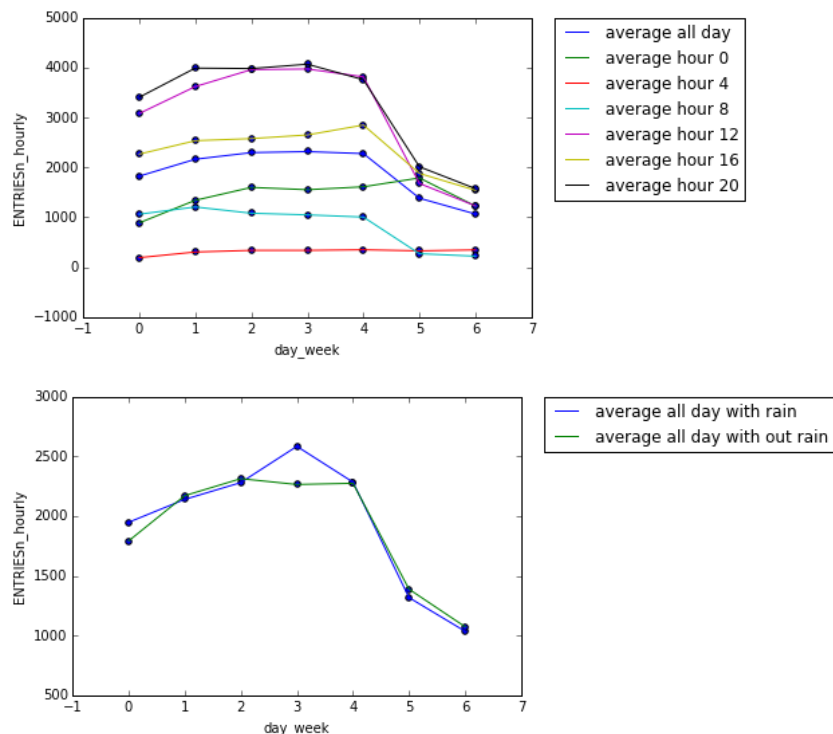
plt.legend(bbox_to_anchor=(1.05, 1), loc=2, borderaxespad=0.)
plt.xlabel("day_week")
plt.ylabel("ENTRIESn_hourly")
plt.show()

agg = turnstile_weather[turnstile_weather.rain == 1].groupby(['day_week'], as_index=False).mean()
x_data = agg['day_week']
y_data = agg['ENTRIESn_hourly']
plt.plot(x_data, y_data, label="average all day with rain")
plt.scatter(x=agg['day_week'], y=agg['ENTRIESn_hourly'])

agg = turnstile_weather[turnstile_weather.rain == 0].groupby(['day_week'], as_index=False).mean()
x_data = agg['day_week']
y_data = agg['ENTRIESn_hourly']
plt.plot(x_data, y_data, label="average all day with out rain")
plt.scatter(x=agg['day_week'], y=agg['ENTRIESn_hourly'])
plt.legend(bbox_to_anchor=(1.05, 1), loc=2, borderaxespad=0.)

plt.xlabel("day_week")
plt.ylabel("ENTRIESn_hourly")

plt.show()
```



Section 4. Conclusion

Please address the following questions in detail. Your answers should be 1-2 paragraphs long.

4.1 From your analysis and interpretation of the data, do more people ride the NYC subway when it is raining or when it is not raining?

From the null hypothesis testing rain has significant effect on Hourly_entries.

From rain's coefficient in linear regression we can conclude that "The people do less ride when it is raining."

4.2 What analyses lead you to this conclusion? You should use results from both your statistical tests and your linear regression to support your analysis.

From the null hypothesis testing results with 2 samples rain and with out rain return p value = it mean that the possibility of different has occurred by chance is about 0.0002741% so we assume that the rain has significantly effect on Hourly_entries.

From the rain coefficient of linear regression is about -71.02 so we assume that the rain has effect on Hourly_entries in negative way.

Section 5. Reflection

Please address the following questions in detail. Your answers should be 1-2 paragraphs long.

5.1 Please discuss potential shortcomings of the methods of your analysis, including:

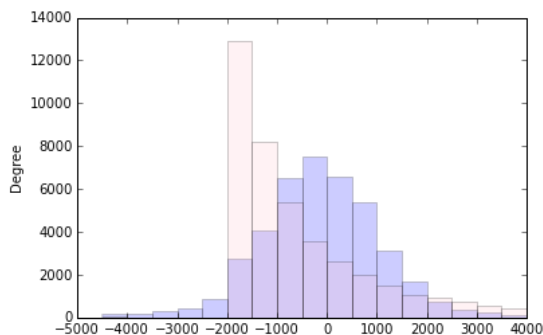
Dataset, Analysis, such as the linear regression model or statistical test.

```
In [20]: def plot_residuals(turnstile_weather, predictions):
plt.figure()
df = (turnstile_weather['ENTRIESn_hourly'] - predictions)
# df.hist()

df.plot(kind='hist',alpha=0.2, y='ENTRIESn_hourly',bins = range(-4500, 4500, 500))

df = (turnstile_weather['ENTRIESn_hourly'] - turnstile_weather['ENTRIESn_hourly'].mean())
df.plot(kind='hist',alpha=0.2, y='ENTRIESn_hourly',bins = range(-4500, 4500, 500) , color="pink")
return plt

predictions = ynewpred
plot_residuals(turnstile_weather, predictions)
plt.show()
print "r^2 is %.6f" % olsres.rsquared
```



r² is 0.546295

By the dataset if we have more parameters that effect "Hourly_entries" we possibly predict the result more accurate such as :

"Does it has nearby special events on the location at that time?" . "Is it holiday?" .

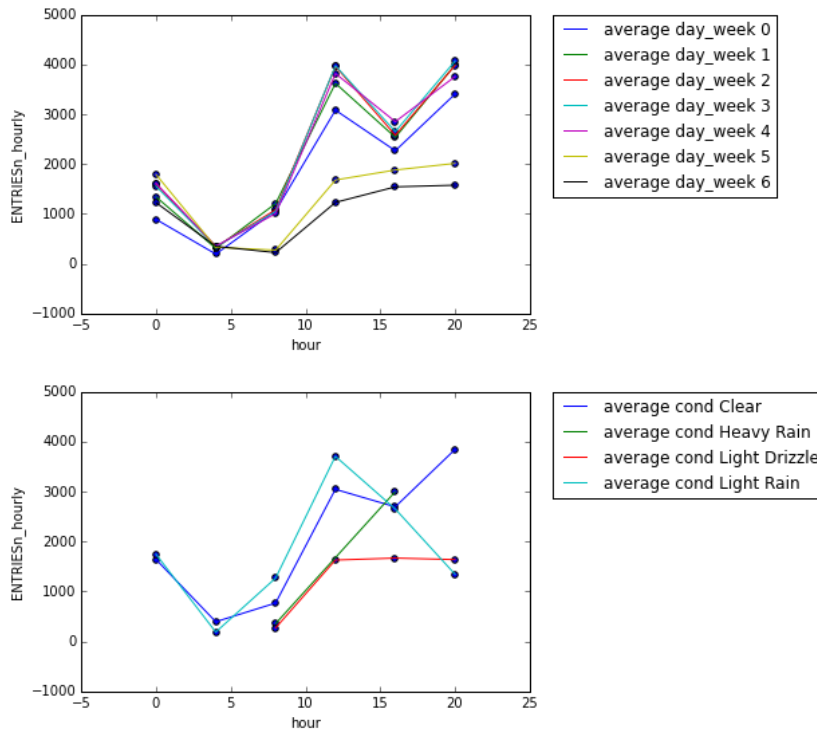
It is possible that the other models such as "Polynomial regression" has a better result.

5.2 (Optional) Do you have any other insight about the dataset that you would like to share with us?

```
In [28]: for day_week in range(0,7):
    agg = turnstile_weather[turnstile_weather.day_week == day_week].groupby(['hour'], as_index=False).mean()
    x_data = agg['hour']
    y_data = agg['ENTRIESn_hourly']
    if np.mean(y_data) > 0 :
        plt.plot(x_data, y_data, label="average day_week %s"%day_week)
        plt.scatter(x=agg['hour'], y=agg['ENTRIESn_hourly'])
plt.legend(bbox_to_anchor=(1.05, 1), loc=2, borderaxespad=0.)
plt.xlabel("hour")
plt.ylabel("ENTRIESn_hourly")
plt.show()

agg2 = turnstile_weather.groupby(['conds'])
for key, cond in agg2['conds']:
    # print key
    if key in ['Clear', 'Light Drizzle', 'Heavy Rain', 'Light Rain']:
        agg = turnstile_weather[turnstile_weather.conds == key].groupby(['hour'], as_index=False).mean()
        x_data = agg['hour']
        y_data = agg['ENTRIESn_hourly']
        # print np.mean(y_data)
        if np.mean(y_data) > 0 :
            plt.plot(x_data, y_data, label="average cond %s"%key)
            plt.scatter(x=agg['hour'], y=agg['ENTRIESn_hourly'])

plt.legend(bbox_to_anchor=(1.05, 1), loc=2, borderaxespad=0.)
plt.xlabel("hour")
plt.ylabel("ENTRIESn_hourly")
plt.show()
```



There are 2 patterns of Hourly_entries by hour first is day_week 0-4 second is dayweek 5-6. Peek Hour for Hourly_entries are about 12:00 and 20:00.

```
In [91]: fig = matplotlib.pyplot.gcf()
fig.set_size_inches(18.5, 10.5)

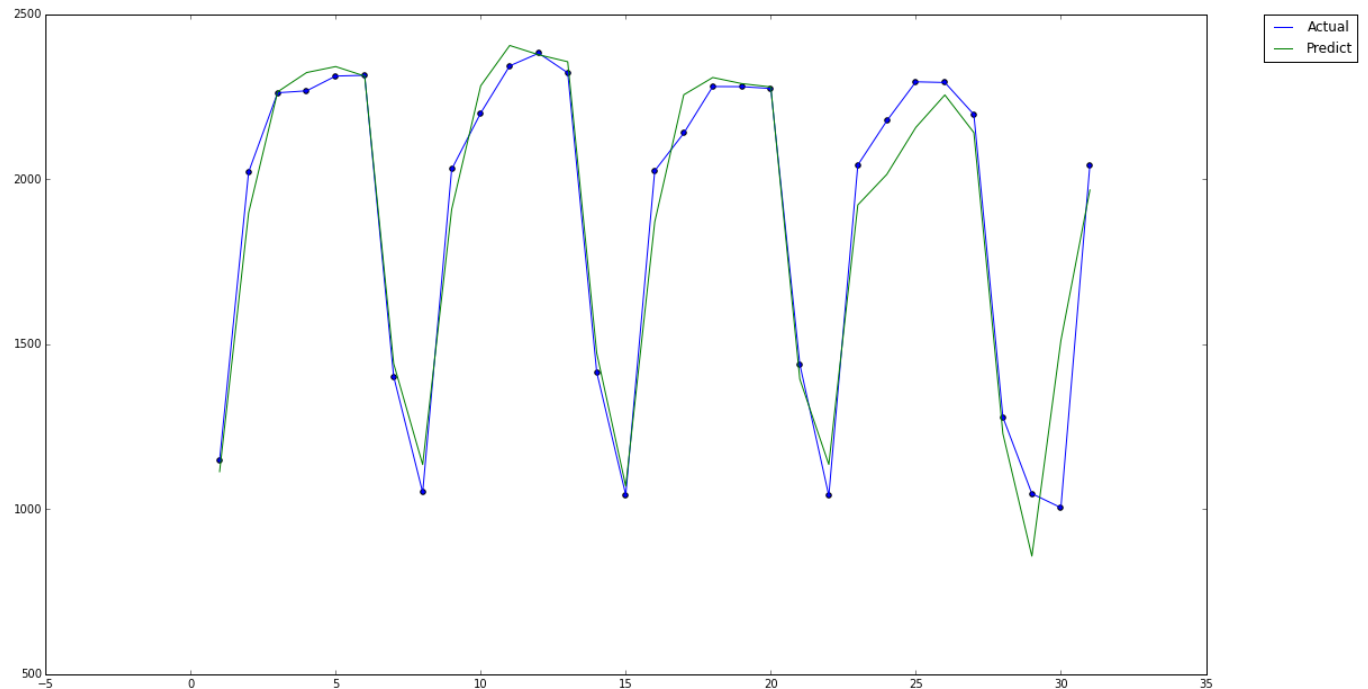
turnstile_weather['predict_ENTRIESn_hourly'] = predictions
agg = turnstile_weather.groupby(['DATEn'], as_index=False).mean()
# agg['datetime_obj'] = pd.to_datetime(agg['DATEn'])
# agg['datetime_obj'] = pd.DatetimeIndex (agg['DATEn'])
# agg['datetime_obj'] = agg['DATEn'].astype('datetime64[ns]')
from time import strftime
# agg = agg[:30]
agg['datetime_obj'] = pd.to_datetime(agg['DATEn'])
# print agg['datetime_obj'][0].day
agg['datetime_num'] = agg['datetime_obj'].apply(lambda x: x.day )

agg = agg.sort(['datetime_num'])
# print len(agg)

import matplotlib.dates as mdates
x_data = agg['datetime_num']
y_data = agg['ENTRIESn_hourly']
y2_data = agg['predict_ENTRIESn_hourly']

plt.scatter(x=agg['datetime_num'], y=agg['ENTRIESn_hourly'])
plt.plot(agg['datetime_num'], agg['ENTRIESn_hourly'],label="Actual")

plt.plot(agg['datetime_num'], agg['predict_ENTRIESn_hourly'],label="Predict")
plt.legend(bbox_to_anchor=(1.05, 1), loc=2, borderaxespad=0.)
plt.show()
```



In []: