

Project Overview

Build a predictive model that identifies the month, hour of day, and appliances (kitchen, laundry, or water heater/ac) that use the most power.

Hypothesis

The most power is used by water heater/ac during months of Dec, Jan, Feb, and July between the hours of 8am - 9pm.

Data Overview/Cleaning

- Drop irrelevant rows
- Change column names for easier readability
- Change date column to month only
- Change time column to hour of day (24 hour clock)
- Remove "I" and null values
- Use MinMaxScaler for Kitchen, LaundryRoom, and WaterHeat_AC columns due to variance in values
- Removed zero values due to imbalanced data

Analysis and Results

- Split data into train and test (80%, 20%)
- Target variable = "Consumer_active_power"
- Run Linear Regression in Scikit Learn and StatsModels
- None of the factors reflected normal distribution!
- Very little correlation for any factors
- Scikit Learn provided score of 0.86
- StatsModels provided R-squared of 0.86

```
In [2]:
import pandas as pd
import numpy as np
import csv
import sys

In [3]:
#df = pd.read_table('household_power_consumption.txt')
with open('household_power_consumption.txt', 'r') as txt_file:
    with open('csv_file.csv', 'w') as csv_file:
        in_txt = csv.reader(txt_file, delimiter = ',')
        out_csv = csv.writer(csv_file)
        out_csv.writerows(in_txt)

In [4]:
df = pd.read_csv('csv_file.csv')
df.head()

/Users/kya/anaconda2/lib/python2.7/site-packages/IPython/core/interactiveshell.py:2717: DtypeWarning: Columns (2,3,4,5,6,7) have mixed types. Specify dtype option on import or set low_memory=False.
Interactivity:interactivity, compiler=compiler, result=result)

Out[4]:
In [5]:
df.shape

Out[5]:
(2075259, 9)

In [6]:
df.drop(df.columns[[3,4,5]], axis=1, inplace=True)

In [7]:
df.head(2)

Out[7]:
In [8]:
df.rename(columns={'Sub_metering_1':'Kitchen',
                  'Sub_metering_2':'LaundryRoom',
                  'Sub_metering_3':'WaterHeat_AC'},
          inplace=True)

In [9]:
df.head(2)

Out[9]:
In [10]:
from datetime import datetime
from datetime.parser import parse

In [11]:
df['Date'] = pd.to_datetime(df['Date'], format='%d/%m/%Y')
df['Date'] = df['Date'].dt.month

In [12]:
df['Time'] = pd.to_datetime(df['Time'], format='%H:%M:%S')
df['Time'] = df['Time'].dt.hour

In [12]:
df.head(2)

Out[12]:
In [13]:
df[['Date', 'Time', 'Global_active_power', 'Kitchen', 'LaundryRoom', 'WaterHeat_AC']]

In [14]:
df['Global_active_power'].value_counts().head()

Out[14]:
In [15]:
df['Kitchen'].value_counts().head(10)

Out[15]:
In [16]:
df['LaundryRoom'].value_counts().head(10)

Out[16]:
In [17]:
df['WaterHeat_AC'].isnull().sum()

Out[17]:
In [18]:
df['Global_active_power'].replace('?', np.nan, inplace=True)

In [19]:
df['Global_active_power'].isnull().value_counts()

Out[19]:
In [20]:
df.dropna(subset=['Global_active_power'], inplace=True)

In [21]:
df['Global_active_power'].isnull().sum()

Out[21]:
In [22]:
df['Kitchen'].replace('?', np.nan, inplace=True)

In [22]:
df['Kitchen'].isnull().value_counts()

Out[22]:
In [24]:
df.dropna(subset=['Kitchen'], inplace=True)

In [25]:
df['Kitchen'].isnull().sum()

Out[25]:
In [26]:
```

```
#df1 = df1.drop(df1[(df1.LaundryRoom == '?').index])
df['LaundryRoom'].replace('?', np.nan, inplace=True)

In [27]:
df['LaundryRoom'].isnull().value_counts()

Out[27]:
False    2049280
Name: LaundryRoom, dtype: int64

In [28]:
df.dropna(subset=['LaundryRoom'], inplace=True)

In [29]:
df['LaundryRoom'].isnull().sum()

Out[29]:
0

In [30]:
df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 2049280 entries, 0 to 2075258
Data columns (total 6 columns):
Date                int64
Time                int64
Global_active_power object
Kitchen             object
LaundryRoom         object
WaterHeat_AC        float64
dtypes: float64(1), int64(2), object(3)
memory usage: 109.4r MB

#

In [31]:
from sklearn.preprocessing import MinMaxScaler
scale = MinMaxScaler()

In [32]:
df[['Kitchen']] = scale.fit_transform(df[['Kitchen']].as_matrix())
/Users/krysl/anaconda2/lib/python2.7/site-packages/sklearn/utils/validation.py:429: DataConversionWarning: Data with input dtype object was converted to float64 by MinMaxScaler.
Warning: warn(msg, _DataConversionWarning)

In [33]:
df['Kitchen'].value_counts()

Out[33]:
0.000000    1880175
0.011364     84936
0.022727     19017
0.431818     14119
0.426455     14892
0.443182      6503
0.408091      5270
0.397727      1359
0.454545       1159
0.363636        802
0.156091        702
0.178455        635
0.375000        603
0.147727        600
0.136364        597
0.352727        586
0.193182        580
0.181818        579
0.113636        576
0.204545        563
0.306818        546
0.125000        541
0.386364        534
0.182727        522
0.238636        510
0.034091        507
0.227273        505
0.215909        487
0.348090        484
0.318182        480
...
0.613636         53
0.818182         46
0.625000         44
0.806818         41
0.875000         34
0.436364         32
0.886364         32
0.794545         29
0.897727         28
0.681818         26
0.908091         19
0.761364         18
0.772727         16
0.715909         16
0.458091         13
0.750000         13
0.478455         13
0.647727         13
0.736364         12
0.727273         12
0.704545         12
0.401182         10
0.784091          9
0.626455          6
0.943182          4
0.986364          3
0.931818          3
1.000000          3
0.954545          2
0.977273          2
Name: Kitchen, dtype: int64

In [34]:
df[['Global_active_power']] = scale.fit_transform(df[['Global_active_power']].as_matrix())

In [35]:
df[['LaundryRoom']] = scale.fit_transform(df[['LaundryRoom']].as_matrix())

In [36]:
df[['WaterHeat_AC']] = scale.fit_transform(df[['WaterHeat_AC']].as_matrix())

In [37]:
df.head(3)

Out[37]:
   Date Time  Global_active_power  Kitchen  LaundryRoom  WaterHeat_AC
0  12  17   0.374796              0.0      0.0125      0.548387
1  12  17   0.478363              0.0      0.0125      0.516129
2  12  17   0.479631              0.0      0.0250      0.548387

In [38]:
df.describe()

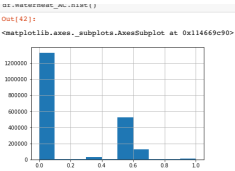
Out[38]:
   Date      Time  Global_active_power  Kitchen  LaundryRoom  WaterHeat_AC
count  2.049280e+06  2.049280e+06  2.049280e+06  2.049280e+06  2.049280e+06  2.049280e+06
mean   6.454453e+00  1.150291e+01  5.194415e-02  1.274813e-02  1.623150e-02  0.853379e-01
std    3.423239e+00  6.829189e+00  8.571738e-02  6.862081e-02  7.277533e-02  2.721663e-01
min    1.000000e+00  0.000000e+00  0.000000e-05  0.000000e+00  0.000000e+00  0.000000e+00
25%    3.000000e+00  2.000000e+00  2.100000e-02  0.000000e+00  0.000000e+00  0.000000e+00
50%    6.000000e+00  2.000000e+01  4.761905e-02  0.000000e+00  0.000000e+00  3.225896e-02
75%    9.000000e+00  1.800000e+01  9.545454e-01  0.000000e+00  0.000000e+00  5.483871e-01
max    1.200000e+01  2.300000e+01  1.000000e+00  1.000000e+00  1.000000e+00  1.000000e+00

In [39]:
import matplotlib.pyplot as plt
matplotlib.rcParams['font.family'] = 'serif'
df.Global_active_power.hist()

Out[39]:
<matplotlib.axes._subplots.AxesSubplot at 0x160cc850>

<matplotlib.axes._subplots.AxesSubplot at 0x160cc850>

```



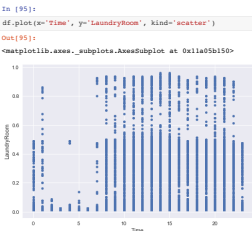
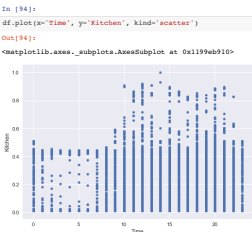
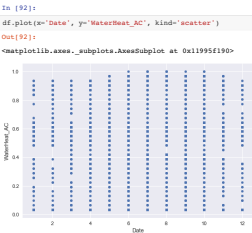
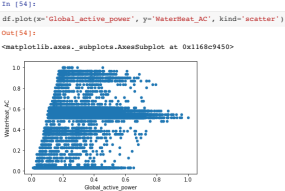
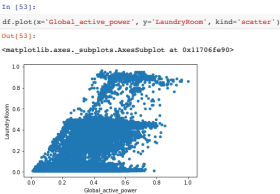
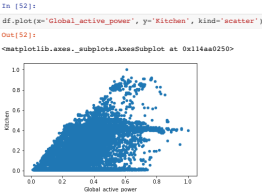
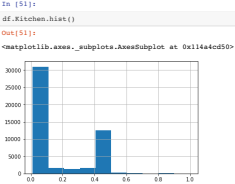
```
In [46]:
print df.groupby('date')['Kitchen']
<pandas.core.groupby.SeriesGroupBy object at 0x1444c590>

In [47]:
df = df[df.Kitchen != 0.0]

In [48]:
df = df[df.LaundryRoom != 0.0]

In [49]:
df = df[df.WaterHeat_AC != 0.0]

In [50]:
df['Kitchen'].value_counts().head()
Out[50]:
0.011364    24533
0.022727     5562
0.420455     4523
0.431818     4188
0.409091     1823
Name: Kitchen, dtype: int64
```



```
In [55]:
df.groupby('Global_active_power')['Kitchen', 'LaundryRoom', 'WaterHeat_AC'].head()
Out[55]:
```

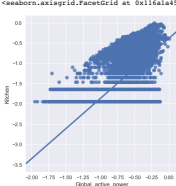
	Kitchen	LaundryRoom	WaterHeat_AC
1019	0.011364	0.0875	0.548387
1020	0.022727	0.4375	0.548387
1021	0.011364	0.5250	0.548387
<v>	0.011364	0.4375	0.548387

	Date	Time	Kitchen	LaundryRoom	WaterHeat_AC
1000	0.011364	0.4375	0.548387		
1024	0.011364	0.4375	0.548387		
1025	0.022727	0.4375	0.548387		
1026	0.011364	0.4375	0.548387		
1027	0.011364	0.4625	0.548387		
1028	0.011364	0.4500	0.548387		
1029	0.011364	0.4625	0.548387		
1030	0.068182	0.4500	0.548387		
1031	0.420455	0.4500	0.548387		
1032	0.420455	0.4375	0.516129		
1033	0.408091	0.4375	0.548387		
1034	0.420455	0.4500	0.516129		
1035	0.408091	0.4375	0.548387		
1036	0.408091	0.4375	0.516129		
1037	0.420455	0.4500	0.548387		
1038	0.408091	0.4375	0.516129		
1039	0.420455	0.4500	0.548387		
1040	0.408091	0.4375	0.548387		
1041	0.425655	0.2875	0.548387		
1042	0.427818	0.2575	0.516129		
1043	0.408091	0.4500	0.548387		
1044	0.284091	0.4375	0.548387		
1045	0.011364	0.3625	0.548387		
1046	0.022727	0.3625	0.516129		
1047	0.011364	0.4500	0.548387		
1048	0.011364	0.4500	0.548387		
...		
2068777	0.011364	0.0875	0.548387		
2068779	0.011364	0.1000	0.548387		
2068812	0.534091	0.0125	0.548387		
2068820	0.227273	0.0125	0.548387		
2068829	0.079545	0.0125	0.548387		
2068901	0.175455	0.0125	0.548387		
2068936	0.188091	0.0125	0.580645		
2068939	0.147727	0.0250	0.548387		
2068940	0.188091	0.0125	0.548387		
2068953	0.011364	0.0125	0.548387		
2068955	0.011364	0.0125	0.548387		
2072999	0.011364	0.0125	0.580645		
2073002	0.011364	0.0250	0.580645		
2073003	0.011364	0.0125	0.580645		
2073013	0.011364	0.0250	0.580645		
2073022	0.420455	0.0125	0.580645		
2073032	0.011364	0.0125	0.580645		
2073040	0.022727	0.0250	0.580645		
2073047	0.011364	0.0250	0.580645		
2073055	0.011364	0.0125	0.580645		
2074471	0.011364	0.0250	0.580645		
2074472	0.022727	0.0125	0.548387		
2074473	0.011364	0.0125	0.548387		
2074478	0.011364	0.0125	0.548387		
2074482	0.022727	0.0250	0.548387		
2074485	0.022727	0.0125	0.548387		
2074488	0.420455	0.0250	0.548387		
2074494	0.431818	0.0250	0.548387		
2074503	0.011364	0.0125	0.580645		
2074504	0.011364	0.0125	0.580645		

16514 rows x 5 columns

```
In [56]:
log_columns = ['Global_active_power', 'Kitchen']
log_df = df.copy()
log_df[log_columns] = log_df[log_columns].apply(np.log10)
```

```
In [57]:
import seaborn as sns
sns.lmplot(log_columns, 'Global_active_power', 'Kitchen', log_df)
Out[57]:
<seaborn.axisgrid.FacetGrid at 0x16a1a150>
```



```
In [58]:
df.describe()
Out[58]:
```

	Date	Time	Global_active_power	Kitchen	LaundryRoom	WaterHeat_AC
Date	1.000000	-0.029794	0.092010	-0.012300	-0.030546	-0.086688
Time	-0.029794	1.000000	0.124053	-0.000787	-0.020223	-0.123444
Global_active_power	0.092010	0.124053	1.000000	0.648848	0.586206	0.305524
Kitchen	-0.012300	-0.000787	0.648848	1.000000	0.006576	0.009362
LaundryRoom	-0.030546	-0.020223	0.586206	0.006576	1.000000	0.056804
WaterHeat_AC	-0.086688	-0.123444	0.305524	0.009362	0.056804	1.000000

```
In [59]:
#From sklearn.model_selection import StratifiedKFold
#X_train, X_test, y_train, y_test = cross_validation.train_test_split(X,y,
#
#
#    test_size=0.20,
#    random_state=0)
```

```
In [60]:
df = df[['Date', 'Time', 'Kitchen', 'LaundryRoom', 'WaterHeat_AC', 'Global_active_power']]
```

```
In [61]:
df.head()
```

```
Out[61]:
```

	Date	Time	Kitchen	LaundryRoom	WaterHeat_AC	Global_active_power
1018	12	10	0.011364	0.0875	0.548387	0.190380
1020	12	10	0.022727	0.4375	0.548387	0.329350
1024	12	10	0.011364	0.3200	0.548387	0.283904
1022	12	10	0.011364	0.4375	0.548387	0.327539
1033	12	10	0.011364	0.3375	0.548387	0.283179

```
In [62]:
#From sklearn.cross_validation import train_test_split
#Users/Users/anaconda2/lib/python2.7/site-packages/sklearn/cross_validation.py:44: DeprecationWarning: This module was deprecated in version 0.18 in favor of the model_selection module into which all the refactored classes and functions were moved. Also note that the interface of the new CV iterators are different from that o
#f this module. This module will be removed in 0.20.
#    "This module will be removed in 0.20.", DeprecationWarning)
```

```
In [63]:
train, test = train_test_split(df, train_size=80, test_size=20)
```

```
In [64]:
from sklearn import linear_model
from sklearn.linear_model import LinearRegression
lin_model = linear_model.LinearRegression()
```

```
In [80]:
feature_cols = ('Date', 'Time', 'Kitchen', 'LaundryRoom', 'WaterHeat_AC')
X = train[feature_cols]
y = train['Global_active_power']
```

```
In [81]:
lin_model.fit(X,y)
Out[81]:
LinearRegression(copy_X=True, fit_intercept=True, n_jobs=1, normalize=False)
```

```
In [86]:
print lin_model.intercept_
print lin_model.coef_
-0.0232587114105
[-1.17121758e-04  4.97507939e-03  4.93865726e-01  4.63246609e-01
 1.97212089e-01]
```

```
In [84]:
pd.DataFrame(zip(train.columns, lin_model.coef_), columns = ['factors', 'est_coef'])
Out[84]:
```

	factors	est_coef
0	Date	-0.000117
1	Time	0.004975
2	Kitchen	0.493866
3	LaundryRoom	0.463247
4	WaterHeat_AC	0.197212

```
In [88]:
```

In () :