

Enron email data set exploration

```
In [1]: # Get better looking pictures
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%config InlineBackend.figure_format = 'retina'

In [2]: df = pd.read_feather('enron.feather')
df['Date'] = pd.to_datetime(df['Date'], format='%Y-%m-%d')
df = df.sort_values(['Date'])
df.tail(5)
```

MailID	Date	From	To	recipients	Subject	filename
603555	2002-07-12	denise.williams	ge_benefits	1	URGENT!!! CUTOVER WEEKEND	fischer-m/discussion_threads/39.
603971	2002-07-12	kurt.anderson	mark.fisher	3	FW: RE: Revised Availability Numbers	fischer-m/discussion_threads/37.
601861	2002-07-12	mark.fisher	tom.nemila	1	WR613 Pitch System performance	fischer-m/all_documents/425.
603131	2002-07-12	mark.fisher	tom.nemila	1	WR613 Pitch System performance	fischer-m/discussion_threads/36.
604569	2002-07-12	kurt.anderson	mark.walker	3	FW: RE: Revised Availability Numbers	fischer-m/notes_inbox/4.

Email traffic over time

Group the data set by Date and MailID, which will get you an index that collects all of the unique mail IDs per date. Then reset the index so that those date and mail identifiers become columns and then select for just those columns; we don't actually care about the counts created by the groupby (that was just to get the index). Create a histogram that shows the amount of traffic per day, then specifically for email sent from richard.shapiro and then john.lavorato . Because some dates are set improperly (to 1980), filter specifically for dates greater than January 1, 1999.

```
In [3]: # filter for dates greater than January 1, 1999
dfQ1 = df[df['Date'] > pd.Timestamp('1999-1-1')]
dfQ1 = dfQ1[['Date','MailID']].reset_index()
dfQ1 = dfQ1.groupby(['Date','MailID']).count().reset_index()
dfQ1 = dfQ1[['Date','MailID']]
```

```
Out[3]: Date MailID
0 1999-01-04 318815
1 1999-01-04 321609
2 1999-01-04 327119
3 1999-01-04 327904
4 1999-01-05 319114
...
347354 2002-07-12 72772
347355 2002-07-12 72905
347356 2002-07-12 73118
347357 2002-07-12 73154
347358 2002-07-12 73164
347359 rows x 2 columns
```

```
In [4]: fig, ax = plt.subplots(figsize=(15,4))
for i in patches:
    rect.set_linewidth(.5)
    rect.set_edgecolor("black")
    ax.spines['top'].set_visible(False)
    ax.spines['bottom'].set_visible(False)
    ax.spines['left'].set_linewidth(0.5)
    ax.spines['right'].set_linewidth(0.5)
    plt.xticks(rotation=45, size = 17)
    plt.yticks(np.arange(0, 12500, 2500), size = 17)
    ax.set_xlabel("Date",size = 17.5)
    ax.set_title("Histogram of emails sent",size = 19.5)
plt.show()
```



