

Exploratory Data Analysis

Introduction

After the data cleaning step where we put our data into a few standard formats, the next step is to take a look at the data and see if what we're looking at makes sense. Before applying any fancy algorithms, it's always important to explore the data first.

When working with numerical data, some of the exploratory data analysis (EDA) techniques we can use include finding the average of the data set, the distribution of the data, the most common values, etc. The idea is the same when working with text data. We are going to find some more obvious patterns with EDA before identifying the hidden patterns with machines learning (ML) techniques. We are going to look at the following for each comedian:

1. **Most common words** - find these and create word clouds
2. **Size of vocabulary** - look number of unique words and also how quickly someone speaks
3. **Amount of profanity** - most common terms

Most Common Words

Analysis

```
In [1]: 1 # Read in the document-term matrix
        2 import pandas as pd
        3
        4 data = pd.read_pickle('dtm.pkl')
        5 data = data.transpose()
        6 data.head()
```

Out[1]:

	ali	anthony	bill	bo	dave	hasan	jim	joe	john	louis	mike	ricky
aaaaah	0	0	1	0	0	0	0	0	0	0	0	0
aaaaahhhhhh	0	0	0	1	0	0	0	0	0	0	0	0
aaaaauugghhhhhh	0	0	0	1	0	0	0	0	0	0	0	0
aaaahhhh	0	0	0	1	0	0	0	0	0	0	0	0
aaah	0	0	0	0	1	0	0	0	0	0	0	0


```
In [3]: 1 # Print the top 15 words said by each comedian
2 for comedian, top_words in top_dict.items():
3     print(comedian)
4     print(', '.join([word for word, count in top_words[0:14]]))
5     print('---')

ali
like, im, know, just, dont, shit, thats, youre, gonna, ok, lot, wanna, gotta, oh
---
anthony
im, like, know, dont, joke, got, thats, said, anthony, day, say, just, guys, people
---
bill
like, just, right, im, know, dont, gonna, got, fucking, yeah, shit, youre, thats, dude
---
bo
know, like, think, im, love, bo, just, stuff, repeat, dont, yeah, want, right, cos
---
dave
like, know, said, just, im, shit, people, didnt, ahah, dont, time, thats, fuck, fucking
---
hasan
like, im, know, dont, dad, youre, just, going, thats, want, got, love, shes, hasan
---
jim
like, im, dont, right, fucking, know, just, went, youre, people, thats, day, oh, think
---
joe
like, people, just, dont, fucking, im, fuck, thats, gonna, theyre, know, youre, think, shit
---
john
like, know, just, dont, said, clinton, im, thats, right, youre, little, hey, got, time
---
louis
like, just, know, dont, thats, im, youre, life, people, thing, gonna, hes, cause, theres
---
mike
like, im, know, said, just, dont, think, thats, says, cause, right, jenny, goes, id
---
ricky
right, like, just, im, dont, know, said, yeah, fucking, got, say, youre, went, id
---
```

NOTE: At this point, we could go on and create word clouds. However, by looking at these top words, you can see that some of them have very little meaning and could be added to a stop words list, so let's do just that.

```
In [4]: 1 # Look at the most common top words --> add them to the stop word list
2 from collections import Counter
3
4 # Let's first pull out the top 30 words for each comedian
5 words = []
6 for comedian in data.columns:
7     top = [word for (word, count) in top_dict[comedian]]
8     for t in top:
9         words.append(t)
10
11 words
```

```
Out[4]: ['like',
'im',
'know',
'just',
'dont',
'shit',
'thats',
'youre',
'gonna',
'ok',
'lot',
'wanna',
'gotta',
'oh',
'husband',
'got',
'right',
'time',
'cause',
.....]
```

```
In [5]: 1 # Let's aggregate this list and identify the most common words along with how many routines they occur in
2 Counter(words).most_common()
```

```
Out[5]: [('like', 12),
('im', 12),
('know', 12),
('just', 12),
('dont', 12),
('thats', 12),
('right', 12),
('people', 12),
('youre', 11),
('got', 10),
('time', 9),
('gonna', 8),
('think', 8),
('oh', 7),
('yeah', 7),
('said', 7),
('cause', 6),
('hes', 6),
('theyre', 6),
.....]
```

```
In [6]: 1 # If more than half of the comedians have it as a top word, exclude it from the list
2 add_stop_words = [word for word, count in Counter(words).most_common() if count > 6]
3 add_stop_words
```

```
Out[6]: ['like',
'im',
'know',
'just',
'dont',
'thats',
'right',
'people',
'youre',
'got',
'time',
'gonna',
'think',
'oh',
'yeah',
'said']
```

```
In [7]: 1 # Let's update our document-term matrix with the new list of stop words
2 from sklearn.feature_extraction import text
3 from sklearn.feature_extraction.text import CountVectorizer
4
5 # Read in cleaned data
6 data_clean = pd.read_pickle('data_clean.pkl')
7
8 # Add new stop words
9 stop_words = text.ENGLISH_STOP_WORDS.union(add_stop_words)
10
11 # Recreate document-term matrix
12 cv = CountVectorizer(stop_words=stop_words)
13 data_cv = cv.fit_transform(data_clean.transcript)
14 data_stop = pd.DataFrame(data_cv.toarray(), columns=cv.get_feature_names())
15 data_stop.index = data_clean.index
16
17 # Pickle it for later use
18 import pickle
19 pickle.dump(cv, open("cv_stop.pkl", "wb"))
20 data_stop.to_pickle("dtm_stop.pkl")
```

```
In [8]: 1 # Let's make some word clouds!
2 # Terminal / Anaconda Prompt: conda install -c conda-forge wordcloud
3 from wordcloud import WordCloud
4
5 wc = WordCloud(stopwords=stop_words, background_color="white", colormap="Dark2",
6               max_font_size=150, random_state=42)
```

```
In [9]: 1 # Reset the output dimensions
2 import matplotlib.pyplot as plt
3
4 plt.rcParams['figure.figsize'] = [16, 6]
5
6 full_names = ['Ali Wong', 'Anthony Jeselnik', 'Bill Burr', 'Bo Burnham', 'Dave Chappelle', 'Hasan Minhaj',
7              'Jim Jefferies', 'Joe Rogan', 'John Mulaney', 'Louis C.K.', 'Mike Birbiglia', 'Ricky Gervais']
8
9 # Create subplots for each comedian
10 for index, comedian in enumerate(data.columns):
11     wc.generate(data_clean.transcript[comedian])
12
13     plt.subplot(3, 4, index+1)
14     plt.imshow(wc, interpolation="bilinear")
15     plt.axis("off")
16     plt.title(full_names[index])
17
18 plt.show()
```



Findings

- Ali Wong says the s-word a lot and talks about her husband. I guess that's funny to me.
- A lot of people use the F-word. Let's dig into that later.

Number of Words

Analysis

```
In [10]: 1 # Find the number of unique words that each comedian uses
2
3 # Identify the non-zero items in the document-term matrix, meaning that the word occurs at least once
4 unique_list = []
5 for comedian in data.columns:
6     uniques = data[comedian].nonzero()[0].size
7     unique_list.append(uniques)
8
9 # Create a new dataframe that contains this unique word count
10 data_words = pd.DataFrame(list(zip(full_names, unique_list)), columns=['comedian', 'unique_words'])
11 data_unique_sort = data_words.sort_values(by='unique_words')
12 data_unique_sort
```

```
<ipython-input-10-ae066bdc1e39>:6: FutureWarning: Series.nonzero() is deprecated and will be removed in a future version. Use Series.to_numpy().nonzero() instead
    uniques = data[comedian].nonzero()[0].size
```

Out[10]:

	comedian	unique_words
1	Anthony Jeselnik	984
9	Louis C.K.	1098
3	Bo Burnham	1272
6	Jim Jefferies	1313
0	Ali Wong	1341
8	John Mulaney	1391
4	Dave Chappelle	1404
7	Joe Rogan	1435
10	Mike Birbiglia	1494
5	Hasan Minhaj	1559
2	Bill Burr	1633
11	Ricky Gervais	1633

```

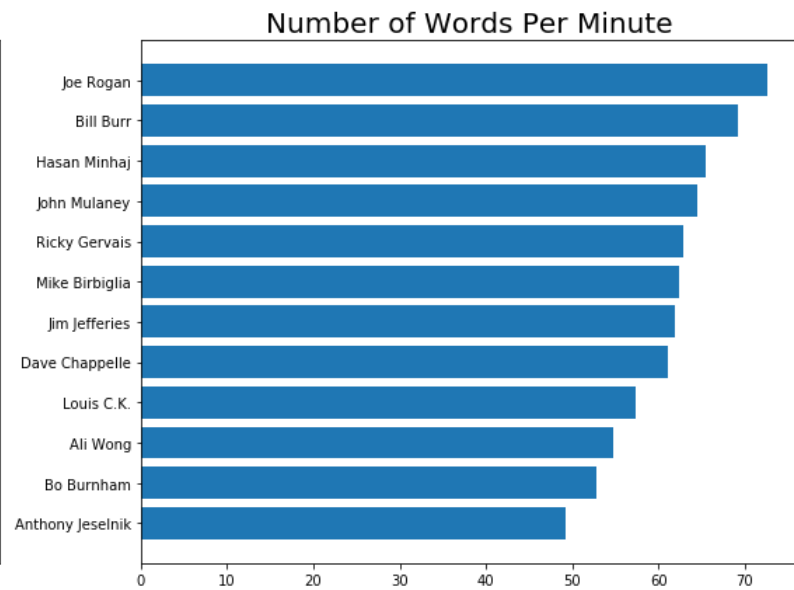
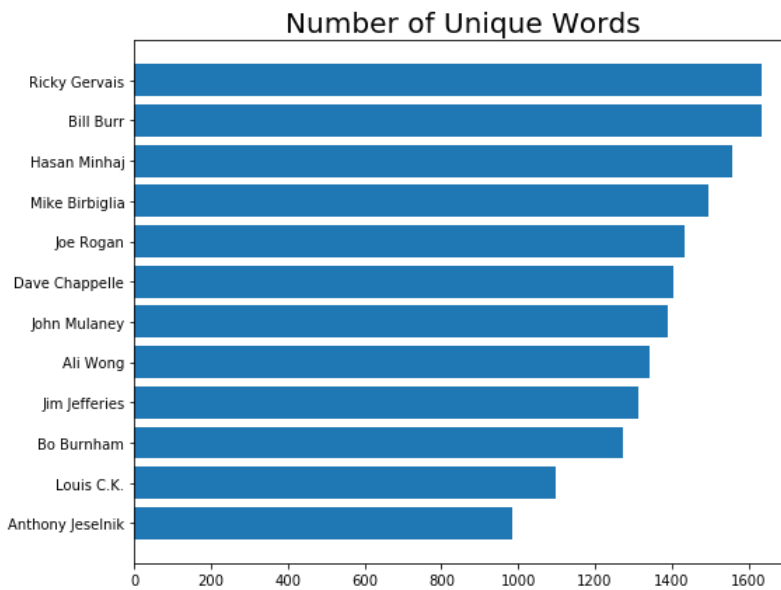
In [11]: 1 # Calculate the words per minute of each comedian
          2
          3 # Find the total number of words that a comedian uses
          4 total_list = []
          5 for comedian in data.columns:
          6     totals = sum(data[comedian])
          7     total_list.append(totals)
          8
          9 # Comedy special run times from IMDB, in minutes
         10 run_times = [60, 59, 80, 60, 67, 73, 77, 63, 62, 58, 76, 79]
         11
         12 # Let's add some columns to our dataframe
         13 data_words['total_words'] = total_list
         14 data_words['run_times'] = run_times
         15 data_words['words_per_minute'] = data_words['total_words'] / data_words['run_times']
         16
         17 # Sort the dataframe by words per minute to see who talks the slowest and fastest
         18 data_wpm_sort = data_words.sort_values(by='words_per_minute')
         19 data_wpm_sort

```

Out[11]:

	comedian	unique_words	total_words	run_times	words_per_minute
1	Anthony Jeselnik	984	2905	59	49.237288
3	Bo Burnham	1272	3165	60	52.750000
0	Ali Wong	1341	3283	60	54.716667
9	Louis C.K.	1098	3332	58	57.448276
4	Dave Chappelle	1404	4094	67	61.104478
6	Jim Jefferies	1313	4764	77	61.870130
10	Mike Birbiglia	1494	4741	76	62.381579
11	Ricky Gervais	1633	4972	79	62.936709
8	John Mulaney	1391	4001	62	64.532258
5	Hasan Minhaj	1559	4777	73	65.438356
2	Bill Burr	1633	5535	80	69.187500
7	Joe Rogan	1435	4579	63	72.682540


```
In [12]: 1 # Let's plot our findings
2 import numpy as np
3
4 y_pos = np.arange(len(data_words))
5
6 plt.subplot(1, 2, 1)
7 plt.barh(y_pos, data_unique_sort.unique_words, align='center')
8 plt.yticks(y_pos, data_unique_sort.comedian)
9 plt.title('Number of Unique Words', fontsize=20)
10
11 plt.subplot(1, 2, 2)
12 plt.barh(y_pos, data_wpm_sort.words_per_minute, align='center')
13 plt.yticks(y_pos, data_wpm_sort.comedian)
14 plt.title('Number of Words Per Minute', fontsize=20)
15
16 plt.tight_layout()
17 plt.show()
```



Findings

- **Vocabulary**
 - Ricky Gervais (British comedy) and Bill Burr (podcast host) use a lot of words in their comedy
 - Louis C.K. (self-deprecating comedy) and Anthony Jeselnik (dark humor) have a smaller vocabulary
- **Talking Speed**
 - Joe Rogan (blue comedy) and Bill Burr (podcast host) talk fast
 - Bo Burnham (musical comedy) and Anthony Jeselnik (dark humor) talk slow

Ali Wong is somewhere in the middle in both cases. Nothing too interesting here.

Amount of Profanity

Analysis

```
In [13]: 1 # Earlier I said we'd revisit profanity. Let's take a look at the most common words again.
          2 Counter(words).most_common()
```

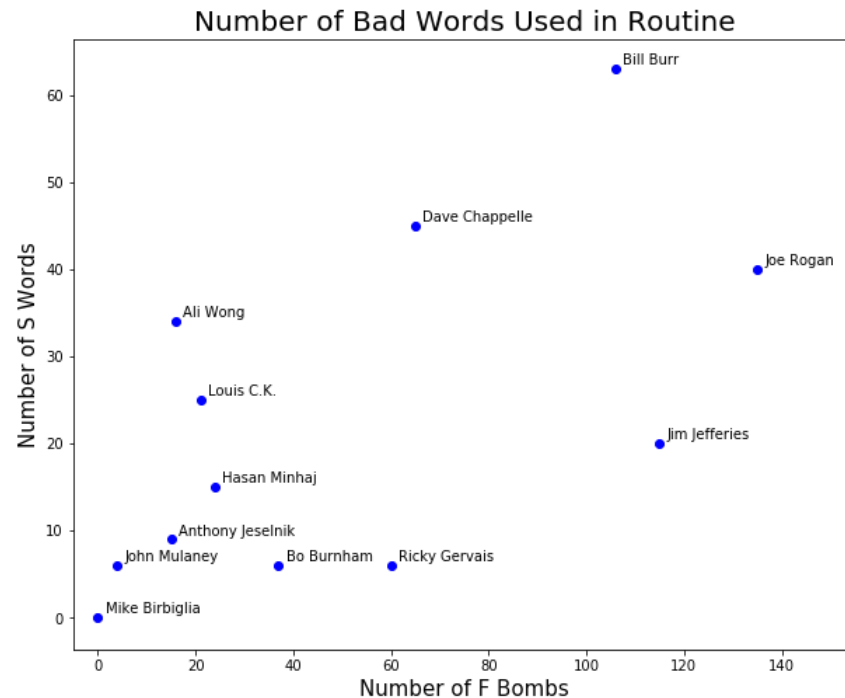
```
Out[13]: [('like', 12),
           ('im', 12),
           ('know', 12),
           ('just', 12),
           ('dont', 12),
           ('thats', 12),
           ('right', 12),
           ('people', 12),
           ('youre', 11),
           ('got', 10),
           ('time', 9),
           ('gonna', 8),
           ('think', 8),
           ('oh', 7),
           ('yeah', 7),
           ('said', 7),
           ('cause', 6),
           ('hes', 6),
           ('theyre', 6),
           ('that', 6)]
```

```
In [14]: 1 # Let's isolate just these bad words
2 data_bad_words = data.transpose()[['fucking', 'fuck', 'shit']]
3 data_profanity = pd.concat([data_bad_words.fucking + data_bad_words.fuck, data_bad_words.shit], axis=1)
4 data_profanity.columns = ['f_word', 's_word']
5 data_profanity
```

Out[14]:

	f_word	s_word
ali	16	34
anthony	15	9
bill	106	63
bo	37	6
dave	65	45
hasan	24	15
jim	115	20
joe	135	40
john	4	6
louis	21	25
mike	0	0
ricky	60	6

```
In [15]: 1 # Let's create a scatter plot of our findings
2 plt.rcParams['figure.figsize'] = [10, 8]
3
4 for i, comedian in enumerate(data_profanity.index):
5     x = data_profanity.f_word.loc[comedian]
6     y = data_profanity.s_word.loc[comedian]
7     plt.scatter(x, y, color='blue')
8     plt.text(x+1.5, y+0.5, full_names[i], fontsize=10)
9     plt.xlim(-5, 155)
10
11 plt.title('Number of Bad Words Used in Routine', fontsize=20)
12 plt.xlabel('Number of F Bombs', fontsize=15)
13 plt.ylabel('Number of S Words', fontsize=15)
14
15 plt.show()
```



Findings

- **Averaging 2 F-Bombs Per Minute!** - I don't like too much swearing, especially the f-word, which is probably why I've never heard of Bill Burr, Joe Rogan and Jim Jefferies.
- **Clean Humor** - It looks like profanity might be a good predictor of the type of comedy I like. Besides Ali Wong, my two other favorite comedians in this group are John Mulaney and Mike Birbiglia.

Side Note

What was our goal for the EDA portion of our journey? **To be able to take an initial look at our data and see if the results of some basic analysis made sense.**

My conclusion - yes, it does, for a first pass. There are definitely some things that could be better cleaned up, such as adding more stop words or including bi-grams. But we can save that for another day. The results, especially the profanity findings, are interesting and make general sense, so we're going to move on.

As a reminder, the data science process is an iterative one. It's better to see some non-perfect but acceptable results to help you quickly decide whether your project is a dud or not, instead of having analysis paralysis and never delivering anything.

Alice's data science (and life) motto: Let go of perfectionism!

Additional Exercises

1. What other word counts do you think would be interesting to compare instead of the f-word and s-word? Create a scatter plot comparing them.

In []:

1
