

## Data Collection and Cleaning

```
1 import pandas as pd
2 import requests
3 from collections import Counter
4 from bs4 import BeautifulSoup
5 import time
6
7 response = requests.get("https://www.gutenberg.org/browse/scores/top#books-last1")#a list of 100 most popular books
8 soup = BeautifulSoup(response.content, "html.parser")
9 ol = soup.find("ol")
10 book_nums = {}
11 for a in ol.find_all("a"):
12     book_nums[str(a).split()[1].split('\"')[1].split("/")[2]] = a.text #key-book number      value-Book name
13 "number of books: "+str(len(book_nums)), list(book_nums.values())[:10]#Only display 10 book names
```

```
➞ ('number of books: 100',
   ['Frankenstein; Or, The Modern Prometheus by Mary Wollstonecraft Shelley (4678)',
    'Pride and Prejudice by Jane Austen (2474)',
    'The Great Gatsby by F. Scott Fitzgerald (1387)',
    'Et dukkehjem. English by Henrik Ibsen (1275)',
    'A Tale of Two Cities by Charles Dickens (1179)',
    "Alice's Adventures in Wonderland by Lewis Carroll (1150)",
    'The Importance of Being Earnest: A Trivial Comedy for Serious People by Oscar Wilde (1050)',
    'The Strange Case of Dr. Jekyll and Mr. Hyde by Robert Louis Stevenson (964)',
    'A Modest Proposal by Jonathan Swift (963)',
    'The Picture of Dorian Gray by Oscar Wilde (959)'])
```

```
1 ### Reads all books into a series
2 ### Series contains-
3 ### index: Book num      Value: full cleaned text
4 books = pd.Series(dtype="object")
5 whitelist = set('abcdefghijklmnopqrstuvwxyz ')
6 for i in book_nums.keys():
7     infile = requests.get("https://www.gutenberg.org/files/%s/%s-0.txt"%(i, i))#opens url for each book num
8     if not infile.text[:15] == "<!DOCTYPE html>": # Checks to make sure url is valid
9         infile.encoding = "UTF-8" # Correct encoding
```

```

10     book = infile.text.replace("\n", " ").replace("\r", " ").lower() #removes type-ins, makes lowercase
11     book = "".join(filter(whitelist.__contains__, book)).split() # Removes all characters not in whitelist
12     books.loc[book_nums[i]] = book[300:len(book)-3000] # Cuts off header and copyright info at top and bottom
13     time.sleep(.25)
14 books

```

```

Frankenstein; Or, The Modern Prometheus by Mary Wollstonecraft Shelley (4678)
Pride and Prejudice by Jane Austen (2474)
The Great Gatsby by F. Scott Fitzgerald (1387)
Et dukkehjem. English by Henrik Ibsen (1275)
A Tale of Two Cities by Charles Dickens (1179)

```

```

[of, those, icy, climes, inspirited, t
[is, let, at, last, mr, bennet, replie
[of, not, a, few, veteran, bores, the,
[burns, in, the, stove, it, is, winter
[the, loadstone, rock, book, the, thir

```

```

The Life and Adventures of Robinson Crusoe by Daniel Defoe (218)
Sense and Sensibility by Jane Austen (214)
David Copperfield by Charles Dickens (208)
Hard Times by Charles Dickens (206)
The Alhambra by Albert Frederick Calvert (203)
Length: 79, dtype: object

```

```

...
[merchandise, and, leaving, off, his,
[owner, of, this, estate, was, a, sing
[a, long, journey, xxxiii, blissful, x
[and, rachael, in, the, sickroom, mr,
[for, this, volume, which, is, humbly,

```

All of the following commented out code works but I opted to not use it because the amount of data it produces is insignificant to the book data

```

1 #books.drop("La Maternelle by Léon Frapié (180)", inplace = True) # That book is in french, only need english
2 #books

```

```

1 # wiki_dir = "https://en.wikipedia.org/wiki/Wikipedia:Most-referenced_articles"
2 # soup = BeautifulSoup(requests.get(wiki_dir).content, "html.parser")
3 # tab = soup.find_all("table")[1]
4 # wiki_articles = set()
5 # for a in tab.find_all("a"):
6 #     a = str(a).split()[1]
7 #     wiki_articles.add(a[6:len(a)-1])
8 #     time.sleep(.25)
9 # wiki_articles.remove('mw-redirect')
10 # "Number of wikipedia articles: " +str(len(wiki_articles)), list(wiki_articles)[10:15]

```

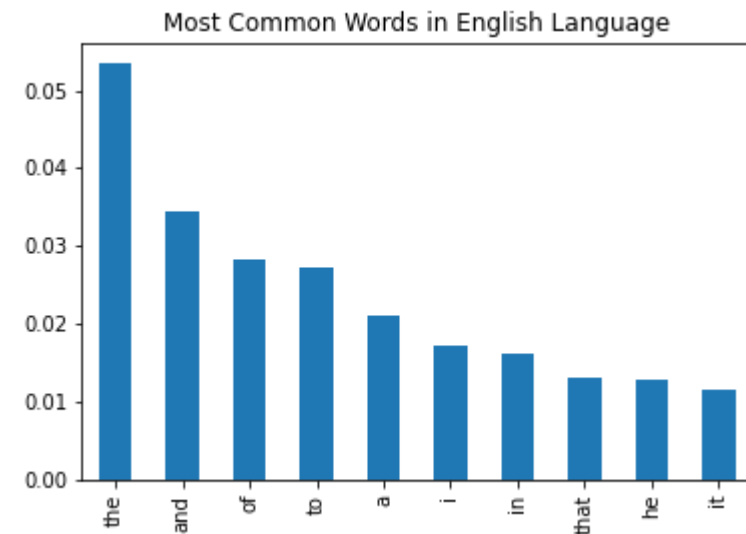
```
1 # for article in wiki_articles:
2 #     soup = BeautifulSoup(requests.get("https://en.wikipedia.org"+article).content, "html.parser")
3 #     txts = soup.find_all("p")
4 #     cleaned = ""
5 #     for txt in txts:
6 #         cleaned += ''.join(filter(whitelist.__contains__, txt.text.lower()))
7 #     time.sleep(.25)
8 # cleaned, len(cleaned.split())
```

```
1 # books.loc["wiki"] = cleaned
2 # books
```

## Data Exploration

```
1 series = pd.Series(books.apply(Counter).sum())
2 (series/series.sum()).sort_values(ascending=False).iloc[:10].plot.bar(title="Most Common Words in English Language")
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fc18e596c50>



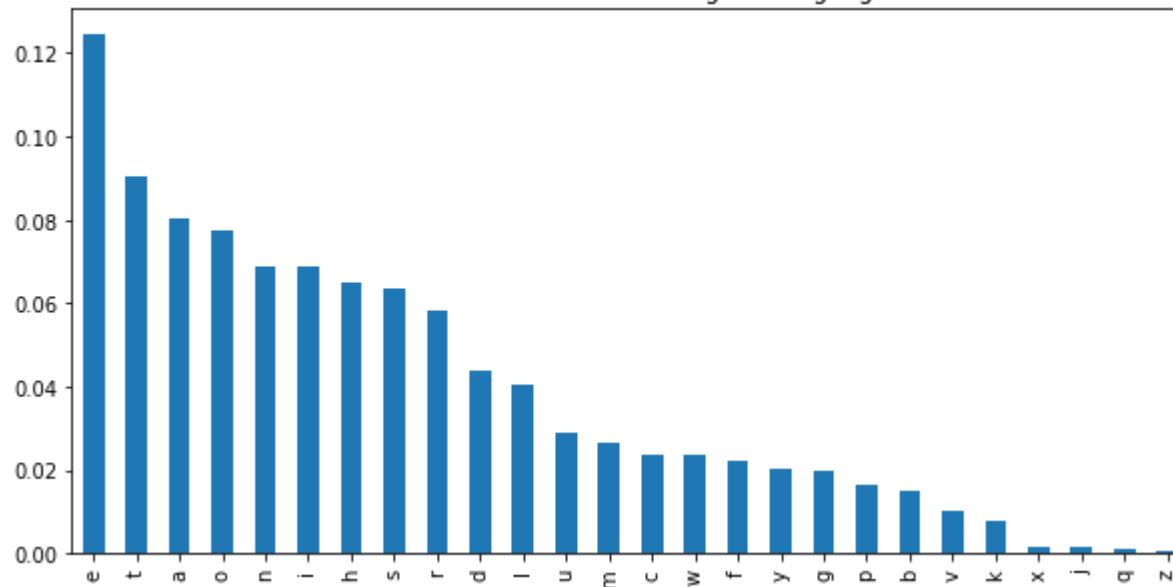
```

1 count = Counter()
2 for i in range(len(books)):
3     count += Counter("".join(books.iloc[i]))
4 series = pd.Series(count)
5 (series/series.sum()).sort_values(ascending=False).plot.bar(figsize=(10, 5), title="Distribution of Letters in English Language")

```

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fc164e0c650>

Distribution of Letters in English Language



```

1
2 ### This block is the majority of the machine learning model. It just lives up here cause
3 ### I used it for visualizations as well
4
5
6
7 ### Breaks every book into bigrams thru 5grams inclusive and counts
8 ### Series contains-
9 ### Index: previous n words      value: count of following words, unsorted
10 grams = {}
11 for n in range(1, 5): #N grams grouped by first n words, N = n+1
12     print("building " +str(n+1)+"-grams")
13     for book in books:# Iterate through all books
14         for i in range(len(book)-n):# Iterate through all words

```

```

14     for i in range(len(book)-n): # Iterate through all words
15         g = " ".join(book[i:i+n])
16         if g in grams: # if first N-1 words already found, append next word to list
17             grams[g].append(book[i+n])
18         else: # Start a new list with following word
19             grams[g]= [book[i+n]]
20 print("counting n-grams")
21 grams = pd.Series(grams).apply(Counter)#Find counts for words following all grams
22 grams #still not sorted

```

```

building 2-grams
building 3-grams
building 4-grams
building 5-grams
counting n-grams
of                {'those': 2056, 'promise': 29, 'frost': 10, 'b...
those             {'icy': 3, 'undiscovered': 1, 'countries': 12,...
icy               {'climes': 1, 'wall': 1, 'and': 2, 'black': 1,...
climes            {'inspired': 1, 'of': 1, 'whiteness': 1, 'th...
inspired          {'by': 3, 'me': 1, 'chapter': 1}
...
for the abandonm of                {'the': 1}
abandonment of the gigantic        {'blocks': 1}
of the gigantic blocks             {'of': 1}
the gigantic blocks of             {'stone': 1}
gigantic blocks of stone           {'which': 1}
Length: 17764391, dtype: object

```

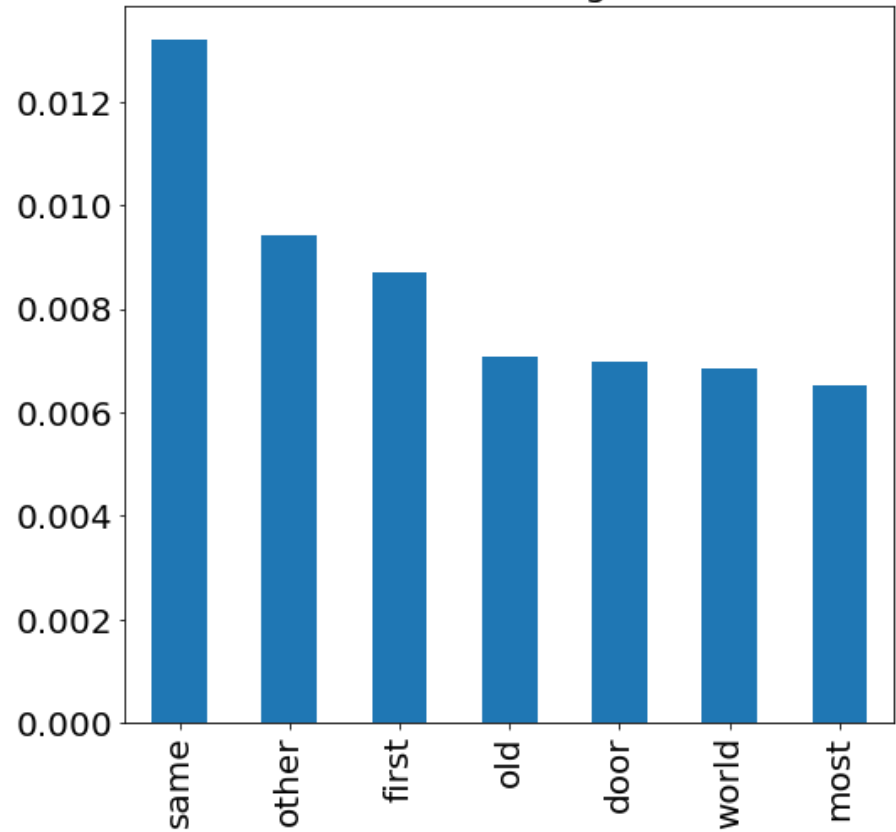
```

1 import matplotlib.pyplot as plt
2
3 fig = plt.figure(figsize=(18, 18))
4 plt.rcParams.update({'font.size': 18}) # must set in top
5 plt.subplots_adjust(hspace=.3)
6 sp = [221, 222, 223, 224]
7 st = ["the", "end of the", "father of the", "heart of the"]
8 m=7
9 for p, s in zip(sp, st):
10     sers = pd.Series(grams[s]).sort_values(ascending=False)
11     (sers/sers.sum()).iloc[:m].plot.bar(title="Words following '"+s+"'", ax=fig.add_subplot(p), fontsize=20)
12

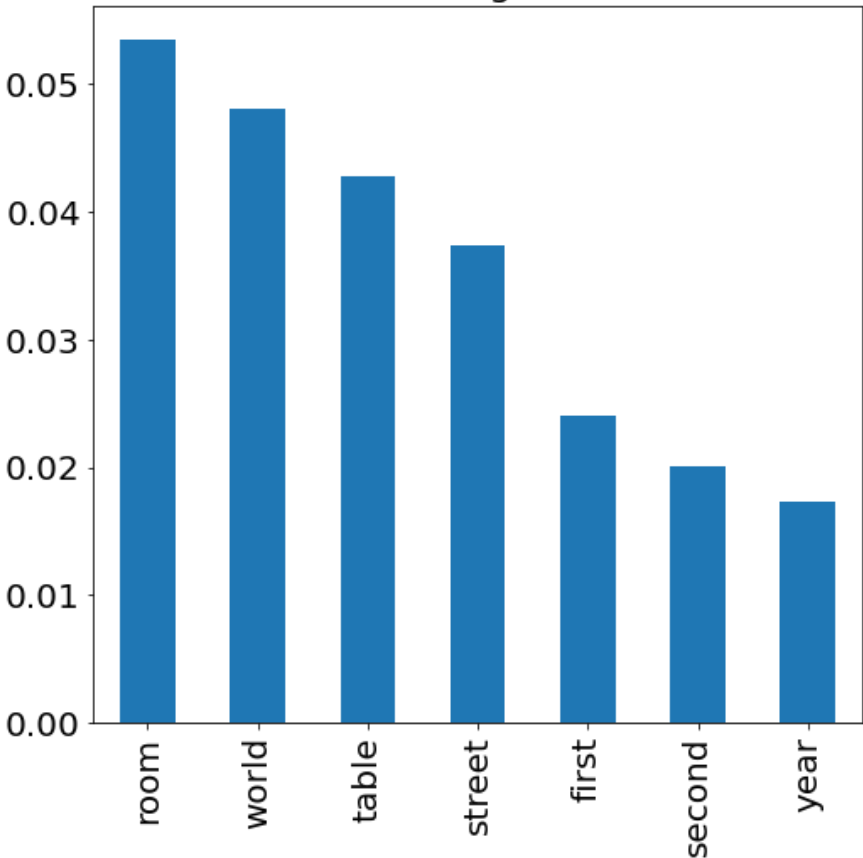
```



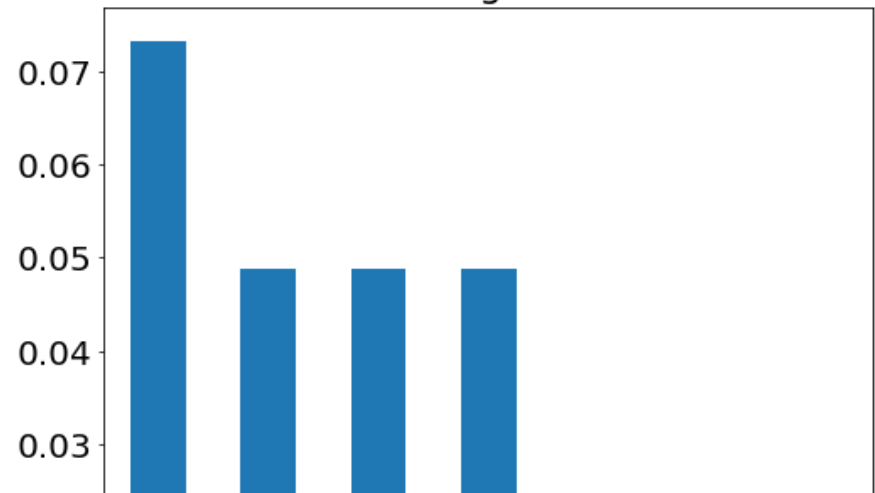
Words following 'the'



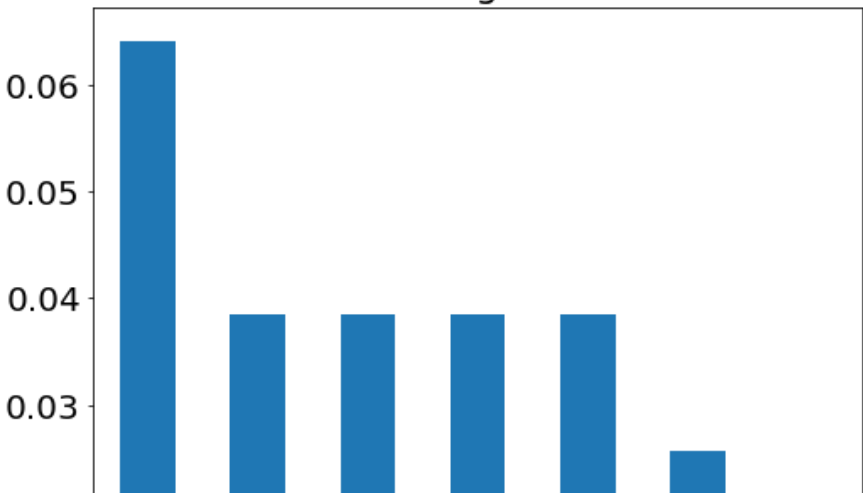
Words following 'end of the'

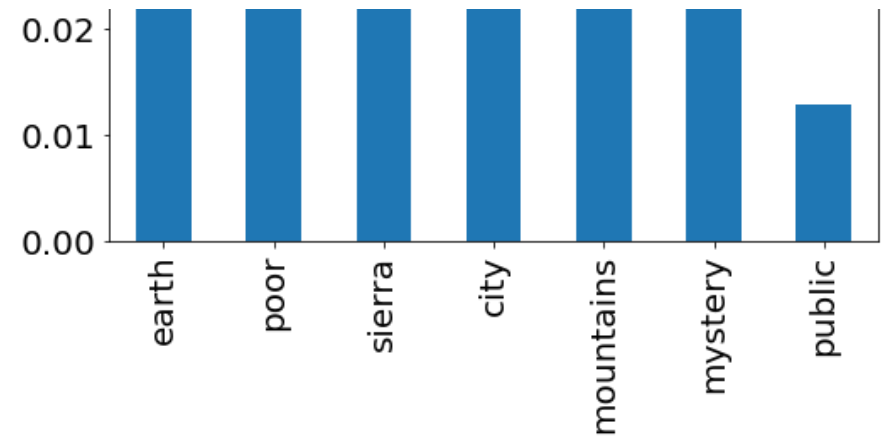
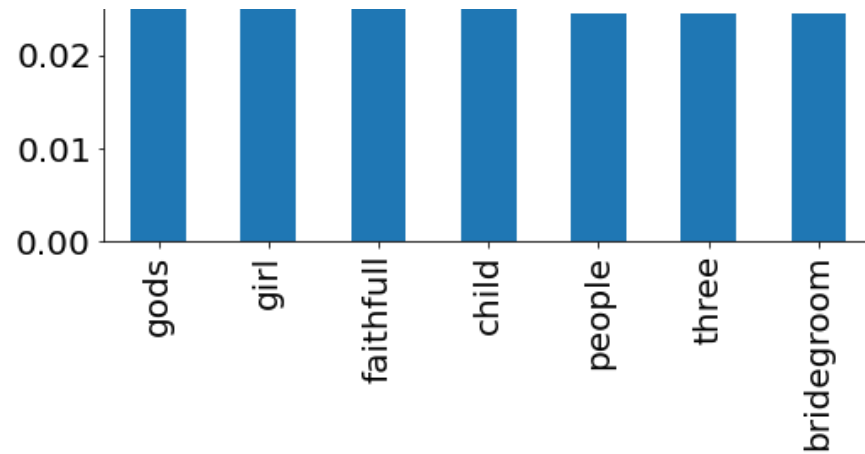


Words following 'father of the'



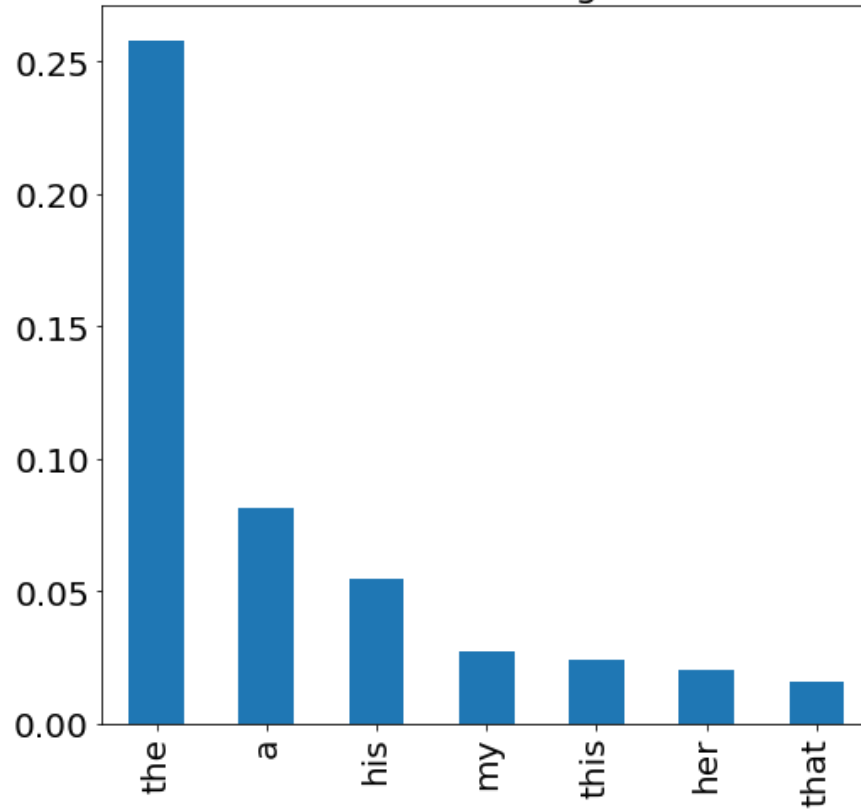
Words following 'heart of the'



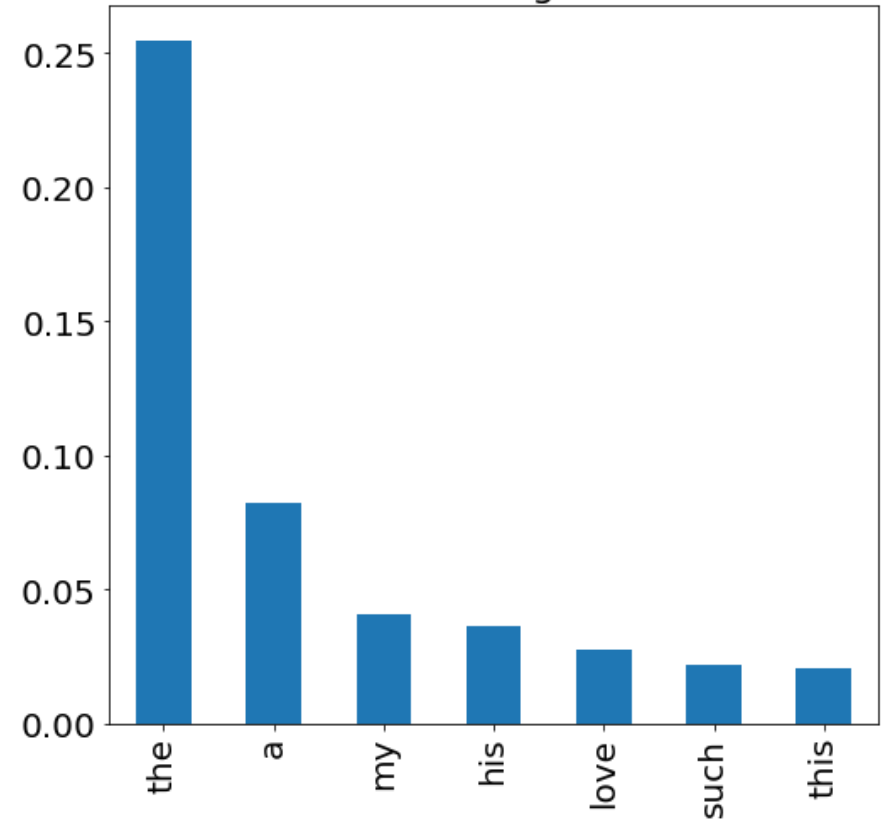


```
1 fig = plt.figure(figsize=(18, 18))
2 plt.rcParams.update({'font.size': 18}) # must set in top
3 st = ["in", "been in", "have been in", "i have been in"]
4 plt.subplots_adjust(hspace=.3)
5 sp = [221, 222, 223, 224]
6 for p, s in zip(sp, st):
7     sers = pd.Series(grams[s]).sort_values(ascending=False)
8     (sers/sers.sum()).iloc[:7].plot.bar(title="Words following '"+s+"'", ax=fig.add_subplot(p), fontsize=20)
9
```

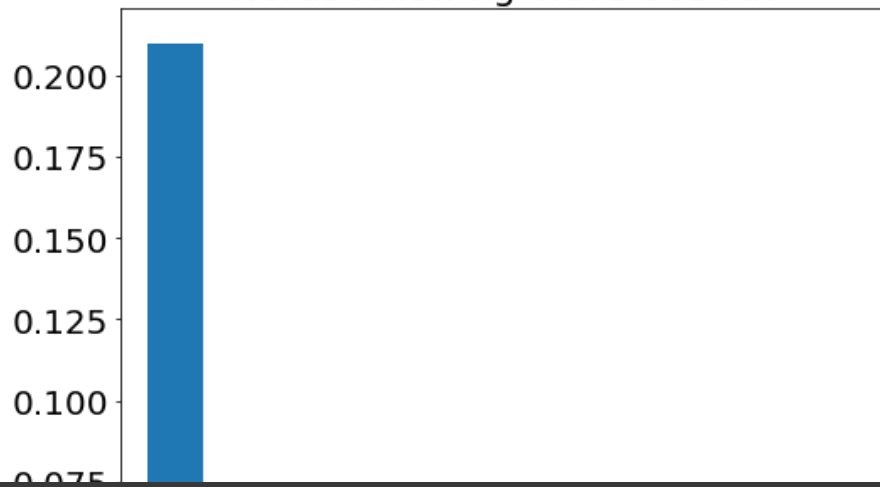
Words following 'in'



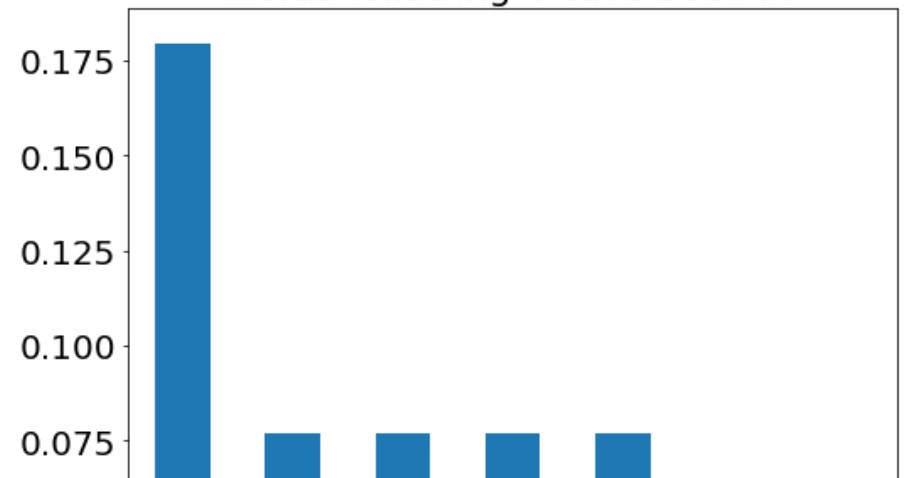
Words following 'been in'



Words following 'have been in'



Words following 'i have been in'





## Machine Learning

```
1 ###Predicts next N words
2 ### String, Int(optional) --> String
3 ### input can be 1 - 3 words long. If longer only
4 ### Considers previous 3 words
5 def predict_text(input):
6     prev = input.lower().split()
7     prev.reverse()
8     result = set()
9     for k in range(4):
10         j = 4-k
11         string = prev[:min(len(prev), j+1)]
12         string.reverse()
13         key = " ".join(string)
14         if key in grams:
15             s = pd.Series(grams[key]).sort_values(ascending=False).index.tolist()
16             z = 0
17             word = s[z]
18             while z < len(s) and z<3:
19                 if len(result) < 3:
20                     result.add(s[z])
21                 z += 1
22
23     return result
```

```
1 predict_text("i want")
```

```
{'a', 'to', 'you'}
```

```
1 predict_text("we are in good")
```

```
{'health', 'repair', 'time'}
```

```
1 predict_text("i want to sail from")
```

```
{'london', 'seville', 'zealand'}
```

```
1 predict_text("do you need")
```

```
{'money', 'not', 'to'}
```

```
1 predict_text("sail the")
```

```
{'sea', 'seas', 'wind'}
```

```
1 predict_text("lets get to")
```

```
{'know', 'the', 'work'}
```

```
1 predict_text("can i have some")
```

```
{'business', 'money', 'of'}
```

```
1 predict_text("would you like a")
```

```
{'cup', 'lift', 'priest'}
```

```
1 predict_text("i want to go to a")
```

```
{'certain', 'court', 'good'}
```

```
1 ###Predicts next N words
2 ### String, Int(optional) --> String
3 ### input can be 1 - 3 words long. If longer only
4 ### Considers previous 3 words
5 def predict_text2(input, num_words=20):
6     prev = input.lower().split()
7     prev.reverse()
8     result = ""
9     for i in range(num_words):
```

```

9     for i in range(num_words):
10         val = 0
11         for j in range(3):
12             string = prev[:min(len(prev), j+1)]
13             string.reverse()
14             key = " ".join(string)
15             if key in grams:
16                 s = pd.Series(grams[key]).sort_values(ascending=False).index.tolist()
17                 z = 0
18                 word = s[z]
19                 while word in prev and z < len(s) and z<5:
20                     word = s[z]
21                     z += 1
22             prev.insert(0, word)
23             result += prev[0] + " "
24     return result

```

```

1 a = open("new.txt", "w") #Creates a text document
2 a.write("test")
3 a.close()

```

```

1 #Can make livetime predictions while editing new.txt
2 #Must hit ctrl-s to save to update predictions
3 import time
4 from IPython.display import clear_output
5
6 for i in range(1):
7     a = open("new.txt", "r")
8     for line in a:#The stupidest way to get to the last line
9         pass
10    a.close()
11    preds = predict_text(line)
12    clear_output()
13    print(preds)
14    time.sleep(1)
15

```

```
{'the', 'and', 'of'}
```

