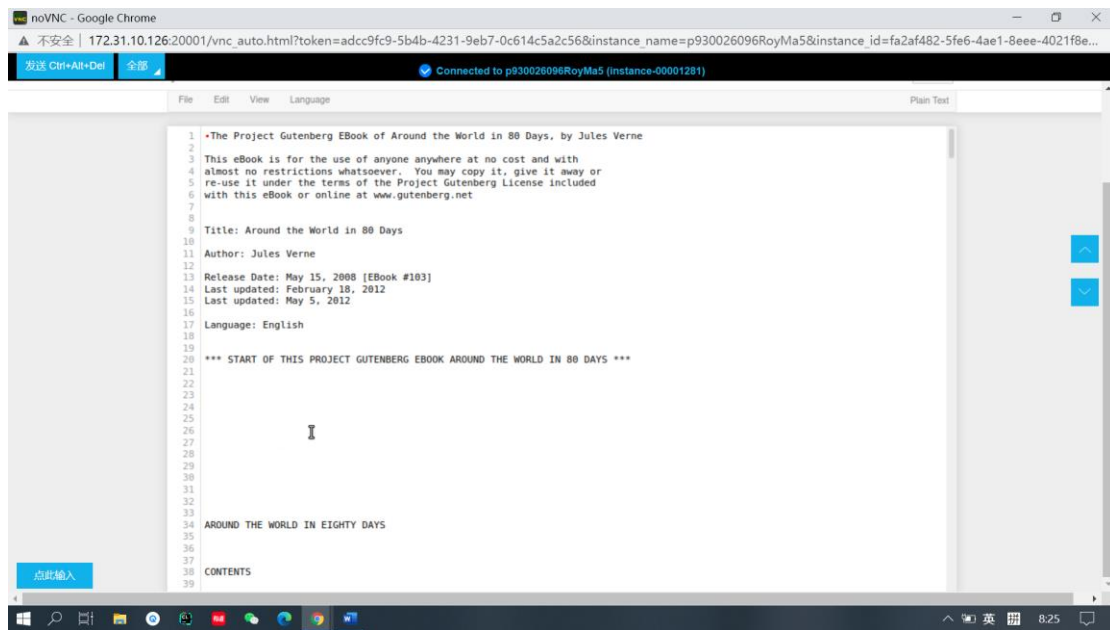
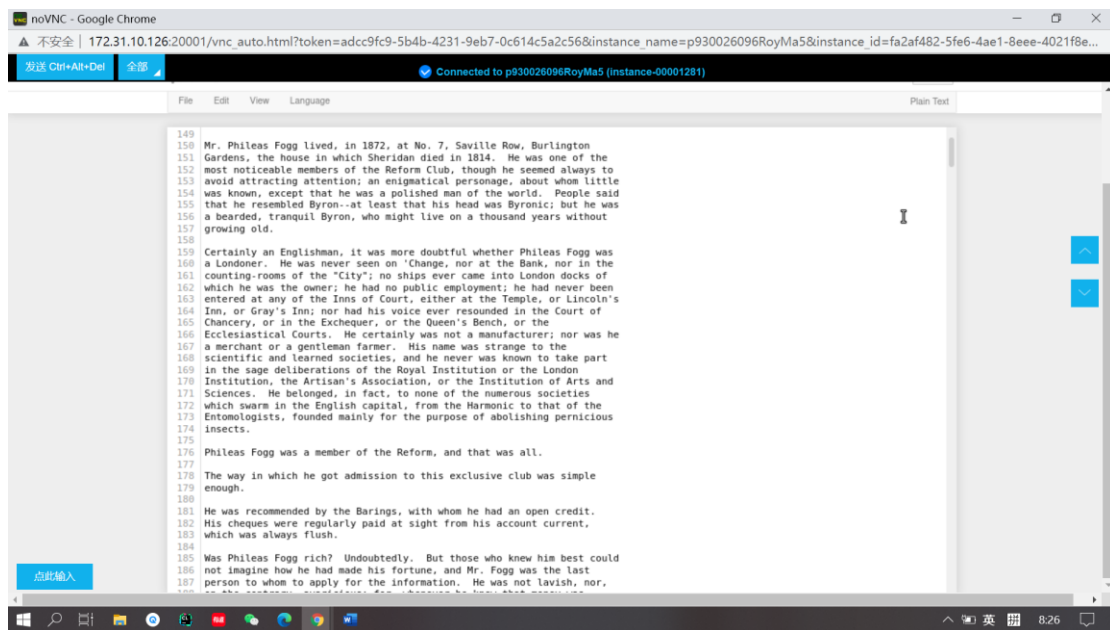


Intermediate



The English text file.



The content of English text file.

Using Spark, Stopwords, RegularExpression, Matplotlib for Word Cloud

```
In [2]: import matplotlib.pyplot as plt
import re
from wordcloud import WordCloud, STOPWORDS
from pyspark import SparkConf, SparkContext
```

Using Spark, stop words and regular expression, matplotlib for word cloud.

Create Spark Context and load text file

```
In [3]: sc = SparkContext.getOrCreate(SparkConf())
text_file = sc.textFile('103.txt')
print(text_file)

103.txt MapPartitionsRDD[1] at textFile at NativeMethodAccessorImpl.java:0
```

Create Spark context and load text file.

Word Count, sort by key

```
In [4]: counts = text_file.flatMap(lambda x: ((v) for v in re.split('\W',x))) \
    .map(lambda word: (word, 1)).reduceByKey(lambda a, b: a + b)
output = counts.sortByKey().collectAsMap()
print(output)

{'': 18494, '0': 1, '000': 3, '1': 48, '103': 4, '10th': 1, '11': 2, '117': 2, '11th': 7, '12th': 3, '13': 2, '13th': 3, '14th': 4, '15': 1, '1500': 1, '158': 1, '15th': 1, '16th': 1, '170': 1, '1756': 1, '17th': 2, '18': 1, '1814': 1, '1825': 1, '1837': 1, '1839': 1, '1842': 1, '1843': 1, '1845': 1, '1849': 1, '1853': 1, '1862': 1, '1867': 1, '1872': 2, '1887': 1, '18th': 1, '19th': 1, '2': 4, '20': 3, '2001': 1, '2008': 1, '2012': 2, '20th': 5, '21st': 1, '22': 1, '22nd': 2, '23': 1, '23rd': 4, '25th': 3, '280': 1, '28th': 1, '2nd': 6, '3': 11, '30': 1, '30th': 2, '35': 1, '3rd': 3, '4': 7, '40': 2, '400': 1, '43': 2, '45': 1, '4557': 1, '48': 1, '4th': 2, '5': 6, '50': 1, '501': 2, '596': 1, '5th': 3, '6': 4, '60': 1, '6221541': 1, '64': 1, '6th': 3, '7': 9, '7th': 3, '8': 6, '80': 6, '801': 1, '809': 1, '80th': 1, '84116': 1, '8th': 1, '9': 4, '90': 2, '99712': 1, '9th': 3, 'A': 186, 'ABOUT': 2, 'ACCEPT': 2, 'ACROBATIC': 1, 'ACROSS': 2, 'ACTUAL': 1, 'AGREE': 2, 'AGREEMENT': 1, 'AIDS': 2, 'AK': 1, 'AMERICAN': 2, 'AN': 2, 'AND': 17, 'ANTIPODES': 2, 'ANY': 3, 'ANYBODY': 2, 'ANYTHING': 1, 'APPEARS': 2, 'ARE': 4, 'AROUND': 5, 'AS': 7, 'ASCII': 2, 'ASTOUNDS': 2, 'AT': 12, 'ATTRACTION': 1, 'Abandoning': 1, 'About': 4, 'Above': 2, 'Abraham': 1, 'Absolute': 1, 'Absolutely': 1, 'According': 2, 'Additional': 1, 'Aden': 8, 'Admitted': 1, 'Admitting': 1, 'Africa': 3, 'African': 1, 'After': 17, 'Ages': 1, 'Agile': 1, 'Agra': 1, 'Agreed': 2, 'Ah': 22, 'Aha': 1, 'Ahmehnnagara': 1, 'Alabama': 2, 'Alas': 1, 'Albamarle': 3, 'All': 15, 'Allahabad': 16, 'Alongside': 1, 'Already': 1, 'Although': 1, 'Always': 2, 'America': 15, 'American': 30, 'Americans': 7, 'Among': 3, 'Amphion': 1, 'An': 15, 'Anam': 1, 'And': 6, 'Andaman': 2, 'Andrew': 14, 'Angelica': 1, 'Anglo': 2, 'Another': 1, 'Any': 1, 'Aouda': 136, 'Apiece': 2, 'Arabic': 2, 'Arabs': 1, 'Archive': 13, 'Are': 7, 'Arkansas': 2, 'Armenian': 1, 'Armenians': 1, 'Around': 4, 'Arrived': 3, 'Articles': 1, 'Artisan': 1, 'Arts': 1, 'As': 42, 'Asia': 4, 'Asian': 1, 'Asphaltite': 1, 'Ass': 1, 'Association': 1, 'Assurghur': 1, 'At': 66, 'Athens': 1, 'Atlantic': 15, 'Attired': 1, 'Auburn': 1, 'Aureng': 1, 'Aurangabad': 1, 'Author': 1, 'B': 1, 'BAD': 1, 'BAG': 1, 'BANKNOTES': 1, 'BATULCAR': 1, 'BE': 1, 'BEAUTIFUL': 1, 'BECOMES': 1, 'BEFORE': 1, 'BETRAYS': 1, 'BRAVE': 1, 'BREACH': 1, 'BUSINESS': 1, 'BUT': 1, 'BY': 1, 'Bab': 1, 'Bah': 1, 'Bank': 1, 'Banks': 1, 'Banyans': 1, 'Banyas': 1, 'Baring': 1, 'Barings': 1, 'Bar num': 1, 'Battery': 1, 'Batulcar': 1, 'Bay': 1, 'Be': 1, 'Because': 1, 'Before': 1, 'Behar': 1, 'Behind': 1, 'Being': 1, 'Below': 1, 'Benares': 1, 'Bench': 1, 'Bengal': 1, 'Benten': 1, 'Besides': 1, 'Bets': 1, 'Between': 1, 'Biblical': 1, 'Birmingham': 1, 'Bitter': 1, 'Blondin': 1, 'Bluffs': 1, 'Bo mbay': 1, 'Bonds': 1, 'Book': 1, 'Boots': 1, 'Bordeaux': 1, 'Boulevard': 1, 'Bow': 1, 'Bradshaw': 1, 'Brahma': 1, 'Brahmin': 1, 'Brahmins': 1, 'Braz
```

Word count and sorted by key.

Get the list of words

```
In [5]: wordlist = output.keys()
print(wordlist)

dict keys(['', '0', '000', '1', '103', '10th', '11', '117', '11th', '12th', '13', '13th', '14th', '15', '1500', '158', '15th', '16th', '170', '1756', '17th', '18', '1814', '1825', '1837', '1839', '1842', '1843', '1845', '1849', '1853', '1862', '1867', '1872', '1887', '18th', '19th', '2', '20', '2001', '2008', '2012', '20th', '21st', '22', '22nd', '23', '23rd', '25th', '280', '28th', '2nd', '3', '30', '30th', '35', '3rd', '4', '40', '400', '43', '45', '4557', '48', '4th', '5', '50', '501', '596', '5th', '6', '60', '6221541', '64', '6th', '7', '7th', '8', '80', '801', '809', '80th', '84116', '8th', '9', '90', '99712', '9th', 'A', 'ABOUT', 'ACCEPT', 'ACROBATIC', 'ACROSS', 'ACTUAL', 'AGREE', 'AGREEMENT', 'AIDS', 'AK', 'AMERICAN', 'AN', 'AND', 'ANTIPODES', 'ANY', 'ANYBODY', 'ANYTHING', 'APPEARS', 'ARE', 'AROUND', 'AS', 'ASCII', 'ASTOUNDS', 'AT', 'ATTRACTION', 'Abandoning', 'About', 'Above', 'Abraham', 'Absolute', 'Absolutely', 'According', 'Additional', 'Aden', 'Admitted', 'Admitting', 'Africa', 'African', 'After', 'Ages', 'Agile', 'Agra', 'Agreed', 'Ah', 'Aha', 'Ahmehnnagara', 'Alabama', 'Alas', 'Albamarle', 'All', 'Allahabad', 'Alongside', 'Already', 'Although', 'Always', 'America', 'American', 'Americans', 'Among', 'Amphion', 'An', 'Anam', 'And', 'Andaman', 'Andrew', 'Angelica', 'Anglo', 'Another', 'Any', 'Aouda', 'Apiece', 'Arabic', 'Arabs', 'Archive', 'Are', 'Arkansas', 'Armenian', 'Armenians', 'Around', 'Arrived', 'Articles', 'Artisan', 'Arts', 'As', 'Asia', 'Asian', 'Asphaltite', 'Ass', 'Association', 'Assurghur', 'At', 'Athens', 'Atlantic', 'Attired', 'Auburn', 'Aureng', 'Aurangabad', 'Author', 'B', 'BAD', 'BAG', 'BANKNOTES', 'BATULCAR', 'BE', 'BEAUTIFUL', 'BECOMES', 'BEFORE', 'BETRAYS', 'BRAVE', 'BREACH', 'BUSINESS', 'BUT', 'BY', 'Bab', 'Bah', 'Bank', 'Banks', 'Banyans', 'Banyas', 'Baring', 'Barings', 'Bar num', 'Battery', 'Batulcar', 'Bay', 'Be', 'Because', 'Before', 'Behar', 'Behind', 'Being', 'Below', 'Benares', 'Bench', 'Bengal', 'Benten', 'Besides', 'Bets', 'Between', 'Biblical', 'Birmingham', 'Bitter', 'Blondin', 'Bluffs', 'Bo mbay', 'Bonds', 'Book', 'Boots', 'Bordeaux', 'Boulevard', 'Bow', 'Bradshaw', 'Brahma', 'Brahmin', 'Brahmins', 'Braz
```

Get the list of words.

Remove stopwords, you can also add your own stopwords

```
In [6]: stopwords = set(STOPWORDS)
#new_words = ['ebook', 'gutenberg', 'https']
#stopwords.update(new_words)
print(stopwords)

{'because', 'all', 'could', 'how's', 'she'd', 'who', 'why', 'he'd', 'own', 'shan't', 'here', 'herself', 'just', 'the', 'm', 'his', 'during', 'doing', 'did', 'didn't', 'shouldn't', 'we', 'about', 'out', 'they'd', 'ought', 'its', 'she's', 'http', 'by', 'when's', 'be', 'haven't', 'for', 'our', 'through', 'at', 'then', 'isn't', 'down', 'than', 'some', 'ha', 'sn't', 'there', 'whom', 'between', 'further', 'you're', 'aren't', 'any', 'a', 'each', 'mustn't', 'wasn't', 'we'll', 'weren't', 'how', 'too', 'what', 'why's', 'shall', 'i'm', 'in', 'an', 'couldn't', 'who's', 'you'd', 'not', 'let's', 'otherwise', 'other', 'was', 'can't', 'from', 'yourself', 'under', 'we're', 'or', 'ourselves', 'i'll', 'being', 'the', 'cannot', 'he'll', 'i've', 'to', 'themselves', 'i'd', 'there's', 'had', 'has', 'what's', 'are', 'doesn't', 'they', 'll', 'would', 'www', 'also', 'should', 'have', 'which', 'against', 'i', 'you', 'having', 'like', 'my', 'of', 'am', 'until', 'get', 'same', 'hers', 'so', 'again', 'such', 'hence', 'you'll', 'few', 'below', 'however', 'myself', 'while', 'don't', 'as', 'off', 'with', 'once', 'where', 'wouldn't', 'do', 'he', 'they', 'you've', 'yourselves', 'since', 'can', 'when', 'they're', 'itself', 'that's', 'therefore', 'here's', 'ours', 'we've', 'that', 'these', 'no', 'r', 'he', 's', 'she'll', 'been', 'if', 'it's', 'were', 'on', 'else', 'does', 'more', 'ever', 'com', 'theirs', 'but', 'k', 'nor', 'both', 'very', 'over', 'only', 'their', 'himself', 'we'd', 'me', 'after', 'above', 'him', 'this', 'most', 'is', 'those', 'your', 'she', 'before', 'into', 'up', 'won't', 'and', 'where's', 'yours', 'it', 'hadn't', 'they've', 'her', '}
```

Remove stop words, and stop words can be updated by using function update.

Generate and show word cloud

```
In [7]: wordcloud = WordCloud(stopwords=stopwords, width=1600, height=800, background_color='white')\
      .generate(" ".join(wordlist))

plt.figure(figsize=(20,10))
plt.imshow(wordcloud)
plt.axis('off')
plt.show()
```



Generate and show word cloud.

You can also include a picture mask for the result

```
In [8]: import numpy as np
        from PIL import Image
        from wordcloud import ImageColorGenerator

        background_image = np.array(Image.open('heart.jpg'))
        #background_image = np.array(Image.open('balloon.jpg'))

        wordcloud = WordCloud(stopwords=stopwords, mask=background_image, width=1600, height=800, background_color='white')\
            .generate(" ".join(wordlist))

        plt.figure(figsize=(20,10))

        image_colors = ImageColorGenerator(background_image)
        plt.imshow(wordcloud.recolor(color_func=image_colors))

        plt.axis('off')
        plt.show()
```

Include a picture mask for the result.

Generate and show word cloud, you may need a Chinese font

```
In [4]: wordcloud = WordCloud(stopwords=stopwords, width=1600, height=800, background_color='white', \
                             font_path='NotoSansCJK-Regular-1.otf').generate(" ".join(wordlist))

plt.figure(figsize=(20,10))
plt.imshow(wordcloud)
plt.axis('off')
plt.show()
```

Building prefix dict from the default dictionary ...
 Dumping model to file cache /tmp/jieba.cache
 Loading model cost 0.858 seconds.
 Prefix dict has been built successfully.

Generate and show word cloud by using a Chinese font.



You can also include a picture mask for the result

```
In [8]: import numpy as np
        from PIL import Image
        from wordcloud import ImageColorGenerator

        background_image = np.array(Image.open('heart.jpg'))
        #background_image = np.array(Image.open('balloon.jpg'))

        wordcloud = WordCloud(stopwords=stopwords, mask=background_image, width=1252, height=1144, background_color='white',
                               font_path='NotoSansCJK-Regular-1.otf').generate(" ".join(wordlist))

        plt.figure(figsize=(20,10))

        image_colors = ImageColorGenerator(background_image)
        plt.imshow(wordcloud.recolor(color_func=image_colors))

        plt.axis('off')
        plt.show()
```

Include a picture for the word cloud.

You can also generate WordCloud from weights/keywords of Jieba

```
In [10]: word_dict = {}
          for i in keywords:
              word_dict[i[0]]=i[1]

          wordcloud = WordCloud(stopwords=stopwords, width=1600, height=800, background_color='white', \
                                font_path='NotoSansCJK-Regular-1.otf').fit_words(word_dict)

          plt.figure(figsize=(20,10))
          plt.imshow(wordcloud)
          plt.axis('off')
          plt.show()
```

Generate word cloud from weights of Jleba.



Word cloud with weights