

CASE STUDY

Frequency Analysis, and Proof of Benford's Law for large text files.

RoyNaaz Banu - 121910302027

Sunanda Tata - 121910302028

Bhargav Kantheti - 121910302029

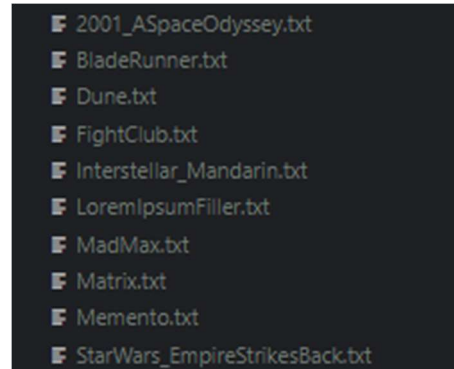
B2 – CSE
Gitam University

Project Repository Link:

<https://github.com/bharxhav/WordFrequencyAnalysis>

Part 0: Collecting Data

Splitting words, punctuation and whitespaces, can be easy, but to test the efficiency of the program, we need to run it with huge data sizes, and the best one is *movie scripts*, since they are thousands of lines.



Here is front page of script of the movie “The Matrix”.



The code has been tested with more than one movie, and with languages such as English and Latin.

Here is Latin filler text called “lorem ipsum”.

```
1 Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt
2
3 In nisl nisi scelerisque eu ultrices vitae auctor eu. Dolor sit amet consectetur adipi
4
5 Urna nunc id cursus metus. Id leo in vitae turpis massa. Blandit turpis cursus in hac
6
7 Egestas egestas fringilla phasellus faucibus scelerisque. Sit amet dictum sit amet jus
8
9 Laoreet suspendisse interdum consectetur libero id faucibus nisl tincidunt. Enim sit a
10
11 Nibh praesent tristique magna sit amet purus gravida quis. Commodó viverra maecenas ac
12
13 Eget mi proin sed libero enim sed faucibus turpis in. Sed odio morbi quis commodo odio
14
15 Nunc lobortis mattis aliquam faucibus purus in massa tempor. Amet consectetur adipisci
16
17 Cras sed felis eget velit aliquet sagittis. Et sollicitudin ac orci phasellus egestas
18
19 Elementum facilisis leo vel fringilla est ullamcorper eget nulla. Sit amet nulla facil
20
21 Mi eget mauris pharetra et. A erat nam at lectus urna duis convallis. Sit amet porttit
22
23 Accumsan lacus vel facilisis volutpat est velit egestas dui id. Diam sit amet nisl sus
```

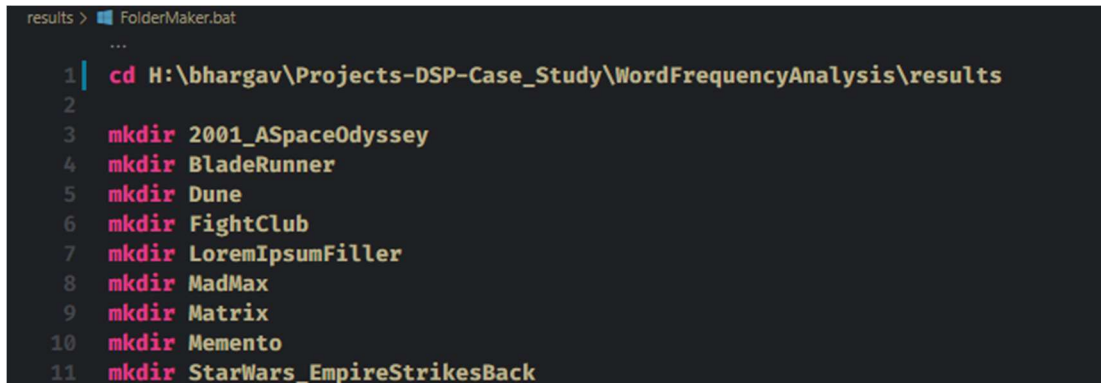
Here is a part of script from movie “Blade Runner”.

```
5673 Then, it dropped into darkness.
5674
5675 The great machine knew that this
5676 tiny scout was reporting back to
5677 its parent; but it was too simple,
5678 too primitive a device to detect
5679 the forces that were gathering
5680 round it now.
5681
5682 Then the pod came, carrying
5683 life. The great machine searched
5684 its memories.
5685
5686 The logic circuits made their
5687 decision when the pod had fallen
5688 beyond the last faint glow of the
5689 reflected Saturnian light.
5690
5691 In a moment of time, too short to
5692 be measured, space turned and
5693 twisted upon itself.
5694
5695 12/9/65 d8
5696 -----
5697 END OF SCREENPLAY
5698 END OF FILE
5699
```

Part 1: Creating Space for Results

To store the frequencies, and also keep them organized, I created a batch file to automate the results for me. By writing onto the batch file, I can create folders by running it.

Inside the Batch File



```
results > FolderMaker.bat
...
1 | cd H:\bhargav\Projects-DSP-Case_Study\WordFrequencyAnalysis\results
2
3 | mkdir 2001_ASpaceOdyssey
4 | mkdir BladeRunner
5 | mkdir Dune
6 | mkdir FightClub
7 | mkdir LoremIpsumFiller
8 | mkdir MadMax
9 | mkdir Matrix
10 | mkdir Memento
11 | mkdir StarWars_EmpireStrikesBack
```

To do all this, I created a class called BatchMaker, which will

- Create a batch file.
- Run the batch file.

BatchMaker constructor

```
class BatchMaker:
    def __init__(self, files, directory="H:\\bhargav\\Projects-
DSP-Case_Study\\WordFrequencyAnalysis\\results"):
        self.files = files
        self.directory = directory
```

Batch File Writer

```
def write_batch(self):
    adj = ["cd " + self.directory, ""]

    for f in self.files:
        adj.append('mkdir ' + f[:-4])

    fp = open('results/FolderMaker.bat', 'w')
    fp.write("\n".join(adj))
    fp.close()

    print("Written The Batch File")
```

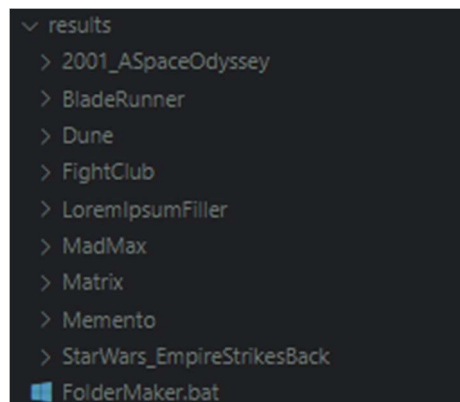
Creating can be done by simply writing onto a file, whereas running it can be problematic, but with the `subprocess` module, it can be done easily.

Batch File Runner

```
def create_folders(self):
    subprocess.run([self.directory+"\\FolderMaker.bat"],
shell=True)

    print()
    print("Created Folders")
```

Result after Running the Batch



Part 2: Purifying the Data

Now, after creating a place for the results to lie, we need to separate the words from punctuation and whitespaces.

TextSplitter constructor

```
class TextSplitter:
    def __init__(self, filename):
        self.fn = filename

        with open("source/" + filename, 'r') as fp:
            self.contents = fp.read()

        self.name = filename[:-4]
        self.path = "results/" + self.name + "/"
```

This can be done using a for loop, and reading the files in source one by one.

I have created a class called TextSplitter which would do this for me, given that I give the file names.

Word Partitioner

```
def partition_words(self):
    words = []
    words_lowercase = []

    for word in self.contents.split():
        if word.isalpha():
            words.append(word.strip())
            words_lowercase.append(word.strip().lower())

    with open(self.path+'words.txt', 'w') as fp:
        fp.write("\n".join(words))
        print("Written words onto a file.")

    with open(self.path+'words_lower.txt', 'w') as fp:
        fp.write("\n".join(words_lowercase))
        print("Written lowercase words onto a file.")
```

Other Partition maker function

```
def partition(self):
    alpha_raw = ""
    alpha_with_whitespace = ""
    binary_raw = ""
    punctuations = ""
    punctuation_with_whitespace = ""
    whitespaces = ""

    for c in self.contents:
        if c in whitespaces:
            whitespaces += c
            punctuation_with_whitespace += c
            alpha_with_whitespace += c

        if c in punctuation:
            punctuations += c
            punctuation_with_whitespace += c

        if c.isalpha():
            alpha_raw += c
            alpha_with_whitespace += c
            binary_raw += str(ord(c)) + " "

    with open(self.path+'alpha_raw.txt', 'w') as fp:
        print('Written raw alphabets onto a file.')
        fp.write(alpha_raw)

    with open(self.path+'alpha_with_whitespace.txt', 'w') as fp:
        print('Written alphabets with whitespaces onto a file.')
        fp.write(alpha_with_whitespace)

    with open(self.path+'binary_raw.txt', 'w') as fp:
        print('Written binary raw onto a file.')
        fp.write(binary_raw)

    with open(self.path+'punctuations.txt', 'w') as fp:
        print('Written raw punctuations onto a file.')
        fp.write(punctuations)

    with open(self.path+'punctuation_with_whitespace.txt', 'w') as fp:
        print('Written punctuations with whitespaces onto a file.')
        fp.write(punctuation_with_whitespace)

    with open(self.path+'whitespaces.txt', 'w') as fp:
        print('Written raw whitespaces onto a file.')
        fp.write(whitespaces)
```

Result

```

  results
  2001_ASpaceOdyssey
    alpha_raw.txt
    alpha_with_whitespace.txt
    binary_raw.txt
    punctuation_with_whitespace.txt
    punctuations.txt
    whitespaces.txt
    words_lower.txt
    words.txt
  BladeRunner
    alpha_raw.txt
    alpha_with_whitespace.txt
    binary_raw.txt
    punctuation_with_whitespace.txt
    punctuations.txt
    whitespaces.txt
    words_lower.txt
    words.txt
  Dune
    alpha_raw.txt
    alpha_with_whitespace.txt
    binary_raw.txt
    punctuation_with_whitespace.txt
    punctuations.txt
    whitespaces.txt
    words_lower.txt
    words.txt
  FightClub
    alpha_raw.txt
    alpha_with_whitespace.txt
    binary_raw.txt
    punctuation_with_whitespace.txt
    punctuations.txt
    whitespaces.txt
    words_lower.txt
    words.txt
  LoremIpsumFiller
    alpha_raw.txt
    alpha_with_whitespace.txt
    binary_raw.txt
    punctuation_with_whitespace.txt
    punctuations.txt
    whitespaces.txt
    words_lower.txt
    words.txt
  MadMax
```


Alpha Raw

```
1  ASpaceOdysseyScreenplaybyStanleyKubrickandArthurCClarkHawkFilmsLtdcoMGMStudios
dlastednowfortenmillionyearsandwouldnotendforanothermillionThereignoftheterrib
rvivalhadreachedanewclimaxofferocityandthevictorwasnotyetinsightInthisdryandba
ERThemanapesofthefieldhadnoneoftheseattributesandtheywereonthelongpatheticroad
byasluggishbrownstreamThetribehadalwaysbeenhungryandnowitwasstarvingAsthefirst
nowtheOldOnewashisfatherforsucharelationshipwasbeyondhisunderstandingbutashest
theroutofthecaveandleaveshimforthehyenasAmonghiskindMoonwatcherisalmostagianth
isquitemanlikeandhisheadisalreadynearerthanapeTheforeheadislowandtherearegr
uponthehostileworldthereisalreadyaACONTINUEDsomethinginhisgazebeyondthegraspof
otfulfillitselfforanothermillionyearsAEXTTHESTREAMTHEOTHERSAsthedawnskybri
hersideeverydaythatdidnotmakeitanylessannoyingThereareeighteenofthemanditisimp
oangrilydanceandshriekontheirsideofthestreamandhisownpeoplereplyInkindTheconfr
hemuddywaterHonorhasbeensatisfiedeachgrouphasstakeditsclaimtoitsownterritoryaA
pangsofhungerwhileallaroundsuchcompetingswiththemforthesamefodderisapotentialso
vannaandthroughthebrushisnotonlybeyondtheirreachtheideaofeatingitisbeyondtheir
beslowlywandersacrossthebareflatcountrysideforagingforrootsandoccasionalberrie
oundSuddenlyMoonwatcherbecomesawareofalionstalkingthemaboutyardsawayDefenceles
DSHONEYithadnotbeenagooddaythoughasMoonwatcherhadnorealremembranceofthepasthec
eeandsoenjoysthefinestdelicacyhispeoplecouldneverknowOfcoursehealsocollectsagoo
sstillhungryheisnotactuallyweakwithhungerThatwasthemostthatanyhominidcouldhope
```

Alphabets with Whitespaces

```
You, 2 days ago | 1 author (You)
1  You, 2 days ago • Add outputs
2
3  A SPACE ODYSSEY
4
5  Screenplay
6
7  by
8
9  Stanley Kubrick and Arthur C Clark
10
11  Hawk Films Ltd
12  co MGM Studios
13  Boreham Wood
14  Herts
15
16
17  TITLE                                PART I
18  AFRICA
19  YEARS AGO
20
21  A
22  VIEWS OF AFRICAN DRYLANDS DROUGHT
23
24  The remorseless drought had lasted now for ten million years
25  and would not end for another million The reign of the ter
26  rible lizards had long since passed but here on the continent
27  which would one day be known as Africa the battle for survival
28  had reached a new climax of ferocity and the victor was not
29  yet in sight In this dry and barren land only the small or
30  the swift or the fierce could flourish or even hope to exist
```

Binary Raw

```
1 65 83 80 65 67 69 79 68 89 83 83 69 89 83 99 114 101 101 110 112 108 97 121 98 121
67 108 97 114 107 72 97 119 107 70 105 108 109 115 76 116 100 99 111 77 71 77 83 11
84 76 69 80 65 82 84 73 65 70 82 73 67 65 89 69 65 82 83 65 71 79 65 86 73 69 87 83
101 109 111 114 115 101 108 101 115 115 100 114 111 117 103 104 116 104 97 100 108
101 97 114 115 97 110 100 119 111 117 108 100 110 111 116 101 110 100 102 111 114 9
111 102 116 104 101 116 101 114 114 105 98 108 101 108 105 122 97 114 100 115 104 9
101 111 110 116 104 101 99 111 110 116 105 110 101 110 116 119 104 105 99 104 119 1
99 97 116 104 101 98 97 116 116 108 101 102 111 114 115 117 114 118 105 118 97 108
101 114 111 99 105 116 121 97 110 100 116 104 101 118 105 99 116 111 114 119 97 115
97 110 100 98 97 114 114 101 110 108 97 110 100 111 110 108 121 116 104 101 115 109
99 101 99 111 117 108 100 102 108 111 117 114 105 115 104 111 114 101 118 101 110 1
79 78 87 65 84 67 72 69 82 84 104 101 109 97 110 97 112 101 115 111 102 116 104 101
114 105 98 117 116 101 115 97 110 100 116 104 101 121 119 101 114 101 111 110 116 1
105 97 108 101 120 116 105 110 99 116 105 111 110 65 98 111 117 116 116 119 101 110
111 102 99 97 118 101 115 111 118 101 114 108 111 111 107 105 110 103 97 115 109 97
98 121 97 115 108 117 103 103 105 115 104 98 114 111 119 110 115 116 114 101 97 109
110 103 114 121 97 110 100 110 111 119 105 116 119 97 115 115 116 97 114 118 105 11
119 110 99 114 101 101 112 115 105 110 116 111 116 104 101 99 97 118 101 77 111 111
105 115 102 97 116 104 101 114 104 97 115 100 105 101 100 100 117 114 105 110 103 1
101 79 108 100 79 110 101 119 97 115 104 105 115 102 97 116 104 101 114 102 111 114
111 110 100 104 105 115 117 110 100 101 114 115 116 97 110 100 105 110 103 98 117 1
116 116 104 101 101 109 97 99 105 97 116 101 100 98 111 100 121 104 101 102 101 101
105 110 116 111 115 97 100 110 101 115 115 84 104 101 110 104 101 99 97 114 114 105
101 99 97 118 101 97 110 100 108 101 97 118 101 115 104 105 109 102 111 114 116 104
110 119 97 116 99 104 101 114 105 115 97 108 109 111 115 116 97 103 105 97 110 116
```

Punctuation with Whitespaces

```
9
0 -----
1
2
3
4
5
6
7
8
9
0
1
2 //
3 -----
4
5 & -
6
7 -
8 ,
9
0 , , ,
1
2 ,
3
4
5
```

Punctuation

```
1  : .. //-----
2  : .. //-----
3  : .. //-----
4  : .. //-----
5  : .. //-----
6  : .. //-----
7  : .. //-----
8  : .. //-----
9  : .. //-----
10 : .. //-----
11 : .. //-----
12 : .. //-----
13 : .. //-----
14 : .. //-----
15 : .. //-----
16 : .. //-----
17 : .. //-----
18 : .. //-----
19 : .. //-----
20 : .. //-----
21 : .. //-----
22 : .. //-----
23 : .. //-----
24 : .. //-----
25 : .. //-----
26 : .. //-----
27 : .. //-----
```

Whitespaces

```
7
8
9  → → → .....
10
11 → → → → → .....
12 → → → → → .....
13 → → → → → .....
14 → → → → → .....
15
16
17 → → → → → .....
18 → → → → → .....
19 → → → → → .....
20
21
22 .....
23
24 .....
25 .....
26 .....
27 .....
```

Words

```
3 ODYSSEY
4 Screenplay
5 by
6 Stanley
7 Kubrick
8 and
9 Arthur
10 Clark
11 Hawk
12 Films
13 Boreham
14 TITLE
15 PART
16 I
17 AFRICA
18 YEARS
```

Lowercase Words

```
1 a You, 2 days ago • Add outputs
2 space
3 odyssey
4 screenplay
5 by
6 stanley
7 kubrick
8 and
9 arthur
10 clark
11 hawk
12 films
13 boreham
14 title
15 part
16 i
17 africa
18 years
19 ago
20 views
21 of
22 african
23 drylands
24 drought
25 the
```


Part 3: Finding Frequencies

After having the required data, I found the frequencies using the Counter class, which will find frequencies of each word and character.

I did all this in a class called FrequencyAnalyser which does this and one more other task which I will discuss in the next part.

FrequencyAnalyser constructor

```
class FrequencyAnalyser:
    def __init__(self, filename):
        self.fn = filename
        self.name = filename[:-4]

        self.path = "results/" + self.name + "/"
```

Analyse Function

```
def analyse(self):
    # Reading Words and Characters and Marking Frequencies
    with open(self.path+'words_lower.txt', 'r') as wl,
    open(self.path+'alpha_raw.txt', 'r') as cl:
        self.words = dict(Counter(wl.read().split()))
        self.chars = dict(Counter(cl.read().lower()))

    print("Analysed the word frequencies.")
    print("Analysed the character frequencies.")
```

Character Frequencies

```
{'t': 10583, 'h': 6882, 'e': 14843, 'm': 2660, 'p': 2540, 'i': 8344, 'r': 7470, 's': 7470, 'f': 2311, 'o': 9621, 'u': 3547, 'x': 306, 'v': 1081, 'z': 143, 'j': 137, 'q': 97}
```

Word Frequencies

```
Script :: StarWars_EmpireStrikesBack.txt  
{'the': 1782, 'empire': 8, 'strikes': 2, 'back': 53, 'written': 1, 'by': 87, 'lawrence': 1, 'galaxy': 5, 'planet': 7, 'hoth': 69, 'star': 69, 'destroyer': 39, 'moves': 60, 'thru': 5, 'probes': 1, 'zooms': 4, 'toward': 74, 'lands': 3, 'on': 210, 'covered': 4, 'an': 77, 'rd': 1, 'mechanical': 5, 'sound': 11, 'rises': 10, 'above': 13, 'whining': 1, 'strange': 7, 'into': 124, 'small': 27, 'figure': 5, 'gallops': 2, 'windswept': 2, 'ice': 48, 'beneath': 4, 'speeding': 1, 'paws': 1, 'up': 83, 'slope': 3, 'reins': 2, 'his': 315, 'something': 18, 'in': 276, 'he': 217, 'takes': 19, 'pair': 3, 'electrobinoculars': 2, 'ut': 11, 'whips': 1, 'at': 144, 'cap': 1, 'activates': 3, 'comlink': 3, 'tauntaun': 8, 'f
```

Part 4: The Plot

Here, in the same class, FrequencyAnalyser I added another subroutine, which will plot all this using pyplot. Here we can find a trend / pattern, which I will explain later.

Plot Subroutine

```
def plot(self):
    def second(x): return x[1]

    self.words = list(self.words.items())
    self.chars = list(self.chars.items())

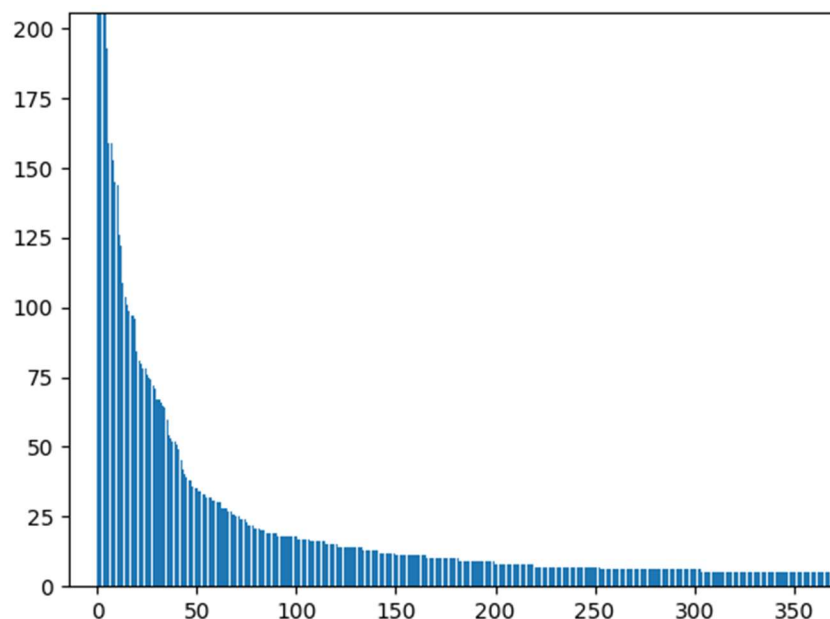
    self.words.sort(key=second, reverse=True)
    self.chars.sort(key=second, reverse=True)

    print("Plotting bar graph for word frequencies.")
    plt.bar(range(len(self.words)), list(
        map(second, self.words)), align='center')
    # plt.xticks(range(len(self.words)), list(
    #     map(lambda x: x[0], self.words)))
    plt.show()

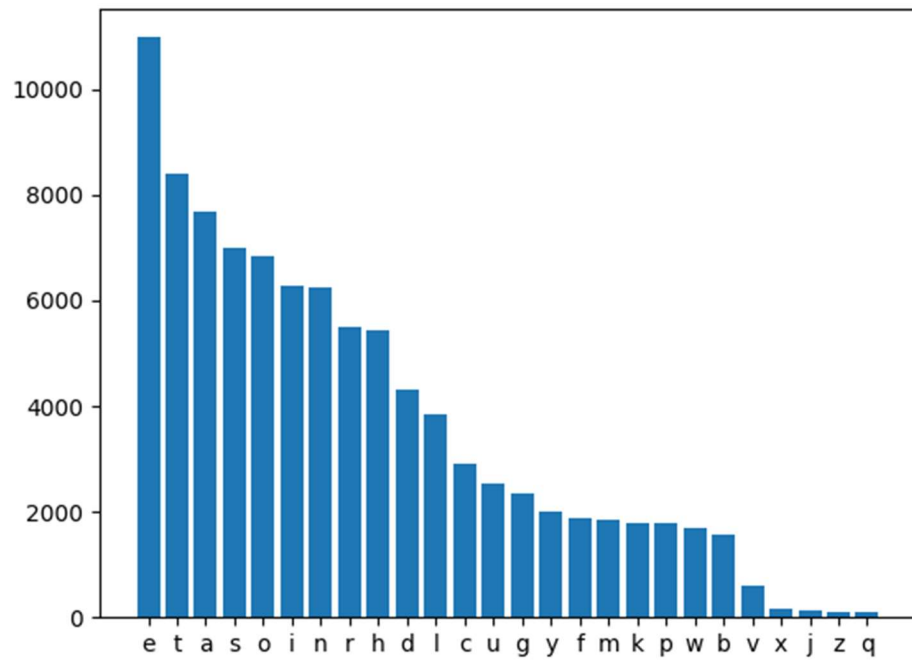
    print("Plotting bar graph for character frequencies.")
    plt.bar(range(len(self.chars)), list(
        map(second, self.chars)), align='center')
    plt.xticks(range(len(self.chars)), list(
        map(lambda x: x[0], self.chars)))
    plt.show()
```

Result

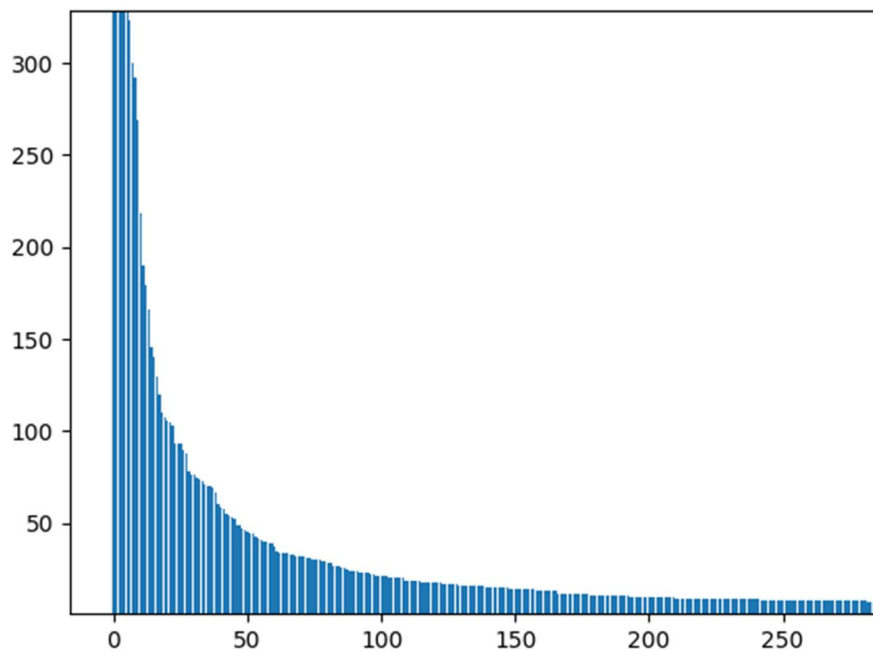
Word Frequencies of Space Odyssey



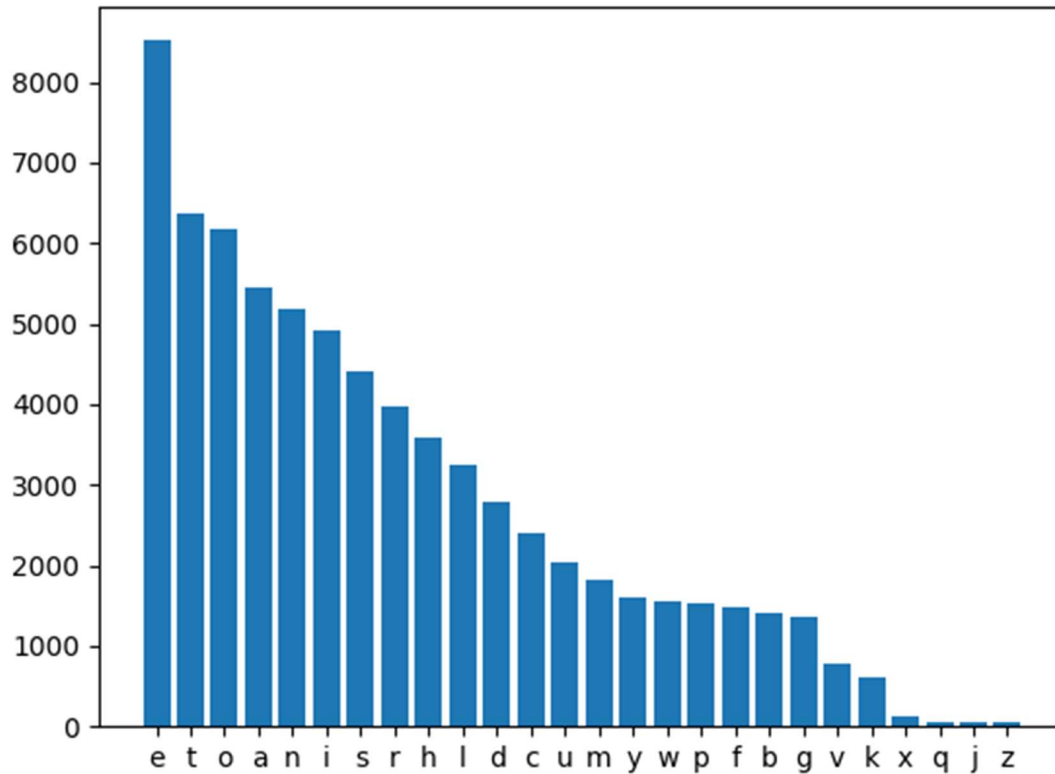
Character Frequencies of Space Odyssey



Word Frequencies of Blade Runner



Character Frequencies of Blade Runner



Part 5: The Proof

In the previous part, the curve that was common in all movies, is called a **Benford's Curve** which is a very common trend observed in every trend aspect.

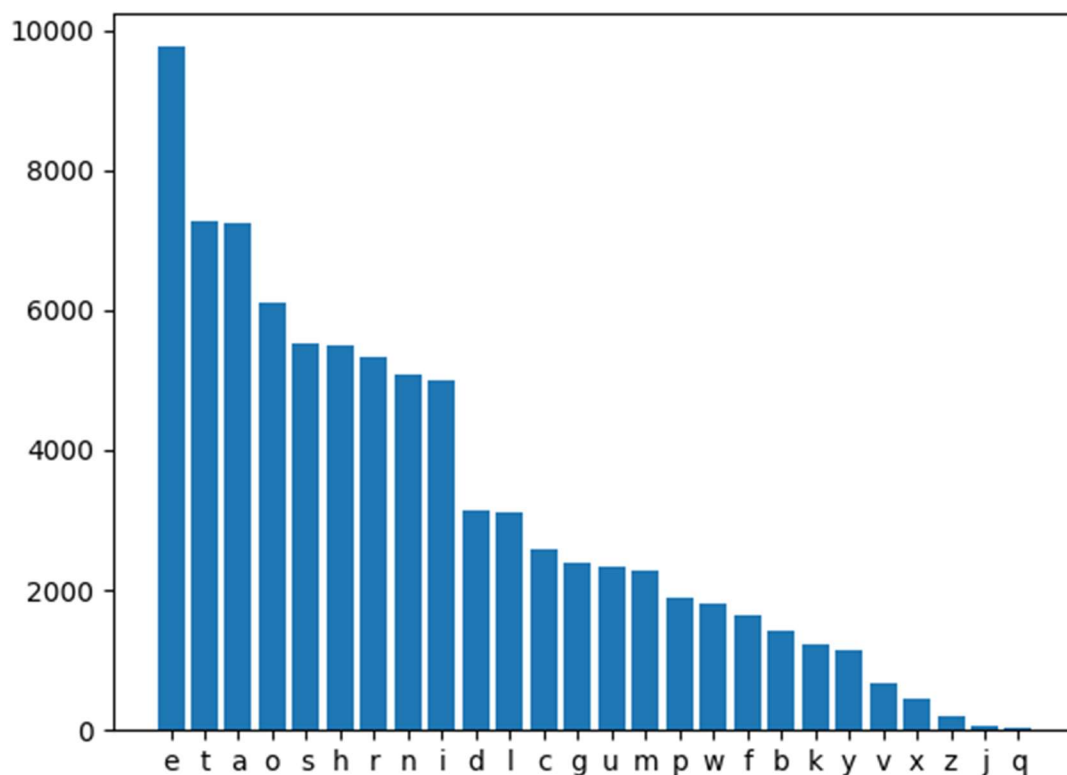
The Benford's law can be used to determine voting fraud, flood detections, etc.

I have proved Benford's law by using these movies, now with this proof, I will be showing how to detect languages of origin.

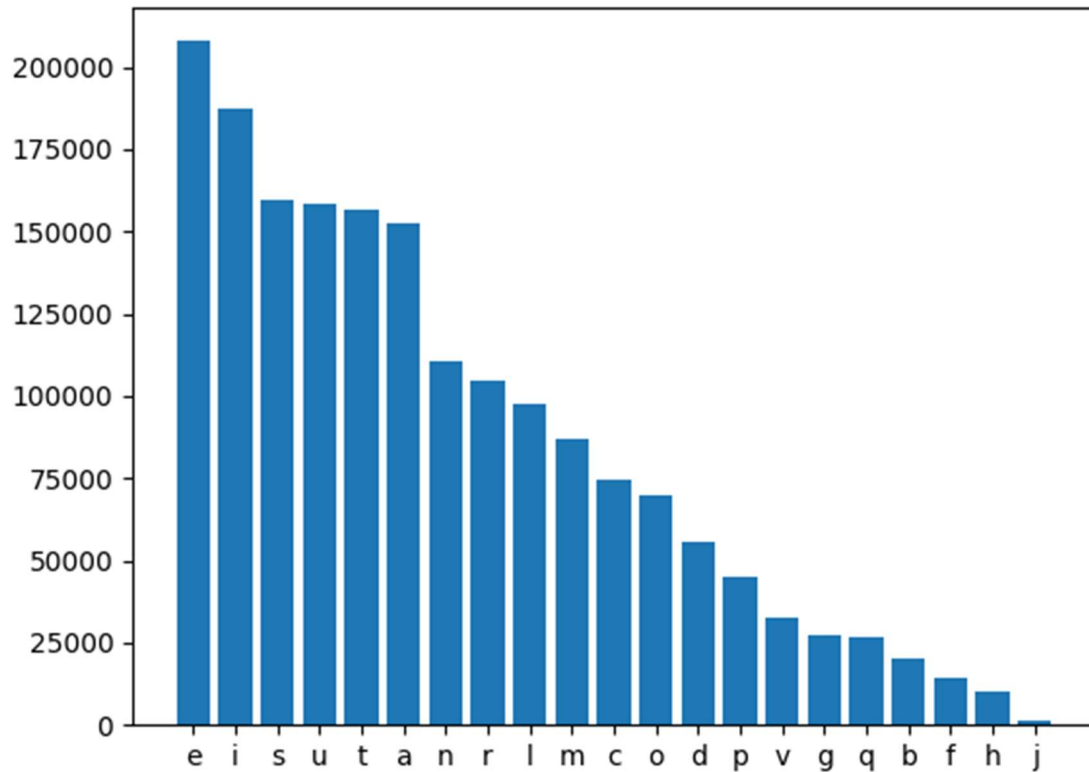
Previously, I have taken the *Lorem Ipsum* file, containing 2MB of filler information. This has a **Latin** origin, and English has **Greek and Latin** origin. Showing us some similarities, as well as **distinct** features that separate them.

Observe the frequency of **q**, **h** and **z** on these.

Here is an English Movie, MadMax.



Here is the Latin, Lorem Ipsum filler text.



We can see very distinct differences in the usage of syllables. The letter h has rarely been used, where as in English, it is part of most common word, which is “The”. Thus, we can find differences like this.

Here is the main, which has been running all these parts

```
# Word Frequency Analysis
from TextSplitter import TextSplitter
from BatchMaker import BatchMaker
from FrequencyAnalyser import FrequencyAnalyser

if __name__ == '__main__':
    files = [
        "2001_ASpaceOdyssey.txt",
        "BladeRunner.txt",
        "Dune.txt",
        "FightClub.txt",
        "LoremIpsumFiller.txt",
        "MadMax.txt",
        "Matrix.txt",
        "Memento.txt",
        "StarWars_EmpireStrikesBack.txt"
    ]

    bm = BatchMaker(files)
    bm.write_batch()
    bm.create_folders()

    for f in files:
        print()
        print(f"Script :: {f}")

        movie = TextSplitter(f)
        movie.partition_words()
        movie.partition()

    for f in files:
        print()
        print(f"Script :: {f}")
        fs = FrequencyAnalyser(f)
        fs.analyse()
        fs.plot()
```

Thank You for reading patiently.

Bhargav Kantheti