# Estimating the distribution of English alphabets in a set of books

Ramkumar

31-10-2024

# Problem definition

Aim is to estimate the distribution of English alphabets in a set of books and check if it is approaching Normal distribution in Probability.

The following books were taken for this work

- ▶ Animal Farm, by George Orwell
- ▶ Around the world in eighty days, by Jules Verne
- ▶ Flow, by Philip Ball
- ▶ For the love of Physics, by Walter Lewin
- ▶ Ikigai, by Hector Garcia and Francesc Miralles
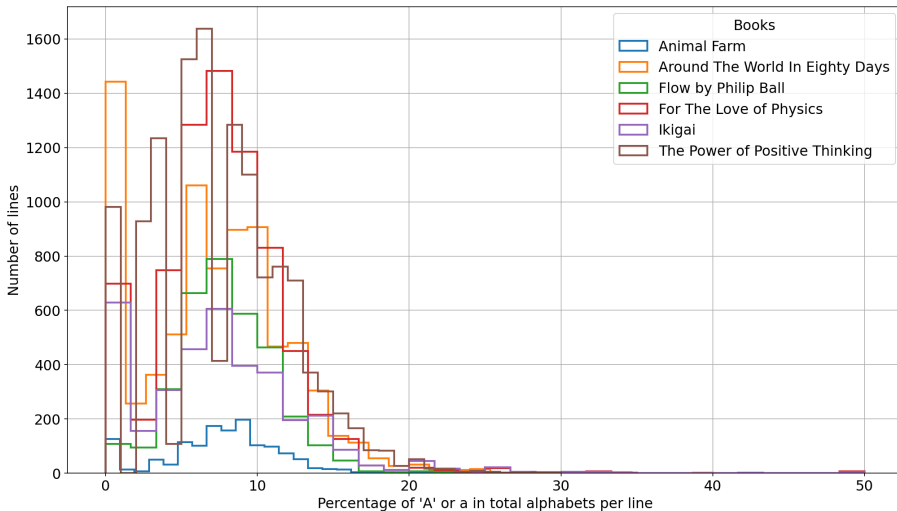- ▶ The power of positive thinking, by Norman Vincent

# Solution approach

A couple of Python codes were generated to read pdf files and process line-wise data with the following Pseudocode.
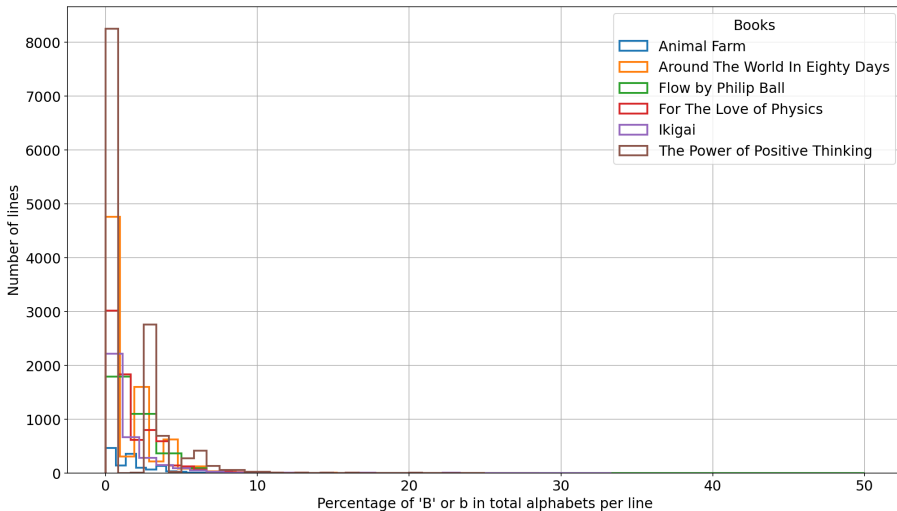
```
1    Start program
2
3    Open pdf file
4    for each page in file:
5        do
6        for each line in current page:
7            do
8            for each alphabet in English Alphabets:
9                do
10                compute no. of occurences of current alphabet in the line
11                compute total number of alphabets in current line
12                compute fraction of above two values and store them
13
14            end for
15        end for
16    end for
17
18    plot histograms for each alphabet
19
20    End program
```

Histograms of all alphabets
compared against
all chosen books
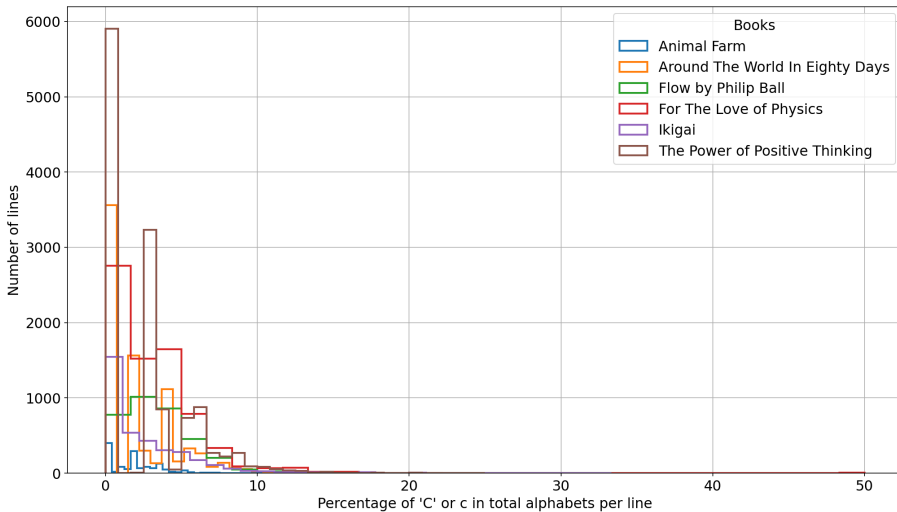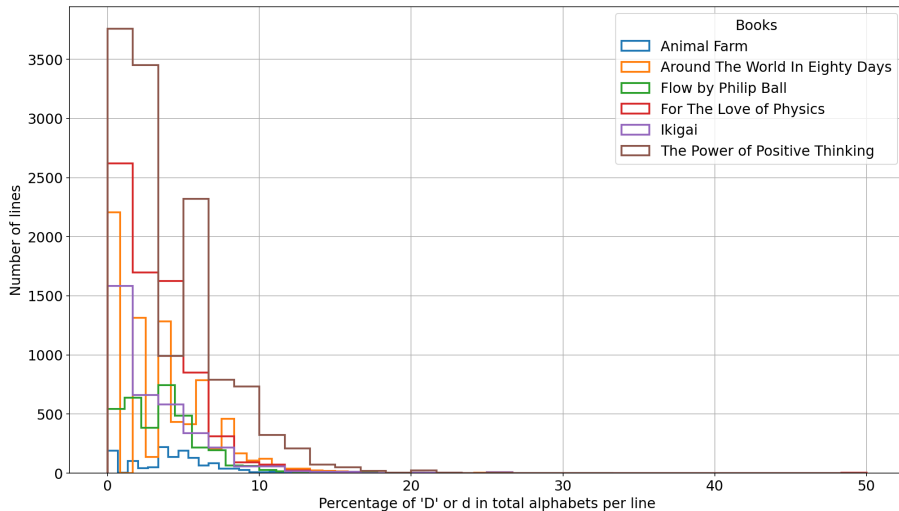
# A or a

# B or b

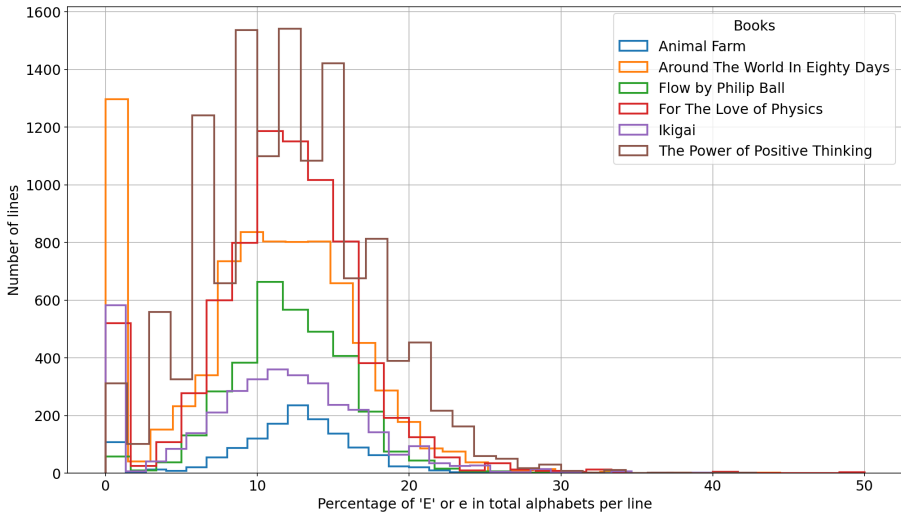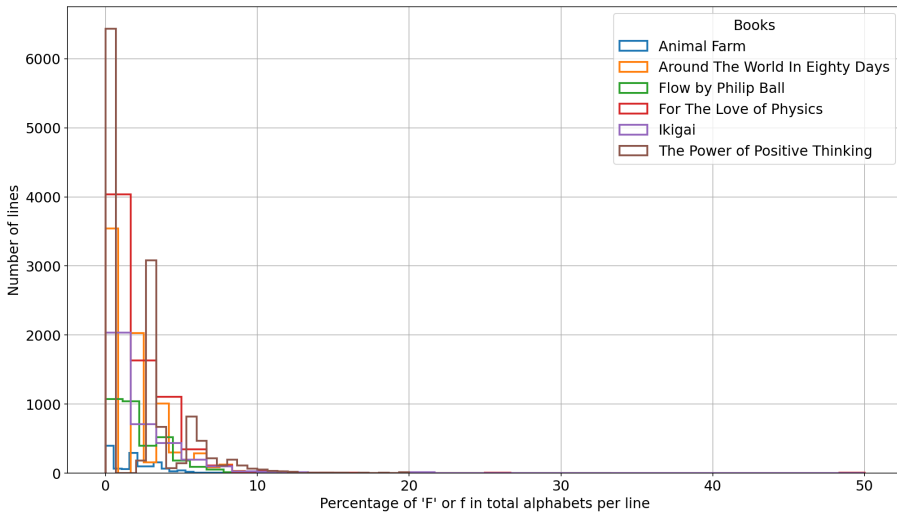# C or c

# D or d

# E or e

# F or f

# G or g

# H or h

# I or i

# J or j

# K or k

# L or l

# M or m

# N or n

# O or o

# P or p

# Q or q

# R or r

# S or s

# T or t

# U or u

# V or v



Alphabets distribution in English books

# W or w

# X or x

# Y or y

# Z or z

# Observations

- Least used letters in the chosen books are X, followed by Q and then Z

- Most used letters in the chosen books are E, followed by A, then O and then T

- Most used letters appear to follow **Normal distribution**

- Least used letters appear to follow **Exponential distribution**

# Python code - Reader script

# Reader script I

```python
#!/bin/python3
"""=========================================================
Alphabet statistics generator from pdf books

Ramkumar
Wed Oct 30 04:55:14 PM IST 2024
=========================================================="""

# importing needed modules
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from pypdf import PdfReader
import os, glob

#==========================================================

# specifying names of books and its contained directory
bookNames = ["For The Love of Physics.pdf",
             "Around The World In Eighty Days.pdf",
             "Flow by Philip Ball.pdf",
             "Animal Farm.pdf",
             "The Power of Positive Thinking.pdf",
             "Ikigai.pdf"]
directory = "books/"

# specifying starting and ending page numbers to exclude title, index etc...
pageStart = [10,4,12,5,5,8]
pageEnd   = [161,320,189,83,700,115]
```

# Reader script II

```python
# creating a directory to store line-wise data files
os.system("rm -rf linewise_data && mkdir linewise_data")

# looping through books
for I in range(len(bookNames)):
    # preparing book name
    bookName = directory + "/" + bookNames[I]

    print("reading book : ",bookNames[I])

    # reading book
    reader = PdfReader(bookName)

    # getting number of pages
    N_pages = len(reader.pages)

    # setting index to start page and end page numbers in python indexing
    start = pageStart[I]
    end   = pageEnd[I]

    # extracting lines
    total_lines = []
    for i in range(start,end):
        # extracting text from current page
        page = reader.pages[i]
        content = page.extract_text()

        # extracting lines from the content
        lines = content.split("\n")
        total_lines.extend(lines)
```

# Reader script III

```
62              print("reading page : ",i-start+1," of ",end-start+1)

64       # counting number of characters in each line
         a_list = []; b_list = []; c_list = []; d_list = []; e_list = []; f_list = []
66       g_list = []; h_list = []; i_list = []; j_list = []; k_list = []; l_list = []
         m_list = []; n_list = []; o_list = []; p_list = []; q_list = []; r_list = []
68       s_list = []; t_list = []; u_list = []; v_list = []; w_list = []; x_list = []
         y_list = []; z_list = []; line_length = [];

70
         idx = 1
72       for line in total_lines:
             # extracting total length of alpha characters
74           length = len([char for char in line if char.isalpha()])

76           # extracting total number of each alphabets
             a_len = len([char for char in line if char == 'a' or char == 'A'])
78           b_len = len([char for char in line if char == 'b' or char == 'B'])
             c_len = len([char for char in line if char == 'c' or char == 'C'])
80           d_len = len([char for char in line if char == 'd' or char == 'D'])
             e_len = len([char for char in line if char == 'e' or char == 'E'])
82           f_len = len([char for char in line if char == 'f' or char == 'F'])
             g_len = len([char for char in line if char == 'g' or char == 'G'])
84           h_len = len([char for char in line if char == 'h' or char == 'H'])
             i_len = len([char for char in line if char == 'i' or char == 'I'])
86           j_len = len([char for char in line if char == 'j' or char == 'J'])
             k_len = len([char for char in line if char == 'k' or char == 'K'])
88           l_len = len([char for char in line if char == 'l' or char == 'L'])
             m_len = len([char for char in line if char == 'm' or char == 'M'])
90           n_len = len([char for char in line if char == 'n' or char == 'N'])
             o_len = len([char for char in line if char == 'o' or char == 'O'])
92           p_len = len([char for char in line if char == 'p' or char == 'P'])
```

# Reader script IV

```python
            q_len = len([char for char in line if char == 'q' or char == 'Q'])
            r_len = len([char for char in line if char == 'r' or char == 'R'])
            s_len = len([char for char in line if char == 's' or char == 'S'])
            t_len = len([char for char in line if char == 't' or char == 'T'])
            u_len = len([char for char in line if char == 'u' or char == 'U'])
            v_len = len([char for char in line if char == 'v' or char == 'V'])
            w_len = len([char for char in line if char == 'w' or char == 'W'])
            x_len = len([char for char in line if char == 'x' or char == 'X'])
            y_len = len([char for char in line if char == 'y' or char == 'Y'])
            z_len = len([char for char in line if char == 'z' or char == 'Z'])

            # appending to the lists
            a_list.append(a_len); b_list.append(b_len); c_list.append(c_len)
            d_list.append(d_len); e_list.append(e_len); f_list.append(f_len)
            g_list.append(g_len); h_list.append(h_len); i_list.append(i_len)
            j_list.append(j_len); k_list.append(k_len); l_list.append(l_len)
            m_list.append(m_len); n_list.append(n_len); o_list.append(o_len)
            p_list.append(p_len); q_list.append(q_len); r_list.append(r_len)
            s_list.append(s_len); t_list.append(t_len); u_list.append(u_len)
            v_list.append(v_len); w_list.append(w_len); x_list.append(x_len)
            y_list.append(y_len); z_list.append(z_len); line_length.append(length)

            print("processing line = ",idx," of ",len(total_lines))
            idx += 1

    # preparing pandas dataframe to store the results
    filename = "alphabetCount_"+bookNames[I].split(".pdf")[0]+".csv"
    fid = pd.DataFrame(np.transpose([
                    a_list,b_list,c_list,d_list,e_list,f_list,g_list,
                    h_list,i_list,j_list,k_list,l_list,m_list,n_list,
                    o_list,p_list,q_list,r_list,s_list,t_list,u_list,
```

# Reader script V

```
124                     v_list , w_list , x_list , y_list , z_list , line_length ]) ,
                        columns = [ "a" , "b" , "c" , "d" , "e" , "f" , "g" ,
126                                  "h" , "i" , "j" , "k" , "l" , "m" , "n" ,
                                   "o" , "p" , "q" , "r" , "s" , "t" , "u" ,
128                                  "v" , "w" , "x" , "y" , "z" , "line_length" ])
         fid . to_csv ( "linewise_data /"+filename , index = None )
130
     print ( "done" )
132
     #════════════════════════════════════════════════════════════
```

# Python code - Post-processing script

# Post-processing script I

```python
#!/bin/python3
"""=====================================================================
post processing for the script_reader.py

Ramkumar
Wed Oct 30 06:12:10 PM IST 2024
====================================================================="""

# importing needed modules
import pandas as pd
import matplotlib.pyplot as plt
import os,glob


#=====================================================================

# reading data files
fileNames = sorted(glob.glob1(os.getcwd()+"/linewise_data/","*.csv"))

# preparing booknames
bookNames = [name.split("alphabetCount_")[1].split(".csv")[0] for name in
        fileNames]

# preparing list to store line-normalized character values
a_list = [];  b_list = [];  c_list = [];  d_list = [];  e_list = [];  f_list = []
g_list = [];  h_list = [];  i_list = [];  j_list = [];  k_list = [];  l_list = []
m_list = [];  n_list = [];  o_list = [];  p_list = [];  q_list = [];  r_list = []
s_list = [];  t_list = [];  u_list = [];  v_list = [];  w_list = [];  x_list = []
y_list = [];  z_list = []

charArray = ["a","b","c","d","e","f","g","h","i","j","k","l","m","n",
```

# Post-processing script II

```
30                      "o","p","q","r","s","t","u","v","w","x","y","z"]

32  # looping through the data files
    for file in fileNames:
34      # reading data
        fid = pd.read_csv("linewise_data/"+file)
36
        # looping through characters, normalizing and storing 'em in the lists
38      for char in charArray:
            fid[char] = fid[char]/(fid["line_length"]+1)*100
40          eval(char+"_list.append(fid[\""+char+"\"])")

42  # preparing directory to store graphs
    os.system("rm −rf histograms && mkdir histograms")
44
    # preparing plots and saving them
46  for char in charArray:

48      # getting the current list
        exec("curr_list = "+char+"_list")
50
        # plotting histogram
52      plt.rcParams.update({"font.size":15})
        plt.figure(figsize=(16,9))
54      for i in range(len(bookNames)):
            plt.hist(curr_list[i],bins=30,histtype="step",label=bookNames[i],
56                  density=False,linewidth=2)
        plt.grid()
58      plt.xlabel("Percentage of \'"+char.upper()+"\' or "+char+" in total
         alphabets per line")
        plt.ylabel("Number of lines")
```

# Post-processing script III

```python
60        #   plt.title("Alphabet : "+char.upper()+" or "+char)
          #   plt.legend(loc=(1.01,0.75))
62        plt.legend(title="Books")
          plt.savefig("histograms/"+char+".png",dpi=150,bbox_inches="tight")
64        plt.close()

66        print("Character : ",char.upper())

68  print("done")

70  #========================================================================
```

End