# Statistical Analysis Of Text Files

**Course Instructor** : Samy S. Soliman

Nada Ismail, 202-001-387

December 17, 2022

# 1    Introduction

This report will discuss project's statiscal analysis, results and the Matlab code.

# 2    Statiscal analysis

let x be an array containing the encoded values (1 to 61) and f(x) be the probability of x.

The equation of mean is:

$$m = \sum_{x=1}^{61} x f(x) \tag{1}$$

The equation of variance is:

$$\sigma^2 = \sum_{x=1}^{61} (x-m)^2 f(x) \tag{2}$$

The equation of skewness is:

$$\frac{1}{\sigma^3} \sum_{x=1}^{61} (x-m)^3 f(x) \tag{3}$$

The equation of kurtosis is:

$$\frac{1}{\sigma^4} \sum_{i=1}^{61} (x-m)^4 f(x) \tag{4}$$

# 3    Code

## 3.1    Data analysis functions

### 3.1.1    filter

this function filters out the unwanted charcters like: $,&,,%. isstrprop() method is used to return an array of ones and zeroes that shows which indexes in the original data are alphanumeric.

```
function filtered_data = filter(app,data)
    filtered_data=[];
    j=1;% counter for the return array
    isAlphaNum=isstrprop(data,"alphanum");%returns boolean array of alphanum
 characters in the data array
    for i=1:length(data)
        if isAlphaNum(i)
            filtered_data(j)= data(i);
            j=j+1;
        end
    end
end
```

### 3.1.2    get_encoded_data

This function returns a 61x1 array where the index-1 is the encoded value and the each element is number of times the corresponding character appeared in the file.

```
function results = get_encoded_data(app,data)%return array where the index is the
 encoded value and the value of each element is the # of occurences
    char_vals=['0','1','2','3','4','5','6','7','8','9','a','A','b','B','c','C','d
','D','e','E','f','F','g','G','h','H','i','I','j','J','k','K','l','L','m','M','n','N
','o','O','p','P','q','Q','r','R','s','S','t','T','u','U','v','V','w','W','x','X','y',
'Y','z','Z'];
    results=zeros(1,62);
```

```
4            for i=1:length(data)
5                for j=1:length(char_vals)
6                    if data(i)==char_vals(j)
7                        results(j)=results(j)+1;
8                    end
9                end
10            end
11        end
```

### 3.1.3 get_sorted_array_of_chars_based_on_occurenceArray

This function Returns a sorted array of character based where the most common character is the first.

```
1        function  [sorted_char_array, sorted_array] =
    get_sorted_array_of_chars_based_on_occurenceArray(app,occurence_array)
2            sorted_char_array=['0','1','2','3','4','5','6','7','8','9','a','A','b','B','c
    ','C','d','D','e','E','f','F','g','G','h','H','i','I','j','J','k','K','l','L','m','M'
    ,'n','N','o','O','p','P','q','Q','r','R','s','S','t','T','u','U','v','V','w','W','x',
    'X','y','Y','z','Z'];
3            sorted_array=occurence_array;
4            [sorted_array,sort_Index] = sort(sorted_array,'descend');
5            sorted_char_array = sorted_char_array(sort_Index);
6        end
```

### 3.1.4 get_statistics

This function return 4x1 array where the elements are the mean, variance, skewness, and kurtosis.

```
1        function [m,v,s,k] = get_statistics(app,data)
2            f_x=data./sum(data);
3            x=1:62;
4            x=x-1;
5            m = sum(x.*f_x);
6            v = sum(((x-m).^2).*f_x);
7            %data_p3=sum((data-m).^3);
8            s = sum(((x-m).^3).*f_x)/(v^(3/2));
9            %data_p4=sum((data-m).^4);
10            k = sum(((x-m).^4).*f_x)/(v^2);
11        end
```

### 3.1.5 get_f_x

This function returns the PMF.

```
1        function f_x = get_f_x(app,data)
2            f_x=data./sum(data);
3        end
```

### 3.1.6 get_F_x

This function returns the CDF.

```
1        function F_x = get_F_x(app,f_x)
2            F_x=[];
3            for i=1:length(f_x)
4                F_x(i)=sum(f_x(1:i));
5            end
6        end
```

## 3.2   GUI callback functions

### 3.2.1   ImportfileButtonPushed

This function promts the user to chose the text file, passes its content to filter() then get occurence array and does the stastistical analysis and plot.

```
1               [filename,path]=uigetfile('*.txt');
2          app.PathEditField.Value=path;
3          app.data=fileread(filename);
4          app.PathEditField.Value=path;
5          %app.draftTextArea.Value=data;
6          app.filter_data=filter(app,app.data);
7          app.filter_data_characters=char(app.filter_data);
8          app.occurences_of_characters=get_encoded_data(app,app.filter_data_characters);
9
10         [m,v,s,k]=get_statistics(app,app.occurences_of_characters);
11         app.ThemeanisTextArea.Value=string(m);
12         app.ThevarianceisTextArea.Value=string(v);
13         app.TheskewdnessisTextArea.Value=string(s);
14         app.ThekurtosisisTextArea.Value=string(k);
15
16         %shift ticks to start at origin and then plot
17
18         app.UIAxes_3.XLim=[1 63];
19         app.UIAxes_3.XTickLabel = {' '; '0'; '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'; '28'; '29'; '30'; '31'; '32'; '33'; '34'; '35'; '36'; '37';
     '38'; '39'; '40'; '41'; '42'; '43'; '44'; '45'; '46'; '47'; '48'; '49'; '50'; '51';
    '52'; '53'; '54'; '55'; '56'; '57'; '58'; '59'; '60'; '61';'62'};
20         app.UIAxes_3.YLim=[0 1];
21         app.f_x=get_f_x(app,app.occurences_of_characters);
22         stairs(app.UIAxes_3,app.f_x);
23
24         app.UIAxes_2.XLim=[1 63];
25         app.UIAxes_2.XTickLabel = {' '; '0'; '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'; '28'; '29'; '30'; '31'; '32'; '33'; '34'; '35'; '36'; '37';
     '38'; '39'; '40'; '41'; '42'; '43'; '44'; '45'; '46'; '47'; '48'; '49'; '50'; '51';
    '52'; '53'; '54'; '55'; '56'; '57'; '58'; '59'; '60'; '61';'62'};
26         app.UIAxes_2.YLim=[0 1];
27         app.F_x=get_F_x(app,app.f_x);
28         stairs(app.UIAxes_2,app.F_x);
```

### 3.2.2   NumberEditFieldValueChanged

This function is triggered when the user enters the most repeated character that he wants. The user cannot choose more than number of diffrent characters in the file.

```
1          function NumberEditFieldValueChanged(app, event)
2             limit= length(nonzeros(app.occurences_of_characters)); %do not allow user to
    ask for more characters than there are in the document if it only has a's and b's
    then he cannot askk for 3 most repeated
3             if (app.NumberEditField.Value>limit)
4                app.NumberEditField.Value=limit;
5                msgbox(sprintf('There are %d diffrent alphanumeric charcters in the
    document,\n you cannot chose more than that',limit));
6             end
7             [uncutmsg,x]=get_sorted_array_of_chars_based_on_occurenceArray(app,app.
    occurences_of_characters);
8             msg=[];
9             for i=1:app.NumberEditField.Value
10                msg(i)=uncutmsg(i);
11            end
12            app.ThemostrepeatedcharctersareTextArea.Value=char(msg);
13
14         end
```

### 3.2.3   EncodingSwitchValueChanged

this function gives the user the option to see the numbers or the characters corresponding to those numbers.

```
1          if app.EncodingSwitch.Value(1:length('Show characters'))=='Show characters'
2                app.UIAxes_3.XTickLabel = {' '; '0'; '1'; '2'; '3'; '4'; '5'; '6'; '7'; '
    8'; '9'; 'a';'A';'b';'B';'c';'C';'d';'D';'e';'E';'f';'F';'g';'G';'h';'H';'i';'I';'j';
    'J';'k';'K';'l';'L';'m';'M';'n';'N';'o';'O';'p';'P';'q';'Q';'r';'R';'s';'S';'t';'T';'
    u';'U';'v';'V';'w';'W';'x';'X';'y';'Y';'z';'Z'};
```

```
3              app.UIAxes_2.XTickLabel = {' '; '0'; '1'; '2'; '3'; '4'; '5'; '6'; '7'; '
      8'; '9'; 'a';'A';'b';'B';'c';'C';'d';'D';'e';'E';'f';'F';'g';'G';'h';'H';'i';'I';'j';
      'J';'k';'K';'l';'L';'m';'M';'n';'N';'o';'O';'p';'P';'q';'Q';'r';'R';'s';'S';'t';'T';'
      u';'U';'v';'V';'w';'W';'x';'X';'y';'Y';'z';'Z'};
4           else
5              app.UIAxes_3.XTickLabel = {' '; '0'; '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'; '28'; '29'; '30'; '31'; '32'; '33'; '34'; '35'; '36';
      '37'; '38'; '39'; '40'; '41'; '42'; '43'; '44'; '45'; '46'; '47'; '48'; '49'; '50'; '
      51'; '52'; '53'; '54'; '55'; '56'; '57'; '58'; '59'; '60'; '61';'62'};
6              app.UIAxes_2.XTickLabel = {' '; '0'; '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'; '28'; '29'; '30'; '31'; '32'; '33'; '34'; '35'; '36';
      '37'; '38'; '39'; '40'; '41'; '42'; '43'; '44'; '45'; '46'; '47'; '48'; '49'; '50'; '
      51'; '52'; '53'; '54'; '55'; '56'; '57'; '58'; '59'; '60'; '61';'62'};
7           end
```
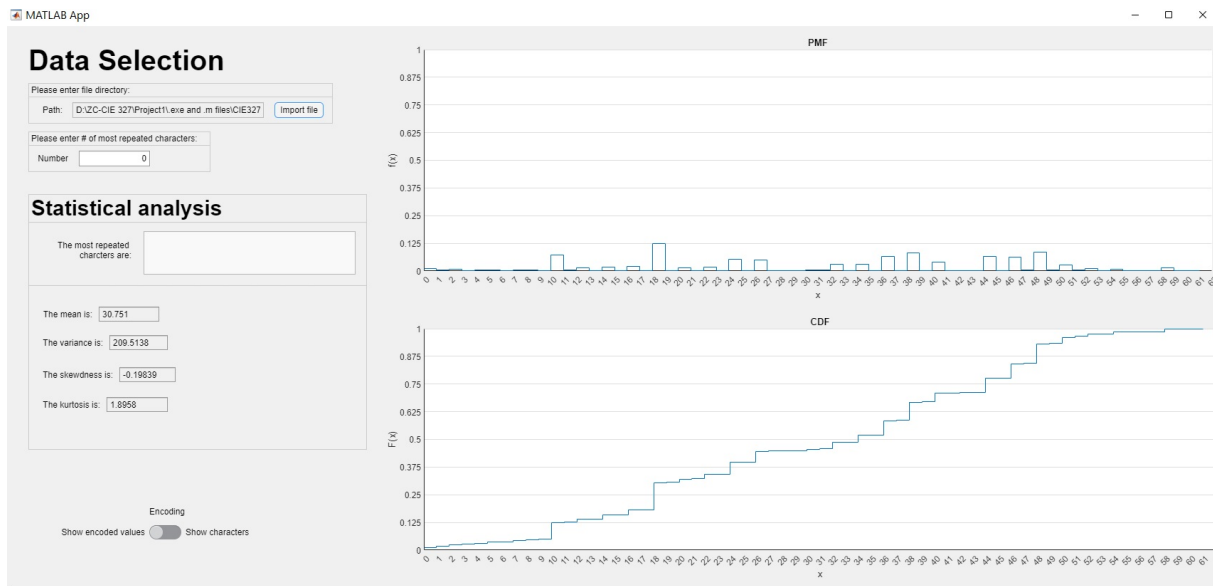
# 4    Results



Figure 1: A screenshot of the results using the provided sample text