

## Experiment A7

### Title- Text Analytics

Aim- 1. Extract sample document & apply following document preprocessing method. Tokenization, PostTagging, stop words removal, Stemming and Lemmatization.

2. Create representation of document by calculating Term frequency & Inverse document frequency

Theory- Text analysis also called text analytics refers to the representation, processing & modeling of textual data to derive useful insights.

Text databases are rapidly growing due to the increasing amount of information available in technique electronic form. Traditional information retrieval techniques becomes inadequate for the increasingly popular & essential part of data mining. The discovery of proper patterns & analyzing the text document from the huge volume of data is a major issue in real-world application.

In text analysis, we take the help of corpus. A corpus is a large collection of text used for various purpose in natural language processing.

Few example of Corpus are Shakespeare Brown Corpus, Google N-grams Corpus.

## Steps in Text Analysis-

- ① Parsing - It is the process that takes unstructured data (text) & imposes a structure for further analysis. The unstructured data could be plain text file, weblog, XML, etc. Parsing deconstructs the provided text and renders it in more structured way.
- ② Search & retrieval - It is the identification of the documents in a Corpus that contains search items like specific words, phrases, topics or entities. Search items are generally called key terms.
- ③ Text mining - It uses the terms & indices produced by the previous two steps to discover meaningful insights pertaining to domains or problems of internet. Clustering & classification techniques can also be applied to derive insights.

## Text Pre-processing techniques

- ① Tokenization - It is the process of splitting a text into a list of tokens. They work by separating the words using punctuation & spaces. It does not discard the punctuation.

Example -

```
from nltk.tokenize import word_tokenize
```

```
text = "Hello world!"
```

```
word_tokenize(text)
```

```
output - ['Hello', 'world', '!']
```

② Pos tagging - Part of speech tagging explains how a word is used in sentence & it is used to extract relationships between words. It tells if a word is noun, pronoun, objective, etc.

Input :- "You gave me a car"

Output - [('You', 'PRP'), ('gave', 'VBD'),  
('me', 'PRP'), ('a', 'DT') ('car', 'NN')]

PRP - Personal Pronoun DT - Determiner

VBD - verb past tense NN - Noun

③ Stop words removal - Stopwords are common words that are mostly irrelevant for text analysis, like 'a', 'or', 'the', 'is', 'for', etc.

Example - Input - I ate an apple

Output - I ate apple

④ Stemming - It is reducing words to their base or root form by removing a few suffices from the words. It is text normalization technique. The Porter stemming algorithm is mostly used for this. However stemming does not provide correct form.

Input :	Aging	Books	Learning	Obesity
Output :	Act	Book	Learn	Obes

\* Term frequency (TF) : It measures the frequency of a word in a given document. It highly depends on the length of the document & the generality

of words:

$$TF = \frac{\text{count of the term in document } (T)}{\text{total number of terms in document}}$$

⑤ Lemmatization - It is similar to stemming technique that aims to get the base words. But unlike stemming that might produce improper words. The lemmatization process does not only turn the suffixes but also use lexical knowledge base.

Input:	learning	caring	causes
Output:	learn	care	cause

\* Inverse document frequency (IDF) - It is measure of the importance of word. IDF provides weightage to each word based on its frequency in the corpus. 0.

$$IDF = \log\left(\frac{\text{Total no of documents in Corpus}}{\text{Number of documents containing t}}\right)$$

Conclusion-

The representation of sample document & text - Preprocessing has been implemented with success.

## Experiment A8

### Title - Data visualization I

Aim - 1. Use the inbuilt dataset 'titanic'. The dataset contains 891 rows & contains information about the passengers who boarded the titanic ship. Use the seaborn library to see if we can find any pattern in data.

2. Write a code to check how the price of ticket (column name : 'fare') for each passenger is distributed by plotting a histogram.

### Theory -

- Data visualization is actually a set of data points & information that are represented graphically to make it easy & quick for user to understand.
- Data visualization is good if it has a clear meaning, purpose & is very easy to interpret without provide an accessible way to see & understand trends, outliers patterns. In data by visual effect or elements such as chart, graph & maps.

### Characteristics of Effective Graphical visual -

- 1) It shows or visualizes data very clearly in an understandable manner.

- 2) It encourages viewers to compose different pieces of data
- 3) It closely integrates statistical & verbal descriptions of dataset.
- 4) It also helps in identifying area that needs more attention & improvement.  
Using graphical representation, a story can be told more efficiently. Also, it requires less time to understand picture than it takes to understand textual data.

### TYPES OF DATA VISUALIZATION-

- ① Comparative plots- They are used for comparing the data points. They are -
  - i. column & bar charts
  - ii. Line chart
  - iii. Area chart
  - iv. Bubble chart
  - v. pie chart
- ② Statistical plots- They are useful to show results of statistical analysis.  
They are -
  - i. histogram
  - ii. treemap
  - iii. Barplot
  - iv. waterfall chart

③ Topology plots - It uses geometric structures to show relationships of connections between datapoints in your dataset.

They are - ① Linear topology  
② Graph

④ Spatial plots - They are logical views for visualizing data which represents a map, location, shape etc.

They are - ① Raster surface  
② Point map  
③ Heat map

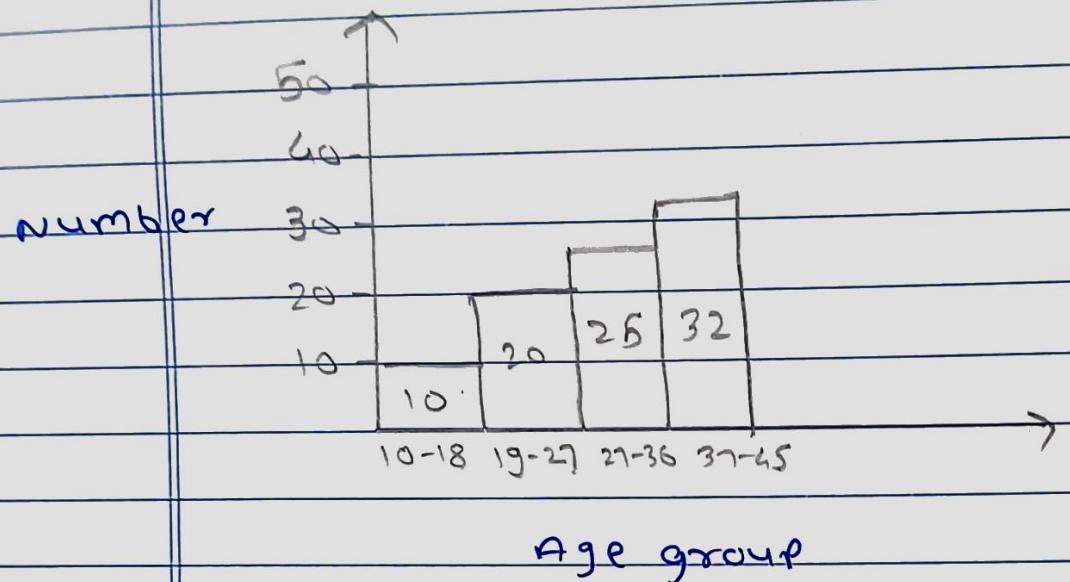
Histogram -

A histogram is used to represent the frequency of data in a dataset as well its distribution. It looks similar to a column chart but is different as it plots the frequency for each distribution rather than the actual value of any datapoint - itself.

Example -

Age group	Number of People
10-18	10
19-27	20
28-36	25
37-45	32

```
Experim  
#Import  
import p
```



### Seaborn Library-

Seaborn is an open-source Python library built on top of matplotlib. It is used for data visualization & exploratory data analysis.

Seaborn work easily with DataFrame & the Pandas library.

### Conclusion -

We use the Seaborn library to see if we can find patterns in data & to check price distribution.

## Experiment - Ag

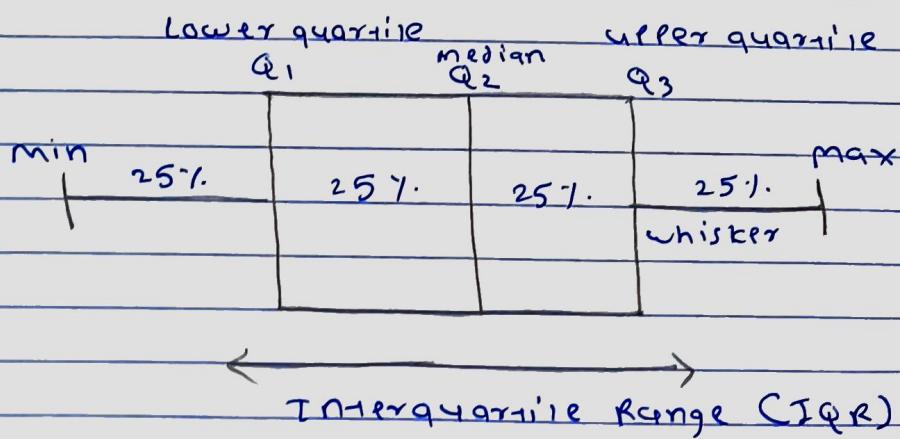
### Title - Data visualization II

Aim - Use the inbuilt dataset 'titanic' as used in the above problem. Plot a bar plot for distribution of age with respect to each gender along with the information about whether they survived or not ('sex' and 'age'). Write observations on the inference from the above statistics.

Theory: Box plot or whisker diagram shows how data points in a dataset are distributed. This plot is quite interesting as it shows several statistical values at once.

It depicts a group of numerical data through their quartiles. It is a simple way to visualize the shape of our data. It makes comparing characteristics of data between categories very easy.

#### Components of Box Plot:



- a) minimum - It is the minimum value in the dataset excluding the outliers.
- b) first quartile ( $Q_1$ ) - 25% of the data lies below the first or lower quartile.
- c) median ( $Q_2$ ) - It is the mid-point of the dataset. Half of the values lie below the median and half above.
- d) third quartile ( $Q_3$ ) - 75% of the data lies below the third or upper quartile.
- e) maximum - It is the maximum value in the dataset excluding the outliers.
- f) Interquartile Range (IQR) - The area inside the box (50% of the data) is known as the IQR.

$$IQR = Q_3 - Q_1$$

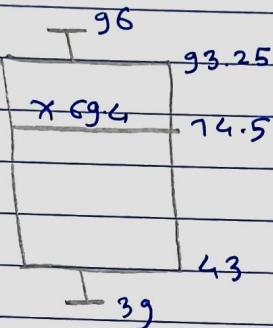
- g) Outliers - They are the data points below and above the lower and upper limit.

The box plots can have skewness and the median might not be at the center of the box.

Example - Below is the distribution of students and their score in science

Student	Ajay	Sunil	Ravi	Pinky	Raju	Shivq	Raj	Rinky	Ramesh	Adam
marks	65	54	96	94	93	85	84	39	66	40

The box plot for the same is:



The values shown by the box plots are:

minimum	39
I quartile	43
mean	69.4
median	74.5
II quartile	74.5
III quartile	93.25
maximum	96

Uses of a Box Plot.

- They provide a visual summary of the data with which we can quickly identify the average value of the data; how dispersed it is, if it is skewed or not.
- ↑ The median gives the average value of the data.
- They show if data is skewed or not
  - i) If median is at center - normal distribution

- ii) If median lies closer to I Quartile: Positive skew
- iii) If median lies closer to III Quartile: Negative Skew

Box Plots give an idea about the outliers which are points numerically distant from the rest of the data.

Conclusion:

Hence, we plot a box plot for the distribution in the given dataset, successfully.

## Experiment - A 10

Title : Data visualization III

Aim: Download the Iris flower dataset or any other dataset into a DataFrame. Scan the dataset and give the inference as

1. How many features are there and what are their types?
2. Create a histogram for each feature in the dataset to illustrate the feature distributions.
3. Create a boxplot for each feature in the dataset. Compare distributions and identify outliers.

Theory: Data visualization is actually a set of data points and information that are represented graphically to make it easy and quick for any user to understand.

It is good if it has a clear meaning, purpose and is very easy to interpret without requiring context.

Tools of data visualization provide an accessible way to see and understand trends, outliers and patterns in data by using visual effects or elements such as chart, graphs and maps.

( characteristics of effective graphical visualizations )

- a) It shows data clearly in an understandable manner.

- b) It encourages viewers to compare different descriptions of data set.
- c) It closely integrates statistical and verbal descriptions of data set.
- d) It helps in identifying areas that need more attention and improvement.

### Types of data / features:

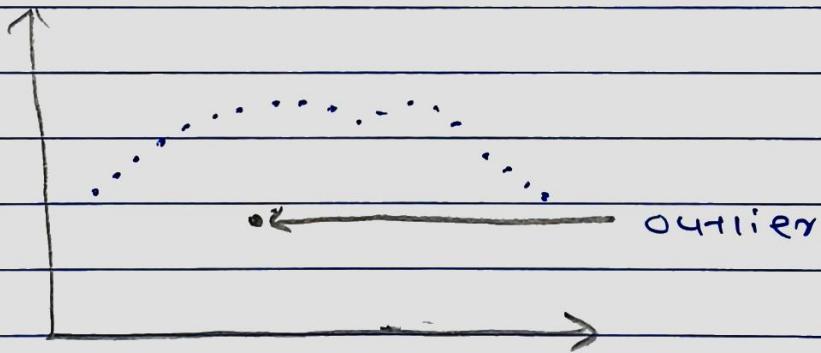
- 1) Numeric data: It consists of all types of numbers. Such features can be measured by any numerical variables. e.g. Age, scores, etc
- 2) Qualitative data: It consists of words, pictures and symbols, basically non-numeric features e.g. colours
- 3) Nominal data: It is used for labelling variables without any type of quantitative value. It is just a name of feature without applying any order e.g. Gender, ethnicity, etc
- 4) Ordinal data - This shows where a numbers is in order. It is the data which is placed into some kind of order by ~~which~~ their position on a scale e.g. letter grades: (A, B, C); Rank (1<sup>st</sup>, 2<sup>nd</sup>, 3<sup>rd</sup>...)

Outlier:

It is a data field that deviates significantly from the rest of the data objects and behaves in a different manner. They can be caused by measurement or execution errors. The analysis of outlier data is referred to as outlier analysis or outlier mining.

An outlier cannot be termed as noise or error.

Instead, they are suspected of not being generated by the same method as the rest of the data objects.



Conclusion: The histograms and box plots for each feature has been plotted and compared

Group B

Assignment - 1

Sahil Gadge

19CO021

TE Comp I

DSBDAL

Aim: write a code in JAVA for a simple WordCount application that counts the number of occurrences of each word in a given inputset using the Hadoop MapReduce framework on local Standalone setup.

Algorithm:

1. Start
2. Take the input data or file to be processed
3. Splits the incoming data into the smaller pieces (split)
4. maps data by performing action like filtering, grouping and sorting
5. Combines the data to improve the performance by reducing amount of data transferred across network
6. shuffles and sorts the output from all the mappers
7. Reduce aggregates the output of mappers using reduce() function
8. End

Steps :

Step 1: Create a text file in your local machine and write some text into it.

```
$ nano data.txt
```

Step 2: Check the text written in the data.txt file

```
$ cat data.txt
```

Step 3: Create a directory in HDFS, where to kept text file.

```
$ hdfs dfs -mkdir /test
```

Step 4: Upload the data.txt file on HDFS in the specific directory.

```
$ hdfs dfs -put /home/data.txt /test
```

Step 5: Write the MapReduce program using eclipse

Result :

The word count application is implemented in hadoop mapreduce

Group B

Assignment - 2

Sahil Gadge

19CO021

TE Comp 1

DSBDAL

Aim: design a distributed application using map Reduce which process a log file of a system.

Algorithm:

1. Start

2. flume collects streaming log data into HDFS from various HTTP sources and application servers

3. Data stored in HDFS which was collected by flume

4. pig parses these log files into the structured format

5. schema will defined by hive to the structured data and schema will be stored in hive metastore

6. TI picks any visualization supporting hive connectivity among Tableau and Hunk

7. End

Steps:

Step 1: Fume collects streaming log data into HDFS from various HTTP sources and application servers.

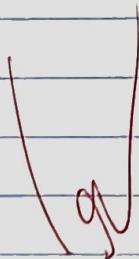
Step 2: The data collected by fume is stored into the HDFS

Step 3: Pig parses these log files into structured format through various UDF's

Step 4: Hive or Metastore will define schema to this structured data and schema will be stored in hive metastore.

Step 5: TI picks any visualization supporting hive connectivity among Tableau and Hunk.

Result: we have designed distributed application using mapReduce which process a log file of a system and resultant output is verified



Group B  
Assignment-3

Sahil Fadge  
19C0021  
TE Comp  
DSBDAL

Aim: Locate dataset (e.g. sample-weather.txt) for working on weather data which reads the text input files and finds average for temperature, dew point and wind speed.

Algorithm:

1. Start
2. It exceeds in 3 stages
  - 2.1 map stage process the input data which is in array or file format
  - 2.2 shuffle stage removes duplicate values and the data will be sorted or shuffle
  - 2.3 reduce stage aggregates the values come from previous stages against the corresponding keys
3. End

Steps:

Step 1: Download the dataset for various cities in different years.

Step 2: Make a project in eclipse and in that project create Java class and put code in it.

Step 3: Start the hadoop daemons  
start-dfs.sh

Step 4: move the dataset to the hadoop HDFS.

Step 5: Run the Jar file with below command  
hadoop jar /jar file location /dataset location  
in HDFS /output filename

Step 6: Move to localhost:50070, under Utilities  
select Browse file system and download it.

Step 7: See the result in downloaded file.

Result: Located the dataset working <sup>on</sup> weather data which reads text input files and finds the average for temperature, dew point and wind speed.