



Title - Text Analytics.

objective - students should able to do the text analytics using TFIDF algorithm.

prerequisite -> Basic of python

2) concept of TFIDF and Text analysis.

Theory -

1. Basic Concepts of Text Analytics.

One of the most frequent types of day-to-day conversion is text communication. In our everyday routine, we chat, message, tweet, share status, email, create blogs, and offer opinions and criticism.

2. Text Analysis operations using natural language Toolkit.

NLTK (natural language Toolkit) is a leading platform for building python programs to work with human language data.

It provides easy-to-use interfaces and lexical resources such as WordNet, along with a suite of text processing libraries for classification. ~~Req~~

Tokenization

Tokenization is the first step in text analytics. The process of breaking down a text paragraph into smaller chunks such as words or sentences is called Tokenization.

Sentence tokenization : split the paragraph into list of sentences using sent_tokenize() method.

- word tokenization - split a sentence into list of words - tokenize() method.

stop words removal -

stopwords considered as noise in the text. Text may contain stop words such as is, am, are, this, a, an, the, etc. In NLTK for removing stopwords, you need to create a list of stopwords, you need to create and filter out your list of token from these words.

Stemming and Lemmatization.

stemming is a normalization technique where list of tokenized words are converted into shortened root words to remove redundancy. stemming is the process of reducing inflected (or sometimes derived) words to their word stem, base or root form.

↳ Lemmatization vs Stemming

Stemming algorithm works by cutting the suffix from the word. In a broader sense cuts either the beginning or end of the word.

On the contrary, Lemmatization is a more powerful operation, and it takes into considering morphological analysis of the words. It returns the lemma which is the base from of all its inflectional forms.



In-depth linguistic knowledge is required to create dictionaries and look for the proper form of the word.

Pos Tagging:

Pos (parts of speech) tell us about grammatical information of words of the sentence by assigning specific token (Determiner, noun, adjective, adverb, verb, personal pronoun etc.) as tag (DT, NN, JJ, RB, PRP etc) to each words.

Text Analysis Model using TF-IDF

Term Frequency - inverse document frequency (TFIDF), is a numerical statistic that is intended to reflect how important a word is to a document in a collection or corpus.

• Term Frequency (TF)

It is measure of the frequency of word (w) in a document (d). TF is defined as the ratio of a word's occurrence in a document to the total number of words in a document. The denominator term in the formula is to normalize since all the corpus documents are of different lengths.

$$TF(w,d) = \frac{\text{occurrences of } w \text{ in document } d}{\text{Total no. of words in document } d}$$

ex.

Documents	Text	Total no. of words in a document
-----------	------	----------------------------------

A Jupiter is the largest planet

5

B Mars is the fourth planet from the sun.

8

The initial steps is to make a vocabulary of unique words and calculate TF for each document.

Inverse Document frequency (IDF)

It is the measure of the word Term frequency (TF) does not consider the importance of words such as of, 'and', etc. can be most frequently present but are of little significance.

Total no. of documents (N)

$$IDF(w,D) = \ln \left(\frac{1}{\text{no. of documents containing } w} \right)$$

In our example, since we have two documents in the corpus, N=2.

Words	TF (For A)	TF (For B)	IDF
Jupiter	115	0	$\ln(2/1) = 0.69$
is	115	118	$\ln(2/2) = 0$
the	115	218	$\ln(2/2) = 0$
largest	115	0	$\ln(2/1) = 0.69$
From	0	118	$\ln(2/1) = 0.69$
Sun	0	118	$\ln(2/1) = 0.69$

Term Frequency - Inverse Document Frequency (TFIDF)

It is the product of TF and IDF.

TFIDF gives more weightage to the word that is more in the corpus (all the documents).

TFIDF provides more importance to the word that is more frequent in the document.

$$TFIDF(w, d, D) = TF(w, d) * IDF(w, D)$$

words	TF (For A)	TF (For B)	IDF	TFIDF(A)	TFIDF(B)
Jupiter	115	0	$\ln(2/1) = 0.69$	0.138	0
is	115	118	$\ln(2/2) = 0$	0	0
The	115	218	$\ln(2/1) = 0.69$	0.138	0
planet	115	0	$(\ln 2/1) = 0$	0.138	0
From	0	118	$(\ln 2/2) = 0$	0.138	0
sun	0	118	$(\ln 2/1) = 0.69$	0	0.086

After applying TFIDF, text in A and B documents can be represented as a TFIDF vector of dimension equal to the vocabulary words. The value corresponding to each word represents.

the importance of that word in particular document.

Disadvantage of TFIDF

It is unable to capture the semantics. For example Funny and humorous as Synonyms, but TFIDF does not capture that. TFIDF can be computationally expensive if the vocabulary is vast.

Bag of words (BOW)

Machine learning algorithms cannot work with raw Text directly. Rather, the text must be converted into vectors of numbers. In natural language processing, a common technique for extracting features from text is to place all of the words that occur in the text in a bucket.

Algorithm for Tokenization, pos Tagging, stop words removal, Stemming and Lemmatization.

Conclusion :-

In this way we have done text data analysis using TF-IDF algorithm.



Assignment No - 8.

Page No.	
Date	

Title - Data visualization.

1. use the inbuilt dataset 'titanic'. The datasets contains 891 rows and contains information about.

objective - students should able to do the data visualization for Titanic dataset using seaborn library.

Prerequisite - 1) Basics of python

2) Concepts of seaborn library.

~~contents For Theory.~~

1) Seaborn library Basics.

2) Know your data.

3) Finding patterns of data.

4) Checking how the price of ticket for each passenger is distributed by plotting histogram.

~~Theory - Data visualization plays a very important role in data mining. To accelerate this process we need to have a well document - ation of all plots.~~

1) Seaborn library Basics.

Seaborn is a python data visualization library based on matplotlib. It provides



Page No.	
Date	



high level interprice for drawing attractive and informative statistical graphics.

2) know your data

The dataset that we are going to are to draw our plots will be titanic dataset which is downloaded by default with the seaborn library.

let's see what the titanic dataset looks like.

Execute the following scripts.

```
import pandas as pd.  
import numpy as np.  
import matplotlib.pyplot as plt.  
import Seaborn as sns.
```

```
dataset = sns.load_dataset("titanic")
```

```
dataset.head()
```

3. Finding patterns of data.

patterns of data can be find out with the help of different types of plots

Types of plots are :-

A. distribution plots.

a. Dist-plot

b. Joint plot

c. Rug plot

A. Distribution plots :

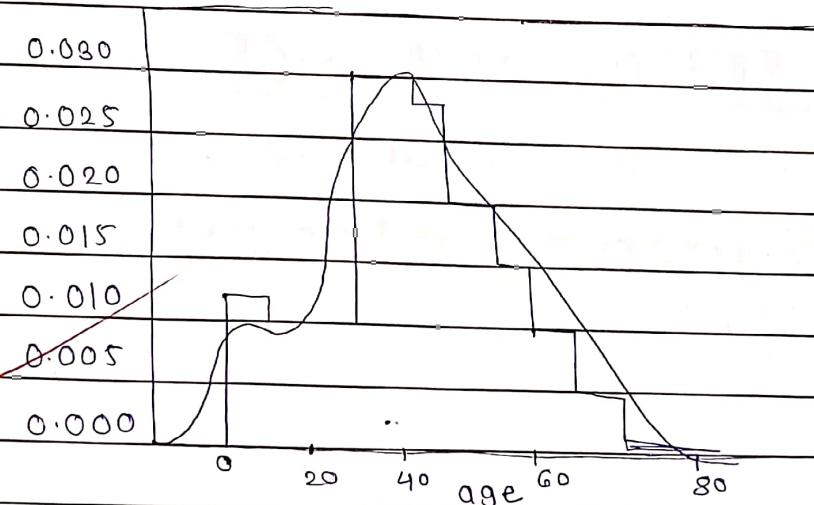
These plots help us to visualize the distribution of data. we can use these plots to understand the mean, median, range, variance, deviation, etc of the data.

a. Displot .

- Dist plot gives us the histogram of the selected continuous variable.
- It is an example of a univariate analysis.
- we can change the number of bins, i.e. no of vertical bars in a histogram.

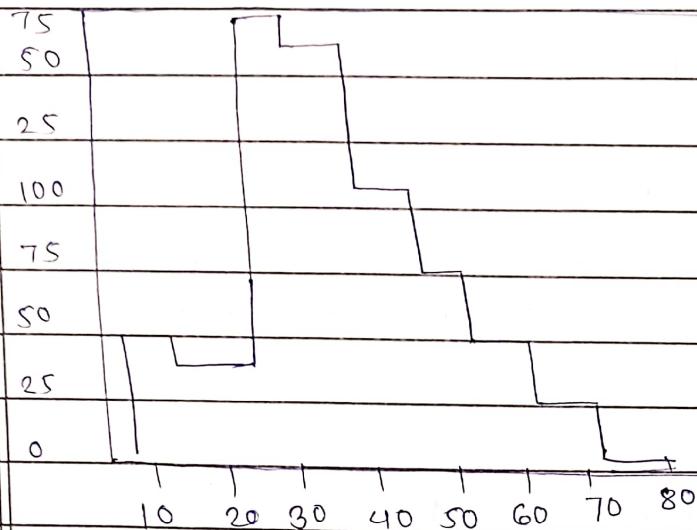
```
import seaborn as sns
```

```
sns.distplot(x = dataset['age'], bins = 10).
```



The line that you see represents the kernel density estimation, you can remove this line by passing False as the parameter for the kde Attribute as shown below.

```
sns.distplot(dataset['age'], bins=10, kde=False)
```



i.b. age.

Joint plot

- It is the combination of the distplot of two variables.
- It is an example of bivariate analysis.

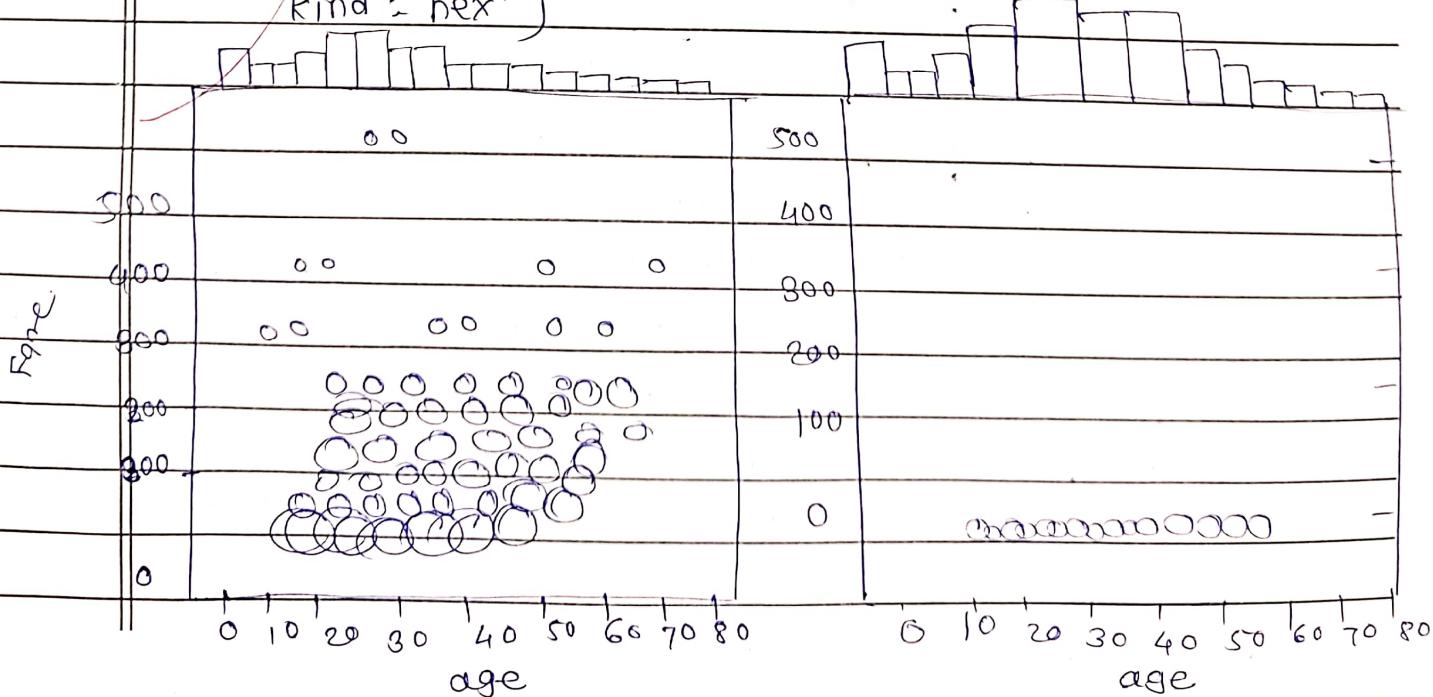
import seaborn as sns.

For plot 1

sns.jointplot(x= dataset['age'], y= dataset['Fare'], kind = 'scatter')

For plot 2.

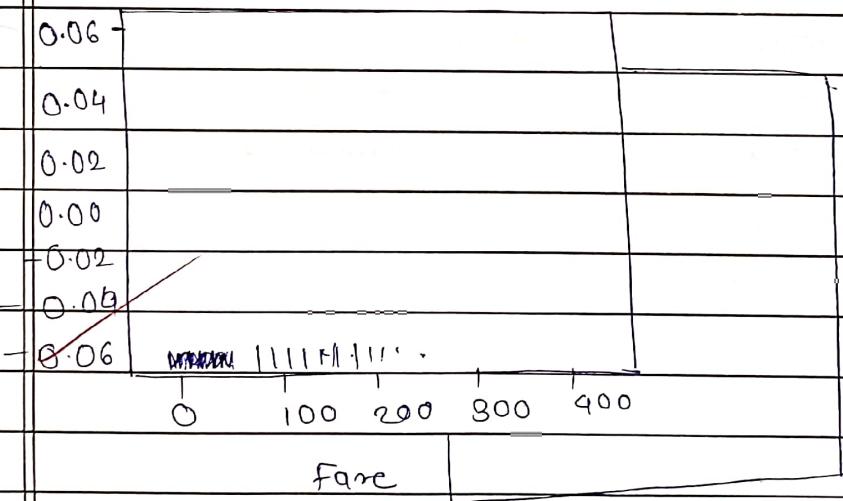
sns.jointplot(x= dataset['age'], y= dataset['Fare'], kind = 'hex')



- from the output, you can see that a joint plot has three parts. A distribution plot at the Top for the column on the x-axis, a distribution plot on the right for the column on the y-axis & a Scatter plot in betn that shows the mutual distribution of data.

~~e.~~ The Rug plot

b. The rugplot () is used to draw small bars along the x-axis for each points in the dataset. To plot a rug plot, you need to pass the most of the p name of the column . Let's plot a rug plot for fare.



From the output, you can see that most of the instances for the fares have values betn 0 and 100.



Conclusion -

Seaborn is an advanced data visualization library built on top of matplotlib library. In this assignment, we looked at how we can draw distributed & categorical plots using the Seaborn library.

TS	PR	VC	VA	RN	Total	Sign
(2)	(2)	(2)	(2)	(2)	Marks (10)	
✓						
✓	✓	✓	✓	✓	✓	L
✓	✓	✓	✓	✓	✓	✓



Assignment No-9.

Title - Data visualization.

- 1> Use the inbuilt dataset 'titanic' The dataset used in the problem plot box plot for distribution of age with respect to each gender
- 2> write observations on the inference from the above statistics.

~~Objectives~~ - students should able to do the data visualization for titanic dataset using Seaborn library.

~~Pre-requisite~~ - 1> Basic of python
2> Concept of seaborn library.

Contents of Theory

- 1> know your data.
- 2> Finding pattern of Data
- 3> checking how price of Ticket.

Theory -

Categorical plots;

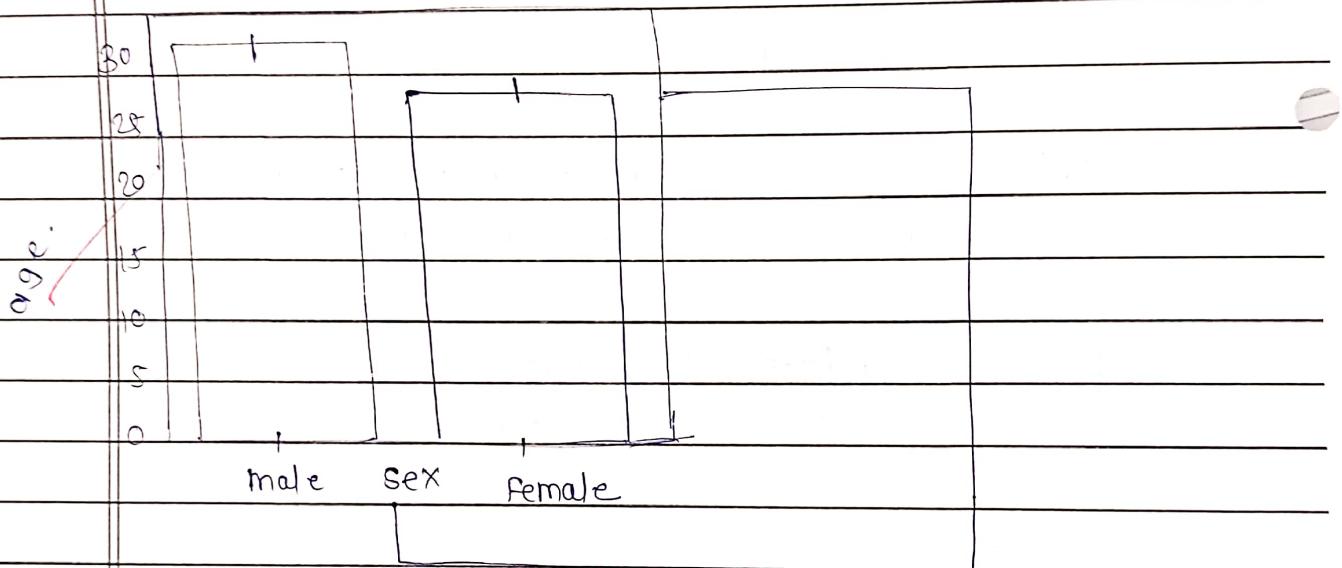
categorical plots, as the name suggests, are normally used to plot categorical data. The categorical plots the values in the categorical columns.

b. The Bar plot.

The barplot () is used to display the mean value in a categorical column, against a numeric column against & while the Third parameter is the dataset.

For instance, if you want to know the mean value the age of the male and female passenger. you can use the bar plot as follows.

`sns.barplot(x='sex', y='age', data=dataset)`



From the output, you can clearly see that the average age of male passengers is just less than 40 while the average age of female passengers is around 33.

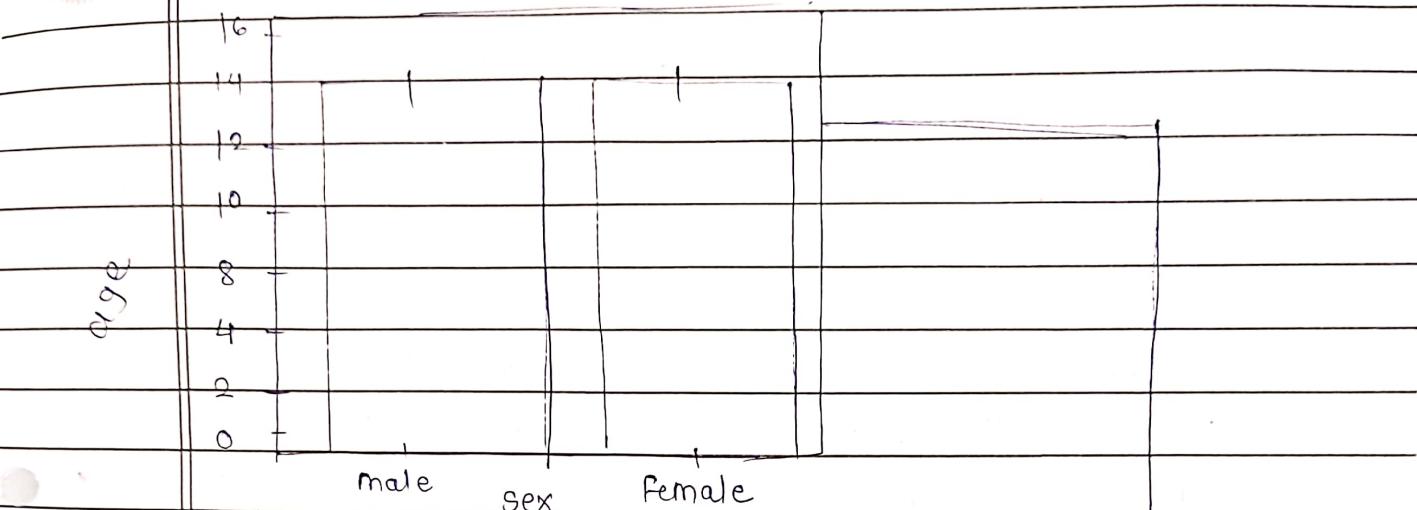
```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns.
```

```
Sns.barplot(x='sex', y='age', data=dataset, estimator=np.std)
```

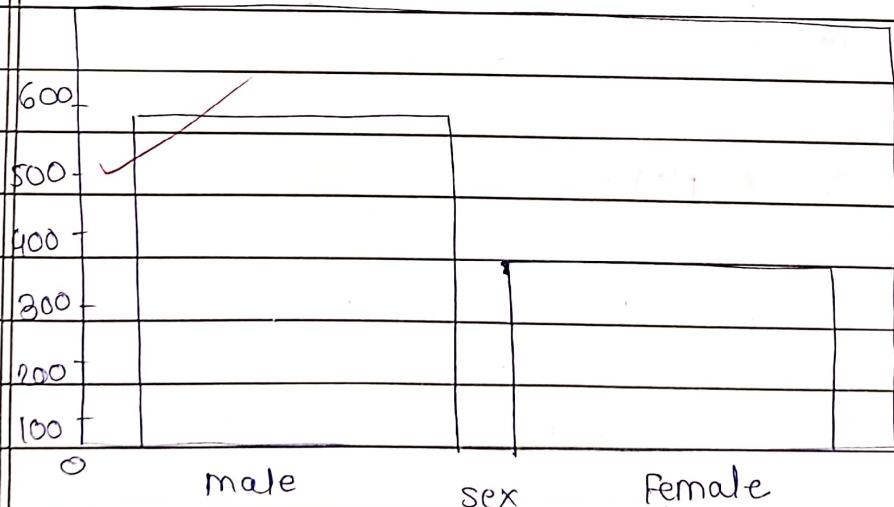
c.



c. The count plot .

The count plot is similar to the bar plot, however it displays the count of the categories in the specific column.

`sns.countplot(x='sex', data = dataset)`



d. The Box plot .

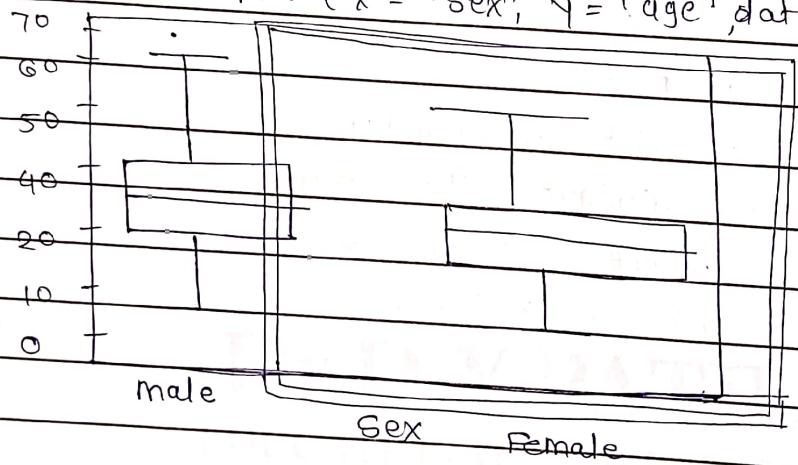
The Box plot is used to display the distribution of the categorical data in the

form of quartiles. The centre of the box shows the median value.

The value from the lower whisker to the bottom of the box shows the first quartile.

Now let's see the box plot that displays, the distribution for the age with respect to each gender.

Sns.boxplot(x = 'sex', y = 'age', data = dataset)

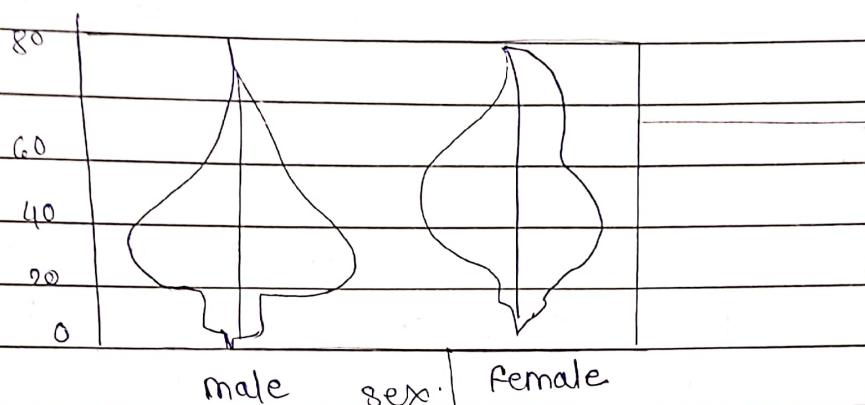


e. The violin plot.

The violin plot is similar to the box plot however, the violin plot allows us to display all the components that actually correspond to the data point.

Let's plot in violin plot.

Sns.violinplot(x = 'sex', y = 'age', data = dataset)



You can see from the figure below the violin plots provide much more information about the data as compared to the box plot.

~~Instead of plotting the quartile, the violin plot allows us to see all the components the actually correspond to the data.~~

b. The swarm plot.

The swarm plot is a combination of strip and violin plot.

In the swarm plot, the points are adjusted in such way that they dont overlaps.

c. A strip plot

~~A strip plot draws a scatter plot where one of the variables is categorical. we have seen scatter plots is the point plot and the pair plot sections where we had two numeric variables.~~

The strip plot is different in a way that one of the variables is categorical in this case.

Sns. strip plot ($x = \text{'sex'}$, $y = \text{'age'}$, data = dataset,
jitter = False)

Conclusion -

Seaborn is an advanced data visualisation library built on top of Matplotlib library. In this assignment, we looked at how we can draw distributional and categorical plots using the seaborn library. We have seen how to plot matrix plots in Seaborn. We also saw how to change plot styles and use grid functions to manipulate subplots.

TS	PR	VC	VA	RN	Total	Sign
(0)	1 2)	(2)	(2)	(2)	marks (10)	
02	02	01	02	02	05	2



Page No.	
Date	

Assignment No-10.

Title - Data visualization.

- Q 1. Download the Iris Flower dataset or any other dataset into a Dataframe (e.g. <https://archive.ics.uci.edu/ml/datasets/Iris>)

Scan the dataset and give the inference as-

2. List down the features and their types.
3. Create a boxplot for each features in the dataset.

objective - students should able to do the data visualization for Titanic dataset using seaborn library.

Theory :-

Matrix plots are the type of plots that show data in the form of rows and columns.

Heat maps are the prime examples of matrix plots.

a. Heat Maps.

Heat maps are normally used to plot correlation betn numeric column in the form of matrix.

It is important to mention here that to draw matrix plots, you need to have meaningful information on rows as well as columns.

import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

import Seaborn as sns

dataset = sns.load_dataset('titanic')

dataset.head()

	Survived	Pclass	Sex	Age	SibSp	Parch	Embarked	Class	Who	adult-male	deck
0	0	3	male	22.0	1	0	S	Third	man	True	NAN
1	1	1	female	38.0	1	0	C	First	woman	False	C
2	1	3	female	26.0	0	0	S	Third	women	False	NAN
3	1	1	female	35.0	1	0	S	First	woman	False	C
4	0	3	male	35.0	0	0	S	Third	man	True	NAN

betw all the numeric columns of the dataset

Execute the following script.

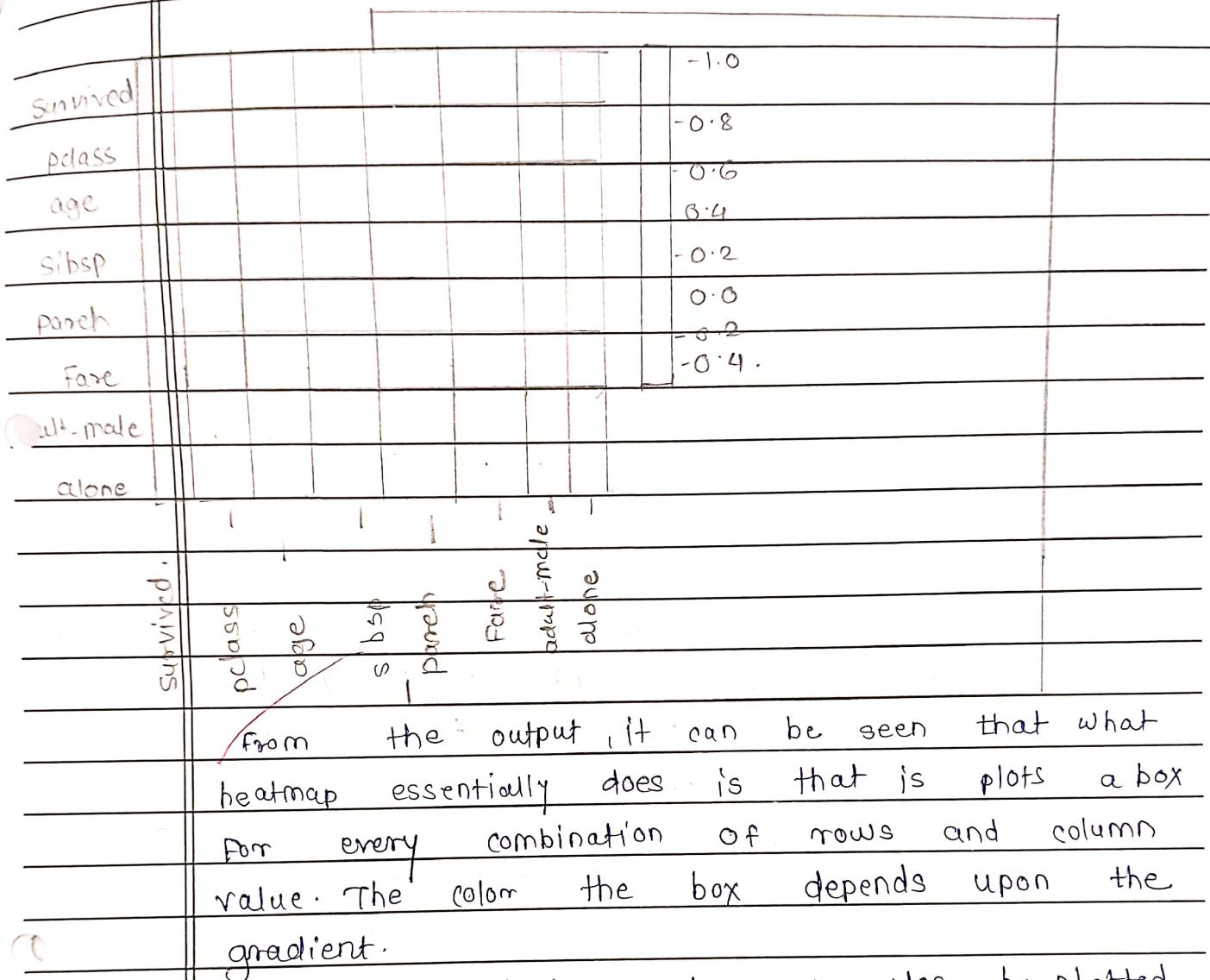
dataset.corr().

In the output, you will see that both the columns and the rows have meaningful header information as shown below.

	survived	Pclass	age	sibsp	parch	alone
survived	1.000000	-0.338481	-0.077221	-0.085322	0.081629	-0.203367
Pclass	-0.338481	1.000000	-0.369226	0.083081	0.018443	0.135207
age	-0.085322	-0.369226	1.000000	-0.308247	-0.189119	0.198270
sibsp	-0.081629	-0.083081	-0.308247	1.000000	-0.414838	0.189471
parch	0.257307	0.018443	-0.189119	0.414838	1.000000	-0.583388
alone	0.257307	-0.549800	0.96067	0.159657	0.216225	-0.271832
adult-male	0.557080	0.094085	0.285328	-0.349943	-0.349943	0.404744
alone.	-0.203367	0.135207	0.19270	-0.389398	-0.271822	1.000000

Now create a heat map with these correlation values, you need to call the heatmap() function and pass it your DataFrame. Look at the following script

corr = dataset.corr()



The correlation values can also be plotted on the heatmap by passing True for the annot parameter.

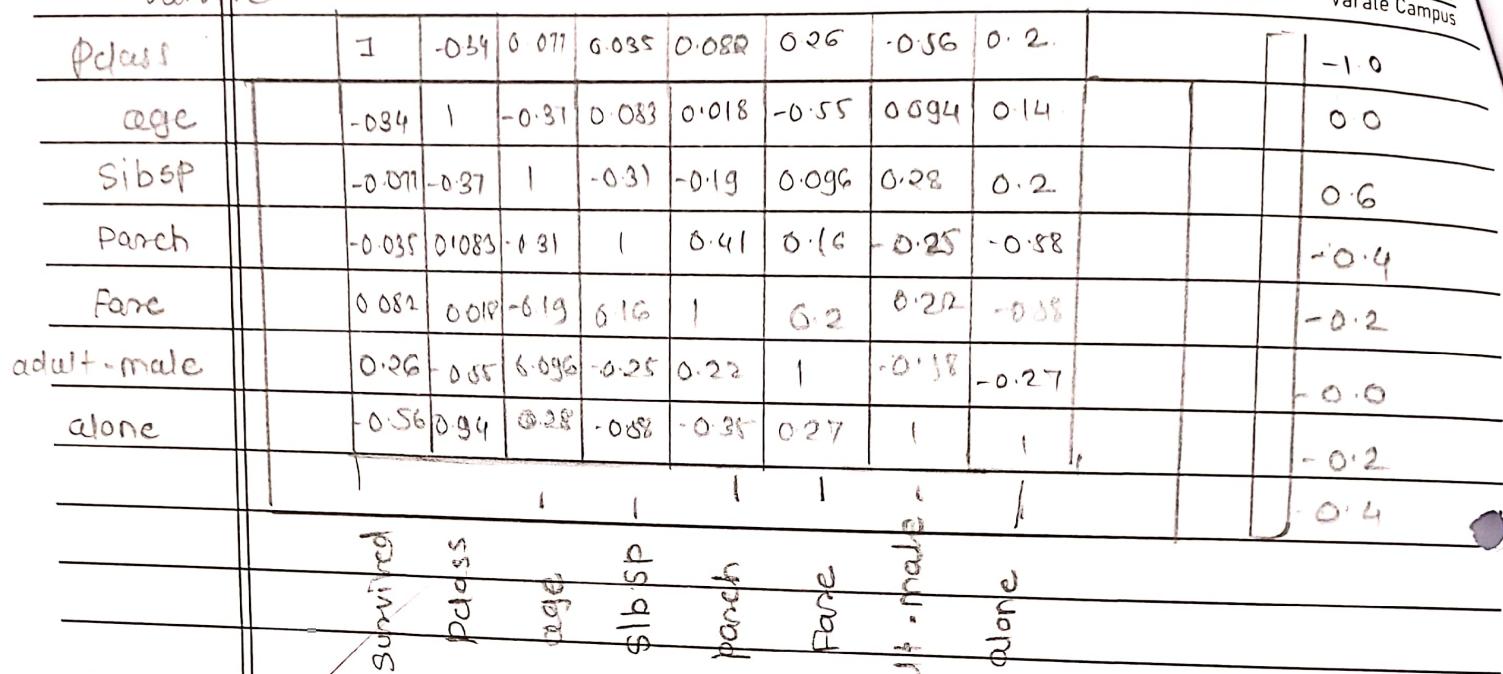
execute the following script to see this in action.

```
corr = dataset.corr()
sns.heatmap(corr, annot=True)
```



Survived

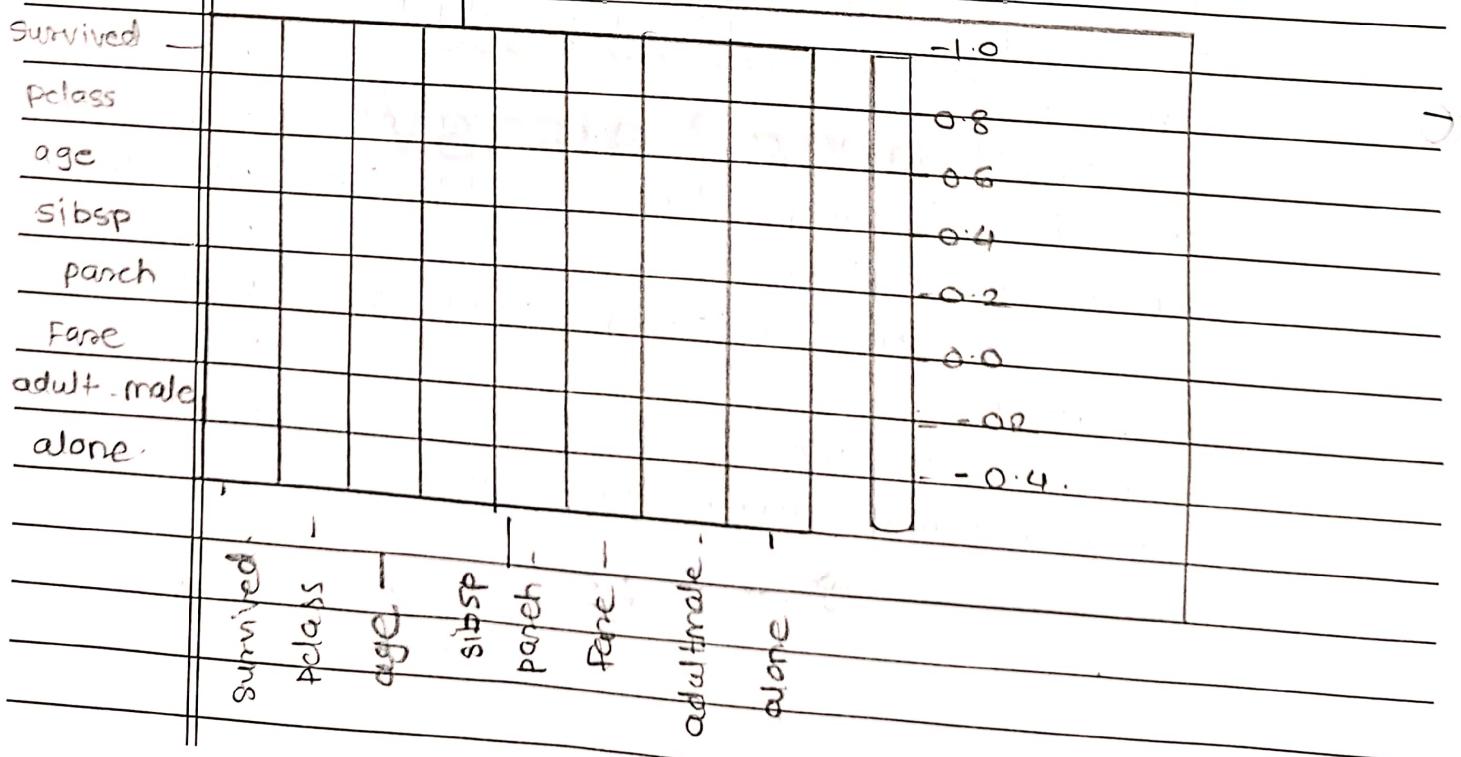
Page No.	
Date	



You can also change the color of the heatmap by passing an argument for the cmap parameter for now, just look at the following script.

```
corr = dataset.corr()
```

```
sns = heatmap(corr)
```



b. cluster map.

In addition to the heat map, another commonly used matrix plot is cluster map.

The cluster map basically uses Hierarchical clustering to cluster the rows and columns of the matrix.

b. clustermap

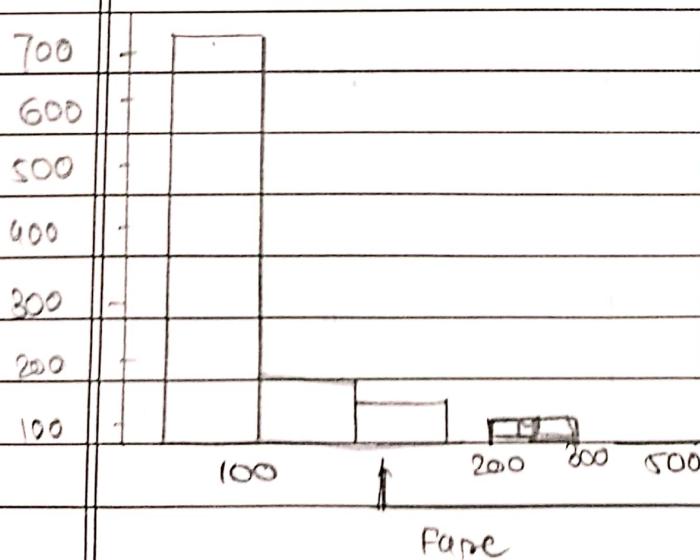
In addition to the heat map, another commonly used matrix plot is the cluster map. The cluster map basically uses Hierarchical clustering to cluster the rows and columns of the matrix.

4. checking how the price of the ticket (column name: 'fare') for each passenger is distributed by plotting a histogram.

import seaborn as sns.

dataset = sns.load_dataset('titanic')

sns.histplot(dataset['fare'], kde=False, bins=10)





from the histogram it is seen that for around 730 passengers the price of the ticket is 50. for 100 passengers the price of the ticket is 100 and so on.

Conclusion -

Seaborn is an advanced data visualisation library built on top of matplotlib library. In this assignment we looked at how we can draw distributional and categorical plots using the seaborn library.

TS	PR	VC	VA	RN	Total marks	Sign
✓	✓	✓	✓	✓	✓	P