

MALLA REDDY COLLEGE OF ENGINEERING & TECHNOLOGY (autonomous institution – ugc, govt. of india)



:S)

# (Emerging Technologies) DATA SCIENCE



# B.TECH III YEAR – I SEM (R20 Regulation) 2023-2024

# **STATISTICAL FOUNDATIONS LAB MANUAL**

# MALLA REDDY COLLEGE OF ENGINEERING & TECHNOLOGY

(Affiliated to JNTUH, Hyderabad, Approved by AICTE - Accredited by NBA & NAAC – 'A' Grade - ISO 9001:2015 Certified) Maisammaguda, Dhulapally (Post Via. Hakimpet), Secunderabad – 500100, Telangana State, India

# **Department of CSE-Emerging Technologies**

# Vision

"To be at the forefront of Emerging Technologies and to evolve as a Centre of Excellence in Research, Learning and Consultancy to foster the students into globally competent professionals useful to the Society."

# Mission

#### The department of CSE (Emerging Technologies) is committed to:

- To offer highest Professional and Academic Standards in terms of Personal growth and satisfaction.
- Make the society as the hub of emerging technologies and thereby capture opportunities in new age technologies.
- To create a benchmark in the areas of Research, Education and Public Outreach.
- To provide students a platform where independent learning and scientific study are encouraged with emphasis on latest engineering techniques.

# **QUALITY POLICY**

- To pursue continual improvement of teaching learning process of Undergraduate and Post Graduate programs in Engineering & Management vigorously.
- To provide state of art infrastructure and expertise to impart the quality education and research environment to students for a complete learning experiences.
- Developing students with a disciplined and integrated personality.
- To offer quality relevant and cost effective programmes to produce engineers as per requirements of the industry need.

For more information: www.mrcet.ac.in

# INDEX

S.No	List of Programs	PageNos.
1	<b>WEEK 1:</b> Study R Languages, Commands, etc Consider 50 observations (dataset), generating random data using functions provided, like rbinom, performing basic statistical computations using built- in functions of R. Discussion of R graphics. Histograms. Stem and leaf plots,Boxplots,Scatterplots. Bar graphs plotting the data using line graph, histograms, multiple graphs, etc. Generate 3D graphs or plots	3-6
2	<b>WEEK 2</b> : Measures of Central Tendency: Given a sample of 50 Observations (from any dataset), use possible functions R or Python and calculate mean, sd, var, min, max, median, range, and quantile. Discuss the properties of this distribution. generate bell curve of a random normal distribution.	7-8
3	<b>WEEK 3</b> : Pragmatic matters tabulating data, transforming a variable. Subsetting vectors and dataframes	9 -13
4	<b>WEEK 4</b> : (i) Consider 100 observations, find out Correlation "cor()" and Covariance "cov()" and programs on Frequencies and Crosstabs. (ii) Finding and analyze the missing data.	14-20
5	<b>WEEK 5</b> : Sorting, transposing and merging data. Reshaping a data frame. Basics of text processing.Reading unusual data files. Basics of variable coercion.	21-23
6	<b>WEEK 6</b> : Hypothesis testing and t-test for any given dataset. Find out null hypothesis, alternate hypothesis, draw the picture (graph) to visualize problem. Test the value of population mean.	24-25
7	<b>WEEK 7</b> : State alpha level and rejection region, estimate the maximum likely hood and inference.	26-28
8	<b>WEEK 8</b> : Binomial simulation: Making the computer flip coins for you. Make use of rbinom function of R to generate samples, and other functions: counts, avgs, mean, sd, sqrt, hist (histogram).	29-32
9	<b>WEEK 9</b> : Bayesian Hypothesis testing on any given dataset or dataframe.	33-36
10	<b>WEEK 10</b> : Use seaborn and combines simple statistical fits with plotting on pandas dataframes	37
11	WEEK 11: Working on Linear Algebra and Linear Systems	38
12	<b>WEEK 12</b> : Working on Monte Carlo Integration (Quasi-random numbers and find out thevariance on any dataframe)	39

003 1.5

# MALLA REDDY COLLEGE OF ENGINEERING AND TECHNOLOGY (R20A6703) Statistical Foundations of Data Science Lab B.Tech. III Year I Sem L T P C

# **COURSE OBJECTIVES:**

• The students are exposed to various experimental skills in data analytics which is very essential forData Science.

- Students are exposed to the Probability distribution using R & Python Programming.
- Students are able to use R and Python Programming and perform all types of operators andfunctions to generate the effective reports.
- To inculcate in students professional and ethical attitude, multidisciplinary approach and an ability to relate statistical analysis in data science by using various statistical methods or principles.
- Students should aware of Hypothetical Tests, Regression Analysis and Monte Carlo Integration.
- To provide student with an academic environment aware of excellence, written ethical codesand guidelines and lifelong learning needed for a successful bright and professional career.

**Required Software Tools:** R and Python (numpy, scipy, matplotlib)Also required: Pandas, Statsmodels and Seaborn.

# **COURSE OUTCOMES:**

- The student learns the concept of R and Python Programming and Statistical analysis and try formulate new solutions or programs.
- Demonstrate an ability to design and develop R and Python programs with this, analysis the data and generate the related report or results.
- Demonstrate an ability to design programming on probability distribution and computed all possible outcomes or required reports.
- Able to do hypothetical analysis and transformation of data into useful manner using Python Programming
- Able to generate the linear algebra, Monte Carlo Integrations etc by using Python programming.

WEEK 1: Study R Languages, Commands, etc

Consider 50 observations (dataset), generating random data using functions provided, like rbinom, performing basic statistical computations using built-in functions of R. Discussion of R graphics. Histograms. Stem and leaf plots. Boxplots. Scatterplots. Bar graphs plotting the data using line graph, histograms, multiple graphs, etc. Generate 3D graphs or plots.

**WEEK 2:** Measures of Central Tendency: Given a sample of 50 Observations (from any dataset), use possible functions R or Python and calculate mean, sd, var, min, max, median, range, and quantile. Discuss the properties of this distribution. generate bell curve of a random normal distribution.

**WEEK 3:** Pragmatic matters. Tabulating data. Transforming a variable. Subsetting vectors anddata frames.

**WEEK 4:** (i) Consider 100 observations, find out Correlation "cor()" and Covariance "cov()" and programs on Frequencies and Crosstabs. (ii) Finding and analyze the missing data.

**WEEK 5:** Sorting, transposing and merging data. Reshaping a data frame. Basics of text processing.Reading unusual data files. Basics of variable coercion.

**WEEK 6:** Hypothesis testing and t-test for any given dataset. Find out null hypothesis, alternate hypothesis, draw the picture (graph) to visualize problem. Test the value of population mean.

**WEEK 7:** State alpha level and rejection region, estimate the maximum likely hood and inference.

**WEEK 8:** Binomial simulation: Making the computer flip coins for you. Make use of rbinomfunction of R to generate samples, and other functions: counts, avgs, mean, sd, sqrt, hist (histogram).

WEEK 9: Bayesian Hypothesis testing on any given dataset or dataframe.

WEEK 10: Use seaborn and combines simple statistical fits with plotting on pandas dataframes.

WEEK 11: Working on Linear Algebra and Linear Systems

**WEEK 12**: Working on Monte Carlo Integration (Quasi-random numbers and find out thevariance on any dataframe)

WEEK 1:Study R Languages, Commands, etc

Consider 50 observations (dataset), generating random data using functions provided, like rbinom, performing basic statistical computations using built-in functions of R. mydata<- sample(1:nrow(Student\_Mark), 10)

> mydata

[1] 17 93 90 33 37 69 62 31 98 15

> Student\_Mark[mydata,]

#### # A tibble: 10 x 3 number\_courses time\_study Marks $\langle dbl \rangle$ *<dbl> <dbl>* 1 5 5.72 30.5 2 4 5.03 23.9 3 7 6.38 40.0 4 8 0.932 15.0 5 4 2.97 13.1 6 4 1.40 8.92 7 4 2.44 10.8 8 8 3.86 24.2 9 4 7.16 41.4 10 3 2.91 11.4

Discussion of R graphics. Histograms. Stem and leaf plots. Boxplots. Scatterplots. Bar graphs plotting the data using line graph, histograms, multiple graphs, etc. Generate 3D graphs or plots.

### Histogram

# **CSE(DATA SCIENCE)**

#### B.Tech-R20



v <- c(9,13,21,8,36,22,12,41,31,33,19)
save.image("C:/Users/MRECT/Music/1.png.RData")
save.image("C:/Users/MRECT/Music/1.png")
hist(v,xlab = "Weight",col = "yellow",border = "blue")
dev.off()</pre>

#### **Stem and Leaf Plots**

[ reached 'max' / getOption("max.print") -- omitted 328 rows ] stem(ChickWeight\$weight)

The decimal point is 1 digit(s) to the right of the |

2 | 599999999

- 32 | 12712
- 34 | 1
- 36 | 13

#### **Scatter plot**

# **CSE(DATA SCIENCE)**



```
input <- mtcars[,c('wt','mpg')]</pre>
```

print(head(input))

input <- mtcars[,c('wt','mpg')]</pre>

png(file = "scatterplot.png")

plot(x = input\$wt,y = input\$mpg,xlab = "Weight", ylab = "Milage", xlim = c(2.5,5),ylim = c(15,30),main = "Weight vs Milage")

#### dev.off()

#### Bar graphs using line graphs

**3D** Graphs or Plots

**WEEK 2:** Measures of Central Tendency: Given a sample of 50 Observations (from any dataset), use possible functions R or Python and calculate mean, sd, var, min, max, median, range, and quantile. Discuss the properties of this distribution. generate bell curve of a random normal distribution.

• 9	II 😅 🗉 🖯 🖯	🗌 🚍 🗌 🍌 Go	to file/function	Addins -				🔋 Project: (None
Stude	nt_Mark ×			-0	Environment History	Connections		- 0
	Filter			Q	💣 🔒 📑 Import D	ataset 🖌 🕑		🗏 List + 🗌 🤅
-	number_courses	time_study ÷	Marks ÷		📑 Global Environment	•		9,
1	3	4.508	19.202	^	Data			
2	4	0.096	7.734		🔘 mydata	100 obs. of 3 vari	ables	
3	4	3.133	13.811		Student_Mark	100 obs. of 3 vari	ables	
4	6	7.909	53.018					
5	8	7.811	55.299					
6	6	3.211	17.822					
7	3	6.063	29.889					
8	5	3.413	17.264		Files Plots Packag	es Help Viewer		- 5
9	4	4.410	20.348		O New Folder O D	elete 📑 Rename 🏻 🍪 More 🗸		
10	3	6.173	30.862		Home			
11	3	7.353	42.036		A Name		Size	Modified
12	7	0.423	12.132	¥	🗌 🔁 1-20RH1A626	0-CS.pdf	89.9 KB	Jul 9, 2021, 12:14 PM
howing	1 to 12 of 100 entries				🗌 🔁 1-20RH1A626	0.CS.pdf	89.9 KB	Jul 9, 2021, 12:23 PM
Console	Terminal ×				🗌 📩 1-60.CS.pdf		83 KB	Jul 9, 2021, 12:25 PM
~100				4	🗌 🔁 1-6260.CS.pdf		92 KB	Jul 9, 2021, 12:30 PM
mean	(mydata \$Marks	)	ener alta anesa		🗆 莒 110.png		42.5 KB	Aug 29, 2019, 4:51 PM
rror libr	in mydata\$Mark arv(readxl)	s : \$ operat	tor 1s 1nva	lid for atomic vectors	🗆 📕 111.jpg		24.3 KB	Aug 29, 2019, 4:51 PM
Stud	ent_Mark <- re	ad_excel("C	/Users/MRE	CW/Music/Datasets/archive (3)/Student_Mark.xls"	🗌 📕 112.jpg		10.7 KB	Aug 29, 2019, 4:52 PM
view	(Student Mark)				🗆 📕 113.jpg		110.8 KB	Aug 29, 2019, 4:52 PM
• myda	ta<-read_excel	("C:/Users/M	ARECW/Music	/Datasets/archive (3)/Student_Mark.xls")	🗆 🗮 114.jpg		30.9 KB	Aug 29, 2019, 4:52 PM
1] 24	(mydata\$Marks) .41769				🗆 📕 115.jpg		39.2 KB	Aug 29, 2019, 4:53 PM
					🗆 📕 116.jpg		5.1 KB	Aug 29, 2019, 4:53 PM
					🗌 📕 117.jpg		Activa2107 KB/in	id (Aug 29, 2019, 4:55 PM
					🗆 📕 118.jpg		Go to S87 KBigs to	) a Aug 29; 2019; 4:55 PM
					2-6260 CS ndf		02.2 VD	1.1 0 2021 12:24 DM

#### 📵 RStudio

ð ×

Stude	nt_Mark ×					Environment History Connections		_
				٩		🞯 📊 📑 Import Dataset 🗸 🖌		≣ List •
	number_courses	time_study ÷	Marks ÷			🐴 Global Environment 👻		Q,
1	3	4.508	19.202		^	Data		
2	4	0.096	7.734			🕐 mydata 100 obs. of 3 vari	ables	
3	4	3.133	13.811			<pre>O Student_Mark 100 obs. of 3 vari</pre>	ables	
4	6	7.909	53.018					
5	8	7.811	55.299					
6	6	3.211	17.822					
7	3	6.063	29.889					
8	5	3.413	17.264			Files Plots Parkages Help Viewer		
9	4	4.410	20.348			Parama Anno Parama Anno Mora -		_
10	3	6.173	30.862					
11	3	7.353	42.036			A Name	Size	Modified
12	7	0.423	12.132		~	1-20RH1A6260-CS.pdf	89.9 KB	Jul 9, 2021, 12:14 PM
ing	to 12 of 100 entries					1-20RH1A6260.CS.pdf	89.9 KB	Jul 9, 2021, 12:23 PM
sole	Terminal ×					1-60.CS.pdf	83 KB	Jul 9, 2021, 12:25 PM
						1-6260.CS.pdf	92 KB	Jul 9, 2021, 12:30 PM
5.1	509				^	🗆 🛤 110.png	42.5 KB	Aug 29, 2019, 4:51 PM
ix(	nydata\$Marks)					🗆 🗮 111.jpg	24.3 KB	Aug 29, 2019, 4:51 PM
55 edi	.299 an(mydata\$Marks	.)				🗆 🛤 112.јрд	10.7 KB	Aug 29, 2019, 4:52 PM
20	0595					🗌 🔜 113.jpg	110.8 KB	Aug 29, 2019, 4:52 PM
49	(range(mydata\$M .69	larks))				🗌 🚨 114.jpg	30.9 KB	Aug 29, 2019, 4:52 PM
uan	tile(mydata)			and the second second second		🗆 💴 115.jpg	39.2 KB	Aug 29, 2019, 4:53 PM
or:	Must use a vec	tor in [ or() to s	, not an o ee a backt	bject of class matrix. race		🗆 본 116.jpg	5.1 KB	Aug 29, 2019, 4:53 PM
uan	ile(mydata\$Mar	ks)	7.50	0.00		🗆 🧮 117.јрд	Activa21:7 KB/in	d Aug 29, 2019, 4:55 PM
	00 12.63300 20.	50% 05950 36.6	7625 55.29	9900		🗆 🗮 118.jpg	Go to S87.KBigs to	a Aug 29; 2019; 4:55 PM
609								

Malla Reddy College of Engineering and Technology

Student\_Mark <- read\_excel("C:/Users/MRECW/Music/Datasets/archive (3)/Student\_Mark.xls")</pre> > View(Student\_Mark) > mydata<-read\_excel("C:/Users/MRECW/Music/Datasets/archive (3)/Student\_Mark.xls") > mean(mydata\$Marks) [1] 24.41769 > sd(mydata\$Marks) [1] 14.3262 > var(mydata\$Marks) [1] 205.24 > min(mydata\$Marks) [1] 5.609 > max(mydata\$Marks) [1] 55.299 > median(mydata\$Marks) [1] 20.0595 > diff(range(mydata\$Marks)) [1] 49.69 > quantile(mydata\$Marks) 75% 25% 50% 100% 0% 5.60900 12.63300 20.05950 36.67625 55.29900

#### <u>WEEK 3 :</u>

#### Tabulating the Data using Python

Visualizing the data in tabular form is easier than visualizing it in a paragraph or commaseparated form. Nicely <u>formatted tables</u> not only provide you with a better way of looking at tables it can also help in <u>understanding each data</u> point clearly with its heading and value.

Tabulate is an open-source <u>python package/module</u> which is used to print tabular data in nicely formatted tables. It is easy to use and contains a variety of formatting functions. It has the following functionalities:

- One function call for all types of formatting
- Can be downloaded in multiple output formats
- Provides a better presentation with text and data.

In this week programs, we will see what are the different types of table formatting we can perform using Tabulate.

#### **Implementation:**

We will start by installing tabulate using **pip install tabulate**.

#### 1. Importing required libraries

We will be using the tabulate function from the tabulate library so we need to import that. Other than this we do not require to import any <u>python module</u>.

from tabulate import tabulate

#### 2. Creating Formatted Tables

Now we will start by creating different types of formatted tables. data = [["Himanshu",1123, 10025], ["Rohit",1126,10029], ["Sha",111178,7355.4]] print(tabulate(data))

Here we can see a plain table which is nicely formatted. Now let's see how we can add a header to this table that we just created.

print(tabulate(data, headers=["Name","User ID", "Roll. No."]))

In order to define the header along with the data, we can set that header='firstrow', let us see it through an example.

data = [['Name','ID'],["Himanshu",1123], ["Rohit",1126], ["Sha",111178]] print(tabulate(data, headers='firstrow'))

We can also display the indices of the rows by using the show index parameter. data = [['Name','ID'],["Himanshu",1123], ["Rohit",1126], ["Sha",111178]] print(tabulate(data, headers='firstrow', showindex='always'))

# tabular data types supported by tabulate

The tabulate function can transform any of the following into an easy to read plain-text table:

- *list of lists or another iterable of iterables*
- list or another iterable of dicts (keys as columns)

- *dict of iterables (keys as columns)*
- two-dimensional NumPy array
- NumPy record arrays (names as columns)
- pandas.DataFrame

# I) list of lists

For example, if we have the following list of lists:

table = [['First Name', 'Last Name', 'Age'], ['John', 'Smith', 39], ['Mary', 'Jane', 25], ['Jennifer', 'Doe', 28]]

We can turn it into into a much more readable plain-text table using the *tabulate* function:

print(tabulate(table))

```
print(tabulate(table))
```

First Name	Last Name	Age
John	Smith	39
Mary	Jane	25
Jennifer	Doe	28

Since the first list in the list of lists contains the names of columns as its elements, we can set it as the column or header names by passing *'firstrow'* as the argument for the *headers* parameter:

print(tabulate(table, headers='firstrow'))

print(tabulat	e(table, head	<pre>lers='firstrow'))</pre>
First Name	Last Name	Age
John	Smith	39
Mary	Jane	25
Jennifer	Doe	28

The *tabulate* function also contains a *tablefmt* parameter, which allows us to improve the appearance of our table using pseudo-graphics:

print(tabulate(table, headers='firstrow', tablefmt='grid'))

print(tabulate(table, headers='firstrow', tablefmt='grid'))

<b>_</b>	L	L
First Name	Last Name	Age
John	Smith	39
Mary	Jane	25
Jennifer	Doe	28
T	T	<b>T</b>

Use the '*fancy\_grid*' argument for *tablefmt*: print(tabulate(table, headers='firstrow', tablefmt='fancy\_grid'))

print(tabulate(table, headers='firstrow', tablefmt='fancy\_grid'))

First Name	Last Name	Age
John	Smith	39
Mary	Jane	25
Jennifer	Doe	28

# II) dictionary of iterables

We can create the same table above using a dictionary:

**info** = {'First Name': ['John', 'Mary', 'Jennifer'], 'Last Name': ['Smith', 'Jane', 'Doe'], 'Age': [39, 25, 28]}

In the case of a dictionary, the *keys* will be the *column headers*, and the **values** will be the **elements of those columns**. We specify that the keys will be the headers by passing '*keys*' as the argument for the *headers* parameter:

print(tabulate(info, headers='keys'))

print(	tabulate	(1nto,	headers='keys'))

First Name	Last Name	Age
John	Smith	39
Mary	Jane	25
Jennifer	Doe	28

And of course we can use the *tablefmt* parameter to improve the table's appearance:

```
print(tabulate(info, headers='keys', tablefmt='fancy_grid'))
```

```
print(tabulate(info, headers='keys', tablefmt='fancy_grid'))
```

First Name	Last Name	Age
John	Smith	39
Mary	Jane	25
Jennifer	Doe	28

# **III**) adding an index

We can also add an **index** to our table with the *showindex* parameter:

```
print(tabulate(info, headers='keys', tablefmt='fancy_grid', showindex=True))
```

	First Name	Last Name	Age
0	John	Smith	39
1	Mary	Jane	25
2	Jennifer	Doe	28

We can add a custom index by passing in an *iterable* to the *showindex* parameter. For example, if we want the index to start at 1, we can pass in a *range object* as the argument:

print(tabulate(info, headers='keys', tablefmt='fancy\_grid', showindex=range(1,4)))

	First Name	Last Name	Age
1	John	Smith	39
2	Mary	Jane	25
3	Jennifer	Doe	28

# IV) missing values

If we remove 'Jennifer' from the above *info* dictionary, our table will contain an empty field:

print(tabulate({'First Name': ['John', 'Mary'], 'Last Name': ['Smith', 'Jane', 'Doe'], 'Age': [39, 25, 28]}, headers="keys", tablefmt='fancy\_grid'))

First Name	Last Name	Age
John	Smith	39
Mary	Jane	25
	Doe	28

If there any **missing values** in our table, we can choose what to fill them in with using the *missingval* parameter. The default value for *missingval* is an empty string. If we change it to  $\frac{N}{A}$ , this is what our table will look like:

print(tabulate({'First Name': ['John', 'Mary'], 'Last Name': ['Smith', 'Jane', 'Doe'], 'Age': [39, 25, 28]}, headers="keys", tablefmt='fancy\_grid', missingval='N/A'))

First Name	Last Name	Age
John	Smith	39
Mary	Jane	25
N/A	Doe	28

For more detailed info about tabulate, visit: <u>https://pypi.org/project/tabulate/</u>

**WEEK 4:** (i) Consider 100 observations, find out Correlation "cor()" and Covariance "cov()" and programs on Frequencies and Crosstabs.

Let  $\Sigma(X)$  and  $\Sigma(Y)$  be the expected values of the variables, the covariance formula can be represented as:

Covariance(x, y) = 
$$\frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})$$

Where,

- xi = data value of x
- yi = data value of y
- $\bar{\mathbf{x}} = \text{mean of } \mathbf{x}$
- $\bar{y} = \text{mean of } y$

N = number of data values

## What Is Correlation?

In statistics, correlation is a measure that determines the degree to which two or more random variables move in sequence. When an equivalent movement of another variable reciprocates the movement of one variable in some way or another during the study of two variables, the variables are said to be correlated. The formula for correlation is:

$$\rho_{xy} = Correlation (x, y) = \frac{cov(x, y)}{\sqrt{var(x)}\sqrt{var(y)}}.$$

where,

var(X) = standard deviation of X

var(Y) = standard deviation of Y

Positive correlation occurs when two variables move in the same direction. When variables move in the opposite direction, they are said to be negatively correlated

Let's we will be working on the well-known Iris dataset.

(Download IRIS Dataset from https://www.kaggle.com/datasets/uciml/iris)

Consider: setosa species field or attribute

>>>import pandas as pd

>>>df=pd.read\_csv('E:/Kamal Files & Notes/Statastical Foundation Lab/Iris.csv')

#### **Output Screen:**



### **To Find Covariance:**

```
def covariance(x, y):
    # Finding the mean of the series x and y
    mean_x = sum(x)/float(len(x))
    mean_y = sum(y)/float(len(y))
    # Subtracting mean from the individual elements
    sub_x = [i - mean_x for i in x]
    sub_y = [i - mean_y for i in y]
    numerator = sum([sub_x[i]*sub_y[i] for i in range(len(sub_x))])
    denominator = len(x)-1
    cov = numerator/denominator
    return cov
```

```
with open('<--file path--->', 'r') as f:
    cov_func = covariance(sep_length, sep_width)
```

print("Covariance from the custom function:", cov\_func)

```
In [13]: def covariance(x, y):
    # Finding the mean of the series x and y
    mean_x = sum(x)/float(len(x))
    mean_y = sum(y)/float(len(y))
    # Subtracting mean from the individual elements
    sub_x = [i - mean_x for i in x]
    sub_y = [i - mean_y for i in y]
    numerator = sum([sub_x[i]*sub_y[i] for i in range(len(sub_x))])
    denominator = len(x)-1
    cov = numerator/denominator
    return cov
with open('E:/Kamal Files & Notes/Statastical Foundation Lab/Iris.csv', 'r') as f:
    cov_func = covariance(sep_length, sep_width)
    print("Covariance from the custom function:", cov_func)
```

Covariance from the custom function: 25.782885906040274

# **To Find Correlation:**

```
def correlation(x, y):
  # Finding the mean of the series x and y
  mean_x = sum(x)/float(len(x))
  mean_y = sum(y)/float(len(y))
  # Subtracting mean from the individual elements
  sub_x = [i-mean_x \text{ for } i \text{ in } x]
  sub_y = [i-mean_y for i in y]
  # covariance for x and y
  numerator = sum([sub_x[i]*sub_y[i] for i in range(len(sub_x))])
  # Standard Deviation of x and y
  std_deviation_x = sum([sub_x[i]**2.0 for i in range(len(sub_x))])
  std_deviation_y = sum([sub_y[i]**2.0 for i in range(len(sub_y))])
  # squaring by 0.5 to find the square root
  denominator = (std\_deviation\_x*std\_deviation\_y)**0.5 \# short but equivalent to
(std_deviation_x**0.5) * (std_deviation_y**0.5)
  cor = numerator/denominator
  return cor
with open(<<Dataset Path>>, 'r') as f:
  cor func
                     correlation(sep_length,
               =
                                                sep_width)
  print("Correlation from the custom function:", cor_func)
```

```
In [12]: def correlation(x, y):
             # Finding the mean of the series x and y
             mean x = sum(x)/float(len(x))
             mean_y = sum(y)/float(len(y))
             # Subtracting mean from the individual elements
             sub x = [i-mean x for i in x]
             sub_y = [i-mean_y for i in y]
             # covariance for x and y
             numerator = sum([sub x[i]*sub y[i] for i in range(len(sub x))])
             # Standard Deviation of x and y
             std_deviation_x = sum([sub_x[i]**2.0 for i in range(len(sub_x))])
             std_deviation_y = sum([sub_y[i]**2.0 for i in range(len(sub_y))])
             # squaring by 0.5 to find the square root
             denominator = (std deviation x*std deviation y)**0.5 # short but equivalent to (std deviation x**0.5) * (std deviation y**0.5)
             cor = numerator/denominator
             return cor
         with open('E:/Kamal Files & Notes/Statastical Foundation Lab/Iris.csv', 'r') as f:
             cor_func = correlation(sep_length, sep_width)
             print("Correlation from the custom function:", cor func)
         Correlation from the custom function: 0.7166762728539001
```

# II Analyse the Missing Data (Using Python)

# **Working with Missing Data in Pandas**

Pandas treat None and NaN as essentially interchangeable for indicating missing or null values. To facilitate this convention, there are several useful functions for detecting, removing, and replacing null values in Pandas DataFrame :

```
• isnull()
```

```
• notnull()
```

<u>Is null() and notnull()</u>

# importing pandas package import pandas as pd # making data frame from csv file data = pd.read\_csv(<<<dataset\_path>>>) # creating boolean series True for NaN values boolean\_series = pd.isnull(data["Col\_Name"]) #Col\_Name is the one where the blank values are there # filtering data # displaying data only with Gender = NaN data[boolean\_series]

### **CSE(DATA SCIENCE)**

```
In [20]: # importing pandas package
import pandas as pd

# making data frame from csv file
data = pd.read_csv('E:/Kamal Files & Notes/Statistical Foundation Lab/Dept_Table.csv')

# creating boolean series True for NaN values
boolean_series = pd.isnull(data["MANAGER_ID"])

# filtering data

# displaying data only with Gender = NaN
data[boolean_series]
```

Out[20]:

DEPARTMENT_ID DEPARTMENT_N	IAME MANAGER_ID	LOCATION_ID
----------------------------	-----------------	-------------

11	120	Treasury	NaN	1700
12	130	Corporate Tax	NaN	1700
13	140	Control And Credit	NaN	1700
14	150	Shareholder Services	NaN	1700
15	160	Benefits	NaN	1700
16	170	Manufacturing	NaN	1700
17	180	Construction	NaN	1700
18	190	Contracting	NaN	1700
19	200	Operations	NaN	1700
20	210	IT Support	NaN	1700
21	220	NOC	NaN	1700
22	220	IT Heindesk	MaM	1700

# NOT NULL()

# importing pandas package import pandas as pd # making data frame from csv file data = pd.read\_csv(<<<dataset\_path>>>) # creating boolean series True for NaN values boolean\_series = pd.notnull(data["Col\_Name"]) #Col\_Name is the one where the blank values are there # filtering data # displaying data only with Gender = NaN data[boolean\_series]

```
In [22]: # importing pandas package
```

import pandas as pd

# making data frame from csv file
data = pd.read\_csv('E:/Kamal Files & Notes/Statistical Foundation Lab/Dept\_Table.csv')

```
# creating boolean series True for NaN values
boolean_series = pd.notnull(data["MANAGER_ID"])
```

```
# filtering data
# displaying data only with Gender = NaN
data[boolean_series]
```

#### Out[22]:

	DEPARTMENT_ID	DEPARTMENT_NAME	MANAGER_ID	LOCATION_ID
0	10	Administration	200.0	1700
1	20	Marketing	201.0	1800
2	30	Purchasing	114.0	1700
3	40	Human Resources	203.0	2400
4	50	Shipping	121.0	1500
5	60	IT	103.0	1400
6	70	Public Relations	204.0	2700
7	80	Sales	145.0	2500
8	90	Executive	100.0	1700
9	100	Finance	108.0	1700
10	110	Accounting	205.0	1700

# To Count the Blank Cells or Missing values in a data\_frame

```
# To count how many blank cell are there in CSV file
data['Col Name'].isna().sum()
```

```
In [24]: # To count how many blank cell are there in CSV file
data['MANAGER_ID'].isna().sum()
Out[24]: 16
```

```
In [26]: # importing pandas package
         import pandas as pd
         # making data frame from csv file
         data = pd.read_csv('E:/Kamal Files & Notes/Statistical Foundation Lab/Dept_Table.csv')
         # creating boolean series True for NaN values
         boolean series = pd.isnull(data["MANAGER ID"])
         # filtering data
         # displaying data only with Gender = NaN
         data[boolean series]
         count_nan = data['MANAGER_ID'].isna().sum()
         print ('Count of NaN: ' + str(count_nan))
         Count of NaN: 16
```

COULD OF MAN, 10

**WEEK 5:** Sorting, transposing and merging data. Reshaping a data frame. Basics of text processing. Reading unusual data files. Basics of variable coercion.

# The various forms of reshaping data in a data frame are:

Transpose of a Matrix, Joining Rows and Columns, Merging of Data Frames

```
# R program to find the transpose of a matrix
```

```
first <- matrix(c(1:12), nrow=4, byrow=TRUE)
print("Original Matrix")
first
```

```
first <- t(first)
print("Transpose of the Matrix")
first
```

**Output:** 

[1] "Original Matrix"

- [,1] [,2] [,3]
- [1,] 1 2 3
- [2,] 4 5 6
- [3,] 7 8 9
- [4,] 10 11 12

[1] "Transpose of the Matrix"

- [,1] [,2] [,3] [,4]
- [1,] 1 4 7 10
- [2,] 2 5 8 11
- [3,] 3 6 9 12

# Joining Rows and Columns in Data Frame

In R, we can join two vectors or merge two data frames using functions. There are basically two functions that perform these tasks:

# cbind():

We can combine vectors, matrix or data frames by columns using **cbind**() function. Syntax: cbind(x1, x2, x3)where x1, x2 and x3 can be vectors or matrices or data frames.

# rbind():

We can combine vectors, matrix or data frames by rows using **rbind**() function. *Syntax: rbind*(*x1, x2, x3*) *where x1, x2 and x3 can be vectors or matrices or data frames.* 

# **CSE(DATA SCIENCE)**

# Cbind and Rbind function in R
name <- c("Shaoni", "esha", "soumitra", "soumi")
age <- c(24, 53, 62, 29)
address <- c("puducherry", "kolkata", "delhi", "bangalore")</pre>

# Cbind function info <- cbind(name, age, address) print("Combining vectors into data frame using cbind ") print(info)

# Rbind function
new.info <- rbind(info, newd)
print("Combining data frames using rbind ")
print(new.info)</pre>

#### **Output:**

[1] "Combining vectors into data frame using cbind "

name age address

- [1,] "Shaoni" "24" "puducherry"
- [2,] "esha" "53" "kolkata"
- [3,] "soumitra" "62" "delhi"
- [4,] "soumi" "29" "bangalore"
- [1] "Combining data frames using rbind "

name age address

- 1 Shaoni 24 puducherry
- 2 esha 53 kolkata
- 3 soumitra 62 delhi
- 4 soumi 29 bangalore
- 5 sounak 28 bangalore
- 6 bhabani 87 kolkata

### Merging two Data Frames

In R, we can merge two data frames using the **merge**() function provided both the data frames should have the same column names. We may merge the two data frames based on a key value.

# Merging two data frames in R

d1 <- data.frame(name=c("shaoni", "soumi", "arjun"),

ID=c("111", "112", "113"))

d2 <- data.frame(name=c("sounak", "esha"),

ID=c("114", "115"))

total <- merge(d1, d2, all=TRUE)

print(total)

## **Output:**

name ID

- 1 arjun 113
- 2 shaoni 111
- 3 soumi 112
- 4 esha 115
- 5 sounak 114

WEEK 6: Hypothesis testing and t-test for any given dataset. Find out null hypothesis, alternate hypothesis, draw the picture (graph) to visualize problem. Test the value of population mean.





Step 1: Collect Data
import pandas as pd
data = pd.read\_csv('diameter.csv')
Step 2: Define Null and Alternative Hypotheses
H0 = 'Data is normal'
Ha = 'Data is not normal'
Step 2: Set the level of significance (α) = 5%
alpha = 0.05
Step 3: Run a test to check the normality
I am using the Shapiro test to check the normality.
from scipy.stats import
shapirop =
round(shapiro(data)[1], 2)

Step 4: Conclude using the p-value from step 3

if p > alpha:

print(f"{p} > {alpha}. We fail to reject Null Hypothesis.

{H0}")else:

print(f"{p} <= {alpha}. We reject Null Hypothesis. {Ha}")</pre>

The above code outputs "0.52 > 0.05. We fail to reject Null Hypothesis.Data is Normal."

**WEEK 7**: State alpha level and rejection re-gion, estimate the maximum likely hood and inference.

import numpy as np import pandas as pd from matplotlib import pyplot as pltimport seaborn as sns from statsmodels import apifrom scipy import stats from scipy.optimize import minimize

Generate some synthetic data based on the assumption of Normal Distribution.

# generate an independent variablex = np.linspace(-10, 30, 100)
# generate a normally distributed residuale = np.random.normal(10, 5, 100)
# generate ground truthy = 10 + 4\*x + e
df = pd.DataFrame({'x':x, 'y':y})df.head()

```
%matplotlib inline
import matplotlib.pyplot as plt
plt.rcParams["figure.figsize"] = (11, 5) #set default figure size
import numpy as np
from numpy import exp
from scipy.special import factorial
import pandas as pd
from mpl_toolkits.mplot3d import Axes3D
import statsmodels.api as sm
from statsmodels.api import Poisson
from scipy import stats
from scipy.stats import norm
from statsmodels.iolib.summary2 import summary col
```

```
poisson pmf = lambda y, \mu: \mu**y / factorial(y)
* \exp(-\mu) y values = range(0, 25)
fig, ax = plt.subplots(figsize=(12, 8))
for \mu in [1, 5, 10]:
   distribution =
   [] for y i in
   y values:
       distribution.append(poisson pmf(y
   i, μ))ax.plot(y values,
          distribution,
          label=f' \ mu = {\mu}'
          , alpha=0.5,
          marker='o',
          markersize=8)
ax.grid()
ax.set xlabel('$y$', fontsize=14)
ax.set ylabel('$f(y \mid \mu)$',
fontsize=14) ax.axis(xmin=0, ymin=0)
ax.legend(fontsize=14)
plt.show()
```



**WEEK 8:** Binomial simulation: Making the computer flip coins for you. Make use of rbinom function of R to generate samples, and other functions: counts, avgs, mean, sd, sqrt,hist (histogram).

```
# r binomial - binomial simulation in r
rbinom(7, 150,.05)
[1] 10 12 10 2 5 5 14
```

We can model individual Bernoulli trials as well. We do this be setting the trials attribute to one. Here is the outcome of 10 coin flips:

```
# bernoulli distribution in r
rbinom(10, 1,.5)
[1] 1 0 1 1 1 0 0 0 0 1
```

Or stepping it up a bit, here's the outcome of 10 flips of 100 coins:



### Adding Mean and Median

mean <- mean(1) # Mean: 16.25</pre>

med <- median(1) # Meadian: 16.5</pre>

# Example 1: Add Mean to Histogram in R

abline(v = mean, col = 'blue')



#### Histogram of I

# Example 2: Add Median to Histogram in R

```
hist(l)
abline(v = med, col = 'red')
```

Histogram of I



# **Computing Average in R Programming**

```
# R program to get average of a list
# Taking a list of elements
list = c(2, 4, 4, 4, 5, 5, 7, 9)
# Calculating average using mean()
print(mean(list))
Output:
[1] 5
# R program to get average of a list
# Taking a list of elements
list = c(2, 40, 2, 502, 177, 7, 9)
# Calculating average using mean()
print(mean(list))
[1] 105.5714
```

# Computing Variance in R Programming

#R program to get variance of a list

# Taking a list of elements

list = c(2, 4, 4, 4, 5, 5, 7, 9)

# Calculating variance using var()

print(var(list))

# Output:

[1] 4.571429

#### **CSE(DATA SCIENCE)**

```
# R program to get variance of a list
# Taking a list of elements
list = c(212, 231, 234, 564, 235)
# Calculating variance using var()
```

[1] 22666.7

print(var(list))

# Standard Deviation in R

```
# R program to get
# standard deviation of a list
# Taking a list of elements
list = c(2, 4, 4, 4, 5, 5, 7, 9)
```

# Calculating standard
# deviation using sd()
print(sd(list))

#### OUTPUT

[1] 2.13809

```
# R program to get
# standard deviation of a list
# Taking a list of elements
list = c(290, 124, 127, 899)
# Calculating standard
# deviation using sd()
```

```
print(sd(list))
```

#### OUTPUT

[1] 367.6076

### WEEK 9: Bayesian Hypothesis testing on any given dataset or dataframe.

#Author DataFlair

t.test(x, mu = 5)	
RStudio	00
<u>File Edit Code View Plots Session Build Debug Profile Tools H</u> elp	
🝳 🔹 🖓 🖆 🖬 🔚 🚺 🛑 🚺 🏕 Go to file function 🔄 🛛 📅 🖌 Addins 🦂	🤱 Project: (None) 🗸
lattice.R* x     in hypothesis.R x	Environment History Connections
🕼 🖉 🔚 🖸 Source on Save I 🔍 🎽 🗐 So	ource 🗸 🧃 🚰 Import Dataset 🛛 🧹 🛛 🗮 List 🗸 🔞 🗸
1 #Author DataFlair	Global Environment 🖌 🔍
2 t.test(x, $mu = 5$ )	Values
	x num [1:10] -0.408 -2.13
	y num [1:10] 1.0347 1.653
2:18 (Top Level) :	R Script :
Console Terminal X Jobs X	-0
10 M	
<pre>&gt; source('~/hypothesis.R', echo=TRUE)</pre>	Files Plots Packages Help Viewer
	🍐 🧑 🖉 Zoom 🎝 Export 🗸 🚺 🧹
> #Author DataFlair	
<pre>&gt; t.test(x, mu = 5)</pre>	
2 2 3 3 7 3	
One Sample t-test	
ii	
t = -20.555, dt = 9, p-Value = 7.151e-09	
alternative hypothesis: true mean is not equal to 5	
95 percent confidence interval:	
-1.019153/ U.1/43669	
sample estimates:	
mean of x	
-0.4223834	

```
#Author DataFlair
```

t.test(y, mu = 5, alternative = 'greater')

# Defining sample vector

#### **CSE(DATA SCIENCE)**

```
x <- rnorm(100)
# One Sample T-Test
t.test(x, mu = 5)
Output:
    One Sample t-test
data: x
t = -49.504, df = 99, p-value < 2.2e-16
alternative hypothesis: true mean is not equal to 5
95 percent confidence interval:
 -0.1910645 0.2090349
sample estimates:
  mean of x
0.008985172
 Two Sample T-Testing
# Defining sample vector
x <- rnorm(100)
y <- rnorm(100)
# Two Sample T-Test
t.test(x, y)
Welch Two Sample t-test
data: x and y
t = -1.0601, df = 197.86, p-value = 0.2904
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -0.4362140 0.1311918
sample estimates:
  mean of x mean of y
-0.05075633 0.10175478
```

# 1. Example Data Frame.

#### Dummy example with 2 groups A/B composed of 20 observations each of Average order value per user.

df <- cbind(group, aov)
df <- as.data.frame(df)</pre>

# 2. Distribution and overlap.

```
ggplot(df, aes(x=aov, fill=group)) +
   stat_count(width = 0.5, alpha=.5, position="dodge")
```



Group B targeting is likely to secure higher AOV.

WEEK 10: Use seaborn and combines simple statistical fits with plotting on pandas dataframes.



**Solution: Juypter Notebook** 

# WEEK 11: Working on Linear Algebra and Linear Systems

## 11 a) Linear Algebra- Juypter Notebook

	er Week 11 a_Linear_algebra_with_Numpy Last Checkpoint: a few seconds ago (autosaved)		n Logou	É.
File Edit	View Insert Cell Kernel Widgets Help		Trusted Python 3	0
<b>B</b> + %	P3 B ↑ ↓ H Run ■ C → Markdown ∨			
	A * Ainv			
Out[1	<pre>i]: matrix([[ 1.00000000e+00, 2.77555756e-17, 3.05311332e-16],</pre>			
			Add tag	
	The above might not look like the identity matrix but if you look closer you see that the diagonals are all 1 and the (which from a computer's point of view is $\theta$ ).	e off diagonals ar	e a <b>very</b> small number	
			Add tag	
	To calculate the determinant:			
In [1	5]: nbval-ignore-output x		Add tag	
	np.linalg.det(A)			
Out[1	J: -128.99999999999999			
Out[1	J: -128.99999999999997		Add tag	

## 11b) Working on Linear Systems – Solution Jupyter Notebook

C Downloads/SDFS/ X C Downloads	/SDFS/ 🗙 🥃 Week 11 b Linear Systems - J 🗙	🧧 Week 11 a_Linear_algebra_w 🗙 📔 🦉 Week 12 Montecarlo - Jupyt	×   +	~ - 0 ×
← → C () localhost:8888/notebooks/Dc	wnloads/SDFS/Week%2011%20b%20Linear%20System	s.ipynb	Ê	A D Paused :
📰 Friday Fun: Directly 🚦 Setting up a Raspb 🧟	Smartphone Contro əs Python Modules ガ How	to take the sa		
💭 jupyter 🛛 Week 11 b	D Linear Systems Last Checkpoint: 27 minutes ago	(autosaved)	6	Logout
File Edit View Insert	Cell Kernel Widgets Help		Trusted	Python 3 O
E + ≫ £ E ↑ ↓	H Run C H Code V	ar rie dagurtres v2) m) v2.		
In [53]: #Actual Prog import nump from scipy.	gram for Week 11 b y as np linalg import solve			
In [54]: A = np.array : : : : : : : 	y( [1, 9, 2, 1, 1], [10, 1, 2, 1, 1], [1, 0, 5, 1, 1], [2, 1, 1, 2, 9], [2, 1, 2, 13, 2], ]			
In [56]: b = np.array	y([170, 180, 140, 180, 350]).reshape((5, 1))			
In [57]: x = solve(A	, b)			
In [58]: x				
Out[58]: array([[10.] [10.] [20.] [10.]	], , , , 1)			
O Type here to search	l D 🔒 🗋 👰 🤅			^ 4 <sub>8</sub> 10:24 □

**WEEK 12**: Working on Monte Carlo Integration (Quasi-random numbers and find out the variance on any data frame)

```
from scipy import stats
import numpy as np
import matplotlib.pyplot as plt
```

```
N = 10000

a, b = (50,50)

x_min, x_max = (0, .55)

randx = np.random.uniform(x_min, x_max, N)

y = stats.beta.pdf(randx, a, b)

print(f'Real value to find: {stats.beta(a,b).cdf(.55)}')

print(f'Integral value: {(x_max-x_min)*y.sum()/N}')

print(f'Calculation error: {np.sqrt((x_max-x_min)*(y*y).sum()/N - (x_max-x_min)*y.mean()**2)/np.sqrt(N)}')
```

Real value to find: 0.8413478010629016 Integral value: 0.8486374584821799 Calculation error: 0.019813380205560064

# Plotting
plt.figure(figsize=(8,5))
x = np.linspace(0, 1, 1000)
plt.plot(x, stats.beta.pdf(x, a, b))
# Then, let's only plot a thousand points for more readability
plt.scatter(randx[:1000], y[:1000], alpha=.08, label='value of f for random uniform x')
plt.xlabel('x')
plt.ylabel('density')
plt.legend()
plt.show()



The blue points correspond to the 10 000 values of  $f(x_i)$  computed from the uniform draws we made over X.

0