# Report on Exploratory Data Analysis (Report)

Mathematical Foundation of Computer Science

Masters in Computer Science (2022-2024)

Department of Computer Science, University of Delhi

**Submitted To:** Dr. Vasudha Bhatnagar

**Submitted By:** Arpit Kumar Mishra, Bhushan, Mukesh Kumar

# Exploratory Data Analysis (EDA)

Exploratory Data Analysis (EDA) is a data analysis approach that utilizes visual techniques to analyze data. EDA is primarily used to discover trends and patterns, and to check assumptions using statistical summary and graphical representations. By visually exploring the data, EDA enables a better understanding of the data's main features, the relationships between variables, and the identification of important variables for problem-solving. EDA often involves the use of descriptive statistics to gain quick insights into the characteristics of the data. Ultimately, EDA serves as a powerful tool for gaining a comprehensive understanding of data and identifying important trends and patterns.

It is a phenomenon under data analysis used for gaining a better understanding of data aspects like:
- main features of data
- variables and relationships that hold between them
- identifying which variables are important for our problem

**Descriptive statistics:** It refers to a set of techniques and measures used to summarize and describe the main characteristics of a dataset. These statistics provide a clear and concise summary of the data, including measures of central tendency.

**Name of the Data set:** Cause-wise Distribution of Suicides by Means/Mode Adopted in India during 2021

**Link to the data set:**

https://data.gov.in/resource/cause-wise-distribution-suicides-meansmode-adopted-india-during-2021

**Source:** Open Government Data(OGD) Platform India

**Released under:** National Data Sharing and Accessibility Policy(NDSAP)

**Contributor:** Ministry of Home Affairs , Department of States , National Crime Records Bureau (NCRB)

**Description:** The Cause-wise Distribution of Suicides by Means/Mode Adopted in India during 2021 dataset provides information on the number of suicides committed in India in 2021, broken down by means/mode of suicide. The data was collected by the National Crime Records Bureau (NCRB), a department of the Indian Ministry of Home Affairs, and released under the National Data Sharing and Accessibility Policy (NDSAP).
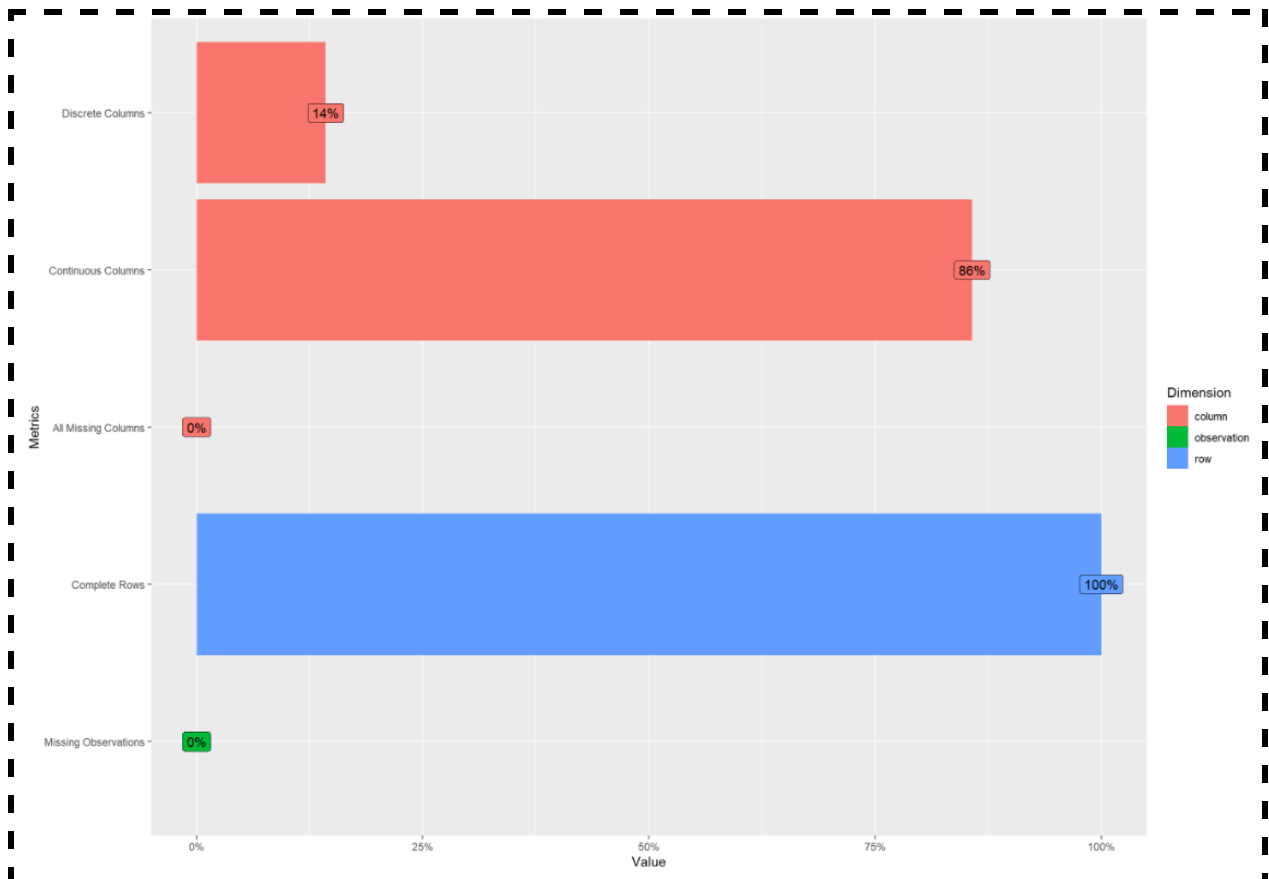
The dataset includes information on the number of suicides committed by various means/modes such as hanging, poisoning, drowning, and others.
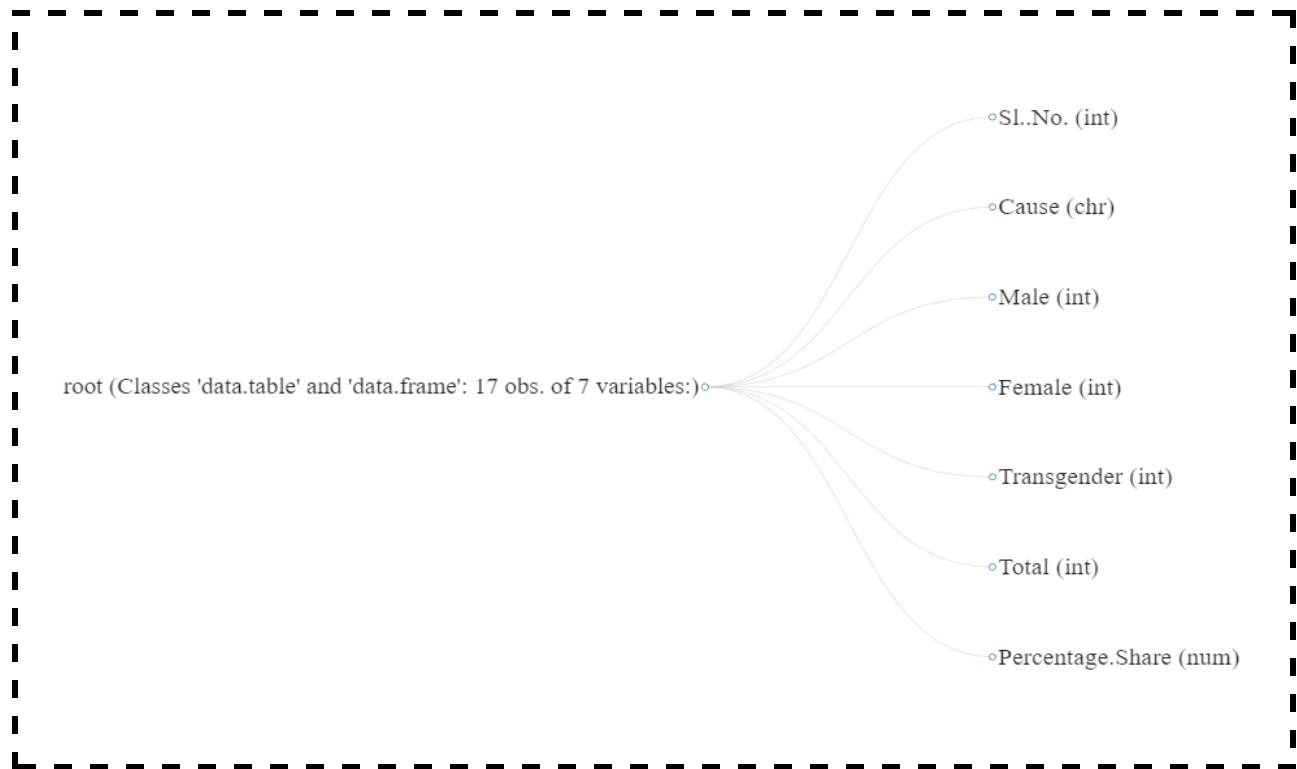
# Dataset Statistics(1/2)

**Code:**

```
 7  #Description of Dataset--> data1_cause
 8  data1_cause <- read.csv("C:/Users/kumar/Documents/Datasets/Cause-wise NCRB_ADSI-2021_Table_2.12.csv")
 9  describe(data1_cause)
10  head(data1_cause)
11  summary(data1_cause)
12
13  install.packages("skimr")
14  library(skimr)
15  skim(data1_cause)
16
17  install.packages("DataExplorer")
18  library(DataExplorer)
19  DataExplorer::create_report(data1_cause)
```

**Output:**

root (Classes 'data.table' and 'data.frame': 17 obs. of 7 variables:)

- Sl..No. (int)
- Cause (chr)
- Male (int)
- Female (int)
- Transgender (int)
- Total (int)
- Percentage.Share (num)

**Name of the Data set:** State/UTs-wise Educational Status Distribution of Suicides during 2021

**Link to the data set:**

https://data.gov.in/resource/stateuts-wise-educational-status-distribution-suicides-during-2021

**Source:** Open Government Data(OGD) Platform India

**Released under:** National Data Sharing and Accessibility Policy(NDSAP)

**Contributor:** Ministry of Home Affairs , Department of States , National Crime Records Bureau (NCRB)

**Description:** The State/UTs-wise Educational Status Distribution of Suicides during 2021 dataset provides information on the number of suicides committed in India in 2021, broken down by state/UT and educational level. The data was collected by the National Crime Records Bureau (NCRB), a department of the Indian Ministry of Home Affairs, and released under the National Data Sharing and Accessibility Policy (NDSAP).
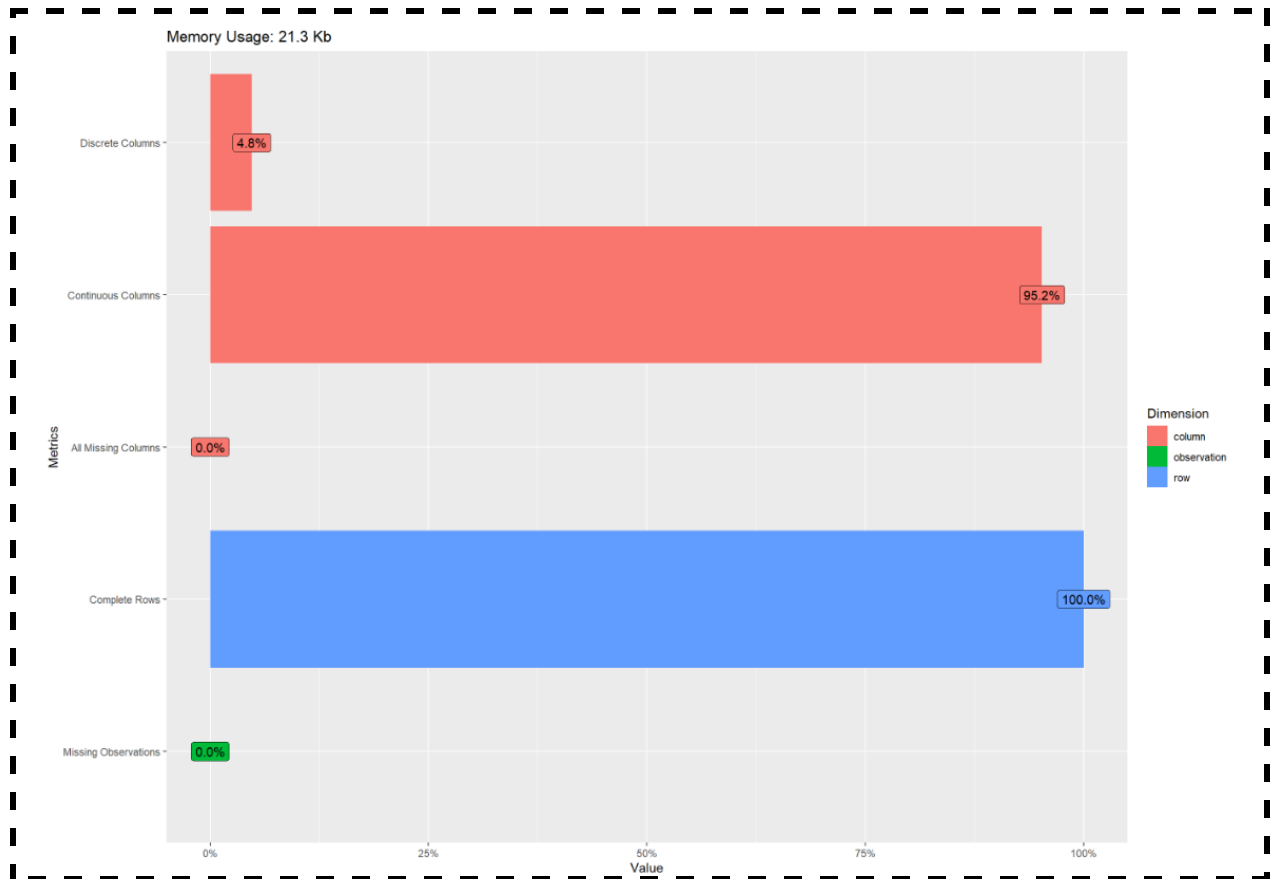
The dataset includes information on the number of suicides committed by individuals with different educational backgrounds, such as those with no education, primary education, secondary education, and others. The data is further classified by state/UT, and gender.

# Dataset Statistics(2/2)

**Code:**

```
22  #Description of Dataset--> data2_edu
23  data2_edu <- read.csv("C:/Users/kumar/Documents/Datasets/Education-wise NCRB_ADSI-2021_Table_2.11.csv")
24  describe(data2_edu)
25  head(data2_edu)
26  summary(data2_edu)
27  skim(data2_edu)
28  DataExplorer::create_report(data2_edu)
```

**Output:**

root (Classes 'data.table' and 'data.frame': 39 obs. of 42 variables:)

- Category (chr)
- State.UT (chr)
- No.Education...Male (int)
- No.Education...Female (int)
- No.Education...Transgender (int)
- No.Education...Total (int)
- Primary..up.to.class...5th....Male (int)
- Primary..up.to.class...5th....Female (int)
- Primary..up.to.class...5th....Transgender (int)
- Primary..up.to.class...5th....Total (int)
- Middle..up.to.class...8th....Male (int)
- Middle..up.to.class...8th....Female (int)
- Middle..up.to.class...8th....Transgender (int)
- Middle..up.to.class...8th....Total (int)
- Matriculate..Secondary..up.to.class...10th....Male (int)
- Matriculate..Secondary..up.to.class...10th....Female (int)
- Matriculate..Secondary..up.to.class...10th....Transgender (int)
- Matriculate..Secondary..up.to.class...10th....Total (int)
- Hr..Secondary..Intermediate..Pre.University..up.to.class...12th....Male (int)
- Hr..Secondary..Intermediate..Pre.University..up.to.class...12th....Female (int)
- Hr..Secondary..Intermediate..Pre.University..up.to.class...12th....Transgender (int)
- Hr..Secondary..Intermediate..Pre.University..up.to.class...12th....Total (int)
- Diploma.Certificate..ITI....Male (int)
- Diploma.Certificate..ITI....Female (int)
- Diploma.Certificate..ITI....Transgender (int)
- Diploma.Certificate..ITI....Total (int)
- Graduate.and.above...Male (int)
- Graduate.and.above...Female (int)
- Graduate.and.above...Transgender (int)
- Graduate.and.above...Total (int)
- Professionals..MBA.etc.....Male (int)
- Professionals..MBA.etc.....Female (int)
- Professionals..MBA.etc.....Transgender (int)
- Professionals..MBA.etc.....Total (int)
- Status.Not.known...Male (int)
- Status.Not.known...Female (int)
- Status.Not.known...Transgender (int)
- Status.Not.known...Total (int)
- Total...Male (int)
- Total...Female (int)
- Total...Transgender (int)
- Total...Total (int)

# Data Collection

## (Dataset 1/2)

**Code:**

```
1   #install.packages("tidyverse")
2   install.packages("gridExtra")
3   library(tidyverse)
4   library(forcats)
5   library(gridExtra)
6   theme_set(theme_bw()+
7     theme(panel.grid = element_blank()))
```

```
7     theme(panel.grid = element_blank()))
8   data1_cause <- read.csv("C:/Users/kumar/Documents/Datasets/Cause-wise NCRB_ADSI-2021_Table_2.12.csv")
9   View(data1_cause)
10  names(data1_cause)
```

**Output:**

```
> names(data1_cause)
[1] "Sl..No."        "Cause"          "Male"           "Female"
[5] "Transgender"    "Total"          "Percentage.Share"
```

**Inferences:**

- **No. of Attributes: 7**
- **No. of Records: 17**
- **No. of "discrete" attributes: 1**
- **No. of "continuous" attributes: 6**
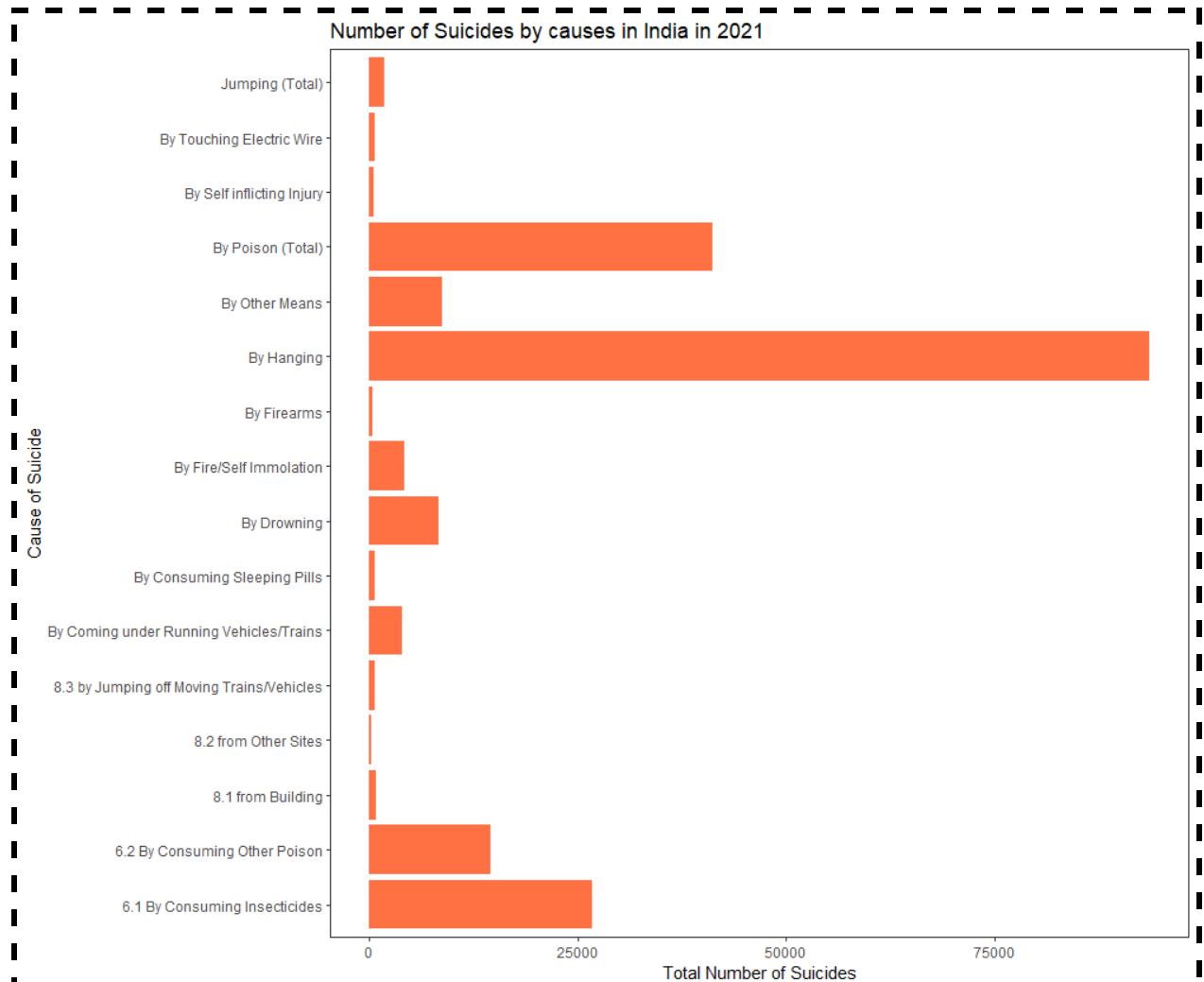- **Total No of observations: 119**

# Data Visualization

Data visualization is the graphical representation of information and data. It uses statistical graphics, plots, information graphics and other tools to communicate information clearly and efficiently.

**************************************************************************************************

**Code:**

```
12  #plotting bar graph to determine Number of Suicides by causes in India in 2021
13  plot1<-data1_cause %>%
14    group_by(Cause)%>%
15    summarize('Total_Number_of_Suicides' = Total) %>%
16    arrange(desc('Total_Number_of_Suicides'))
17  View(plot1)
18  names((plot1))
19  plot1[1:nrow(plot1)-1,] %>%
20    ggplot(aes(Cause,Total_Number_of_Suicides))+
21    geom_bar(stat = "identity", fill = "#ff7043")+
22    #geom_label(aes(label = Total_Number_of_Suicides),size = 1.5)+
23    ggtitle("Number of Suicides by causes in India in 2021")+
24    xlab("Cause of Suicide")+
25    ylab("Total Number of Suicides")+
26    coord_flip()
27    #theme_bw()
```
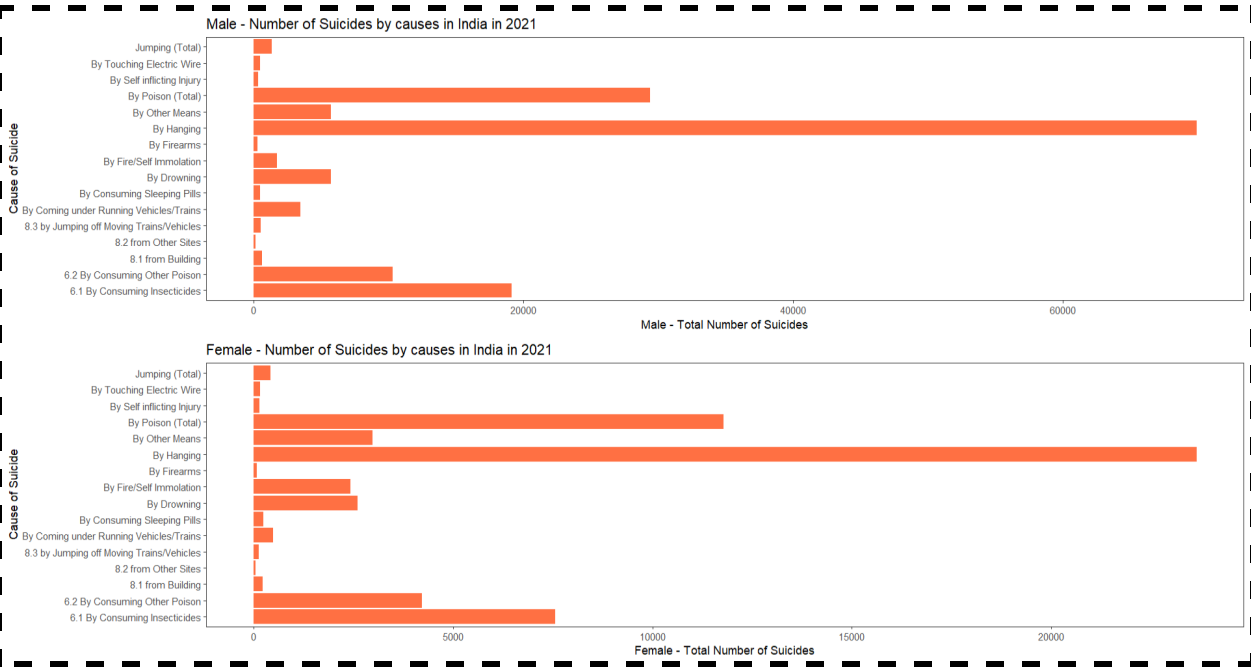
**Output:**

**Inferences:**

The above graph represents the number of suicides in India in 2021 categorised by the cause by which the suicide was comminted. By the bar graph we can make out that "Hanging to death" is the most practiced mode of suicide. And "jumping off Moving Trains/Vehicles", "suicide from building", "suicide from other sites" are the least practiced mode of suicide.

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*
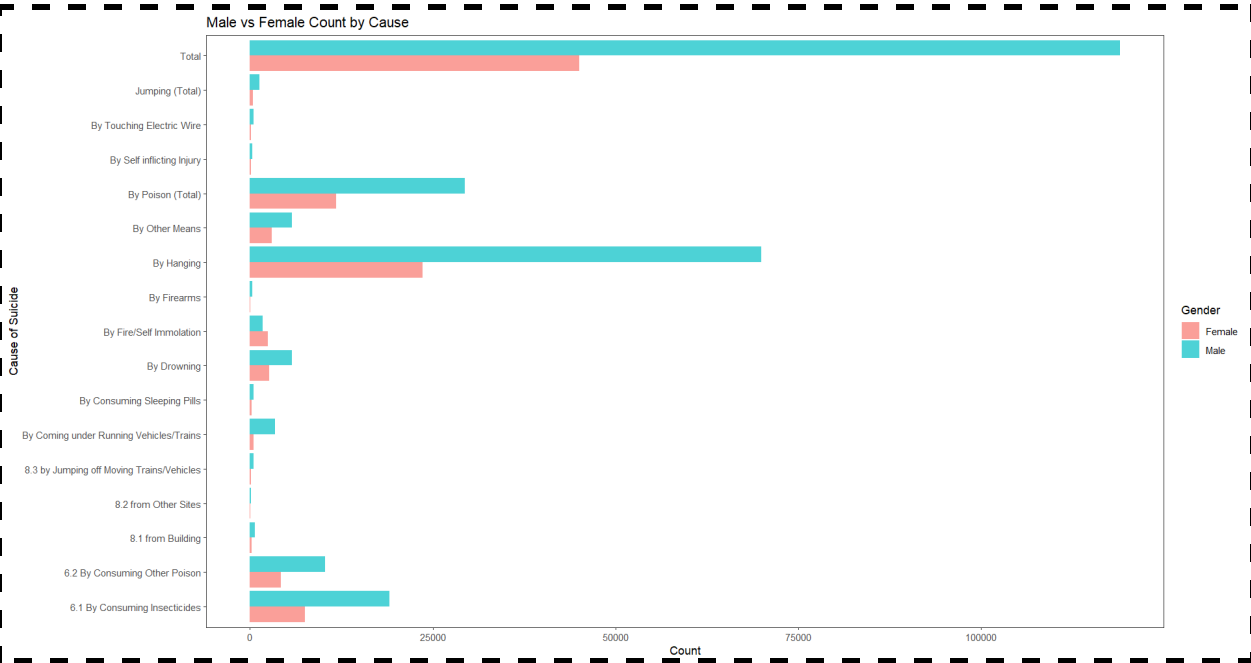
**Code:**

```
30  # plotting bar graph to determine Number of Suicides by causes in India in 2021 Male v/s Female
31
32  #graph1 for male
33  plot2<-data1_cause %>%
34    group_by(Cause)%>%
35    summarize('Total_Number_of_Suicides_Male' = Male) %>%
36    arrange(desc('Total_Number_of_Suicides_Male'))
37  View(plot2)
38  names((plot2))
39  graph1 <- plot2[1:nrow(plot2)-1,] %>%
40    ggplot(aes(Cause,Total_Number_of_Suicides_Male))+
41    geom_bar(stat = "identity", fill = "#ff7043")+
42    #geom_label(aes(label = Total_Number_of_Suicides_Male),size = 1.5)+
43    ggtitle("Male - Number of Suicides by causes in India in 2021")+
44    xlab("Cause of Suicide")+
45    ylab("Male - Total Number of Suicides")+
46    coord_flip()
47  #theme_bw()
48
49
50  #graph2 for female
51  plot3<-data1_cause %>%
52    group_by(Cause)%>%
53    summarize('Total_Number_of_Suicides_Female' = Female) %>%
54    arrange(desc('Total_Number_of_Suicides_Female'))
55  View(plot3)
56  names((plot3))
57  graph2 <- plot3[1:nrow(plot3)-1,] %>%
58    ggplot(aes(Cause,Total_Number_of_Suicides_Female))+
59    geom_bar(stat = "identity", fill = "#ff7043")+
60    #geom_label(aes(label = Total_Number_of_Suicides_Female),size = 1.5)+
61    ggtitle("Female - Number of Suicides by causes in India in 2021")+
62    xlab("Cause of Suicide")+
63    ylab("Female - Total Number of Suicides")+
64    coord_flip()
65  #theme_bw()
66
67  # Arrange the plots in a 2x1 grid layout
68  grid.arrange(graph1, graph2, nrow = 2)
```

## Output:



Male - Number of Suicides by causes in India in 2021

Female - Number of Suicides by causes in India in 2021
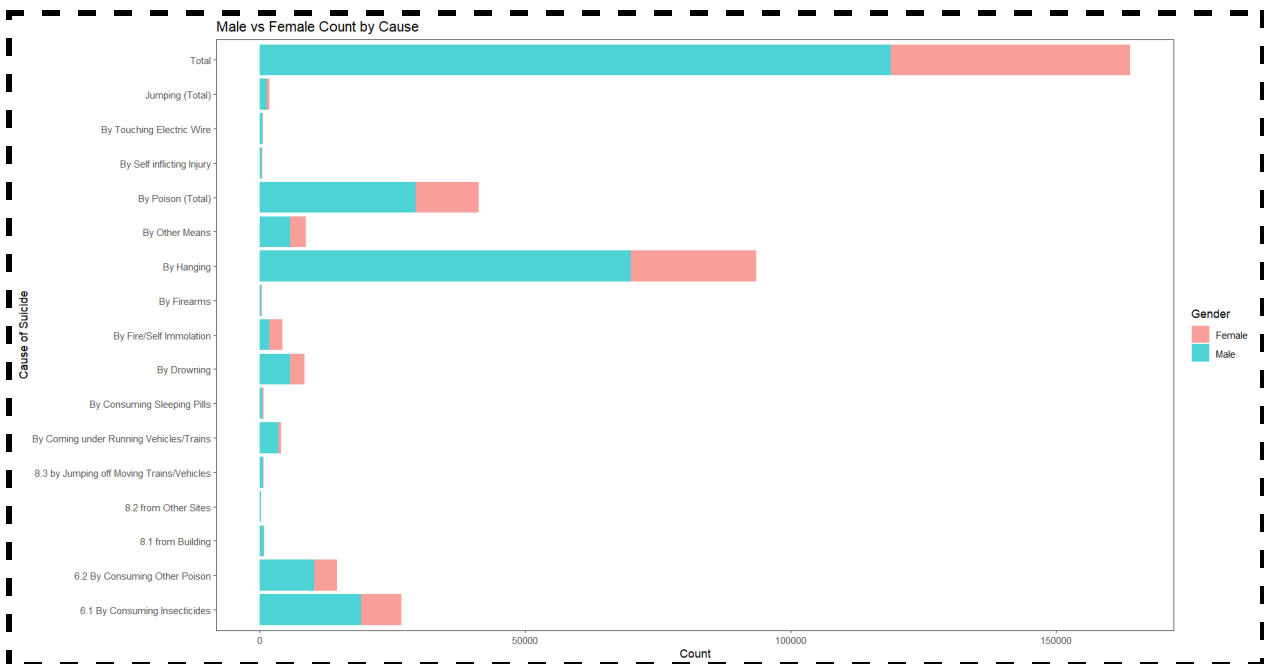
## Code:

```
72  # Male v/s Female Part 1
73  data1_cause %>%
74    pivot_longer(cols = c(Male, Female), names_to = "Gender", values_to = "Count") %>%
75    ggplot(aes(Cause, Count, fill = Gender)) +
76    geom_bar(stat = "identity", position = "dodge", alpha = 0.7) +
77    coord_flip()+
78    labs(x = "Cause of Suicide", y = "Count", title = "Male vs Female Count by Cause")
```



Male vs Female Count by Cause

**Code:**

```
80  # Male v/s Female Part 2
81  data1_cause %>%
82    pivot_longer(cols = c(Male, Female), names_to = "Gender", values_to = "Count") %>%
83    ggplot(aes(Cause, Count, fill = Gender)) +
84    geom_bar(stat = "identity", position = "stack", alpha = 0.7) +
85    coord_flip()+
86    labs(x = "Cause of Suicide", y = "Count", title = "Male vs Female Count by Cause")
```

**Output:**



**Inferences:**

The above three graphs compares the Male and Female suicides with three different perspectives. The first graph shows Male v/s Female Suicide cause side by side, the second graph shows Male v/s Female Suicide cause compared individually and the last graph shows the comparison of Male v/s Female suicide cause, Male and Female count stacked over eachother.
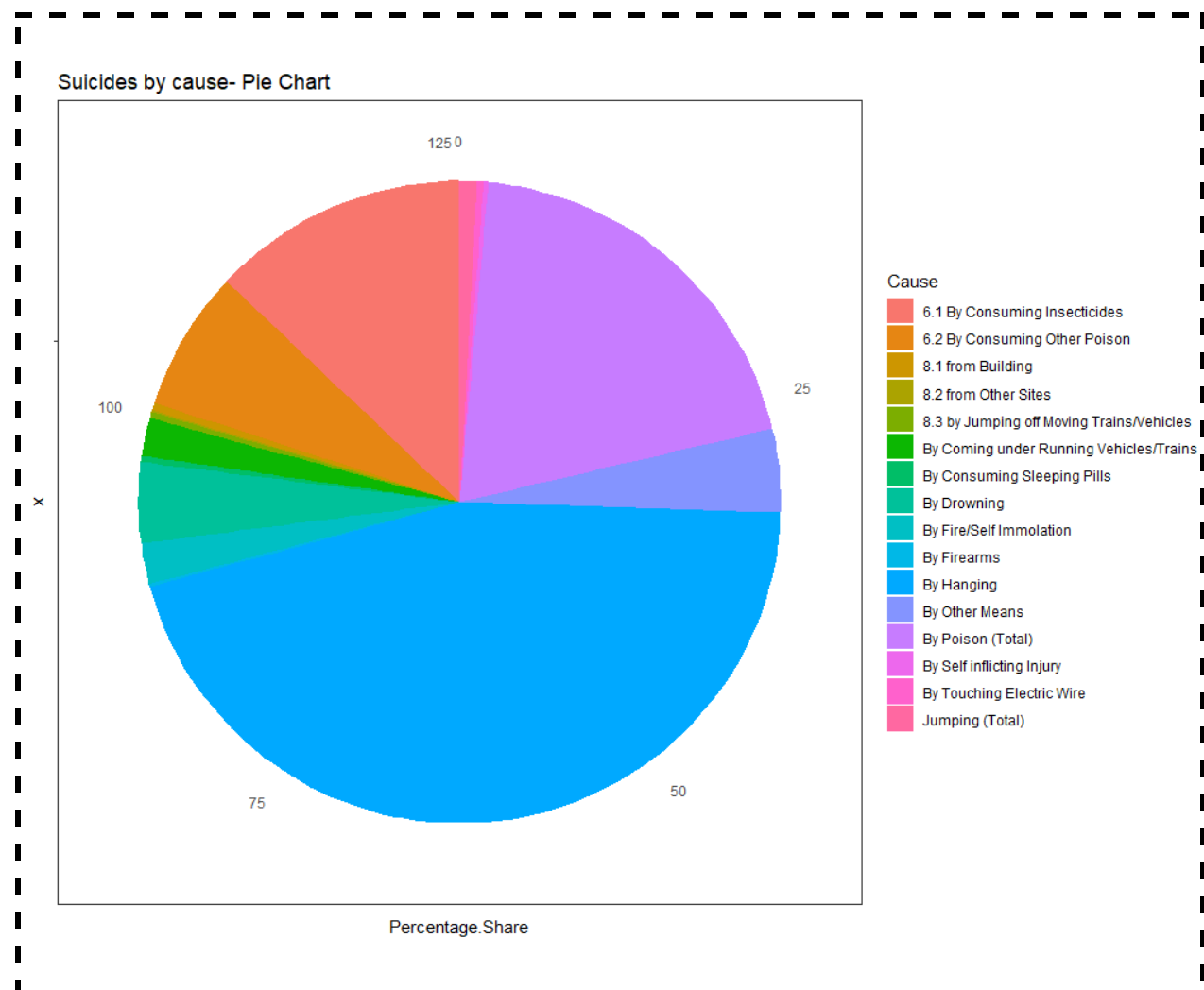
We can make out the following inferences from the graphs:

- Both male and female have committed the most suicides by "hanging".
- In males, "suicide by jumping off Moving Trains/Vehicles", "suicide from building", "suicide from other sites" are the least practiced mode of suicide.
- Whereas in females, "suicide by firearms" and "suicide by touching electric wires" are the least practiced mode of suicide.
- There are less half females suicides than male suicides.

*************************************************************************************************

**Code:**

```
90   # Create a pie chart with cause of suicide
91   plot4<-data1_cause[1:nrow(data1_cause)-1,] %>%
92     select(Cause,Percentage.Share)
93   View(plot4)
94   ggplot(plot4, aes(x = "", y = Percentage.Share, fill = Cause)) +
95     geom_bar(stat = "identity") +
96     coord_polar(theta = "y") +
97     ggtitle("Suicides by cause- Pie Chart")
98
```
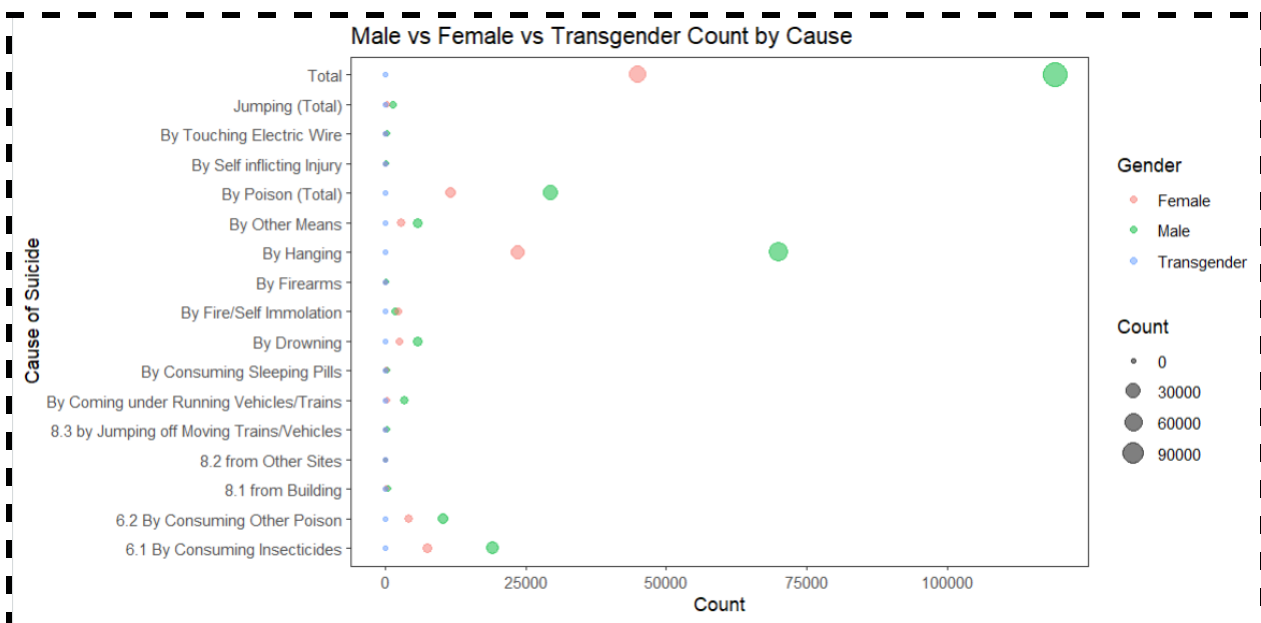
**Output:**



**Inferences:**

The above figure is the Pie Chart of the Suicides Causes in India in 2021.

We can make out from the above Pie chart that "suicides by hanging" is the most practiced mode of suicide followed by "suicide by poision" and "suicide by insecticides".

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

**Code:**

```
104  #scatter plot
105  data1_cause %>%
106    pivot_longer(cols = c(Male, Female,Transgender), names_to = "Gender", values_to = "Count") %>%
107    ggplot(aes(Cause, Count, colour = Gender, size = Count)) +
108    geom_point(alpha = 0.5) +
109    coord_flip()+
110    labs(x = "Cause of Suicide", y = "Count", title = "Male vs Female vs Transgender Count by Cause")
```
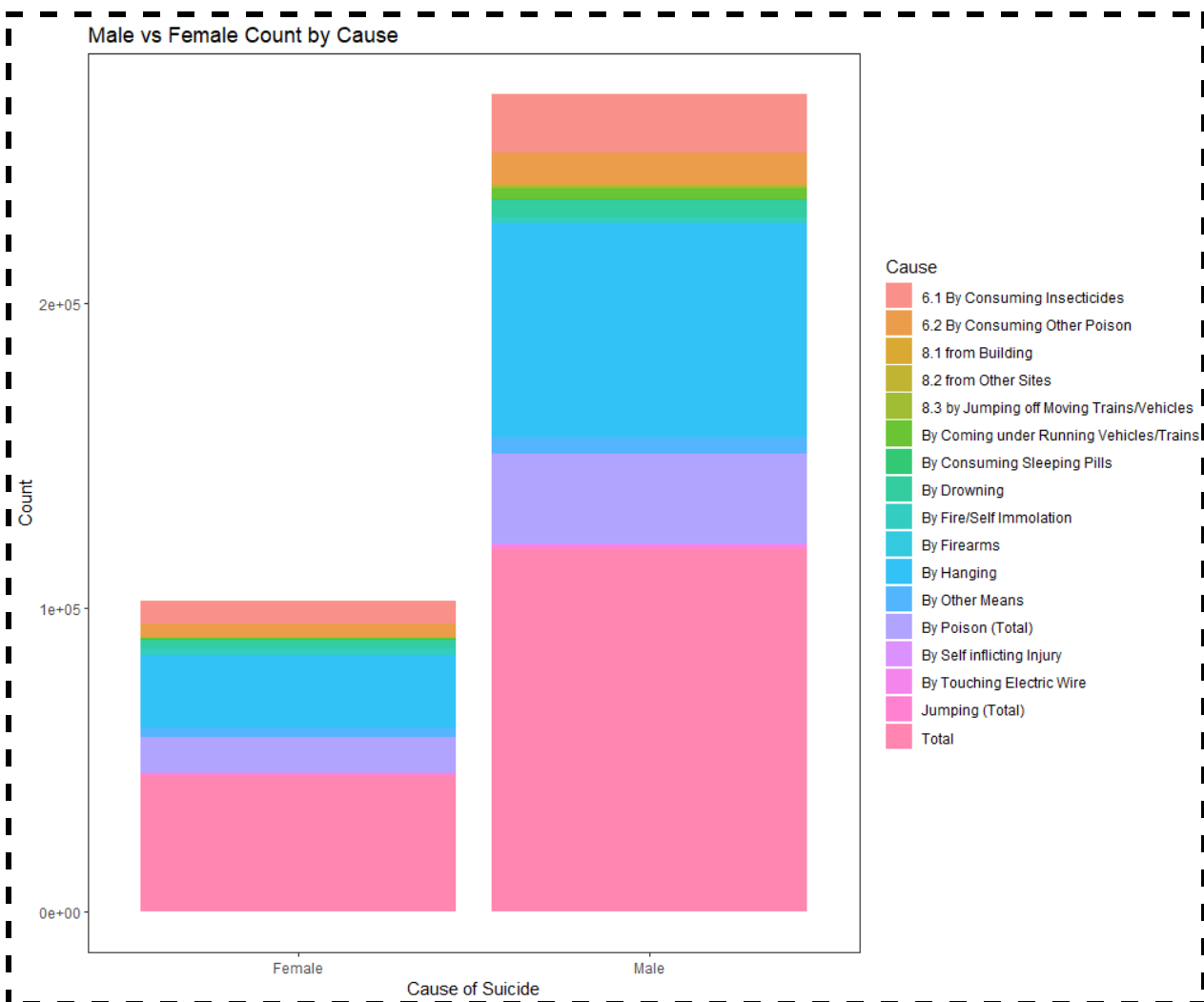
**Output:**



The above scatter-plot graph represents the number of cases in each category of cause of the suicide. The comparison in the above plot is done among the Male, female and transgender. The Point Size of the points are determined by the Count of cases under each category of the 'Cause of suicide'.

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

**Code:**

```
112  #histogram Part(1/2)
113  graph3 <- data1_cause %>%
114    pivot_longer(cols = c(Male, Female), names_to = "Gender", values_to = "Count")
115    ggplot(graph3, aes(Gender, Count, fill = Cause)) +
116    geom_histogram(stat = "identity",position = "stack", alpha = 0.8) +
117    #coord_flip()+
118    labs(x = "Cause of Suicide", y = "Count", title = "Male vs Female Count by Cause")
```
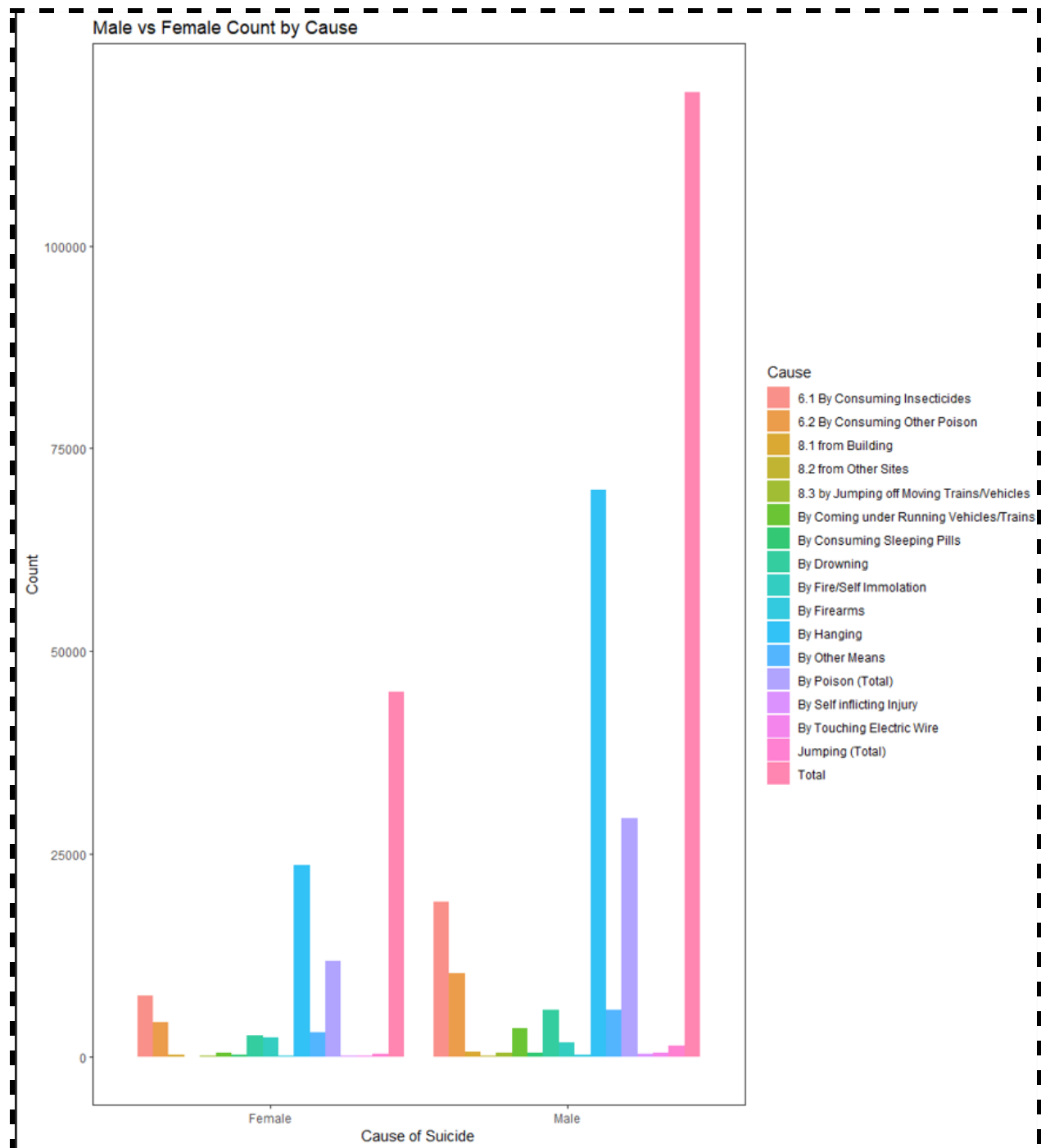
**Output:**



**Inferences:**

The above graph shows the histogram plot of the suicide cases in India in 2021 Male v/s Female. In the graph the cause of the suicide are differentiated by the colour of the stacked bars. Comparison can made through the graph

b/w the Male and Female suicide cases. The graph clearly shows that male total suicides are more than double the female suicides. We can also compare the other category of cause of suicide.

**********************************************************************************************

**Code:**

```
120  #histogram Part(2/2)
121  graph3 <- data1_cause %>%
122    pivot_longer(cols = c(Male, Female), names_to = "Gender", values_to = "Count")
123  ggplot(graph3, aes(Gender, Count, fill = Cause)) +
124    geom_histogram(stat = "identity",position = "dodge", alpha = 0.8) +
125    #coord_flip()+
126    labs(x = "Cause of Suicide", y = "Count", title = "Male vs Female Count by Cause")
```

The above graph is the 'dodge' version of the previous histogram graph.

**********************************************************************************************

# Data Collection

## (Dataset 2/2)

```
1   #install.packages("tidyverse")
2   install.packages("gridExtra")
3   library(tidyverse)
4   library(forcats)
5   library(gridExtra)
6   theme_set(theme_bw()+
7           theme(panel.grid = element_blank()))
8   data2_edu <- read.csv("C:/Users/kumar/Documents/Datasets/Education-wise NCRB_ADSI-2021_Table_2.11.csv")
9   View(data2_edu)
10  names(data2_edu)
```

```
> names(data2_edu)
 [1] "Category"
 [2] "State.UT"
 [3] "No.Education...Male"
 [4] "No.Education...Female"
 [5] "No.Education...Transgender"
 [6] "No.Education...Total"
 [7] "Primary..up.to.class...5th....Male"
 [8] "Primary..up.to.class...5th....Female"
 [9] "Primary..up.to.class...5th....Transgender"
[10] "Primary..up.to.class...5th....Total"
[11] "Middle..up.to.class...8th....Male"
[12] "Middle..up.to.class...8th....Female"
[13] "Middle..up.to.class...8th....Transgender"
[14] "Middle..up.to.class...8th....Total"
[15] "Matriculate..Secondary..up.to.class...10th....Male"
[16] "Matriculate..Secondary..up.to.class...10th....Female"
[17] "Matriculate..Secondary..up.to.class...10th....Transgender"
[18] "Matriculate..Secondary..up.to.class...10th....Total"
[19] "Hr..Secondary..Intermediate..Pre.University..up.to.class...12th....Male"
[20] "Hr..Secondary..Intermediate..Pre.University..up.to.class...12th....Female"
[21] "Hr..Secondary..Intermediate..Pre.University..up.to.class...12th....Transgender"
[22] "Hr..Secondary..Intermediate..Pre.University..up.to.class...12th....Total"
[23] "Diploma.Certificate..ITI...Male"
[24] "Diploma.Certificate..ITI...Female"
[25] "Diploma.Certificate..ITI...Transgender"
[26] "Diploma.Certificate..ITI...Total"
[27] "Graduate.and.above...Male"
[28] "Graduate.and.above...Female"
[29] "Graduate.and.above...Transgender"
[30] "Graduate.and.above...Total"
[31] "Professionals..MBA.etc.....Male"
[32] "Professionals..MBA.etc.....Female"
[33] "Professionals..MBA.etc.....Transgender"
[34] "Professionals..MBA.etc.....Total"
[35] "Status.Not.known...Male"
[36] "Status.Not.known...Female"
[37] "Status.Not.known...Transgender"
[38] "Status.Not.known...Total"
[39] "Total...Male"
[40] "Total...Female"
```

```
[40] "Total...Female"
[41] "Total...Transgender"
[42] "Total...Total"
```
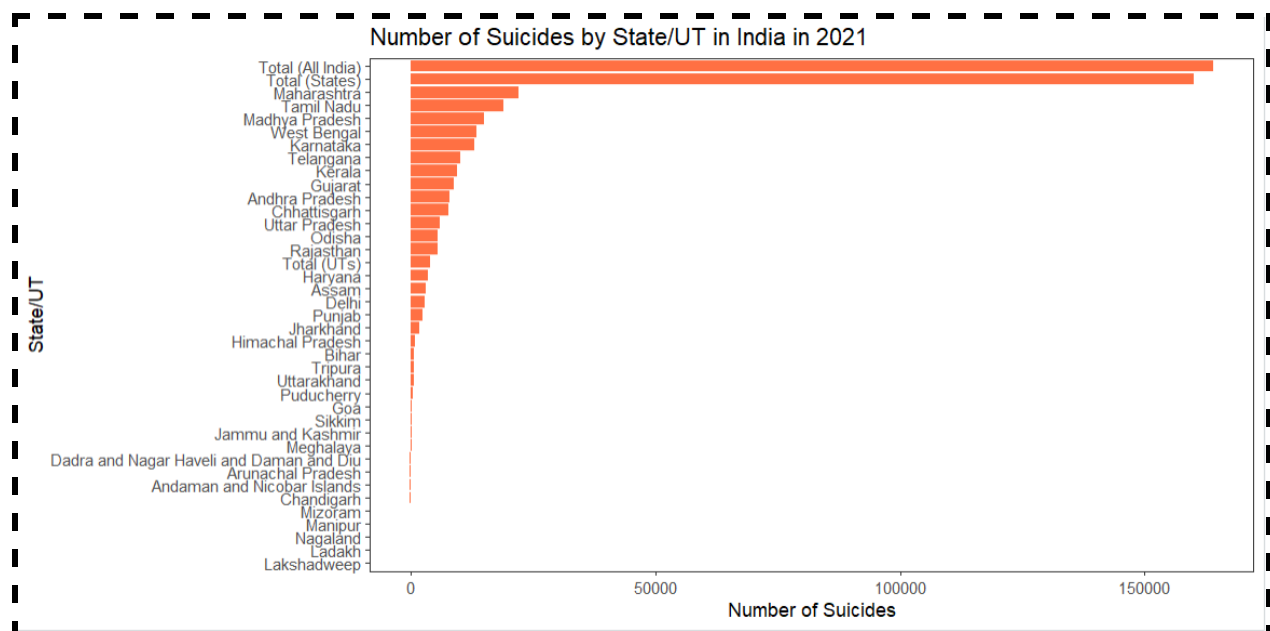
**Inferences:**

- **No. of Attributes: 42**
- **No. of Records: 39**
- **No. of "discrete" attributes: 2**
- **No. of "continuous" attributes: 2**

- **Total No of observations: 1638**

*********************************************************************************************

**Code:**

```
11  # plotting bar graph to determine the number of sucides in each state/UT
12  data2_edu[1:nrow(data2_edu),] %>%
13    ggplot(aes(reorder(State.UT,Total...Total),Total...Total))+
14    geom_bar(stat = "identity", fill = "#ff7043")+
15    ggtitle("Number of Suicides by State/UT in India in 2021")+
16    xlab("State/UT")+
17    ylab("Number of Suicides")+
18    coord_flip()
19  #theme_bw()
```

**Output:**



**Inferences:**

The above graph shows the number of suicides against each states/union territory in India in 2021. From the bar graph we can make out following inferences:

- Maharashtra state has the most number of suicide cases in 2021, followed by Tamil Nadu, Madhya Pradesh, West Bengal, Karnataka...
- Among the states Nagaland has the least number of suicides, followed by Mizoram.
- Among the UTs Delhi has the most number of suicides.
- Among the UTs Lakshadweep has the least number of suicides, followed by Ladakh and Chandigarh.

*********************************************************************************************

```
21  graph1 <- data2_edu %>%
22    pivot_longer(cols = c(No.Education...Male,Primary..up.to.class...5th....Male,Middle..up.to.class...8th....Male,Matriculate..Secor
23
24  View(graph1)
25  graph2<- data2_edu %>%
26    pivot_longer(cols = c(No.Education...Female,Primary..up.to.class...5th....Female,Middle..up.to.class...8th....Female,Matriculate.
27  graph3<- data2_edu %>%
28    pivot_longer(cols = c(No.Education...Transgender,Primary..up.to.class...5th....Transgender,Middle..up.to.class...8th....Transgenc
29  graph4<- data2_edu %>%
30    pivot_longer(cols = c(No.Education...Total,Primary..up.to.class...5th....Total,Middle..up.to.class...8th....Total,Matriculate..Se
31
```
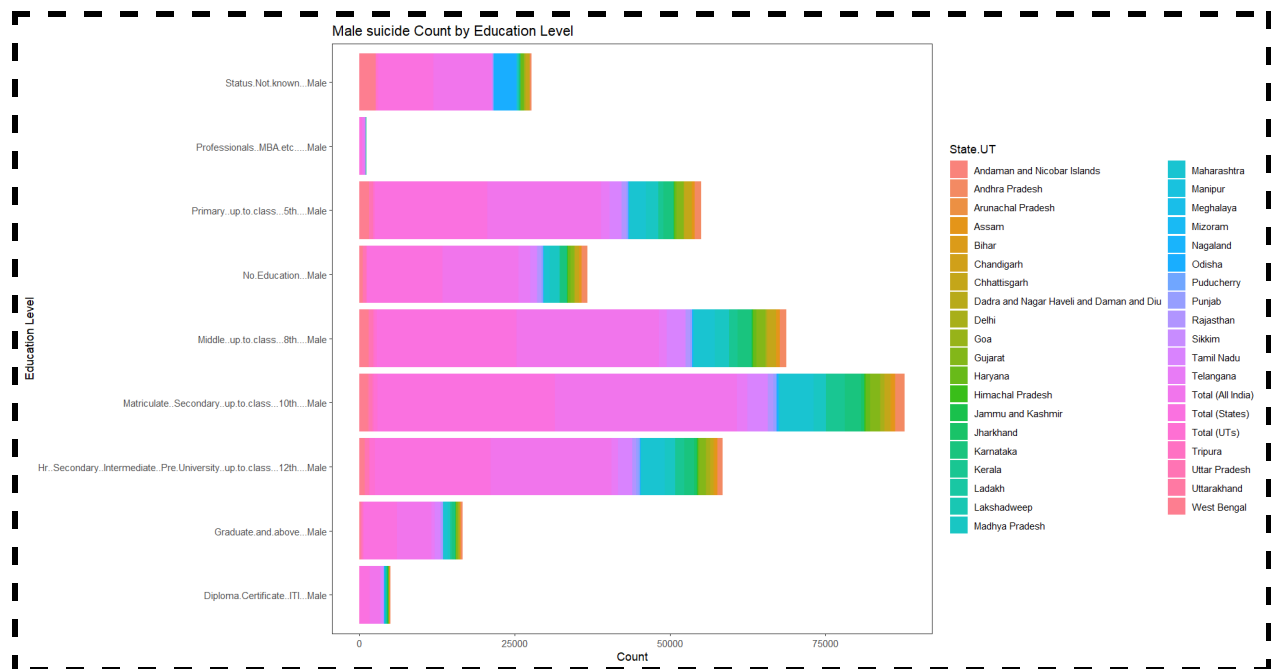
## Code:

```
32  #plot for male
33  ggplot(graph1, aes(Edu_Male, Count_Male, fill = State.UT)) +
34    geom_histogram(stat = "identity",position = "stack", alpha = 0.9) +
35    coord_flip()+
36    labs(x = "Education Level", y = "Count", title = "Male suicide Count by Education Level")
```

## Output:



## Inferences:

The graph depicting suicide counts among males by education level suggests that those who have completed their secondary education up to class 10th are at a greater risk of suicide than males with higher levels of education. Conversely, males with professional degrees, such as an MBA, appear to have the lowest risk of suicide.
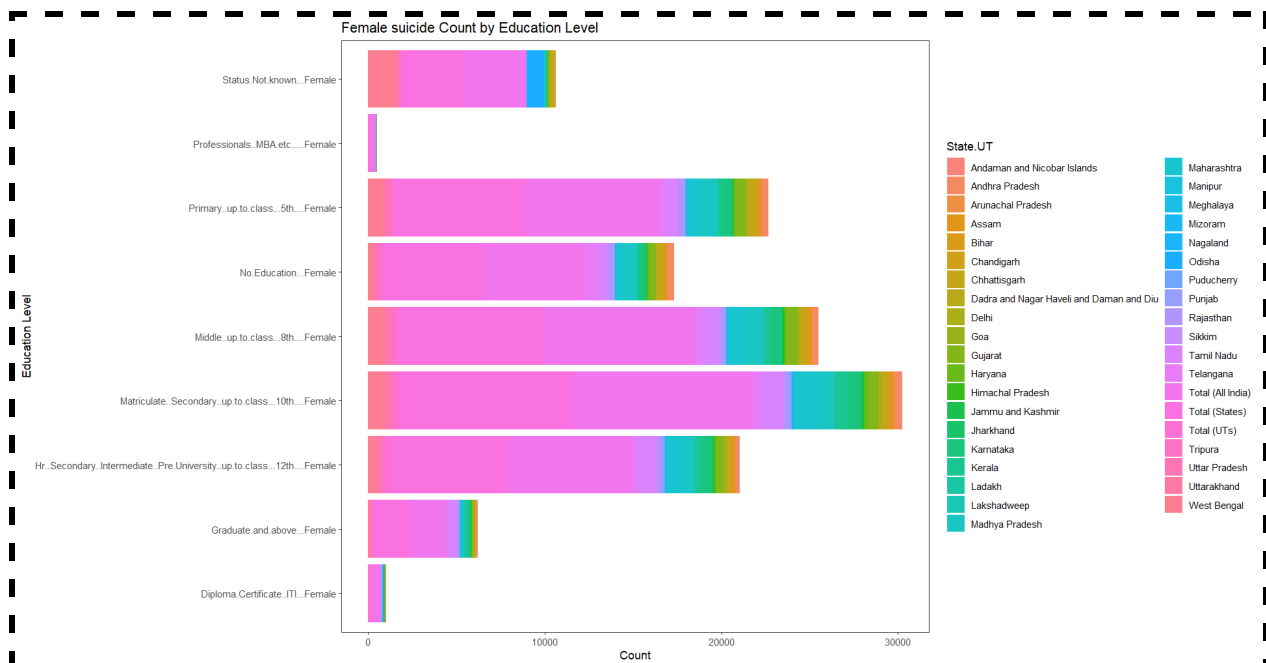
It is important to note that the graph displays variations in suicide rates among states, which are represented by different colors. These variations may reflect underlying socio-economic, cultural, and other contextual factors that

influence suicidal behavior in different regions. Further research is necessary to explore the complex interplay of these factors and their relationship with male suicide rates by education level.

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

**Code:**

```
37  #plot for female
38  ggplot(graph2, aes(Edu_Female, Count_Female, fill = State.UT)) +
39    geom_histogram(stat = "identity",position = "stack", alpha = 0.9) +
40    coord_flip()+
41    labs(x = "Education Level", y = "Count", title = "Female suicide Count by Education Level")
```

**Output:**



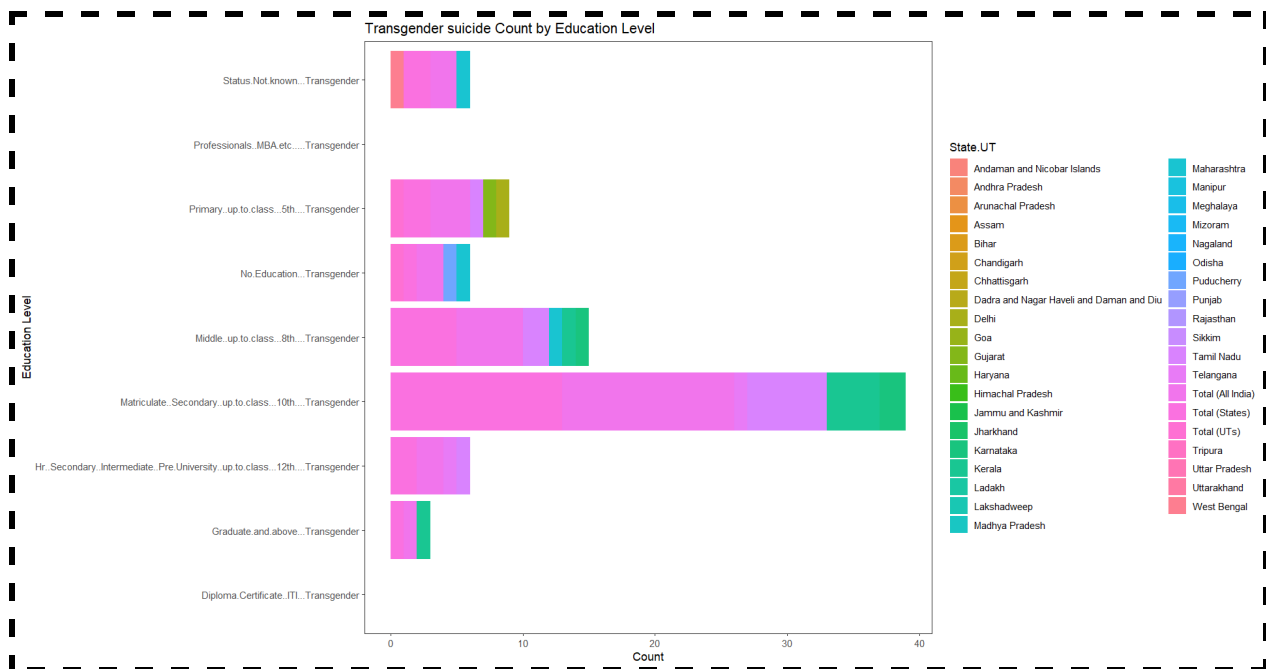Female suicide Count by Education Level

**Inferences:**

The above graph represents the Female suicide count by Eductaion level. In the graph the various states suicides are represented by the various colours. The following inferences can be made from the above data:

- Most suicides are with those females whose education level is 'Mariculate Seconday up to class 10th'.
- Least suicides are by the 'Professionals MBA etc…' Females.
- We can observe a decreasing trends in suicides as a female goes for higher education or have education level more than 'Mariculate Seconday up to class 10th'.

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

**************************************************************************************************

## Code:

```
43  #plot for transgender
44  ggplot(graph3, aes(Edu_Transgender, Count_Transgender, fill = State.UT)) +
45    geom_histogram(stat = "identity",position = "stack", alpha = 0.9) +
46    coord_flip()+
47    labs(x = "Education Level", y = "Count", title = "Transgender suicide Count by Education Level")
```
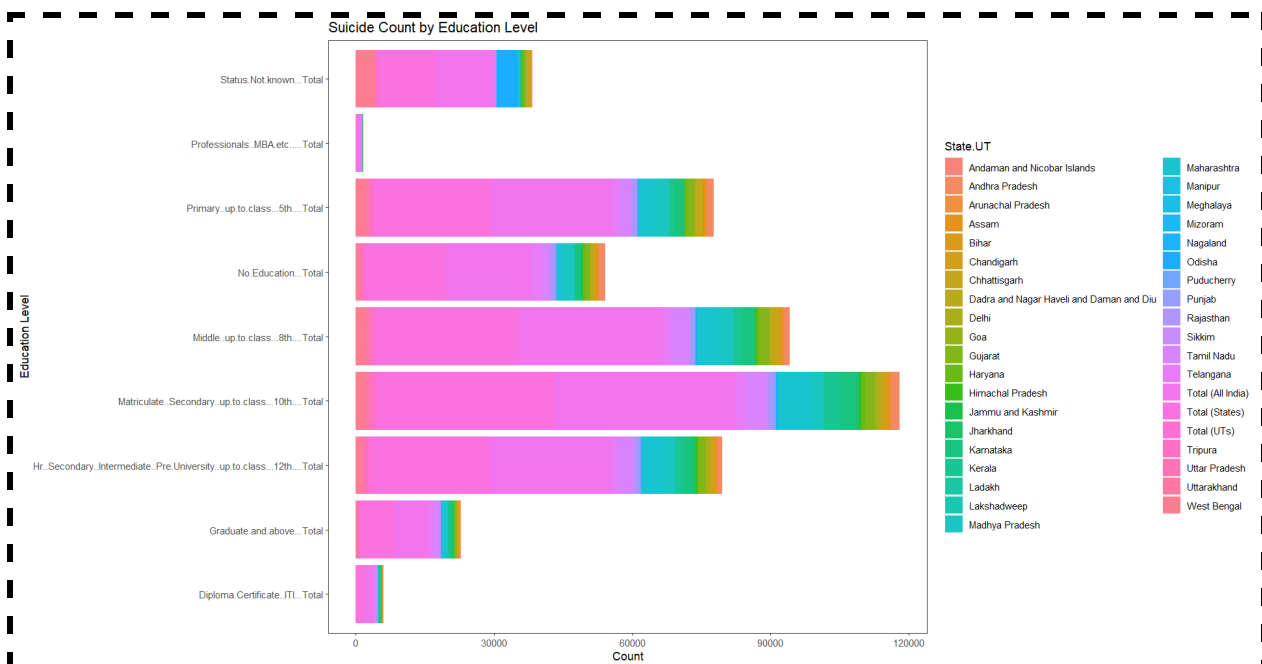
## Output:



## Inferences:

Based on the graph depicting suicide counts among transgender individuals by education level, it can be inferred that those with a secondary education level up to class 10th have the highest incidence of suicide. Conversely, transgender individuals with graduate-level education or higher have the lowest incidence of suicide.

The graph also indicates that suicide rates vary among states, with each state represented by a different color.

**************************************************************************************************

********************************************************************************************

**Code:**

```
48  #plot by Education
49  ggplot(graph4, aes(Edu_Total, Count_Total, fill = State.UT)) +
50    geom_histogram(stat = "identity",position = "stack", alpha = 0.9) +
51    coord_flip()+
52    labs(x = "Education Level", y = "Count", title = "Suicide Count by Education Level")
```

**Output:**



Suicide Count by Education Level

**Inferences:**

According to the graph representing suicide counts by education level, it can be inferred that individuals who have completed their secondary education up to class 10th have the highest suicide rates among all education levels. Conversely, individuals with professional degrees such as an MBA have the lowest suicide rates.

Additionally, the graph indicates that suicide rates vary by state, with each state having a different color. Further analysis is required to understand the reasons for these variations.

********************************************************************************************