**Name:** Shrey Pendurkar
**Class:** D15C
**Batch:** C
**Roll No:** 64

# DMBI - Experiment 3

**Aim: To perform exploratory data analysis and visualization on the dataset using python.**

## Theory:

Exploratory Data Analysis (EDA) is the process of exploring, summarizing, and visualizing data to understand its main characteristics before applying statistical models or predictive analytics. For the **Global YouTube Statistics 2023** dataset (top ~1000 YouTubers worldwide), EDA helps uncover patterns in subscribers, views, uploads, categories, countries, and earnings.

### 1. Descriptive Statistics

- Numerical summaries (mean, median, variance, min/max, std) are computed for key features like `subscribers`, `video_views`, `uploads`, and `earnings`.
- Helps detect skewness — e.g., subscriber counts are highly **right-skewed** with a few mega-channels dominating.
- Outliers (channels with extremely high views or uploads) are identified, showing the power-law nature of YouTube popularity.

### 2. Target Variable Analysis

- Since this is not a binary classification problem, the **main dependent variables** can be considered as `subscribers` or `video_views`.
- Visualizing top 10 YouTubers by subscribers and by views reveals dominance of entertainment, music, and kids' content.
- Distributions show that most channels have fewer subscribers/views compared to a handful of global leaders.

### 3. Correlation Analysis

- Pearson correlation is calculated between numeric features (`subscribers`, `video_views`, `uploads`, `earnings`).
- Heatmap highlights strong positive correlation between **subscribers and views**, confirming loyal audiences drive viewership.
- Uploads show weaker correlation with subscribers/views, indicating that **quantity doesn't guarantee popularity**.
- Earnings are moderately correlated with views but influenced by other factors (CPM, content type).

### 4. Feature Distribution Analysis

- Histograms/KDE plots reveal that features like `subscribers` and `video_views` follow a **long-tailed distribution** (few channels dominate).
- Boxplots of `views per subscriber` highlight engagement efficiency — some categories deliver disproportionately high views relative to their subscriber base.
- `Uploads` distribution shows many small creators with large video counts but relatively fewer views.

### 5. Categorical Feature Analysis

- Countplots of `category` show **Music, Entertainment, People & Blogs, and Shows** as dominant genres.
- Bar plots by `country` reveal the U.S. and India as leaders, followed by Brazil and other high-population nations.
- Category vs. subscribers/earnings analysis shows that **Show and Music channels** achieve higher engagement than other genres.
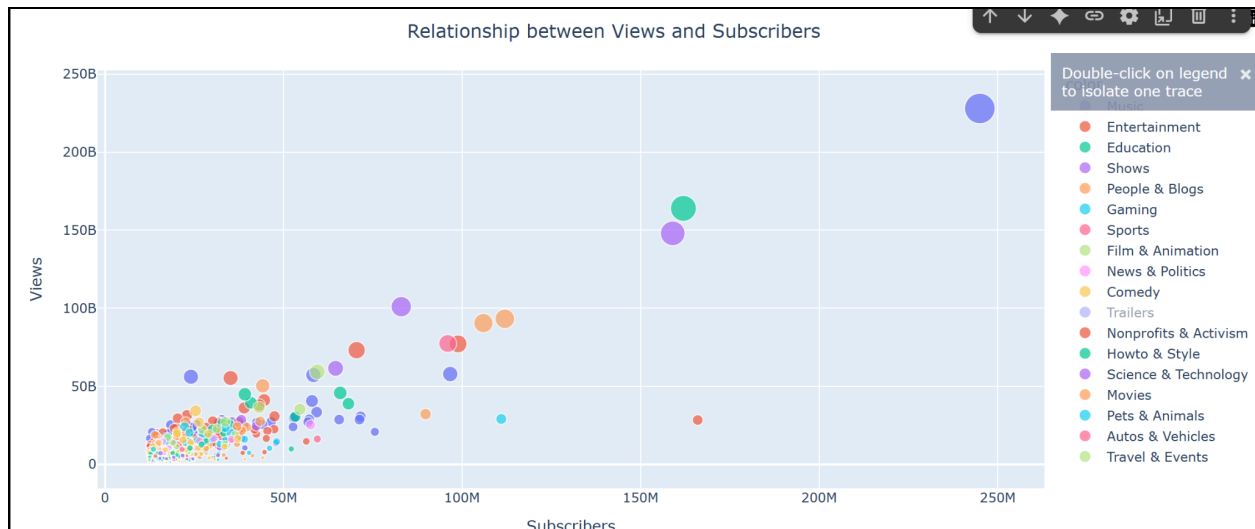
### 6. Multivariate Visualization

- Scatter plots (`subscribers vs. video_views`) show clusters of highly correlated channels, with outliers representing viral mega-channels.
- Pairplots across (`subscribers`, `views`, `uploads`, `earnings`) highlight differences across categories.
- Grouped boxplots (`category vs. views per subscriber`) illustrate niches where content achieves stronger audience retention.

## Conclusion:

The Exploratory Data Analysis (EDA) and visualization of the Global YouTube Statistics 2023 dataset revealed key insights about top YouTube channels worldwide. Subscribers and video views showed a strong positive correlation, indicating audience size drives engagement, while uploads had a weak influence, emphasizing content quality over quantity. Popular categories like Music and Shows achieved higher engagement and earnings, and countries like the U.S. and India led channel counts and viewership. Visualizations such as scatter plots, pie charts, bar charts, box plots, and correlation heatmaps effectively summarized these relationships.
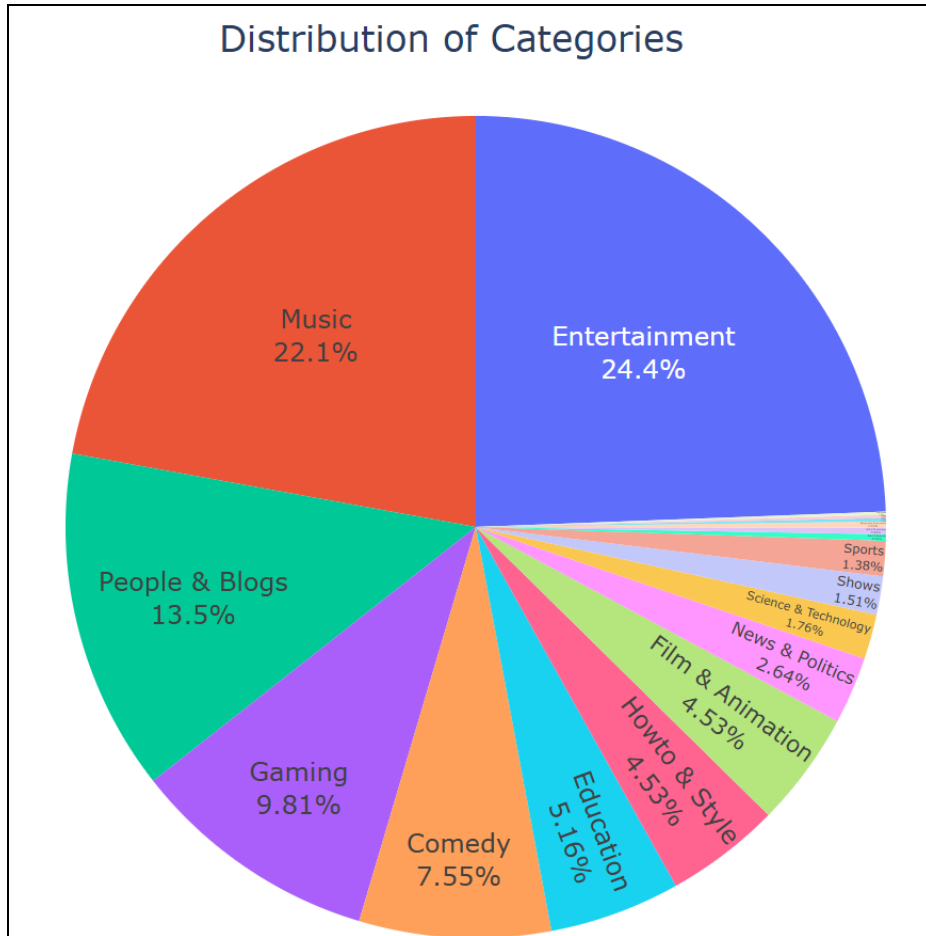
### 1) Scatter Plot:

```python
# Creating a scatter plot using plotly express (px)
fig1 = px.scatter(
    x=df["subscribers"],
    y=df["video views"],
    color=df["category"],
    size=df["video views"],
    hover_name=df["Title"]
)
# Updating the layout of the plot
fig1.update_layout(
    xaxis=dict(title="Subscribers"),
    yaxis=dict(title="Views"),
    title="Relationship between Views and Subscribers",
    title_x=0.48
)
# Displaying the plot
fig1.show()
```
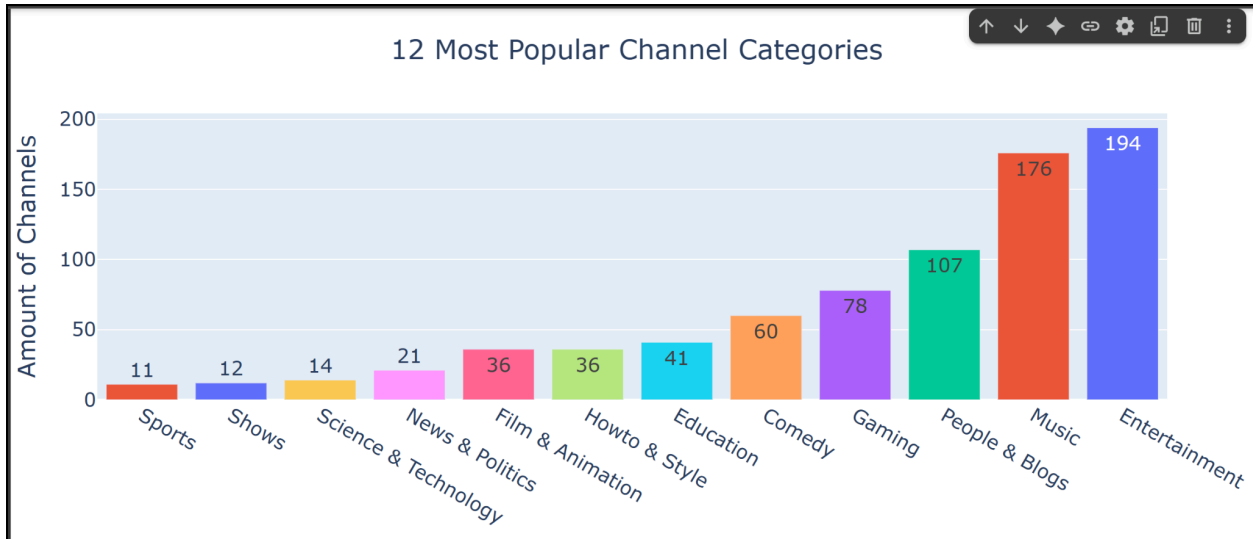
## 2) Pie Chart:

```python
# Calculating the count of each category and sorting them in ascending order
category = df['category'].value_counts().sort_values(ascending=True)
# Creating a pie chart using plotly express (px)
fig3 = px.pie(
    values=category.values,
    names=category.index
)
# Updating trace properties for the pie chart
fig3.update_traces(
    textposition='inside',
    textinfo='percent+label'
)
# Updating the layout of the pie chart
fig3.update_layout(
    title_text="Distribution of Categories",
    title_x=0.49,
    uniformtext_minsize=10,
    showlegend=False
)
# Displaying the pie chart
fig3.show()
```
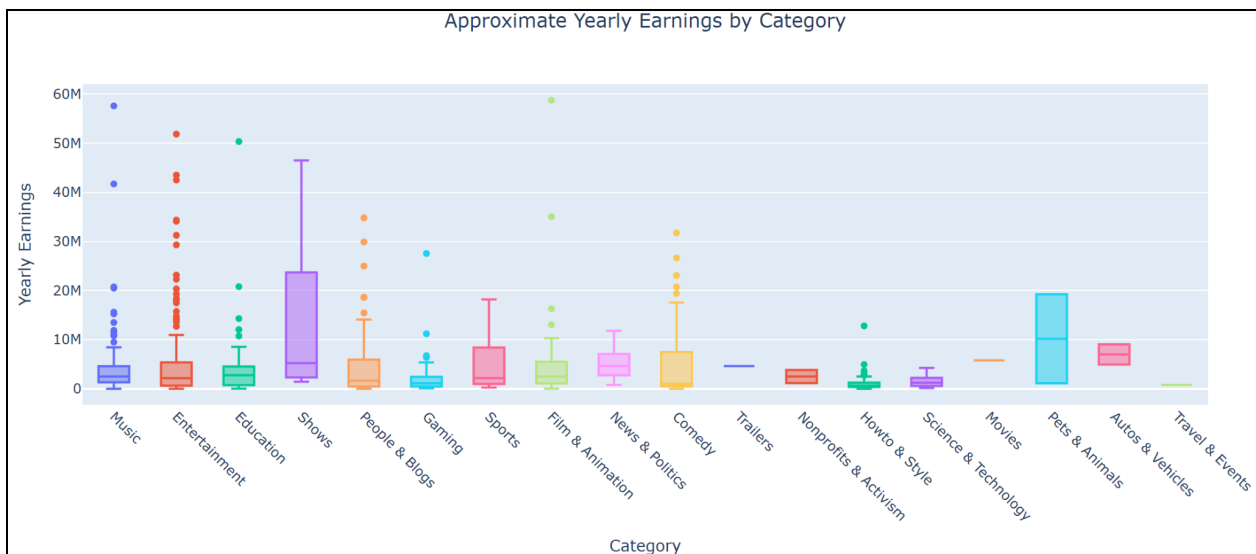
## Distribution of Categories



Entertainment 24.4%
Music 22.1%
People & Blogs 13.5%
Gaming 9.81%
Comedy 7.55%
Education 5.16%
Howto & Style 4.53%
Film & Animation 4.53%
News & Politics 2.64%
Science & Technology 1.76%
Shows 1.51%
Sports 1.38%

**3) Bar Chart:**

```python
# Get the top 12 channel categories
top12_channel_categories = df['category'].value_counts().head(12)
# Plotting the bar chart using plotly.express
fig4 = px.bar(
    x=top12_channel_categories.index,
    y=top12_channel_categories.values,
    color=top12_channel_categories.index,
    text=top12_channel_categories.values,
    title='12 Most Popular Channel Categories',
    labels={'x': 'Category', 'y': 'Amount of Channels'}
)
# Customize the layout
fig4.update_layout(
    title_x=0.5,
    font=dict(size=18),
    showlegend=False,
    xaxis={'categoryorder': 'total ascending'}
)
# Display the plot
fig4.show()
```

12 Most Popular Channel Categories

## 4) Box Plot:

```
# Create a box plot using plotly.express
fig7 = px.box(df, x='category', y='yearly_earnings', color='category',
            labels={'yearly_earnings': 'Yearly Earnings', 'category': 'Category'},
            title='Approximate Yearly Earnings by Category')
# Center the title above the bar chart
fig7.update_layout(title=dict(x=0.5), showlegend=False)
# Customize the layout for better readability
fig7.update_layout(xaxis=dict(tickangle=45))
# Show the plot
fig7.show()
```



Approximate Yearly Earnings by Category

## 5) Correlation Heatmap:

```python
import matplotlib.pyplot as plt
import seaborn as sns
plt.style.use("seaborn")
plt.rcParams['figure.figsize'] = (16, 8)
# Filter numeric columns only
numeric_df = df.select_dtypes(include=[float, int])
# Compute the correlation matrix
correlation_matrix = numeric_df.corr()
# Plot the heatmap
title = "Correlation Heatmap"
plt.title(title, fontsize=18, weight='bold')
sns.heatmap(correlation_matrix, cmap="BuPu", annot=True)
plt.show()
```



Correlation Heatmap