

Data Visualization

Contents

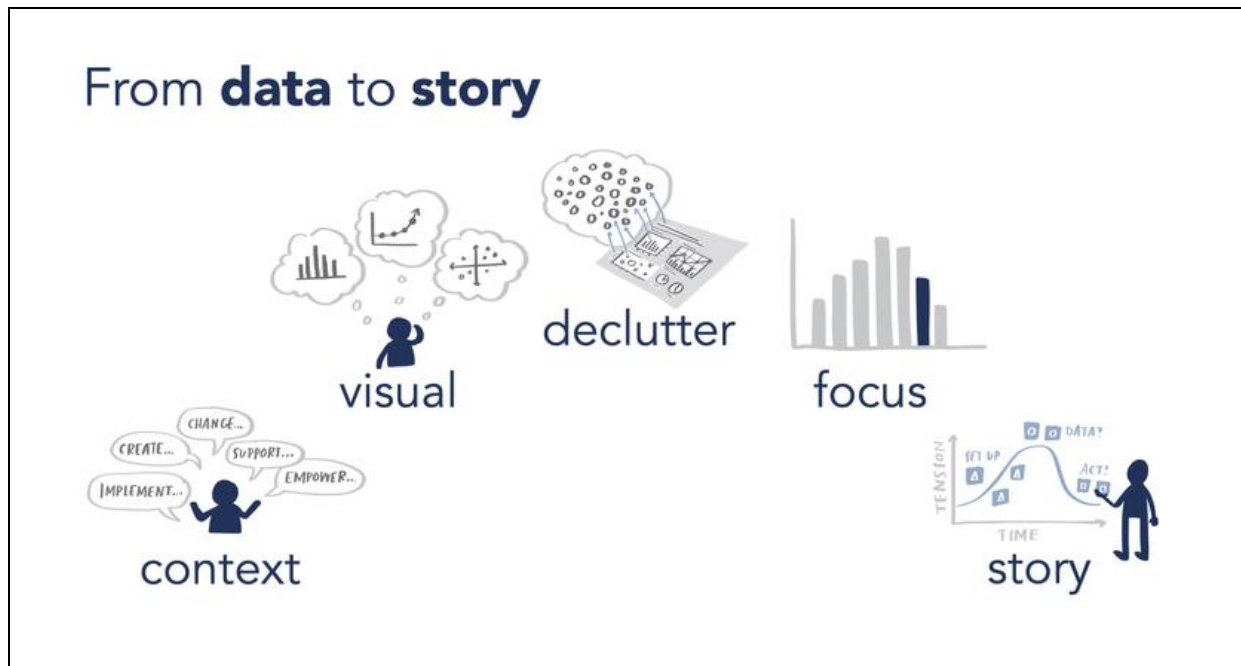
| | |
|--|----|
| 1. Data Introduction | 4 |
| 2. What is Data? | 4 |
| 3. Data Visualization | 5 |
| 1.1. Example in words | 5 |
| 4. Advantages | 5 |
| 4.1. Common data visualization techniques | 6 |
| 4.2. With data visualization | 6 |
| 5. Real time Examples | 7 |
| 5.1. Data & Visual: Netflix Subscribers | 7 |
| 5.2. Data & Visual: Top Social Media Usage Statistics 2024 | 9 |
| 5.3. Data & Visual: Top ChatGPT Statistics (2024) | 10 |
| 6. Process behind the Data Visualization | 14 |
| 6.1. Data collection | 14 |
| 6.2. Data cleaning | 15 |
| 6.3. Data analysis | 16 |
| 6.4. Chose the right visualization | 17 |
| 6.5. Creating visual representation | 18 |
| 6.6. Review and Iterative | 18 |
| 7. Types of Data Visualization | 19 |
| 7.1. Basic Charts | 19 |
| 7.2. Statistical Visualizations | 20 |
| 7.3. Advanced Charts | 21 |
| 7.4. Interactive Visualizations | 21 |
| 7.5. Textual Visualizations | 21 |
| 8. Gestalt principles for Data Visualization | 22 |
| 8.1. Proximity | 22 |
| 8.2. Similarity | 22 |
| 8.3. Closure | 22 |
| 8.4. Continuity | 22 |
| 9. Visualization reference model | 23 |
| 9.1. Data Layer | 23 |
| 9.2. Processing Layer | 23 |

| | |
|---|-----------|
| 9.3. Visualization Layer | 24 |
| 9.4. Interaction Layer | 24 |
| 9.5. Presentation Layer | 24 |
| 9.6. Feedback Layer | 25 |
| 9.7. Deployment Layer | 25 |
| 10. Data visualizations by the number of variables | 26 |
| 10.1. Univariate Visualizations..... | 26 |
| 10.2. Bivariate Visualizations | 26 |
| 10.3. Multivariate Visualizations..... | 27 |
| 11. Matplotlib | 28 |
| 12. Line chart | 29 |
| 12.1. Labelling the axes..... | 33 |
| 13. Bar Chart | 35 |
| 14. Histogram | 39 |
| 15. Pie Chart | 40 |
| 15.1. Attributes | 42 |
| 16. Scatter Plot | 43 |
| 17. Box Plots | 45 |
| 17.1. Use Box plots..... | 45 |
| 17.2. Box plot explanation | 46 |
| 18. Heatmap | 47 |
| 18.1. How to understand? | 47 |

Data Visualization



Data Visualization



1. Data Introduction

- ✓ Currently we are all living in the data world.
- ✓ Everyone is communicating by using devices and social networks, due to this huge amount of data is generating.
- ✓ All applications are generating data.
 - Ecommerce applications.
 - Banking applications.
 - Social network etc.

2. What is Data?

- ✓ Data is a collection of Facts.
- ✓ Facts can be,
 - Numbers
 - Alphabets
 - Alphanumeric
 - Symbols
 - Images
 - Audio
 - Video & etc

3. Data Visualization

- ✓ **Data Visualization** is the process of converting data into a graphical representation.
- ✓ If we visualize the data then it is very easy to understand.

Best quote

- ✓ A picture gives more meaningful information than thousand words

1.1. Example in words

- ✓ Reaching to target



4. Advantages

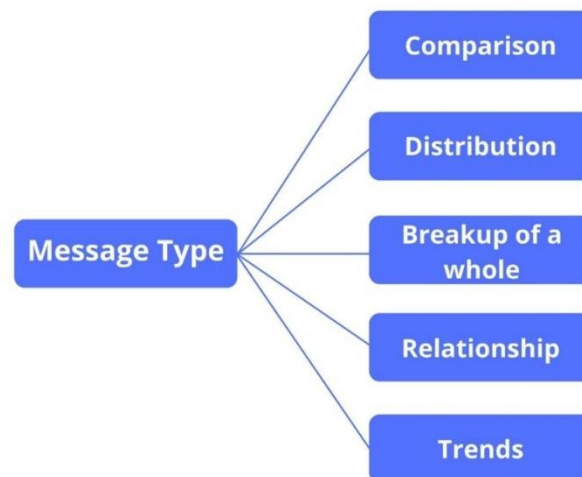
- ✓ To identify **trends**, such as whether sales increasing or decreasing.
- ✓ To identify **patterns**, such as during weekend more sales.
- ✓ To identify **relationships**, such as if we study more hours then we will get good marks.
- ✓ To identify **frequency**, such as how often a product is purchased in a specific area & etc

4.1. Common data visualization techniques

- ✓ Bar charts
- ✓ Pie charts
- ✓ Line graphs
- ✓ Box plot
- ✓ Scatter plot & etc

4.2. With data visualization

- ✓ We are sharing message/information to end users.



5. Real time Examples

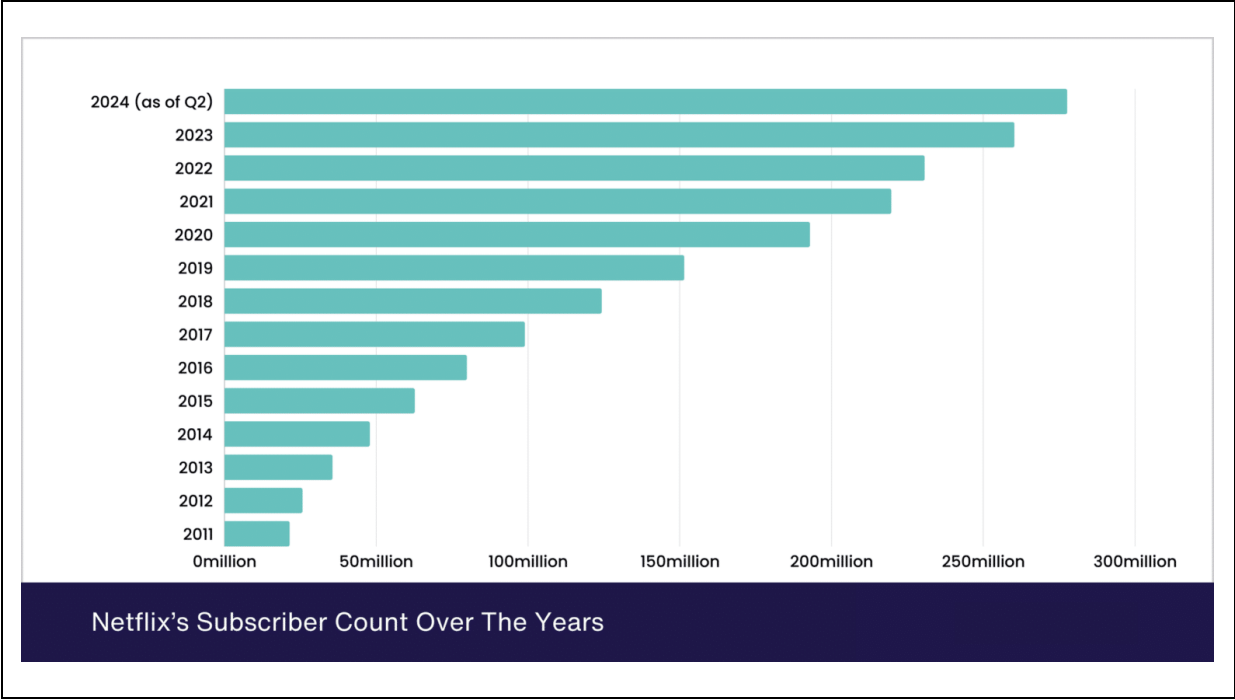
5.1. Data & Visual: Netflix Subscribers

Ref Link: <https://www.demandsage.com/netflix-subscribers/>

Top Netflix Statistics At A Glance

- Netflix has 277.65 million subscribers as of 2024.
- Netflix generated \$18.93 billion in revenue in the first half of 2023.
- Women make up 51%, while Males make up 49% of all Netflix users.
- Netflix is preferred by 47% of Americans over other streaming platforms and is responsible for 8.4% of the screen time in the country.
- Around 65% of Netflix consumers are from outside of the United States of America & Canada.
- Netflix customers spend 62.1 minutes each day on average consuming content.

| Year | Netflix Subscribers |
|-----------------|---------------------|
| 2024 (as of Q2) | 277.65 million |
| 2023 | 260.28 million |
| 2022 | 230.7 million |
| 2021 | 219.7 million |
| 2020 | 192.9 million |
| 2019 | 151.5 million |
| 2018 | 124.3 million |
| 2017 | 99 million |
| 2016 | 79.9 million |
| 2015 | 62.7 million |
| 2014 | 47.9 million |
| 2013 | 35.6 million |
| 2012 | 25.7 million |
| 2011 | 21.5 million |



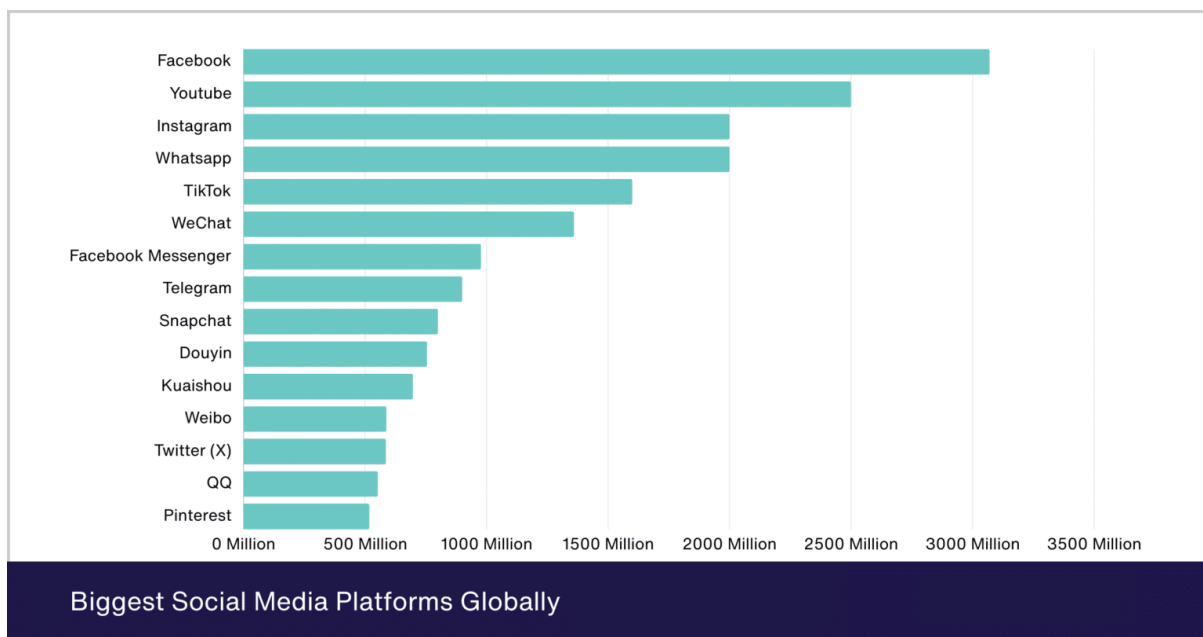
5.2. Data & Visual: Top Social Media Usage Statistics 2024

Ref Link: <https://www.demandsage.com/social-media-users/>

Top Social Media Usage Statistics 2024

- There are 5.17 billion social media users globally.
- 68% of the people in the United States use social media, approximately 308 million people.
- Facebook is the biggest social media platform, with over 3.07 billion users.
- A typical social media user interacts with 6.7 social media platforms.
- On average, users spend 2 hours and 20 minutes daily on Social media platforms.
- China has the highest number of social media users, with 1.07 billion users in the country.

Most Popular Social Media Platforms



5.3. Data & Visual: Top ChatGPT Statistics (2024)

Ref Link: <https://www.demandsage.com/chatgpt-statistics/>

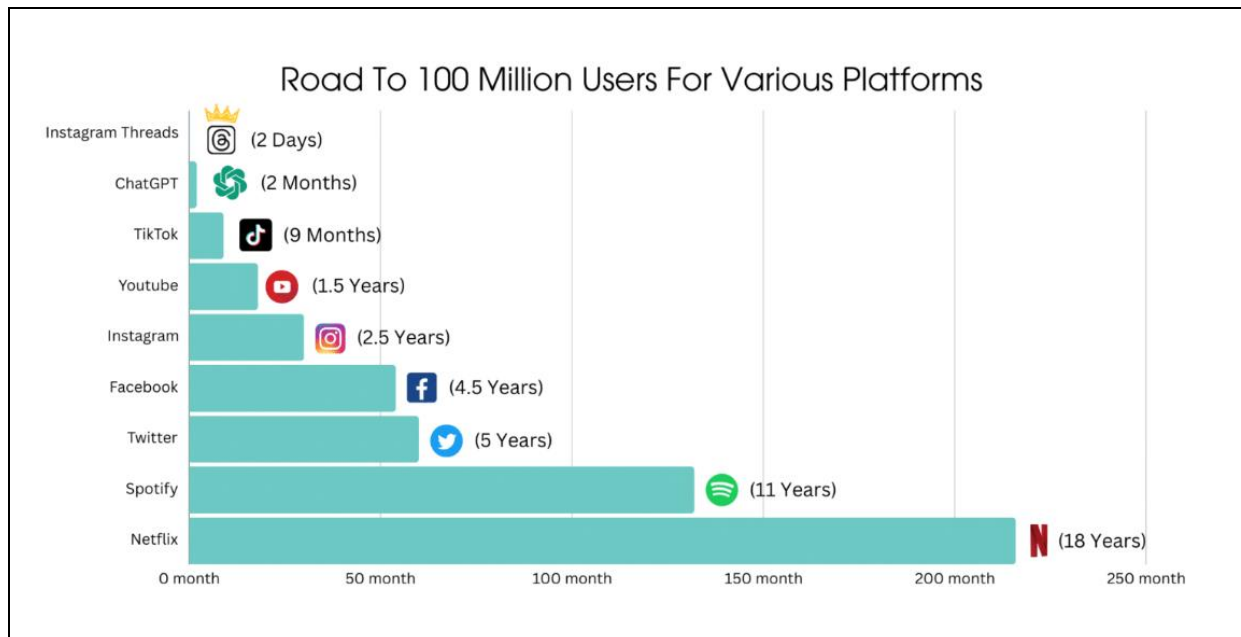
Top ChatGPT Statistics (2024)

- ChatGPT has over 200 million weekly active users as of September 2024.
- Around 77.2 million monthly active users in the US.
- ChatGPT Plus is used by 7.7 million people worldwide.
- ChatGPT reached 1 million users in just five days after its launch.
- More than 92% of Fortune 500 companies are using ChatGPT.
- ChatGPT is forecasted to generate a revenue of \$1 billion in 2024.
- OpenAI spends approximately \$700,000 every day to operate ChatGPT.
- ChatGPT gets over 1.54 billion page visits every month on average.

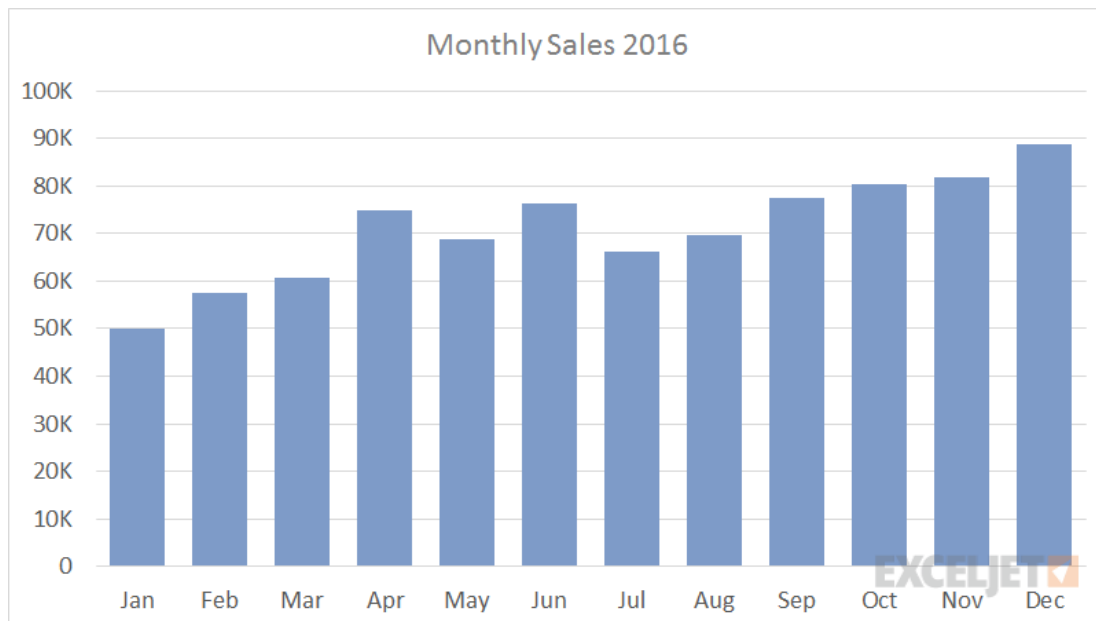
ChatGPT User Demographics

ChatGPT Gender Split

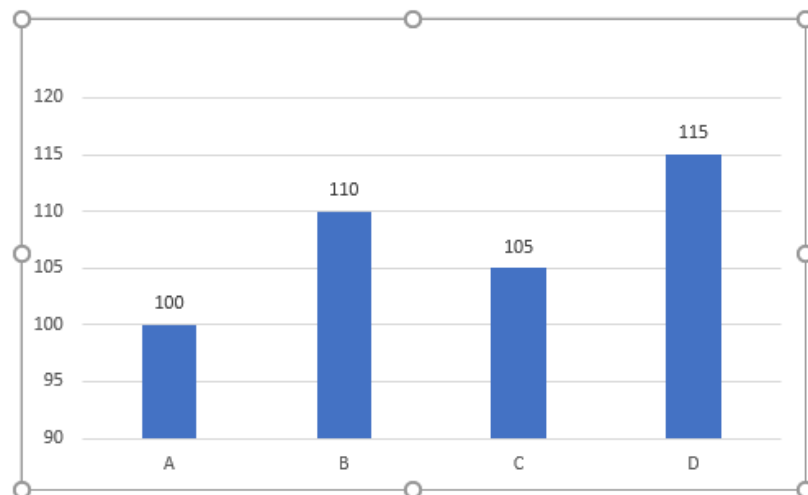




Bar chart



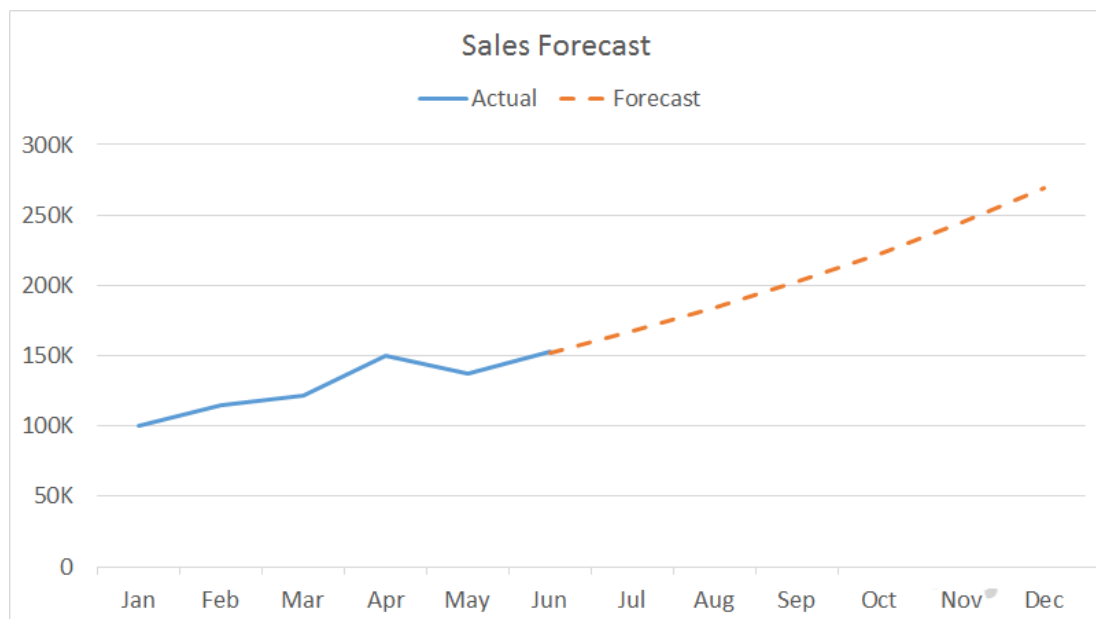
| Group | Value |
|-------|-------|
| A | 100 |
| B | 110 |
| C | 105 |
| D | 115 |



Sales Data

| | Actual | Forecast |
|-----|--------|----------|
| Jan | 100K | |
| Feb | 115K | |
| Mar | 121K | |
| Apr | 150K | |
| May | 137K | |
| Jun | 152K | 152K |
| Jul | | 167K |
| Aug | | 184K |
| Sep | | 202K |
| Oct | | 223K |
| Nov | | 245K |
| Dec | | 269K |

Line Chart



6. Process behind the Data Visualization

Steps

- ✓ Data collection
- ✓ Data cleaning
- ✓ Data analysis
- ✓ Chose the right visualization
- ✓ Creating visual representation
- ✓ Review and Iterative

6.1. Data collection

- ✓ Gather/collect the relevant data from difference sources.



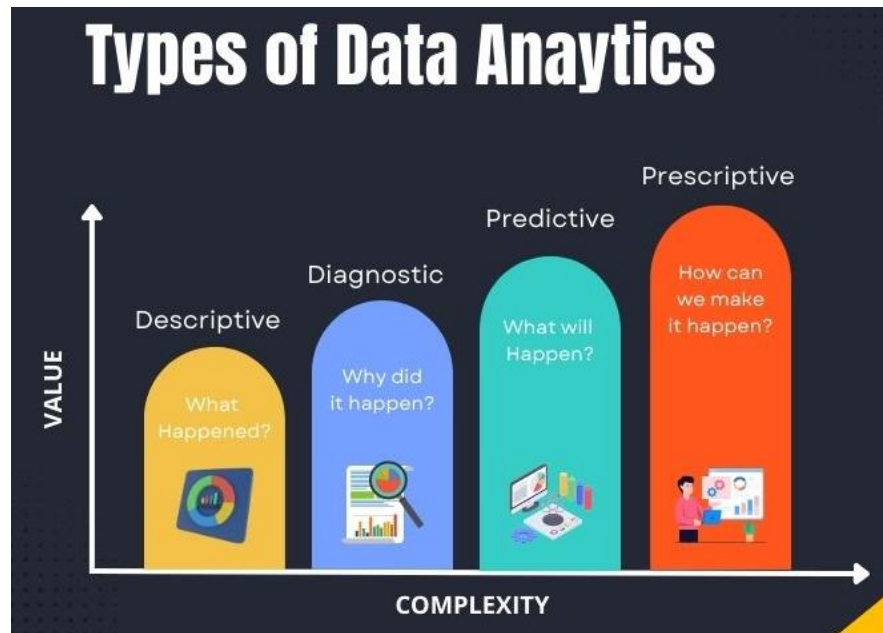
6.2. Data cleaning

- ✓ Ensuring accuracy and consistency.



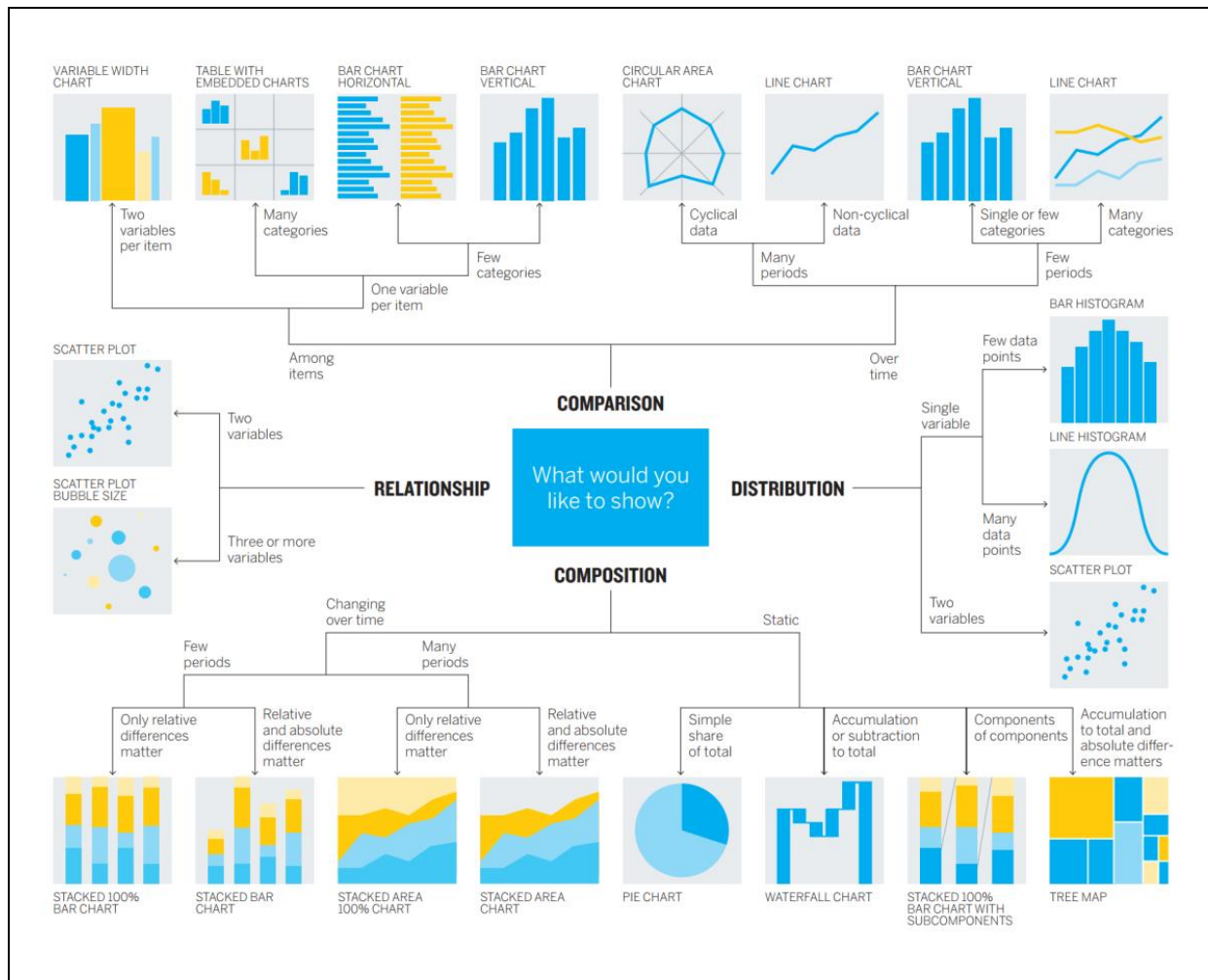
6.3. Data analysis

- ✓ Exploring data to identify trends and patterns.



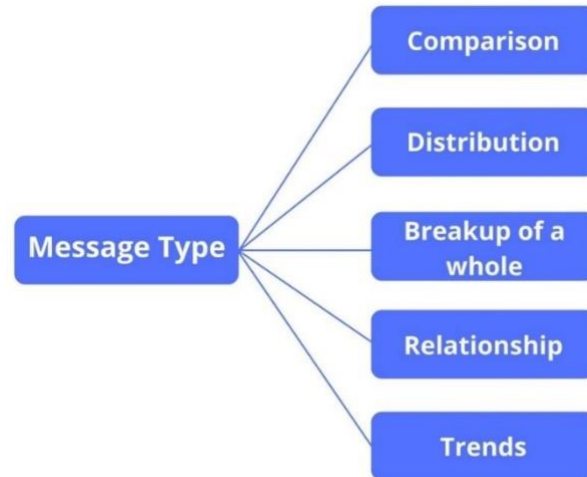
6.4. Chose the right visualization

- ✓ Selecting appropriate charts, graphs, or maps based on requirement.



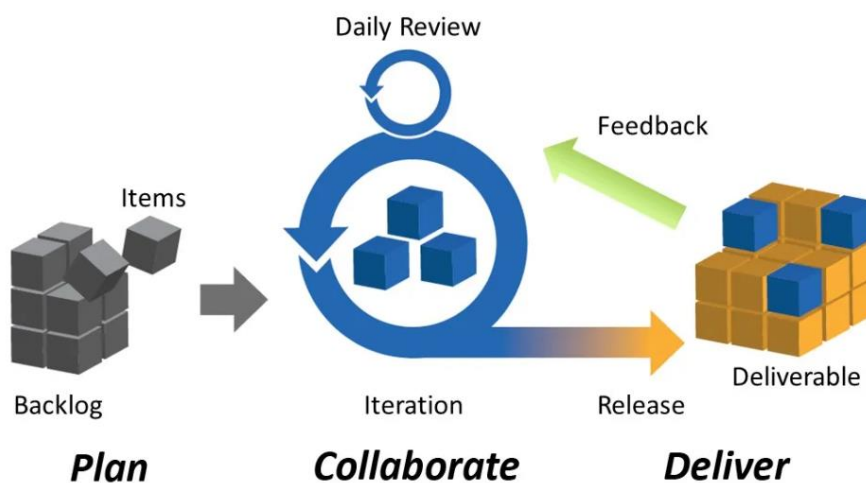
6.5. Creating visual representation

- ✓ Create data visualization to share information effectively



6.6. Review and Iterative

- ✓ Testing the visualization for clarity and effectiveness, making adjustments if needed.



7. Types of Data Visualization

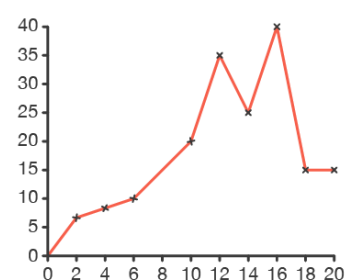
7.1. Basic Charts

- ✓ **Bar Chart:** Compares quantities across different categories.
 - **Column Chart:** Vertical version of the bar chart.
- ✓ **Line Chart:** Shows trends over time or continuous data.
- ✓ **Pie Chart:** Displays proportions of a whole; best for limited categories.

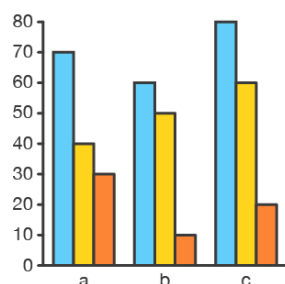
pie chart



line graph



bar chart

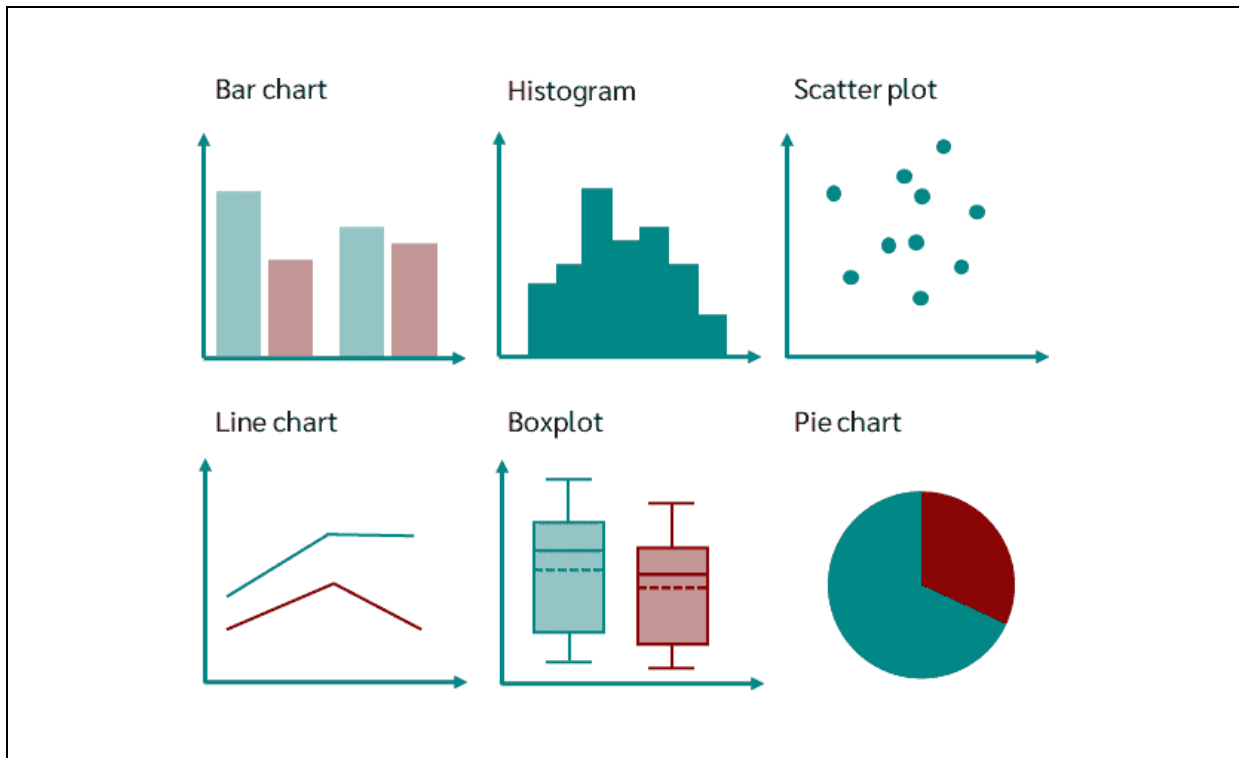


pictogram

| category | frequency |
|----------|-----------|
| A | ■ ■ ■ ■ |
| B | ■ ■ ■ |
| C | ■ |
| D | ■ ■ ■ ■ |

7.2. Statistical Visualizations

- ✓ **Histogram:** Represents the distribution of numerical data.
- ✓ **Box Plot:** Summarizes data distribution.
- ✓ **Scatter Plot:** Shows the relationship between two variables.



7.3. Advanced Charts

- ✓ **Heatmap:** Uses color to represent data values in a matrix format.
- ✓ **Bubble Chart:** A scatter plot with an added dimension represented by the size of the bubbles.

7.4. Interactive Visualizations

- ✓ **Dashboards:** Combines multiple visualizations to provide an overview of key metrics.
- ✓ **Data Explorer:** Allows users to interact with data, filtering and adjusting parameters.

7.5. Textual Visualizations

- ✓ **Word Cloud:** Displays text data, highlighting frequently used terms.
- ✓ **Tag Cloud:** Similar to a word cloud, often used for categorizing data.

8. Gestalt principles for Data Visualization

8.1. Proximity

- ✓ Grouping related items together.
 - Example: In a scatter plot, showing sales data for different products, grouping all electronics together.

8.2. Similarity

- ✓ Using similar shapes or colours to indicate relationships.
 - Example: Using same colors for regions in a bar chart (e.g., blue for East, green for West) helps viewers easily compare sales.

8.3. Closure

- ✓ Completing incomplete shapes to create a whole.
 - Example: A line graph, showing monthly temperature changes can use dotted lines for predictions, allowing viewers to intuitively connect the dots and grasp trends.

8.4. Continuity

- ✓ Following lines and patterns to guide the viewer's eye.
 - Example: A line graph, a stock price graph clearly shows trends, allowing viewers to easily spot increases and decreases over time.

9. Visualization reference model

- ✓ A visualization reference model works as a framework for understanding the components and process in Data Visualization.

Steps

- ✓ Data Collection
- ✓ Data Processing
- ✓ Visualization Design
- ✓ User Interaction
- ✓ Presentation
- ✓ Feedback
- ✓ Deployment

9.1. Data Layer

- ✓ **Data Sources**
 - Identify where the data is coming from (databases, APIs, spreadsheets).
- ✓ **Data Preparation**
 - Data Cleaning, transforming, and aggregating data to ensure it's ready for visualization.

9.2. Processing Layer

- ✓ **Data Analysis:**
 - Apply statistical methods or algorithms to extract insights from the data.
- ✓ **Data Reduction:**
 - Simplifying data by selecting key variables or filtering out noise to focus on important information.

9.3. Visualization Layer

- ✓ **Visual Encoding:**
 - Choosing how to represent data visually (e.g., using colors, shapes, sizes).
- ✓ **Chart Types:**
 - Selecting appropriate visualization types based on the data and insights (e.g., bar charts, line graphs, scatter plots).
- ✓ **Design Principles:**
 - Applying best practices for layout, color schemes, labelling, and accessibility.

9.4. Interaction Layer

- ✓ **Interactivity:**
 - Incorporating features like tooltips, filters, and zooming to allow users to explore data dynamically.
- ✓ **User Experience:**
 - Ensuring that the interaction is intuitive and enhances the understanding of the data.

9.5. Presentation Layer

- ✓ **Contextual Information:**
 - Providing background, legends, and annotations to help interpret the visualizations.
- ✓ **Storytelling:**
 - Structuring the visualizations to convey a narrative and guide the viewer through the insights.

9.6. Feedback Layer

✓ **User Testing:**

- Gathering input from users to assess clarity, effectiveness, and engagement.

✓ **Iteration:**

- Refining the visualizations based on feedback to improve understanding and impact.

9.7. Deployment Layer

✓ **Distribution:**

- Sharing the visualizations through dashboards, reports, or web applications.

✓ **Maintenance:**

- Regularly updating the visualizations to reflect new data and insights.

10. Data visualizations by the number of variables

- ✓ We can divide the data visualization based on the number of variables.

Types

- ✓ Univariate Visualizations
- ✓ Bivariate Visualizations
- ✓ Multivariate Visualizations

10.1. Univariate Visualizations

- ✓ These visualizations focus on a single variable, allowing you to explore its distribution and key statistics.
 - **Histograms:** Show the distribution of a continuous variable.
 - **Bar Charts:** Represent count of categories in a categorical variable.
 - **Box Plots:** Summarize the distribution, median, quartiles, and outliers of a single variable.
 - **Pie Charts:** Explains the proportions of categories in a categorical variable.
 - **Density Plots:** Display the distribution of a continuous variable in a smoothed format.

10.2. Bivariate Visualizations

- ✓ These visualizations explore the relationship between two variables, helping to identify correlations or patterns.
 - **Scatter Plots:** Show the relationship between two continuous variables.
 - **Line Graphs:** By using this we can display how one variable changes in relation to another variable
 - **Grouped Bar Charts:** Compare the values of a categorical variable across different groups.
 - **Heatmaps:** Represent the intensity of a variable across two dimensions (e.g., correlation matrices).
 - **Bubble Charts:** Extend scatter plots by adding a third variable represented by the size of the bubbles.

10.3. Multivariate Visualizations

- ✓ These involve three or more variables and can include
 - **3D Scatter Plots:** Visualize relationships among three continuous variables.
 - **Parallel Coordinates:** By using this we can understand how several variables relate to one another.
 - **Facet Grids:** Display multiple plots in a grid, each representing a subset of data based on one or more categorical variables.

11. Matplotlib

- ✓ Matplotlib is a powerful and widely-used plotting library for Python
- ✓ Using matplotlib we can plot the data.

Environment

- ✓ We can install this library by using pip command.

matplotlib installation

```
pip install matplotlib
```

12. Line chart

- ✓ A line chart or line graph is a type of chart which displays information as a series of data points connected by straight line
- ✓ A line chart is often used to visualize a trend in data over intervals of time.

Program Name Create a simple line chart
demo1.py

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

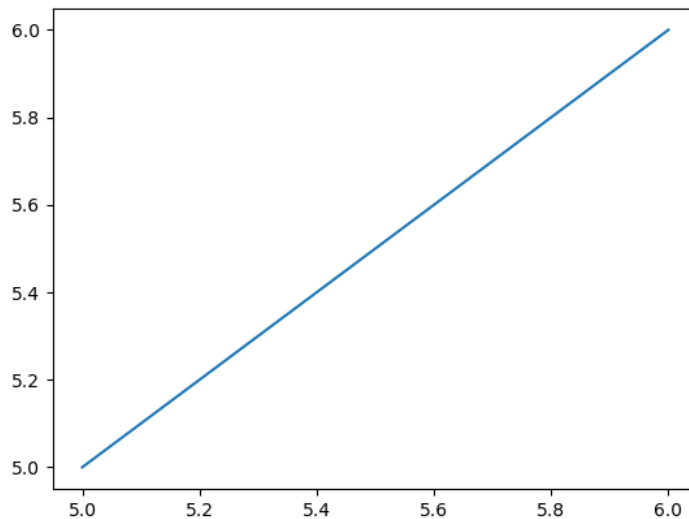
```
x = [5, 6]
```

```
y = [5, 6]
```

```
plt.plot(x, y)
```

```
plt.show()
```

Output



Program Name Create a simple line chart
demo2.py

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

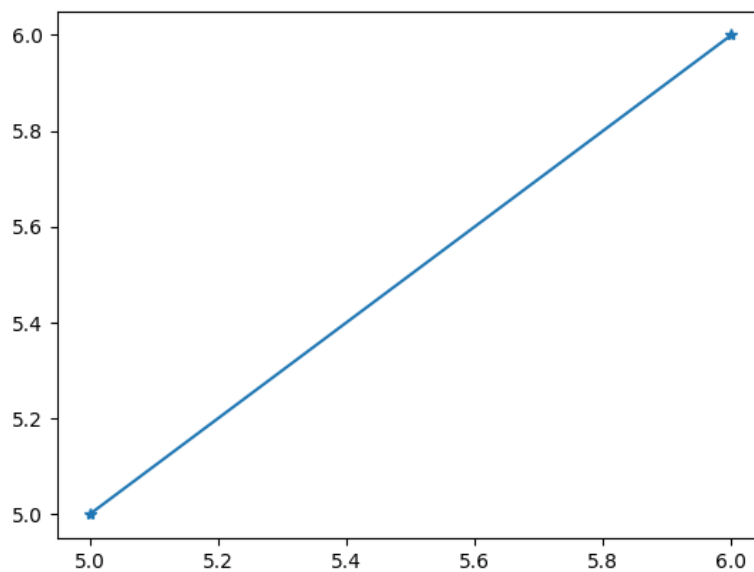
```
x = [5, 6]
```

```
y = [5, 6]
```

```
plt.plot(x, y, marker='*')
```

```
plt.show()
```

Output



Program Name Create a simple line chart
demo3.py

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

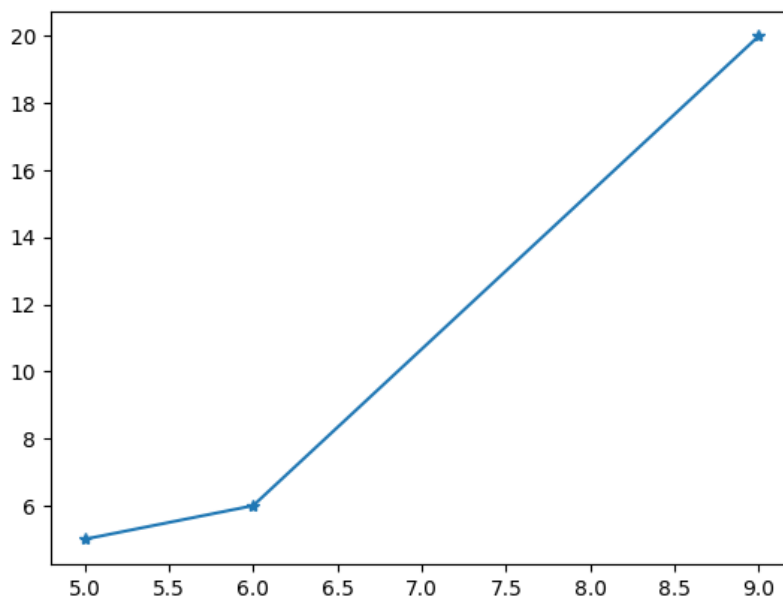
```
x = [5, 6, 9]
```

```
y = [5, 6, 20]
```

```
plt.plot(x, y, marker='*')
```

```
plt.show()
```

Output



Program Name Create a simple line chart and title demo4.py

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

```
x = [5, 6, 9]
```

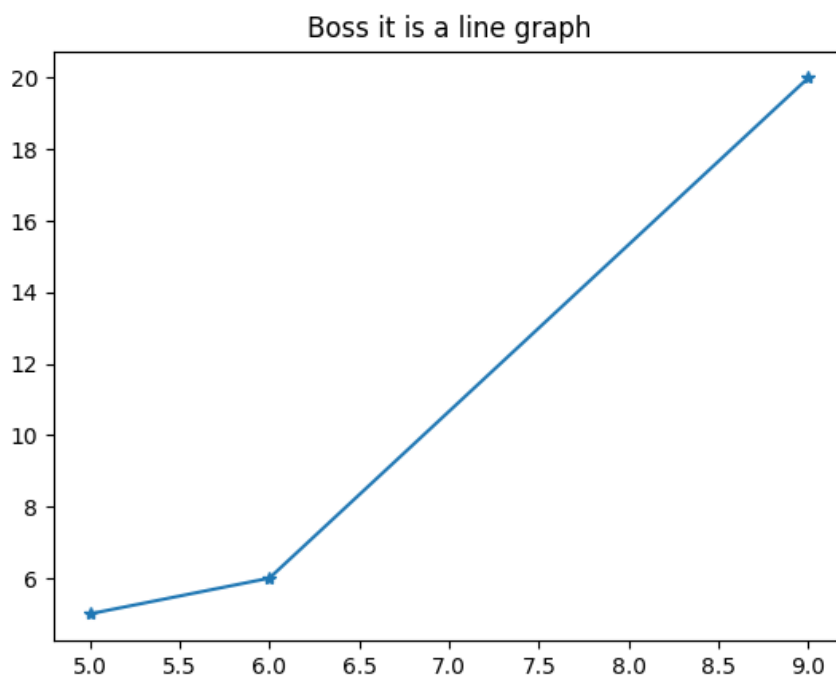
```
y = [5, 6, 20]
```

```
plt.title("Boss it is a line graph")
```

```
plt.plot(x, y, marker='*')
```

```
plt.show()
```

Output



12.1. Labelling the axes

- ✓ We can label **x axis** and **y axis** by using `xlabel` and `ylabel`

Program Name Create a simple line chart and giving title and labelling
demo5.py

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

```
x = [5, 6, 9]
```

```
y = [5, 6, 20]
```

```
plt.title("A line graph")
```

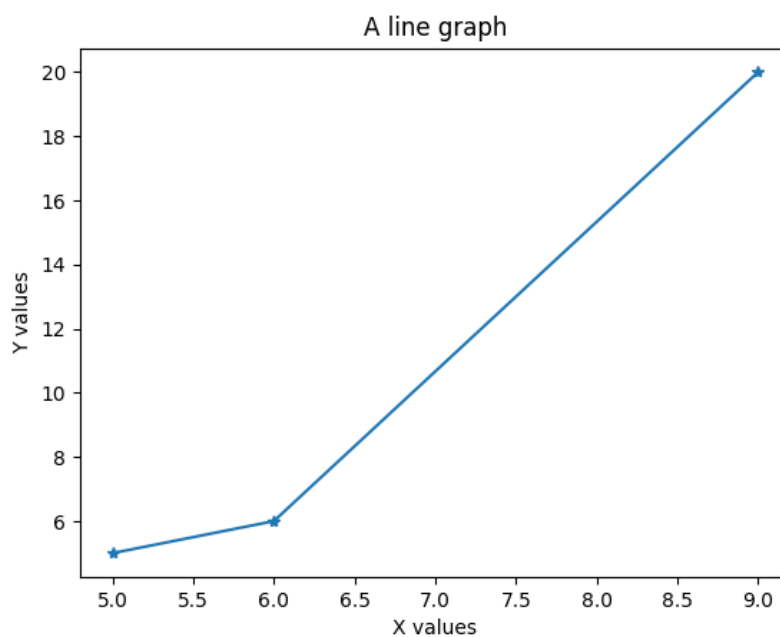
```
plt.xlabel("X values")
```

```
plt.ylabel("Y values")
```

```
plt.plot(x, y, marker = '*')
```

```
plt.show()
```

Output



Program Name Create two lines in single chart
demo6.py

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

```
x = [5, 6, 9]  
y = [5, 6, 20]  
p = [10, 20, 25]
```

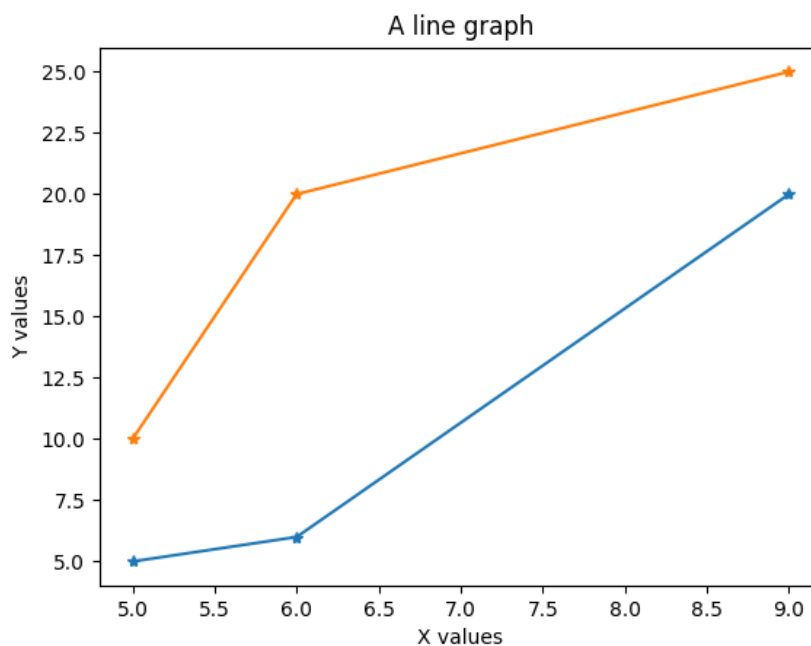
```
plt.title("A line graph")
```

```
plt.xlabel("X values")  
plt.ylabel("Y values")
```

```
plt.plot(x, y, marker = '*')  
plt.plot(x, p, marker = '*')
```

```
plt.show()
```

Output



13. Bar Chart

- ✓ The bar graph is the graphical representation of categorical data.

Program Creating bar chart
Name demo7.py

```
import matplotlib.pyplot as plt

months = ["Jan", "Feb", "Mar", "Apr", "May", "June", "July", "Aug",
"Sep", "Oct", "Nov", "Dec"]
sales = [23, 45, 56, 78, 213, 45, 78, 89, 99, 100, 101, 130]

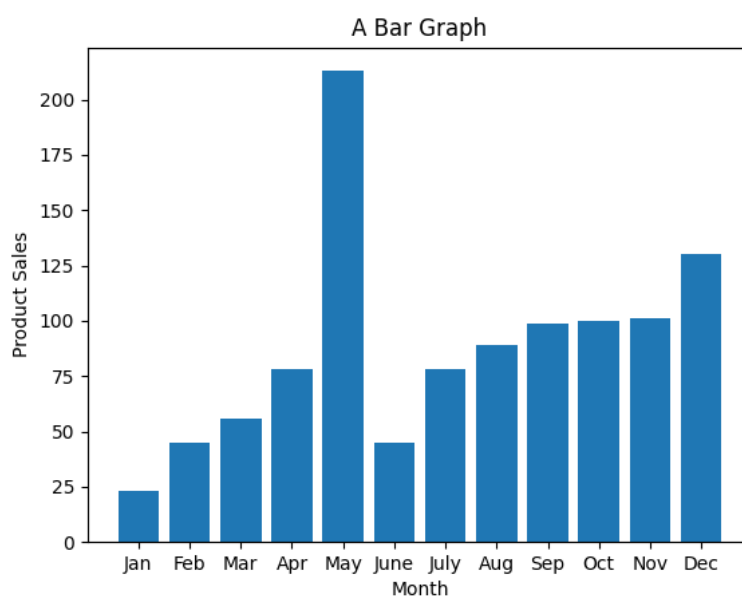
plt.title('A Bar Graph')

plt.xlabel('Month')
plt.ylabel('Product Sales')

plt.bar(months, sales)

plt.show()
```

Output



Program Name Creating horizontal bar chart
demo8.py

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

```
months = ["Jan", "Feb", "Mar", "Apr", "May", "June", "July", "Aug",  
"Sep", "Oct", "Nov", "Dec"]  
sales = [23, 45, 56, 78, 213, 45, 78, 89, 99, 100, 101, 130]
```

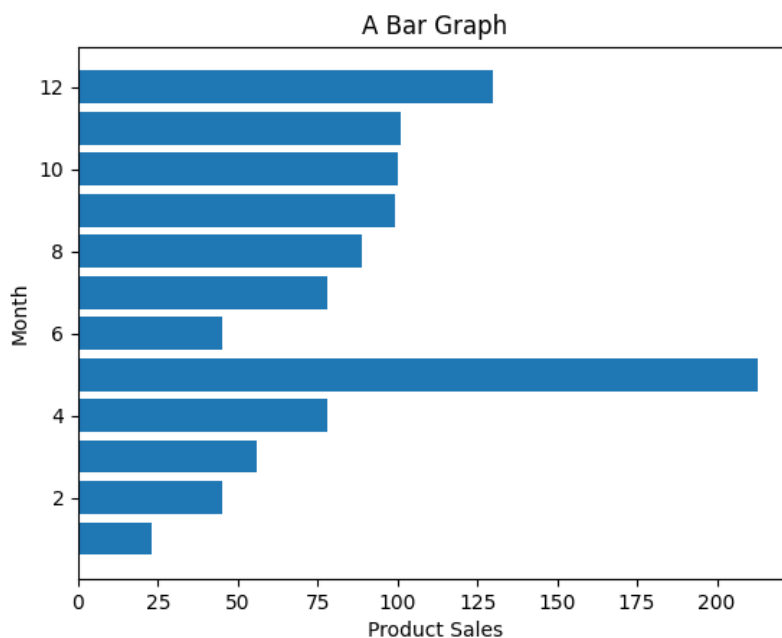
```
plt.title('A Bar Graph')
```

```
plt.xlabel('Product Sales')  
plt.ylabel('Month')
```

```
plt.barh(months, sales)
```

```
plt.show()
```

Output



Program Creating horizontal bar chart
Name demo9.py
File name sales11.csv

```
import matplotlib.pyplot as plt
import pandas as pd

df = pd.read_csv("sales11.csv")

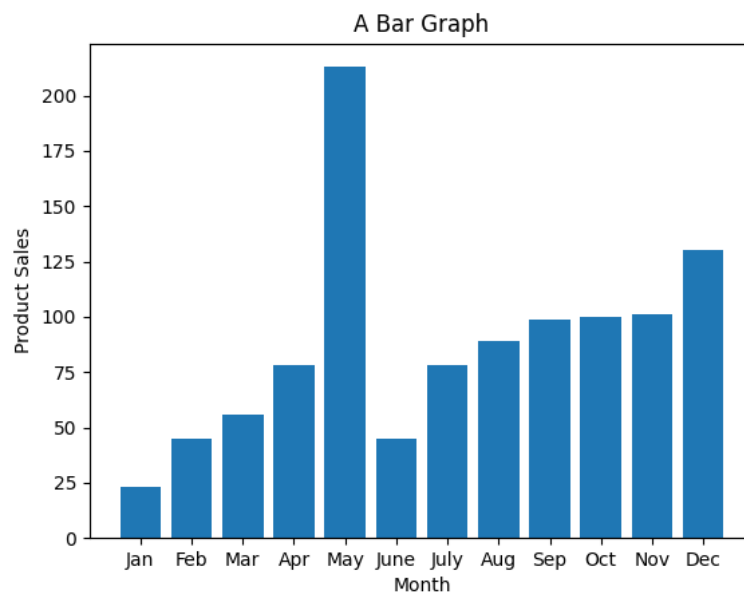
plt.title('A Bar Graph')

plt.xlabel('Month')
plt.ylabel('Product Sales')

plt.bar(df.month, df.sales)

plt.show()
```

Output



Program Name Creating bar chart
demo10.py

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

```
months = ["Jan", "Feb", "Mar", "Apr", "May", "June", "July", "Aug",  
"Sep", "Oct", "Nov", "Dec"]  
sales = [23, 45, 56, 78, 213, 45, 78, 89, 99, 100, 101, 130]
```

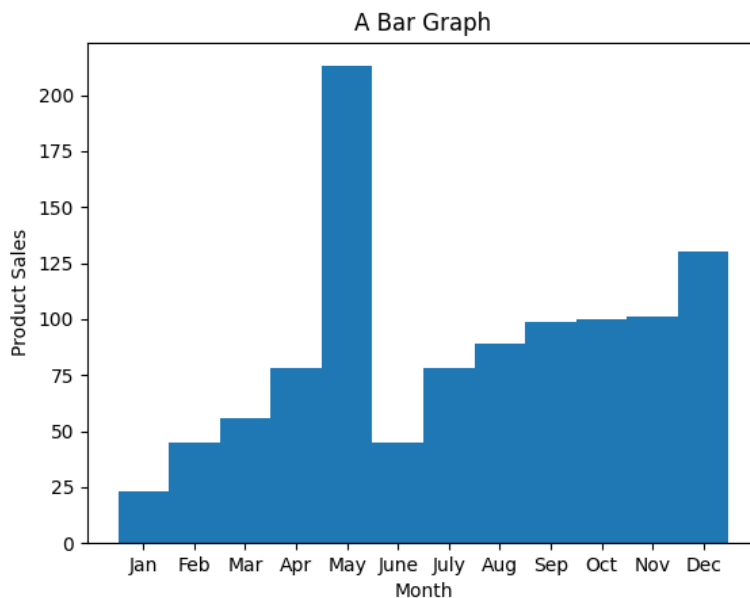
```
plt.title('A Bar Graph')
```

```
plt.xlabel('Month')  
plt.ylabel('Product Sales')
```

```
plt.bar(months, sales, width = 1.0)
```

```
plt.show()
```

Output



14. Histogram

- ✓ A histogram is the graphical representation of quantitative data.
- ✓ This displays the frequency/count of numerical data in bars.

Program Name Creating histogram
demo11.py

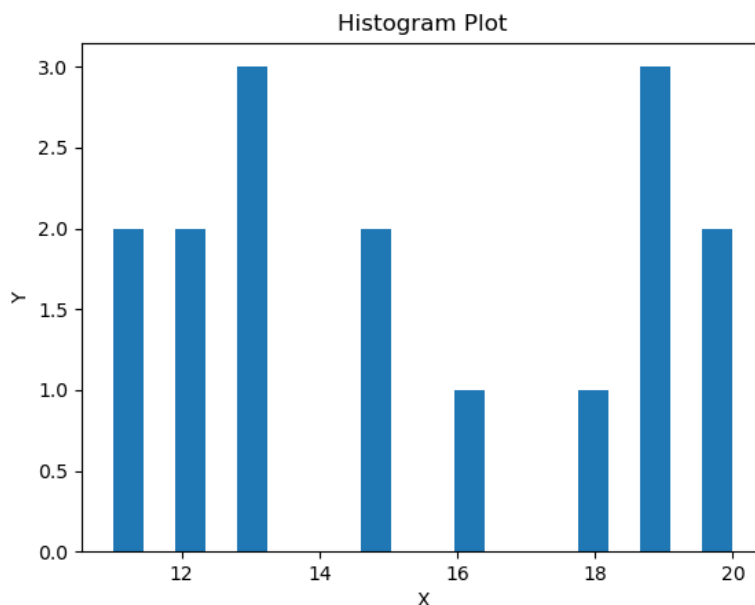
```
import matplotlib.pyplot as plt

data = [12, 15, 13, 20, 19, 20, 11, 19, 11, 12, 19, 13, 15, 16, 18, 13]

plt.xlabel("X")
plt.ylabel("Y")
plt.title("Histogram Plot")

plt.hist(data, bins = 20)
plt.show()
```

Output



15. Pie Chart

- ✓ This is a circular plot that has been divided into slices displaying numerical proportions.
- ✓ Every slice in the pie chart shows the proportion of the element to the whole.
- ✓ A large category means that it will occupy a larger portion of the pie chart.

Program Name Creating pie chart
demo12.py

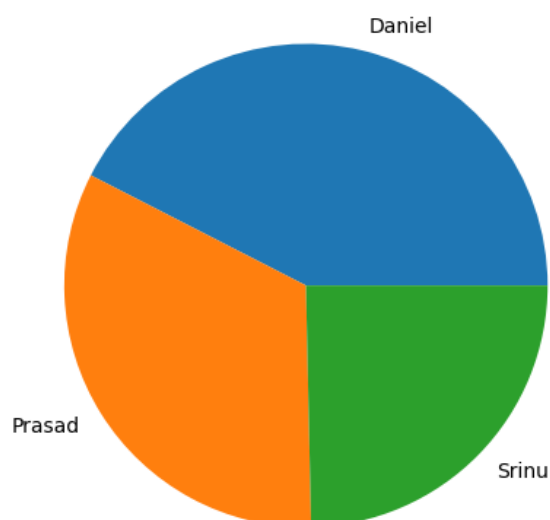
```
import matplotlib.pyplot as plt

students = ['Daniel', 'Prasad', 'Srinu']
points = [62, 48, 36]

plt.pie(points, labels = students)

plt.axis('equal')
plt.show()
```

Output



Program Name Creating pie chart
demo13.py

```
import matplotlib.pyplot as plt

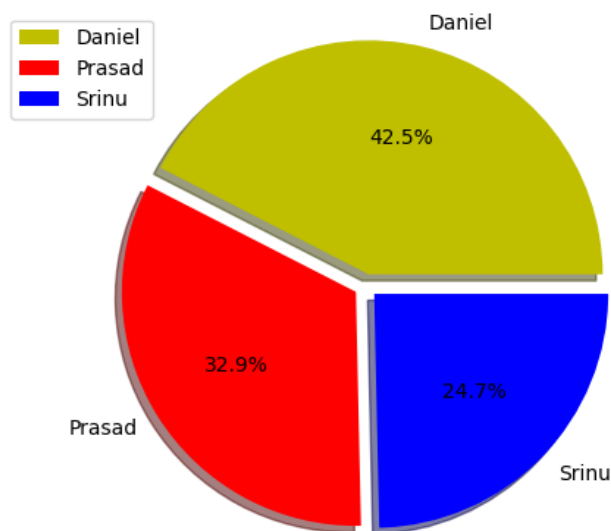
students = ['Daniel', 'Prasad', 'Srinu']
points = [62, 48, 36]

c = ['y', 'r', 'b']

plt.pie(points, labels = students, colors = c , shadow = True,
explode = (0.05, 0.05, 0.05), autopct = '%1.1f%%')

plt.axis('equal')
plt.legend()
plt.show()
```

Output



15.1. Attributes

- ✓ The first parameter to the function is the list of numbers for every category.
 - labels attribute:
 - A list of categories separated by commas is then passed as the argument to labels attribute.
 - colors attribute:
 - To provide the color for every category.
 - To create shadows around the various categories in pie chart.
 - To split each slice of the pie chart into its own.

16. Scatter Plot

- ✓ In scatter plot each value in the data set is represented by a dot.
- ✓ By using this plot we can understand the relationship between two variables.

Program Name Creating Scatter plot
demo14.py

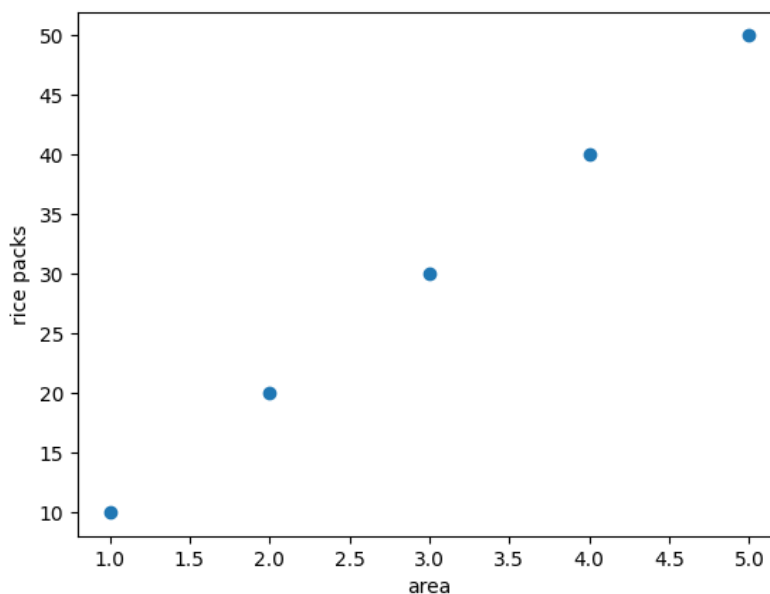
```
import matplotlib.pyplot as plt

area = [1, 2, 3, 4, 5]
rice_packs = [10, 20, 30, 40, 50]

plt.xlabel('area')
plt.ylabel('rice packs')

plt.scatter(area, rice_packs)
plt.show()
```

Output



Program Name Creating Scatter plot
demo15.py

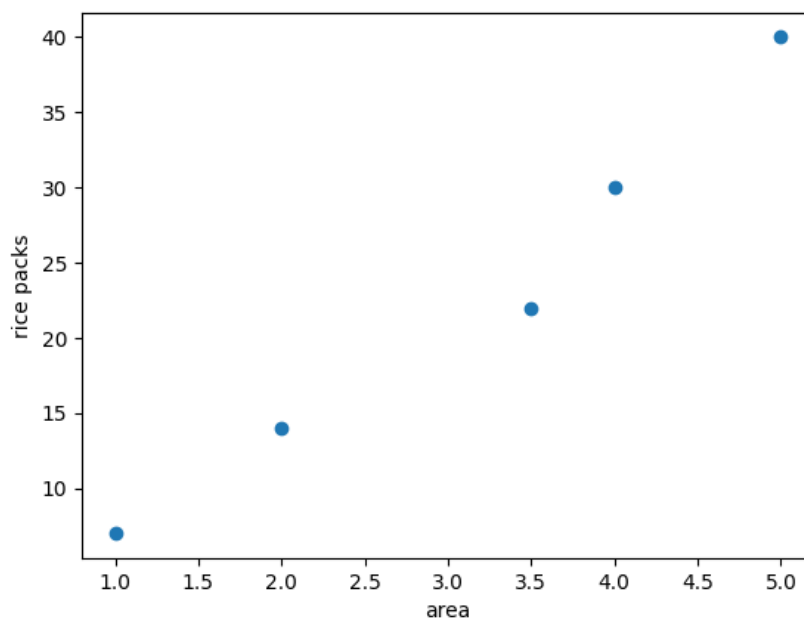
```
import matplotlib.pyplot as plt

area = [1, 2, 3.5, 4, 5]
rice_packs = [7, 14, 22, 30, 40]

plt.xlabel('area')
plt.ylabel('rice packs')

plt.scatter(area, rice_packs)
plt.show()
```

Output



17. Box Plots

- ✓ Box plots help us measure how well data in a dataset is distributed.
- ✓ The graph shows the maximum, minimum, median, first quartile and third quartiles of the dataset.

17.1. Use Box plots

- ✓ Use a boxplot when you need to get the overall statistical information about the data distribution.
- ✓ It is a good tool for detecting outliers in a dataset.

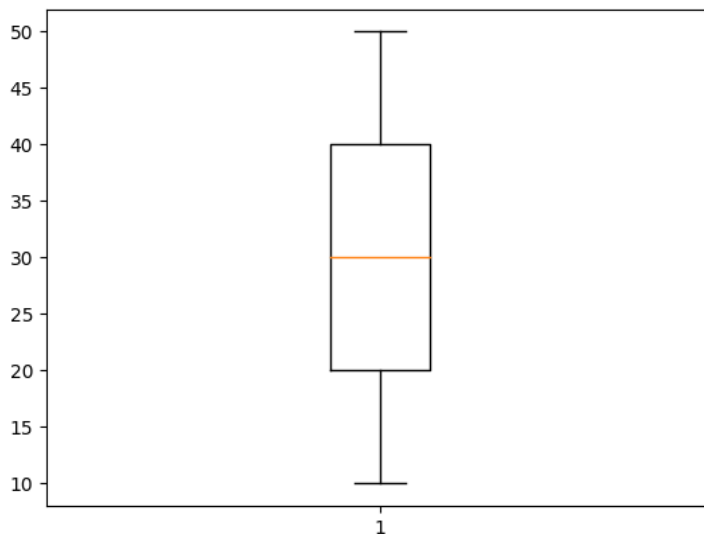
Program Name Creating box plot
demo16.py

```
import matplotlib.pyplot as plt

data = [10, 20, 30, 40, 50]

plt.boxplot(data)
plt.show()
```

Output



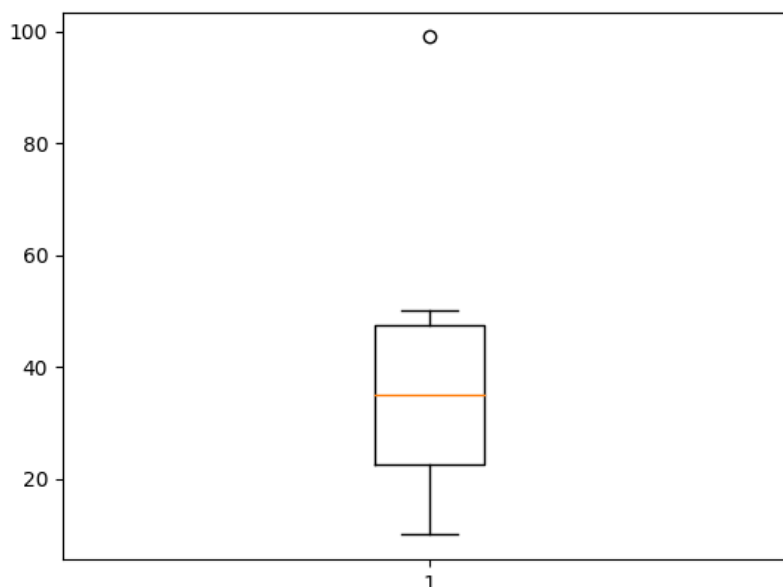
17.2. Box plot explanation

- ✓ The line dividing the box into two shows the median of the data.
- ✓ The end of the box represents the upper quartile (75%) while the start of the box represents the lower quartile (25%).
- ✓ The part between the upper quartile and the lower quartile is known as the Inter Quartile Range (IQR) and helps in approximating 50% of the middle data.

Program Name Creating box plot with outlier
demo17.py

```
import matplotlib.pyplot as plt  
  
data = [10, 20, 30, 40, 50, 90]  
  
plt.boxplot(data)  
plt.show()
```

Output



18. Heatmap

- ✓ A heatmap is a method of data visualization that plots data by replacing numbers with colours.
- ✓ If it is representing with color then it is very easy to understand patterns between different values in the dataset.
- ✓ It is used to visualize data in a two-dimensional format as a coloured map so that different colour variations represent different patterns between features.

18.1. How to understand?

- ✓ A heatmap visualizes the relationship between features as a colour palette.
- ✓ While analysing a heatmap, always remember that **dark shades** represent a **high degree** of linear relationship between features and **light shades** represent a **low degree** of linear relationship between features.

Program Name Creating box plot
demo18.py

```
import matplotlib.pyplot as plt
import pandas as pd

d = {
    "Apple": [10, 20, 30, 40],
    "Orange": [7, 14, 21, 28],
    "Banana": [55, 15, 8, 12],
    "Pear": [15, 14, 1, 8]
}

i = ['Basket1', 'Basket2', 'Basket3', 'Basket4']

df = pd.DataFrame(d, index = i)

plt.imshow(df, cmap = "YlGnBu")
plt.colorbar()

plt.xticks(range(len(df)), df.columns)
plt.yticks(range(len(df)), df.index)

plt.show()
```

Output

