# Visualizing Data

## Lecture 3

Centre for Data Science, ITER
Siksha 'O' Anusandhan (Deemed to be University), Bhubaneswar, Odisha, India.

# Contents

## Introduction

- A fundamental part of the data scientist's toolkit is data visualization.
- Data visualization is the graphical representation of information and data.
- By using visual elements like charts, graphs, and maps, data visualization tools provide an accessible way to see and understand trends, outliers, and patterns in data.
- There are two primary uses for data visualization:
  - To explore data
  - To communicate data

# Matplotlib

- For simple bar charts, line charts, and scatterplots Matplotlib is useful.
- Module **matplotlib.pyplot** of matplotlib is used for visualization of data.
- **pyplot** maintains an internal state in which you build up a visualization **step by step.**
- Once you're done, you can save it with **savefig** or display it with **show**.
- A simple line plot is showen in figure(1).

# Example

- Code for a simple line plot of year v/s gdp

```
1  from matplotlib import pyplot as plt
2  years = [1950, 1960, 1970, 1980, 1990, 2000, 2010]
3  gdp = [300.2, 543.3, 1075.9, 2862.5, 5979.6, 10289.7,
       14958.3]
4  '''create a line chart, years on x-axis, gdp on y-axis'''
5  plt.plot(years, gdp, color='green', marker='o', linestyle='
       solid')
6  '''add a title'''
7  plt.title("Nominal GDP")
8  '''add a label to the y-axis'''
9  plt.ylabel("Billions of $")
10 plt.xlabel("Year")
11 plt.show()
12 plt.savefig()
```
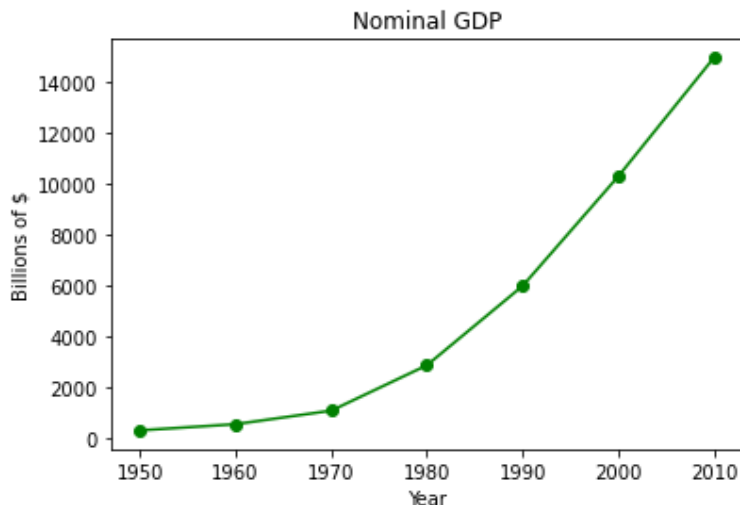
# Simple Line Plot

- Output of above code.



Figure 1: A Simple Line Plot of year v/s gdp

# Bar Charts

- **A bar chart** is a good choice when you want to show how some quantity varies among some discrete set of items.
- Academy Awards were won by each of a variety of movies is shown in figure(2) and code is given below.

```python
1  from matplotlib import pyplot as plt
2  movies = ["Annie Hall", "Ben-Hur", "Casablanca", "Gandhi",
       "West Side Story"]
3  num_oscars = [5, 11, 3, 8, 10]
4  ''' plot bars with left x-coordinates [0, 1, 2, 3, 4],
       heights [num_oscars] '''
5  plt.bar(range(len(movies)), num_oscars)
6  plt.title("My Favorite Movies")
7  plt.ylabel("# of Academy Awards")
8  ''' label x-axis with movie names at bar centers '''
9  plt.xticks(range(len(movies)), movies)
10 plt.show()
```
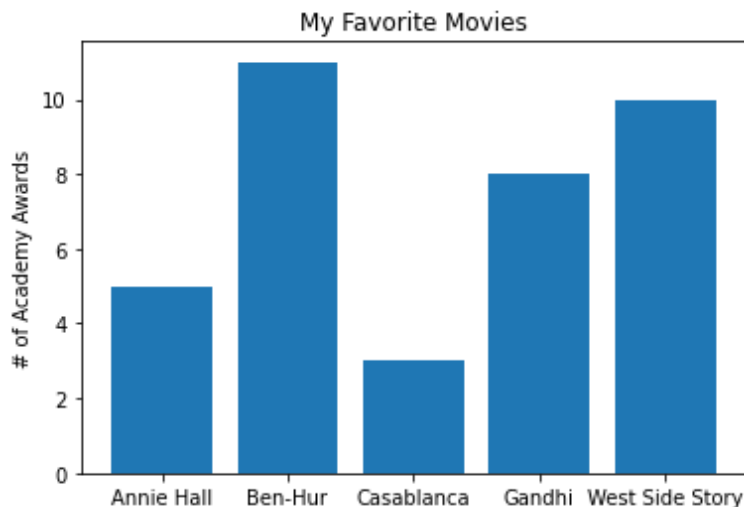
# Bar Charts Cont...

- Output of above code.



Figure 2: A simple bar chart

## Bar Charts Cont...

- A bar chart can also be a good choice for plotting **histograms** of bucketed numeric values as shown in figure(3) and code is given below.
- The third argument to **plt.bar** specifies the bar width.
- The call to **plt.axis** indicates that we want the x-axis to range from –5 to 105, and that the y-axis should range from 0 to 5.
- And the call to **plt.xticks** puts x-axis labels at 0, 10, 20, . . . , 100.

# Bar Charts Cont...

- Code for histogram using **plt.bar** .

```
1  from collections import Counter
2  from matplotlib import pyplot as plt
3  grades=[83, 95, 91, 87, 70, 0, 85, 82, 100, 67, 73, 77, 0]
4  '''Bucket grades by decile, but put 100 in with the 90s'''
5  histogram = Counter(min(grade // 10 * 10, 90) for grade in
      grades)
6  plt.bar([x + 5 for x in histogram.keys()], '''Shift bars
      right by 5'''
7  histogram.values(), '''Give each bar its correct height '''
8  10, '''Give each bar a width of 10 '''
9  edgecolor=(0, 0, 0)) '''Black edges for each bar '''
10 plt.axis([-5, 105, 0, 5]) '''x-axis from -5 to 105,'''
11 plt.xticks([10 * i for i in range(11)]) '''x-axis labels at
       0, 10,...,100'''
12 plt.xlabel("Decile")
13 plt.ylabel("# of Students")
14 plt.title("Distribution of Exam 1 Grades")
15 plt.show()
```

# Bar Charts Cont...
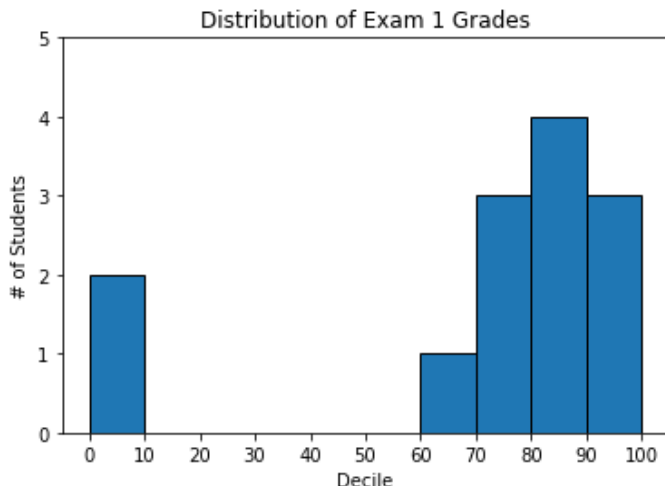
- Output of above code.



Figure 3: Using a bar chart for a histogram

# Line Charts

- **Line charts** can be obtained using **plt.plot**.
- These are a good choice for showing trends.

# Line Charts Cont...

- Several line charts with a legend.

```
1  from matplotlib import pyplot as plt
2  variance = [1, 2, 4, 8, 16, 32, 64, 128, 256]
3  bias_squared = [256, 128, 64, 32, 16, 8, 4, 2, 1]
4  total_error=[x + y for x, y in zip(variance, bias_squared)]
5  xs = [i for i, _ in enumerate(variance)]
6  ''' We can make multiple calls to plt.plot to show multiple
       series on the same chart '''
7  plt.plot(xs, variance, 'g-', label='variance') ''' green
       solid line '''
8  plt.plot(xs, bias_squared, 'r-.', label='bias^2') ''' red
       dot-dashed line '''
9  plt.plot(xs, total_error, 'b:', label='total error') '''
       blue dotted line '''
10 ''' Because we have assigned labels to each series, we can
       get a legend for free (loc=9 means "top center") '''
11 plt.legend(loc=9)
12 plt.xlabel("model complexity")
13 plt.xticks([])
14 plt.title("The Bias-Variance Tradeoff")
15 plt.show()
```
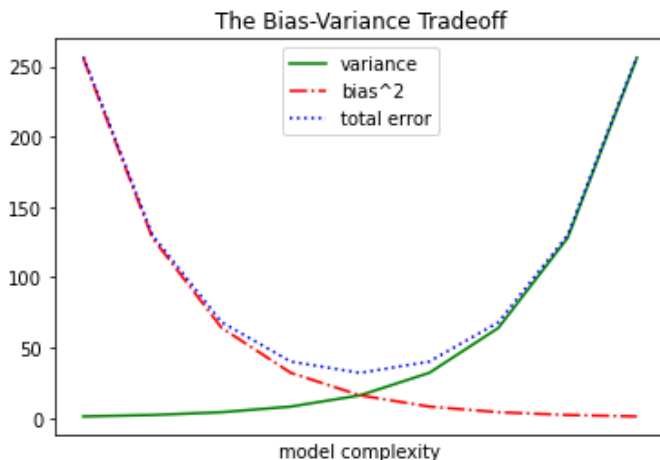
# Line Charts Cont...

- Output of above code.



Figure 4: Several line charts with a legend

# Scatter Plots

- **A scatterplot** is the right choice for visualizing the relationship between two paired sets of data.
- The relationship between the number of friends your users have and the number of minutes they spend on the site every day can be obtained using given code and its visualization is shown in figure(5).

# Scatter Plots Cont...

- Code

```
1  from matplotlib import pyplot as plt
2  friends = [ 70, 65, 72, 63, 71, 64, 60, 64, 67]
3  minutes = [175, 170, 205, 120, 220, 130, 105, 145, 190]
4  labels = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i']
5  plt.scatter(friends, minutes)
6  for label, friend_count, minute_count in zip(labels,
       friends, minutes):
7      plt.annotate(label,
8      xy=(friend_count, minute_count), ''' Put the label with
        its point '''
9      xytext=(5, -5), ''' but slightly offset '''
10     textcoords='offset points')
11 plt.title("Daily Minutes vs. Number of Friends")
12 plt.xlabel("# of friends")
13 plt.ylabel("daily minutes spent on the site")
14 plt.show()
```
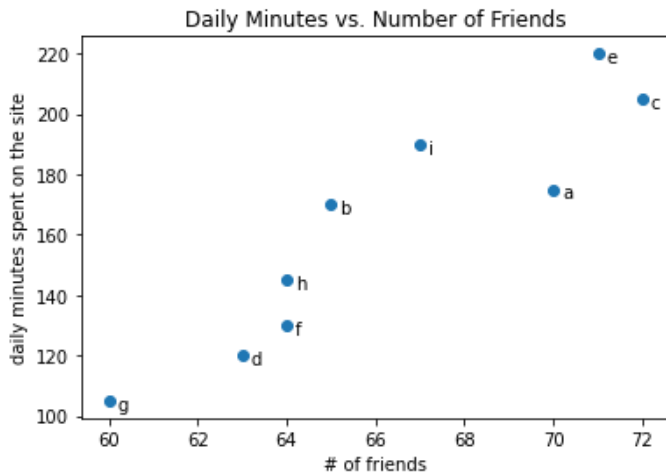
# Scatter Plot Cont...

- Output of above code.



Figure 5: A scatterplot of friends and time on the site

# References

[1] Data Science from Scratch Joel Grus, Shroff/O'reilly, Second Edition

Thank You
Any Questions?