



# Class 9 - Intro to Data Science

MISE Summer Programming Camp 2023



# Motivation



## Why Data Science?

- Data is continuously being collected all the time, both about individuals and collective entities
- This is extremely helpful for making informed decisions about the world
- Data science offers the tools to extract this information from data



# What is Data?

Data is any information that can be collected, stored, and analyzed

- Temperature measurements
- Age of a collection of patients
- Price of an item over time
- GPS coordinates

A list of **floats**

- Subscription status of users
- Binary classification (loan/no loan)

A list of **booleans**

- Product ratings
- Website traffic
- Product inventory
- Exam scores

A list of **ints**



# What is Data Science

- Data science is an interdisciplinary field that combines **statistics, mathematics, programming**, and domain expertise to extract knowledge and insights from data
- It involves various processes, including **data collection, cleaning, analysis**, and visualization.
- Data scientists use techniques such as **machine learning, data mining**, and **statistical modeling**
- The goal of data science is to derive **actionable insights** and **drive data-informed decision-making**



# How does Data Science work?

- **Data Collection:** Gathering relevant and reliable data from various sources, such as databases, APIs, surveys, and sensors.
- **Data Cleaning and Preprocessing:** Ensuring the data is accurate, consistent, and free from errors or missing values. This step involves data validation, transformation, and handling outliers.
- **Exploratory Data Analysis (EDA):** Understanding the data through visualizations, descriptive statistics, and uncovering patterns or relationships.
- **Machine Learning:** Utilizing algorithms to build models that can learn from the data and make predictions or classifications.
- **Data Visualization:** Presenting data and insights in a visual format, such as charts, graphs, and interactive dashboards, to facilitate understanding and communication.



## Where is Data Science used?

- **Finance:** Fraud detection, risk assessment, algorithmic trading.
- **Healthcare:** Disease prediction, personalized medicine, electronic health records analysis.
- **Marketing:** Customer segmentation, targeted advertising, sentiment analysis.
- **Transportation:** Route optimization, demand forecasting, autonomous vehicles.
- **Retail:** Recommender systems, inventory management, pricing optimization.
- **Manufacturing:** Predictive maintenance, quality control, supply chain optimization.

---

# Visualizing Data



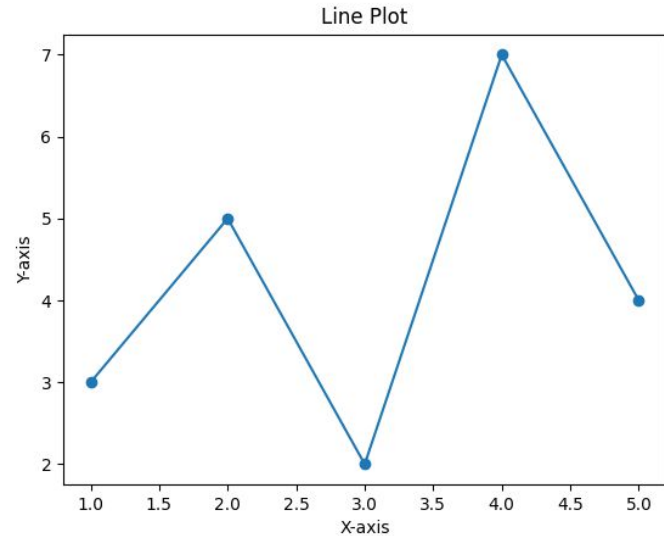




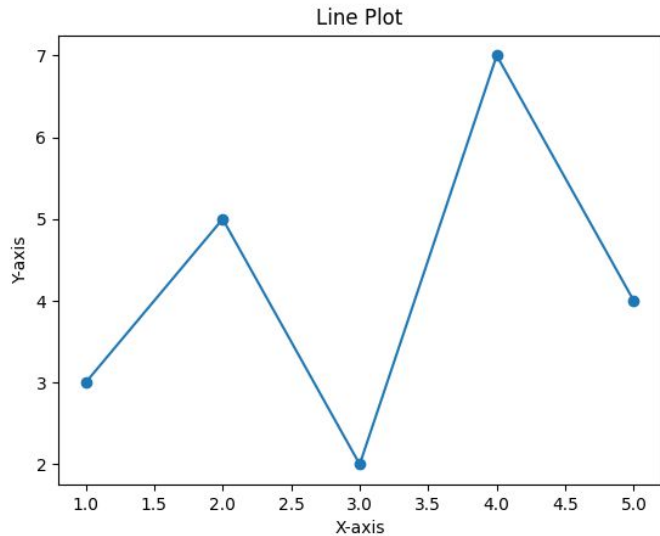
# In Python: Matplotlib

- Matplotlib is a widely used data visualization **library** in Python
- It provides a comprehensive set of tools for creating various types of plots, charts, and graphs

Note: we won't go into how to install any libraries today

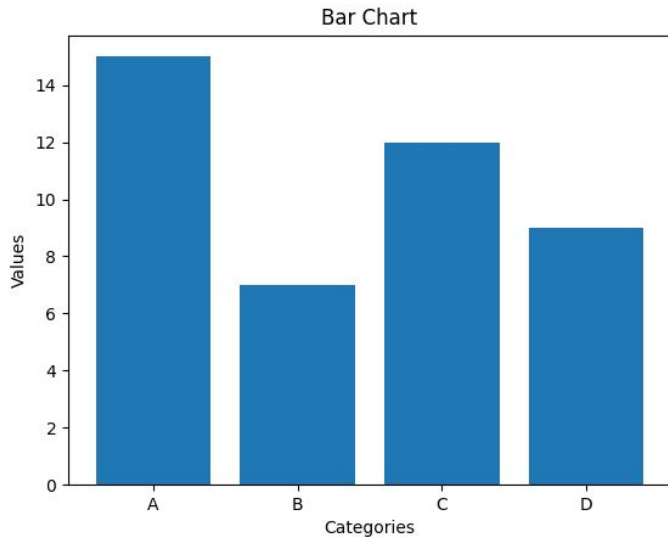


## Example: Line Plot



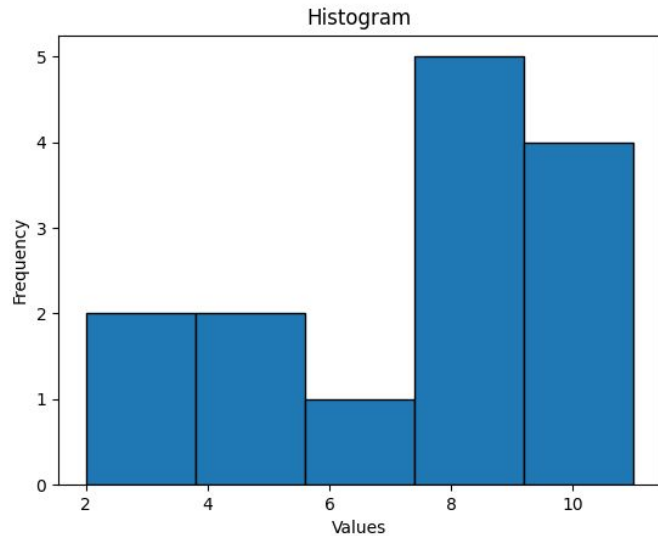
```
1 import matplotlib.pyplot as plt
2
3 # Data
4 x = [1, 2, 3, 4, 5]
5 y = [3, 5, 2, 7, 4]
6
7 # Create a line plot
8 plt.plot(x, y, marker='o')
9
10 # Customize labels and title
11 plt.xlabel('X-axis')
12 plt.ylabel('Y-axis')
13 plt.title('Line Plot')
14
15 # Display the plot
16 plt.show()
```

# Example: Bar Chart



```
1 import matplotlib.pyplot as plt
2
3 # Data
4 categories = ['A', 'B', 'C', 'D']
5 values = [15, 7, 12, 9]
6
7 # Create a bar chart
8 plt.bar(categories, values)
9
10 # Customize labels and title
11 plt.xlabel('Categories')
12 plt.ylabel('Values')
13 plt.title('Bar Chart')
14
15 # Display the plot
16 plt.show()
```

# Example: Bar Chart



```
1 import matplotlib.pyplot as plt
2
3 # Data
4 values = [2, 3, 5, 5, 7, 8, 8, 8, 9, 9, 10, 11, 11, 11]
5
6 # Create a histogram
7 plt.hist(values, bins=5, edgecolor='black')
8
9 # Customize labels and title
10 plt.xlabel('Values')
11 plt.ylabel('Frequency')
12 plt.title('Histogram')
13
14 # Display the plot
15 plt.show()
```

---

# Analyzing Data

# The numpy library

- NumPy (Numerical Python) is a library for scientific computing
- Powerful tools for working with arrays and performing mathematical operations efficiently

**NumPy** 

## Quantum Computing



QuTIP  
PyQuil  
Qiskit  
PennyLane

## Statistical Computing



Pandas  
statsmodels  
Xarray  
Seaborn

## Signal Processing



SciPy  
PyWavelets  
python-control

## Image Processing



Scikit-image  
OpenCV  
Mahotas

## Graphs and Networks



NetworkX  
graph-tool  
igraph  
PyGSP

## Astronomy Processes



AstroPy  
SunPy  
SpacePy

## Cognitive Psychology



PsychoPy

## Bioinformatics



BioPython  
Scikit-Bio  
PyEnsembl  
ETE

## Bayesian Inference



PyStan  
PyMC3  
ArviZ  
emcee

## Mathematical Analysis



SciPy  
SymPy  
cvxpy  
FEniCS

## Chemistry



Cantera  
MDAnalysis  
RDKit  
PyBaMM

## Geoscience



Pangeo  
Simpeg  
ObsPy  
Fatiando a Terra

## Geographic Processing



Shapely  
GeoPandas  
Folium

## Architecture & Engineering



COMPAS  
City Energy Analyst  
Sverchok



# Basic Data Analysis

- **Descriptive Statistics:** statistical measures such as mean, median, standard deviation, variance, and percentiles.
- **Array Operations:** sorting, filtering, searching, and element-wise computations.
- **Linear Algebra:** matrix multiplication, decomposition, and solving linear equations.

```
1 import numpy as np
2
3 # Data
4 dataset = np.array([1, 2, 3, 4, 5, 6, 7, 8, 9])
5
6 # Calculate descriptive statistics
7 mean = np.mean(dataset)
8 median = np.median(dataset)
9 std_dev = np.std(dataset)
10 variance = np.var(dataset)
11
12 # Perform array operations
13 sorted_data = np.sort(dataset)
14 filtered_data = dataset[dataset > 5]
15 search_index = np.where(dataset == 3)
16
17 # Display the analysis results
18 print("Mean:", mean)
19 print("Median:", median)
20 print("Standard Deviation:", std_dev)
21 print("Variance:", variance)
22 print("Sorted Data:", sorted_data)
23 print("Filtered Data:", filtered_data)
24 print("Index of 3:", search_index)
```





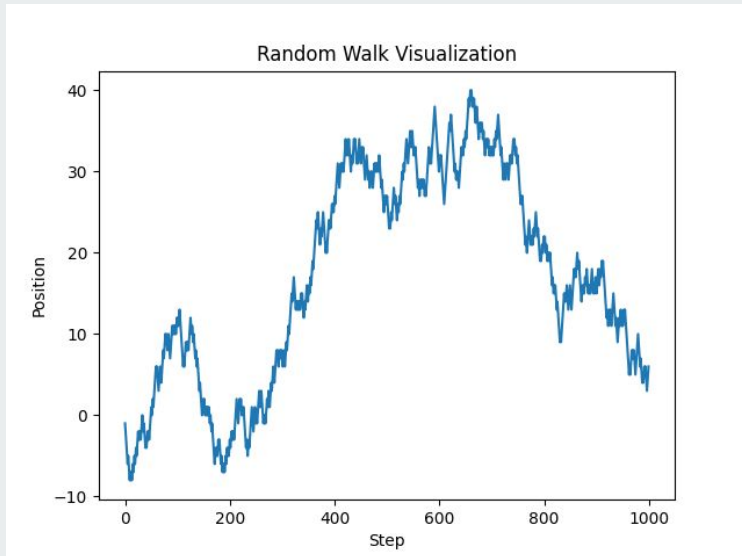
# Basic Data Synthesis

- **Random Number Generation:** generate random numbers from different distributions
- **Synthetic Dataset Creation:** create datasets with specific patterns, or characteristics.
- **Data Reshaping:** transforming data arrays to match specific dimensions

```
1 import numpy as np
2
3 # Generate random numbers from a uniform distribution
4 uniform_numbers = np.random.random(size=10)
5
6 # Generate random integers
7 integer_numbers = np.random.randint(low=1, high=100, size=5)
8
9 # Generate random numbers from a normal distribution
10 mean = 0
11 std_dev = 1
12 normal_numbers = np.random.normal(loc=mean, scale=std_dev, size=10)
13
14 # Display the generated random numbers
15 print("Uniform Numbers:", uniform_numbers)
16 print("Integer Numbers:", integer_numbers)
17 print("Normal Numbers:", normal_numbers)
```



# Example: Random Walk



```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # Set the number of steps
5 num_steps = 1000
6
7 # Generate random steps
8 steps = np.random.choice([-1, 1], size=num_steps)
9
10 # Calculate the accumulated position
11 position = np.cumsum(steps)
12
13 # Plot the random walk
14 plt.plot(position)
15
16 # Customize the plot
17 plt.title('Random Walk Visualization')
18 plt.xlabel('Step')
19 plt.ylabel('Position')
20
21 # Display the plot
22 plt.show()
```

---

# Predicting with Data

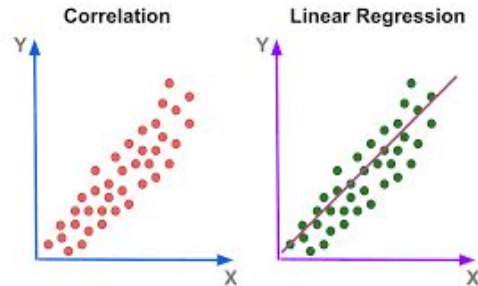


# Prediction Models

- Mathematical algorithms that learn patterns from existing data to make predictions
- **Classification:** label data into discrete types, e.g. “Is this a fraudulent transaction?”
- **Clustering:** Group data together by common attributes, e.g. “Which of these people have similar health conditions?”
- **Regression:** estimating the relationships between dependent variables, e.g. “How much should this house cost?”

# The simplest model: Linear Regression

- Model linear relationship between variable (target) and independent variables (features)
- Assumes a linear equation of the form:  $y = mx + b$ , where  $y$  is the target variable,  $x$  is the feature,  $m$  is the slope, and  $b$  is the intercept
- The goal of is to find the best-fitting line that minimizes the difference between the predicted values and the actual values

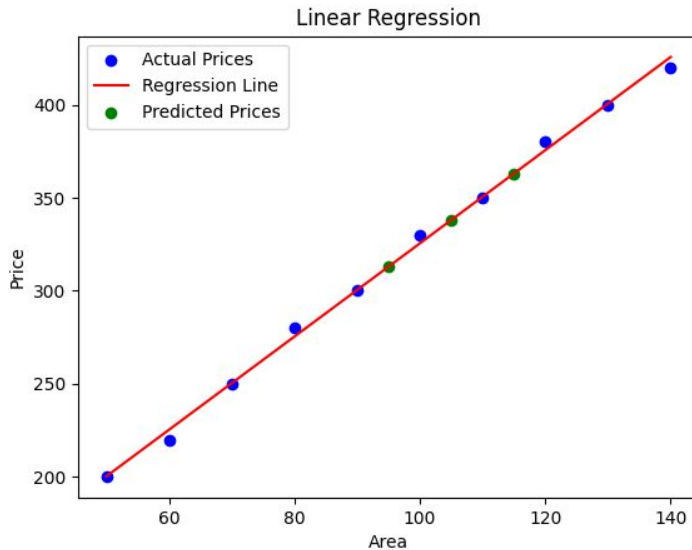




## Example: Linear Regression Equation

- Consider a simple example of predicting house prices based on the area of the house
- Model the relationship between house area ( $x$ ) and house price ( $y$ ):  $y = mx + b$
- The slope ( $m$ ) represents the change in house price for every unit increase in area, and the intercept ( $b$ ) is the house price when the area is zero

# Linear Regression in Python



```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.linear_model import LinearRegression
4
5 # Sample data
6 area = np.array([50, 60, 70, 80, 90, 100, 110, 120, 130, 140])
7 price = np.array([200, 220, 250, 280, 300, 330, 350, 380, 400, 420])
8
9 # Reshape the data
10 area = area.reshape((-1, 1))
11
12 # Create a linear regression model
13 model = LinearRegression()
14
15 # Fit the model to the data
16 model.fit(area, price)
17
18 # Predict house prices for new areas
19 new_areas = np.array([95, 105, 115]).reshape((-1, 1))
20 predicted_prices = model.predict(new_areas)
21
22 # Plot the data points and the regression line
23 plt.scatter(area, price, color='blue', label='Actual Prices')
24 plt.plot(area, model.predict(area), color='red', label='Regression Line')
25
26 # Plot the predicted prices for new areas
27 plt.scatter(new_areas, predicted_prices, color='green', label='Predicted Prices')
28
29 # Customize the plot
30 plt.xlabel('Area')
31 plt.ylabel('Price')
32 plt.title('Linear Regression')
33 plt.legend()
34
35 # Display the plot
36 plt.show()
```

---

# Case Study: Temperature Forecast





# Forecasting temperature

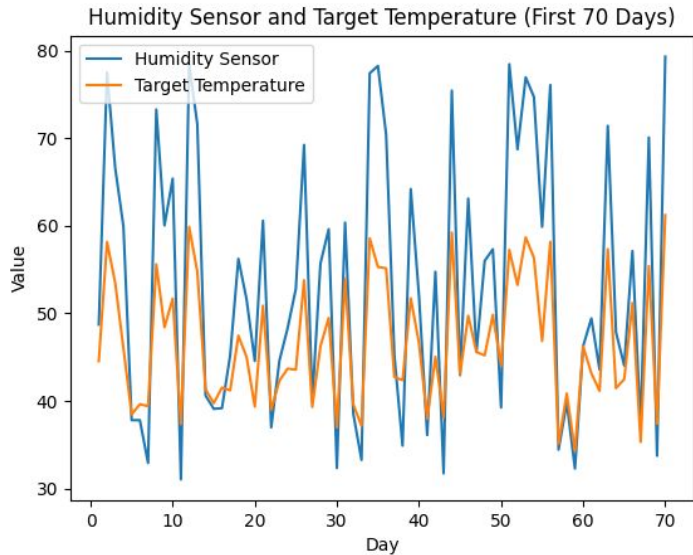
1. **Introduction:**
  - **Objective:** Predicting temperature based on sensor measurements
  - **Dataset:** Generated data simulating temperature and humidity sensor readings
  - **Methodology:** Utilizing linear regression for temperature prediction
2. **Data Generation:**
  - **Humidity sensor data:** Randomly generated values for 100 days
  - **Target temperature:** Linearly dependent on humidity with Gaussian error
  - **Data Exploration:** Displaying the first 70 days of data
3. **Model Training:**
  - **Linear regression:** Training the model using the first 70 days of data
  - Learning the relationship between humidity and temperature
4. **Temperature Prediction:**
  - Applying the trained model to predict temperatures for the last 30 days
  - Comparing the predicted values with the true temperatures
5. **Visualization:**
  - Plotting the true temperature and predicted values for the last 30 days
  - Assessing the accuracy of the temperature predictions



# Data Generation

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.linear_model import LinearRegression
4
5 # Generate random sensor data
6 np.random.seed(42)
7 num_days = 100
8 humidity_sensor = np.random.uniform(30, 80, num_days)
9 error = np.random.normal(0, 2, num_days) # Gaussian error with mean 0 and standard
    deviation 2
10
11 # Generate target temperature data
12 slope = 0.5
13 intercept = 20
14 target_temperature = slope * humidity_sensor + intercept + error
```

# Visualizing the Data



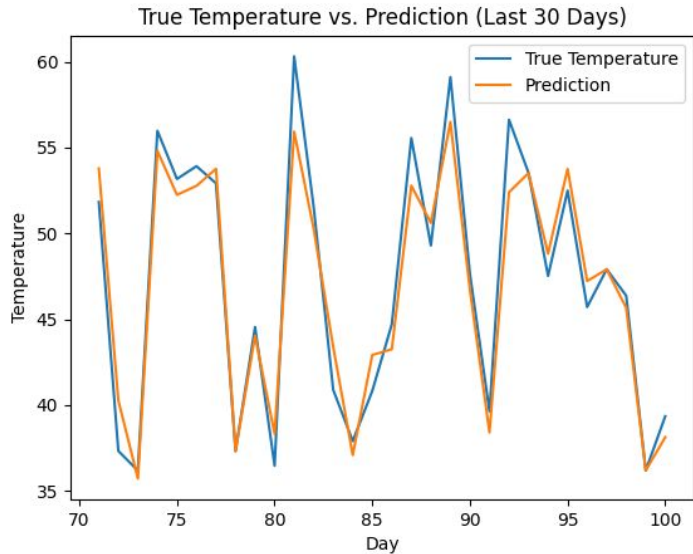
```
16 # Display first 70 days of data
17 days = np.arange(1, 71)
18 plt.plot(days, humidity_sensor[:70], label='Humidity Sensor')
19 plt.plot(days, target_temperature[:70], label='Target Temperature')
20 plt.xlabel('Day')
21 plt.ylabel('Value')
22 plt.title('Humidity Sensor and Target Temperature (First 70 Days)')
23 plt.legend()
24 plt.show()
```



## Model Training

```
26 # Data Analysis and Forecasting
27 X_train = humidity_sensor[:70].reshape(-1, 1)
28 y_train = target_temperature[:70]
29
30 model = LinearRegression()
31 model.fit(X_train, y_train)
32
33 # Predict temperature for the last 30 days
34 X_test = humidity_sensor[70:].reshape(-1, 1)
35 y_test = target_temperature[70:]
36 predicted_temperature = model.predict(X_test)
```

# Visualizing the Results



```
38 # Visualization of true temperature vs. prediction for the last 30 days
39 days = np.arange(71, 101)
40 plt.plot(days, y_test, label='True Temperature')
41 plt.plot(days, predicted_temperature, label='Prediction')
42 plt.xlabel('Day')
43 plt.ylabel('Temperature')
44 plt.title('True Temperature vs. Prediction (Last 30 Days)')
45 plt.legend()
46 plt.show()
47
48 print("Predicted Temperature:", predicted_temperature)
```

---

**Some parting thoughts**

**This is the end of the program!**

**Thank you all for participating!**

