

Intro to Data Science and ML

We acknowledge that all SFU Surge initiatives, including the opportunity to host this workshop today, takes place on x^wməθk^wəyəm (Musqueam), səlilwətaʔ (Tsleil-Waututh), and Sk̓wx̓wú7mesh (Squamish) nations.

As build the tech community of tomorrow, its important to understand and respect indigenous histories!



Access This File.

It is also available online in the
StormHacks discord Server!

Hi 🖐️🖐️

I'm Matt!

👤 I've judged at over 10 hackathons

🏆 Won over 20 hackathons globally

💉 Worked in the healthcare industry

🚀 Organized 6 hackathons



 @ermergesh

 /in/MatthewWong1129

What we'll cover

Introduction to Data Science

What it is and why it matters

Data Types & Structure

Understanding the different kinds of data

Data Cleaning & Preprocessing

Making messy data ready for use

Exploratory Data Analysis

Finding insights through visualization

Machine Learning Basics

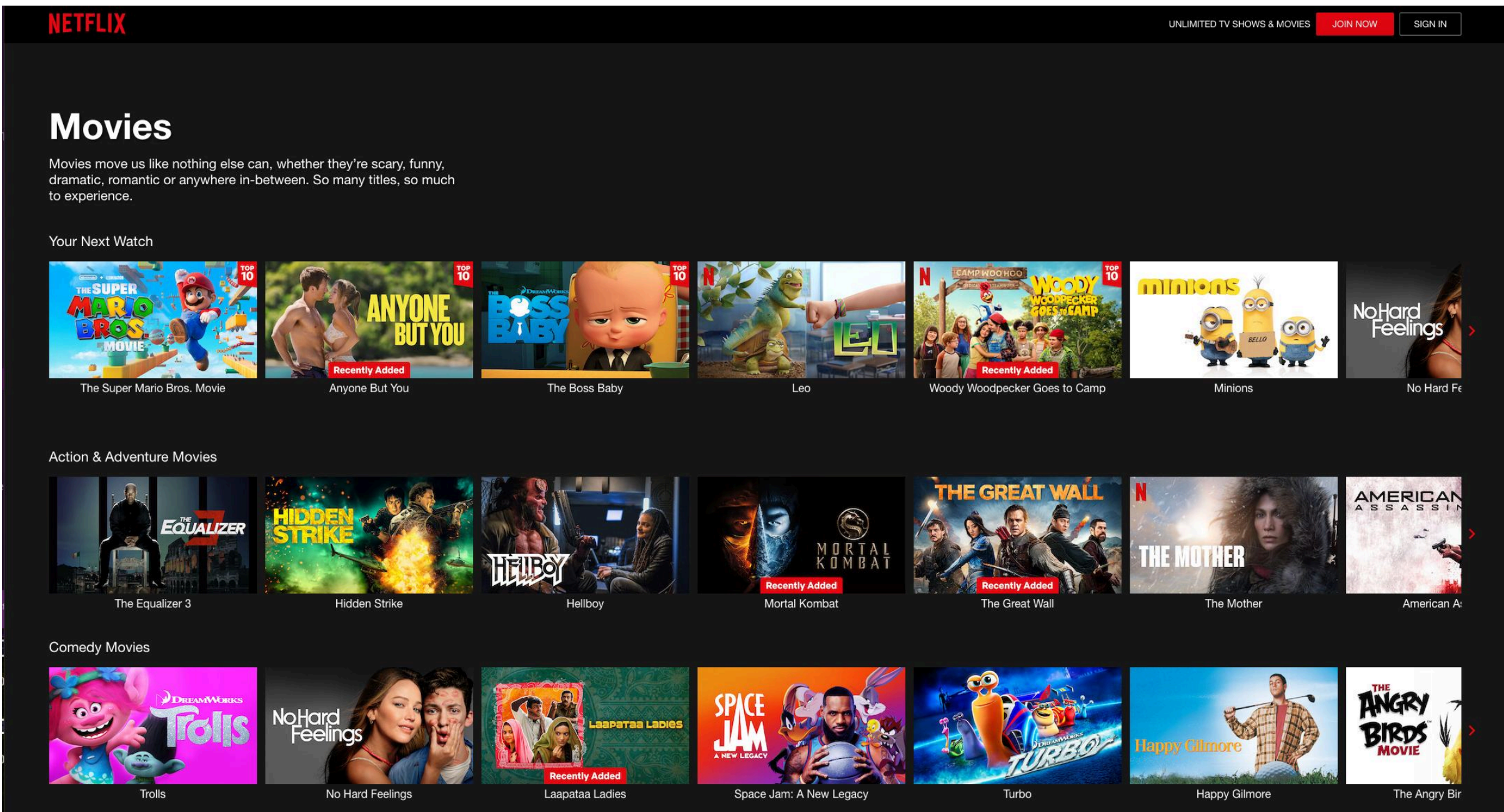
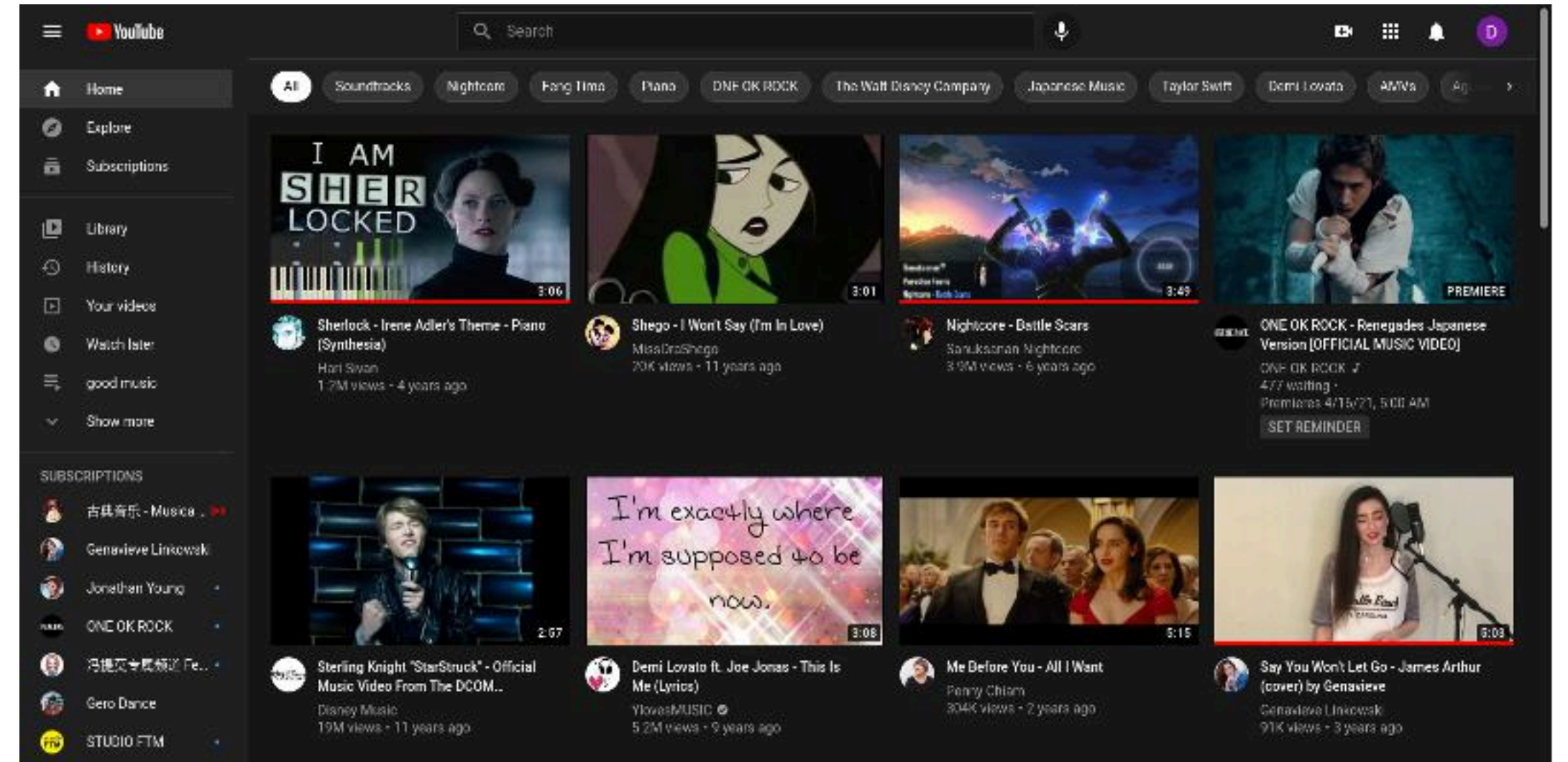
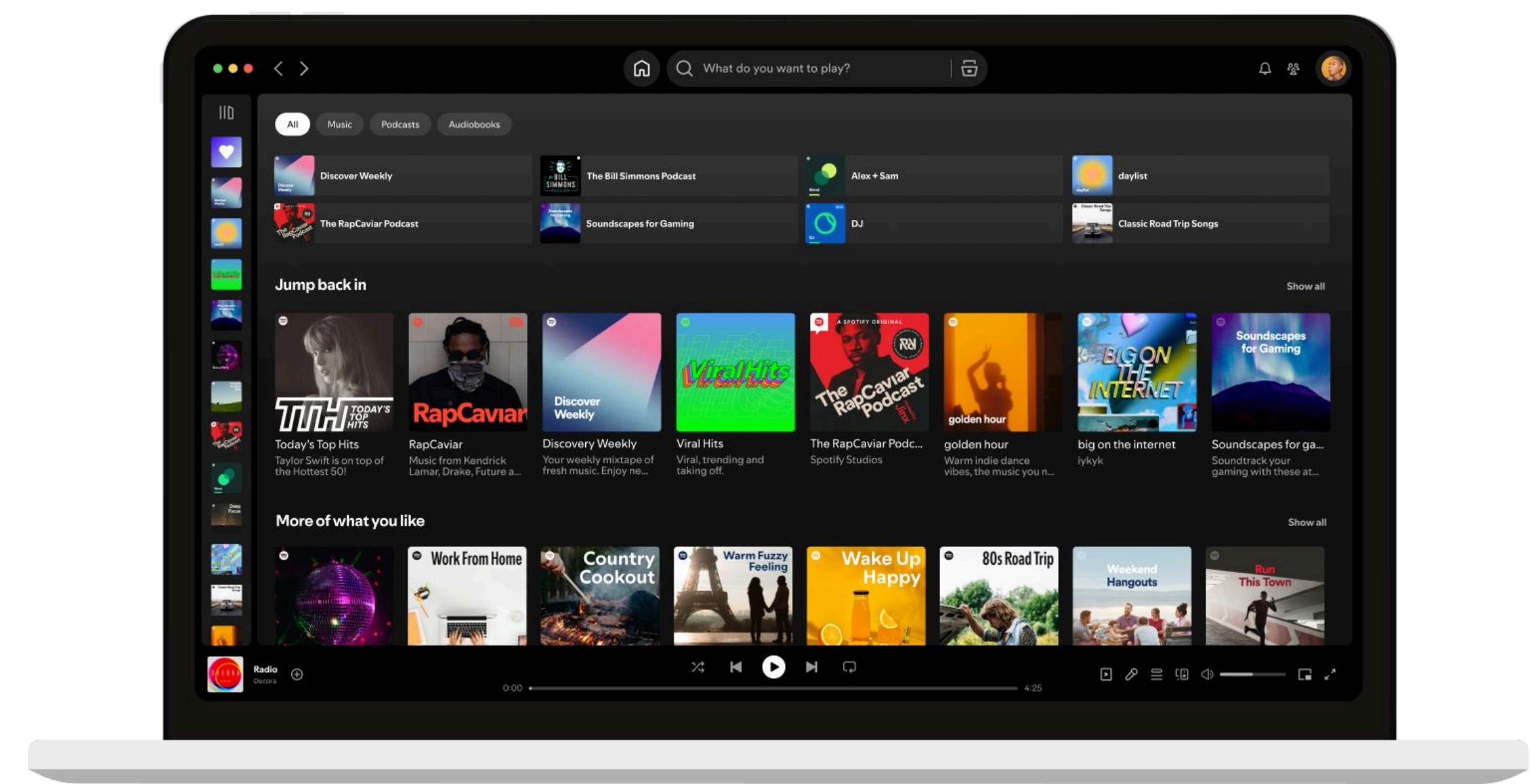
Core ML concepts for a hackathon

Demo: Kaggle Titanic Dataset

Applying what we learned



**Why is it so
important?**





Data Science

Turning raw data into actionable insights
using math, coding and storytelling



Data Cleaning

Removing duplicate/missing values and
standardizing information



Exploratory Data Analysis

Summarizing and visualizing data



Feature Engineering

Creating useful new variables for models



Modeling

Building mathematical representations of data to make predictions

Types of Data

Structured

- Tables/arrays

Numerical

- numbers
 - stock prices

Unstructured

- Raw text or images

Categorical

- types
 - product categories
(clothing, tech, etc)

Data Science Tools

Pandas

- Data manipulation

Scikit-Learn, TensorFlow, Pytorch

- Machine learning models

Seaborn + Matplot

- Data Visualization

NumPy

- Numerical Calculations

Kaggle

- Database of datasets

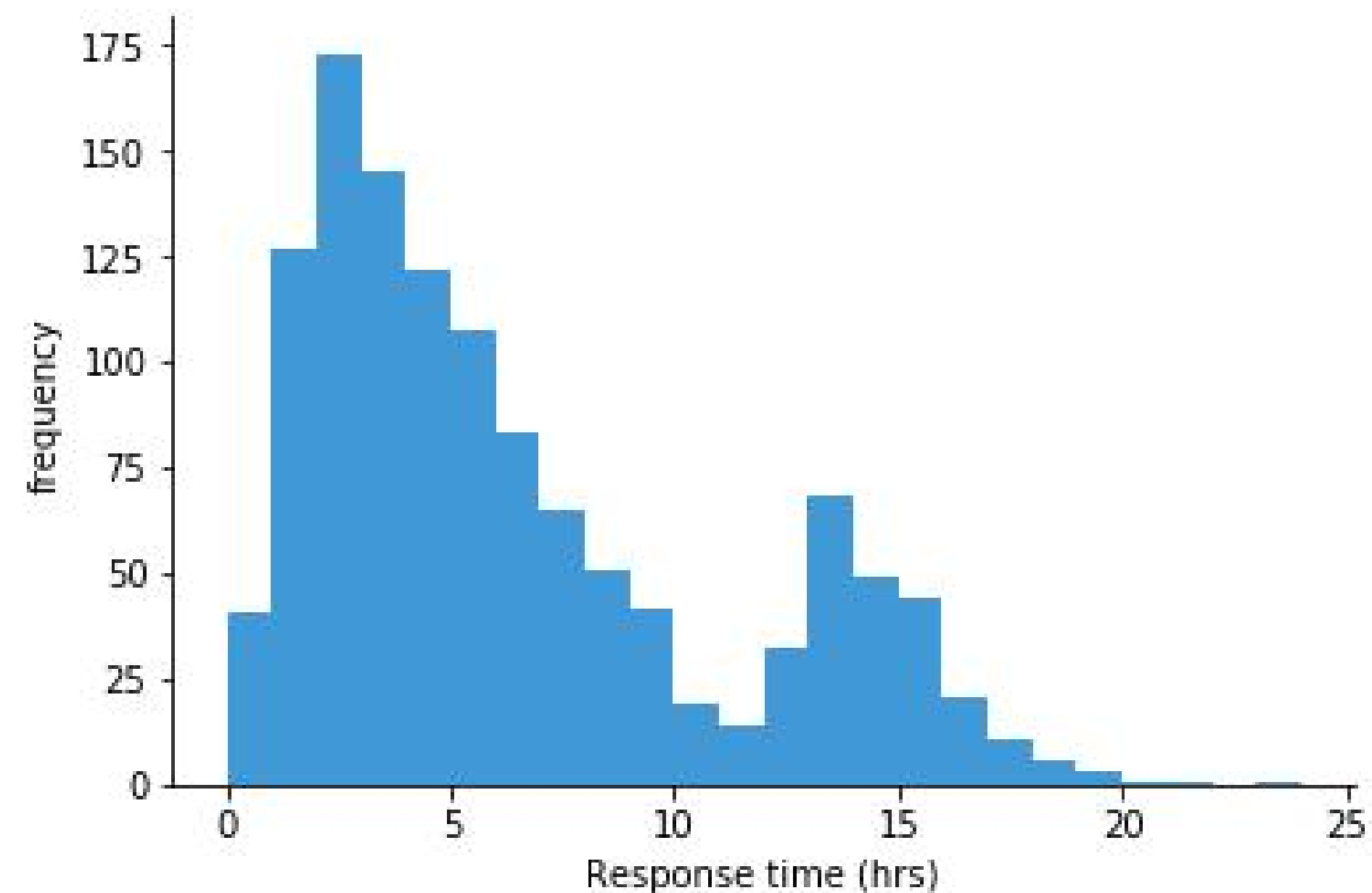
Jupyter Notebook

- IDE

Data Visualization Basics

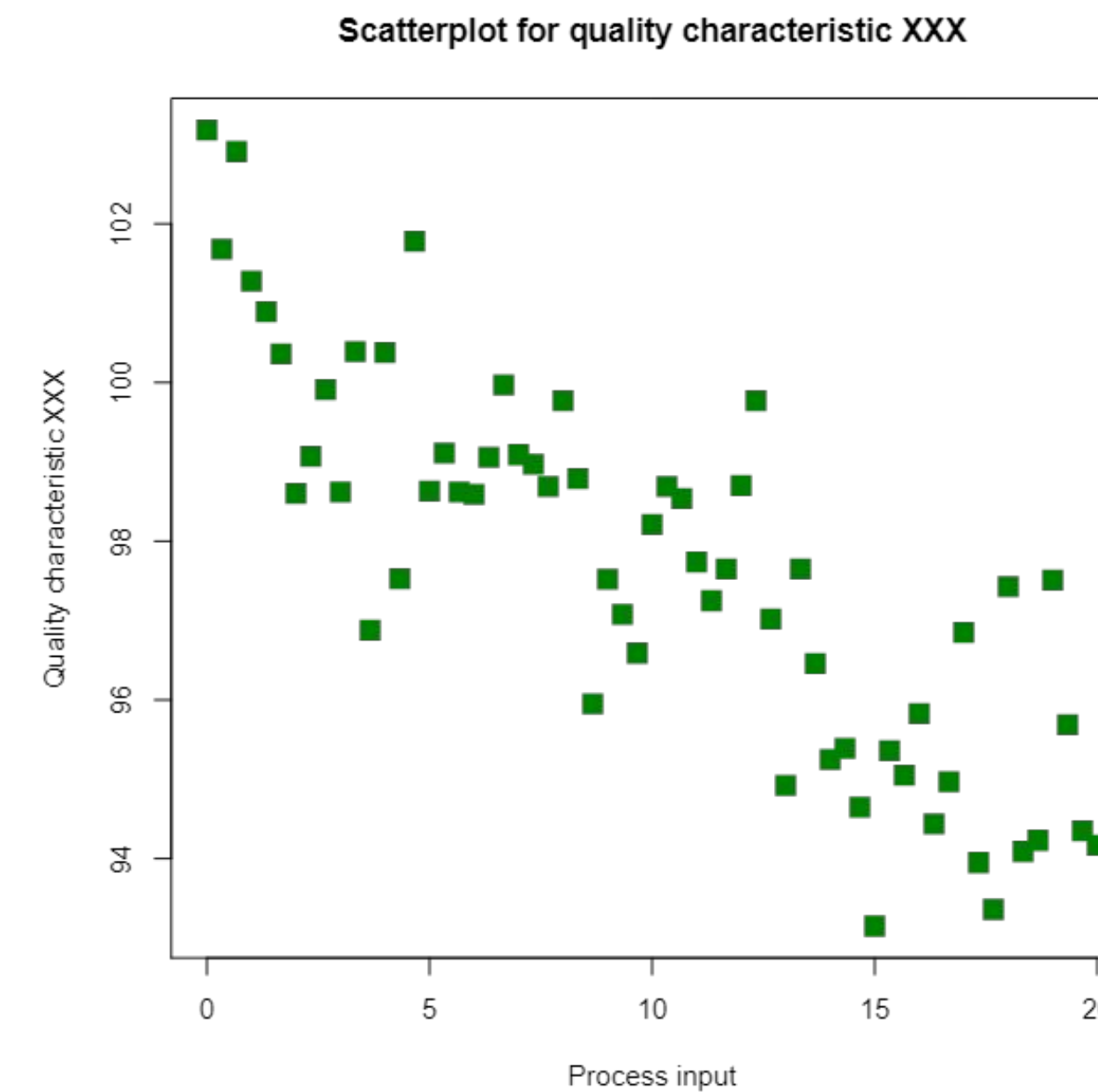
Histogram

- Shows data distribution



Scatterplot

- Identifies relationships



Data Visualization Basics

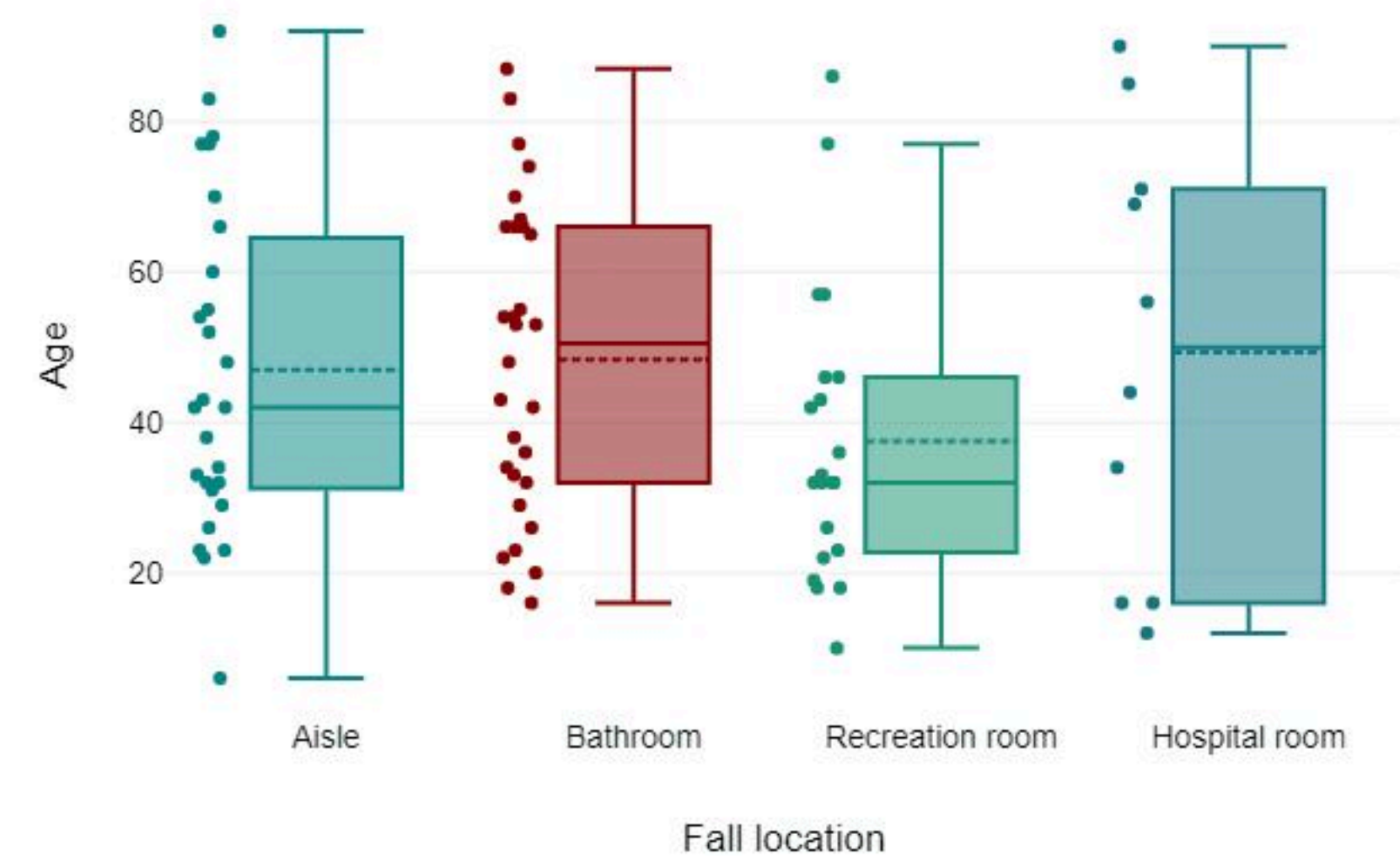
Heatmap

- Shows correlations



Box Plot

- Highlights outliers

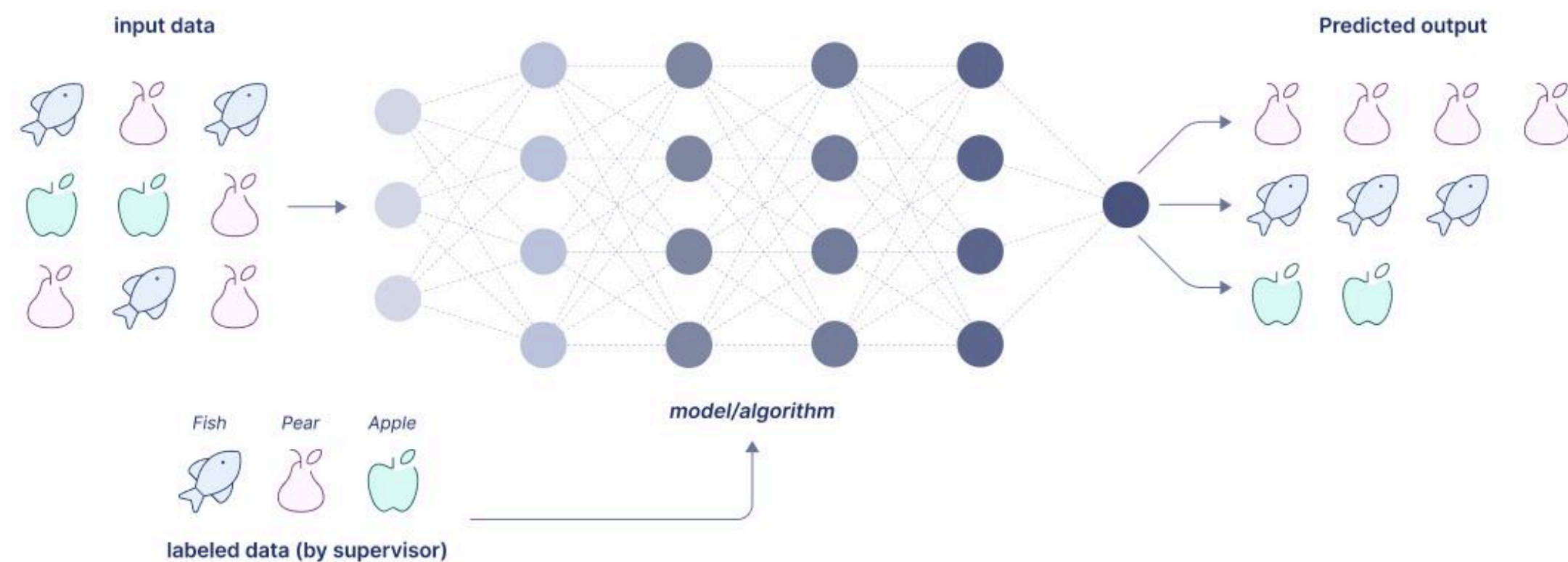


What's wrong?

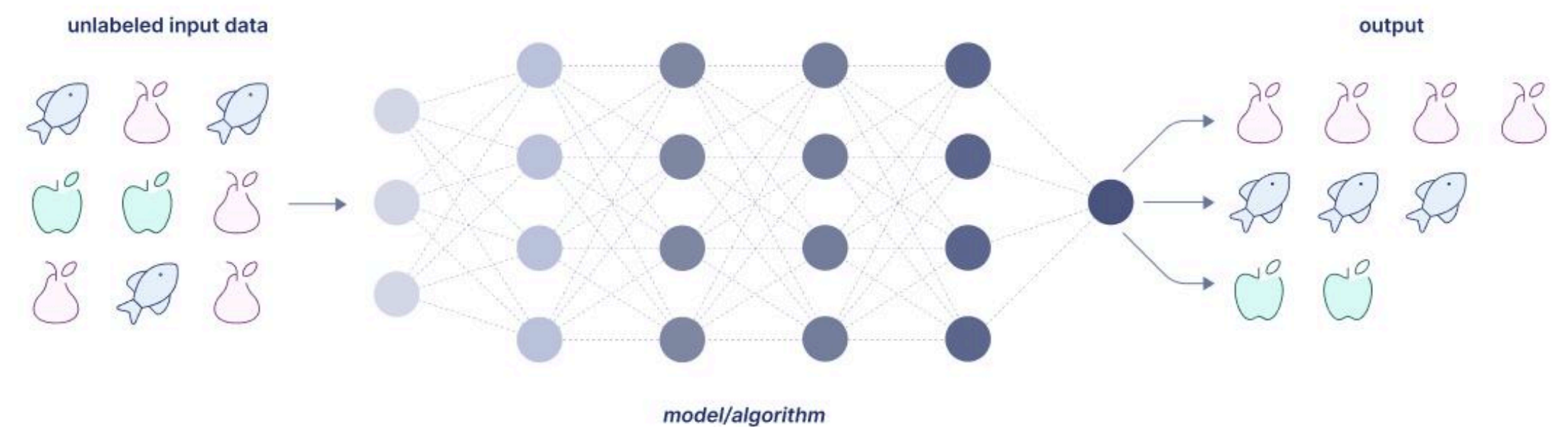
Global Search...													
<input type="checkbox"/>	Name ↑	Status	Temporary Status	Hacker ID	Application Date	Email	Student Number	Major	Enrollment Year	Participant Type	Team Member Names	Dietary Restrictions	PI
<input type="checkbox"/>	Aaron Wong	Accepted	Accepted	200	Wed Feb 12 2025 19:40...	[REDACTED]	[REDACTED]	Computing Scie...	2nd year	Team (4 people...	Alex Zeng, Brayden Ch...		Ye
<input type="checkbox"/>	Abhijot Singh S...	Accepted	Accepted	236	Wed Feb 12 2025 19:40...	[REDACTED]	[REDACTED]	Computing Scie...	2nd year	Team (4 people...	Abhijot Singh Sandhu, ...		Ye
<input type="checkbox"/>	Aditya Kulkarni	Accepted	Accepted	88	Wed Feb 12 2025 19:40...	[REDACTED]	[REDACTED]	Computing Scie...	4th	Individual			Ye
<input type="checkbox"/>	Adrian Crusius	Accepted	Accepted	251	Wed Feb 12 2025 19:41...	[REDACTED]	[REDACTED]	Data Science	4th year	Individual looki...		Vegetarian	Ye
<input type="checkbox"/>	Adriel Adasa	Withdrawn	Accepted	202	Thu Feb 13 2025 22:02...	[REDACTED]	[REDACTED]	Computing Scie...	3rd year?	Individual looki...			No
<input type="checkbox"/>	Ajay Unnikrishn...	Accepted	Accepted	188	Wed Feb 12 2025 19:41...	[REDACTED]	[REDACTED]	Engineering	1st	Team (4 people...	Shameer Khan, Rushee...		Ye
<input type="checkbox"/>	Alex Oliver Reyes	Accepted	Accepted	274	Wed Feb 12 2025 19:41...	[REDACTED]	[REDACTED]	Computing Scie...	1st year	Individual			No
<input type="checkbox"/>	Alex Zeng	Accepted	Accepted	164	Wed Feb 12 2025 19:41...	[REDACTED]	[REDACTED]	Computing Scie...	2nd	Team (4 people...	Arron Wong, Brayden ...		Ye
<input type="checkbox"/>	Alexander Chen	Accepted	Accepted	199	Wed Feb 12 2025 19:41...	[REDACTED]	[REDACTED]	Computing Scie...	2nd	Individual looki...			Ye
<input type="checkbox"/>	Alexander Potia...	Accepted	Accepted	115	Wed Feb 12 2025 19:41...	[REDACTED]	[REDACTED]	Computing Scie...	2nd year	Team (4 people...	Manan Mehta and Khali...		Ye

Machine Learning in 5 minutes

Supervised learning



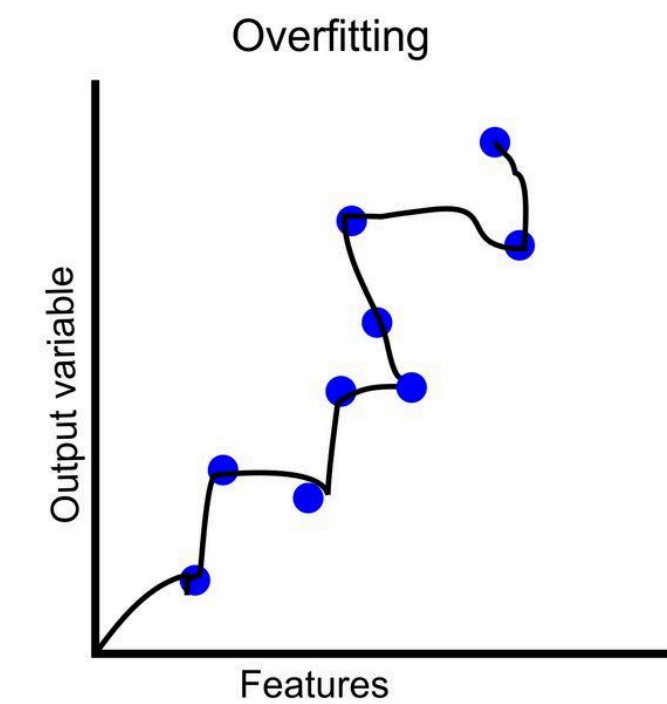
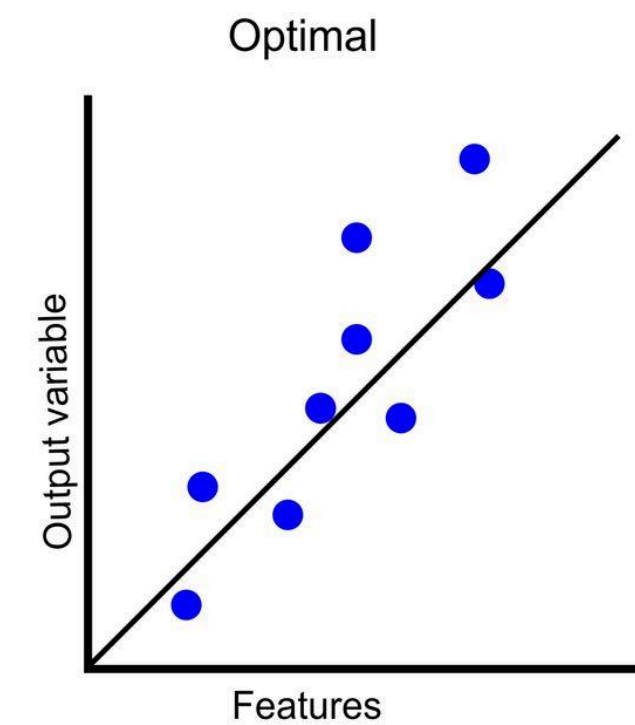
Unsupervised learning



Machine Learning in 5 minutes

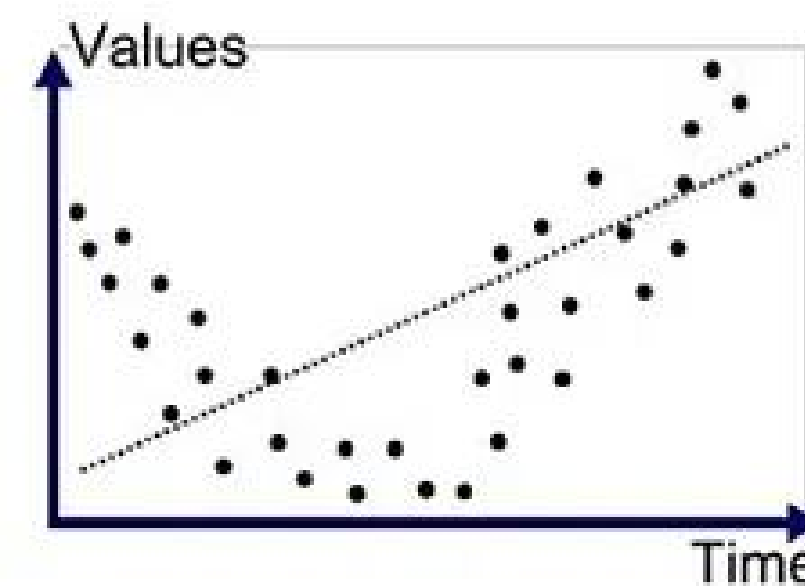
Overfitting

Model memorizes the data instead of generalizing it

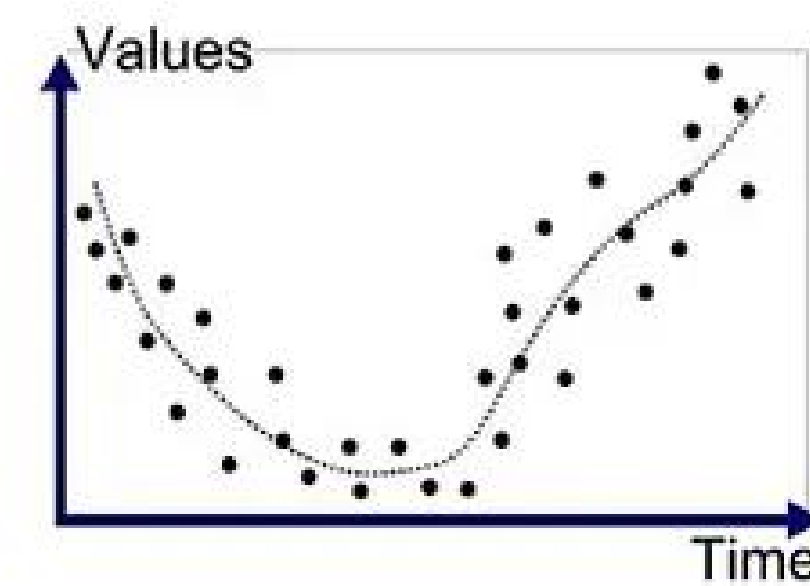


Underfitting

Model is too simple and misses out on key patterns



Underfitted



Good Fit/Robust



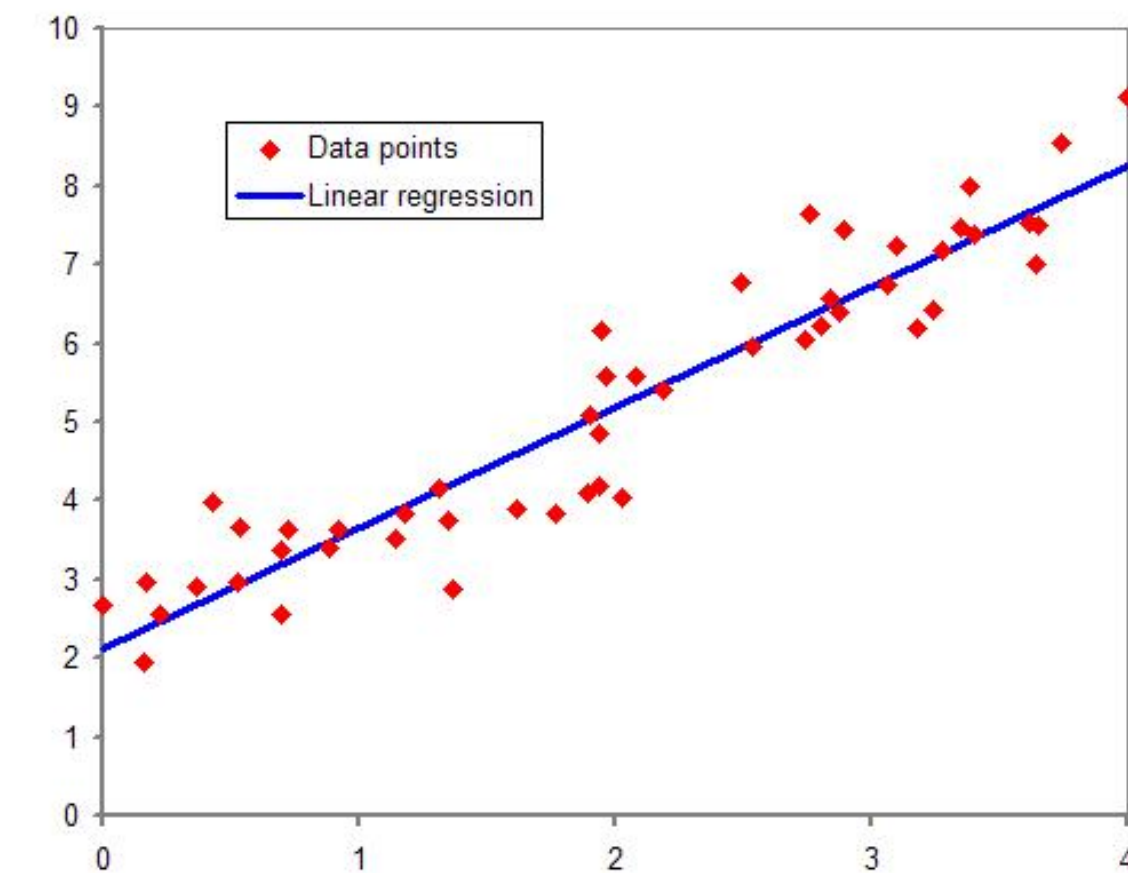
Pro hacker tip:

Garbage in, garbage out

Machine Learning in 5 minutes

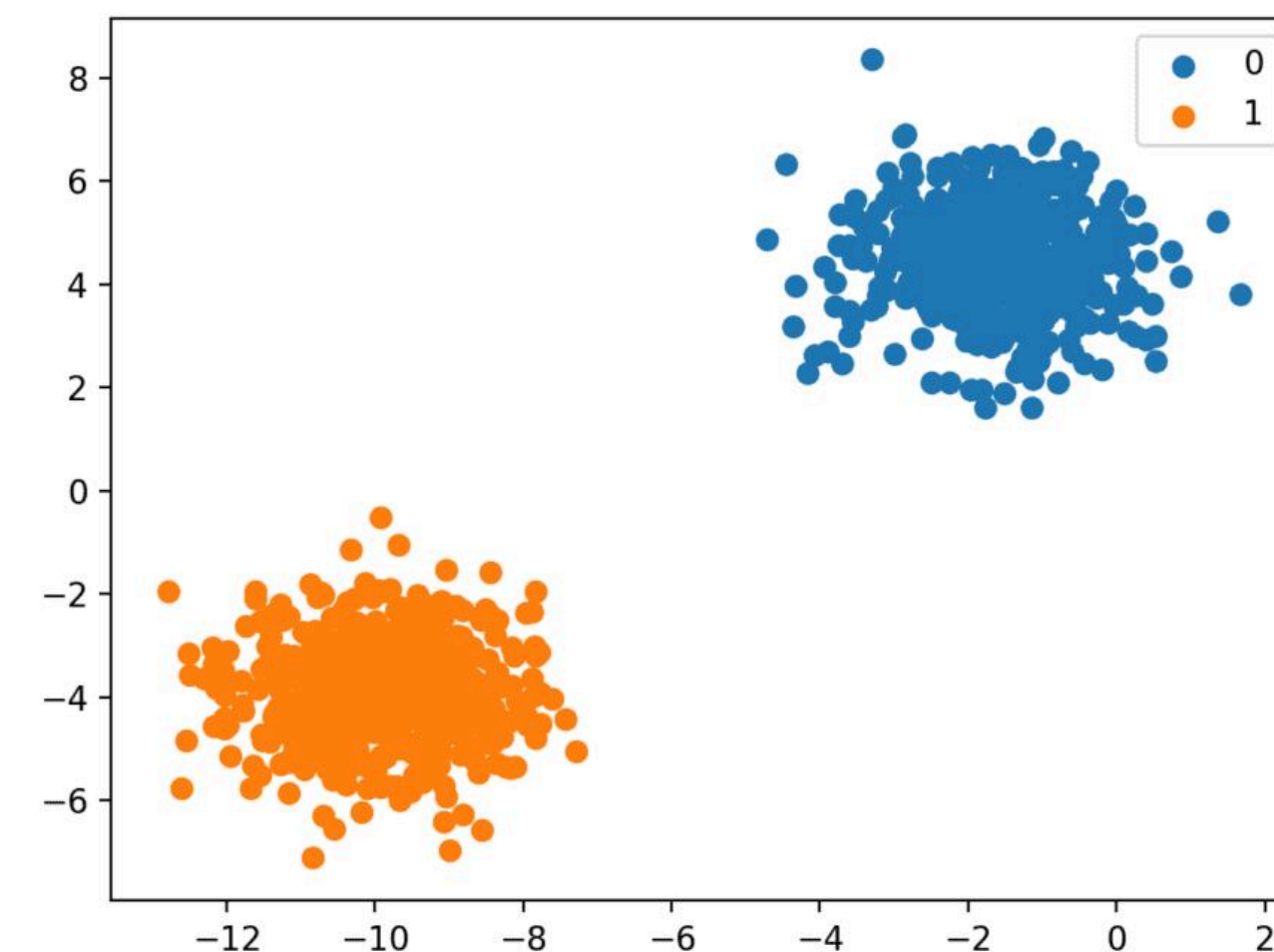
Regression

Predicts recurring values like housing prices



Classification

Predicts categories such as spam emails



Regression Models

Linear Regression

Line of best fit

$$y = mx + b \rightarrow y = \beta_0 + \beta_1 x + \varepsilon$$

$\beta_0 = m$, $\beta_1 = b$, $\varepsilon = \text{error}$

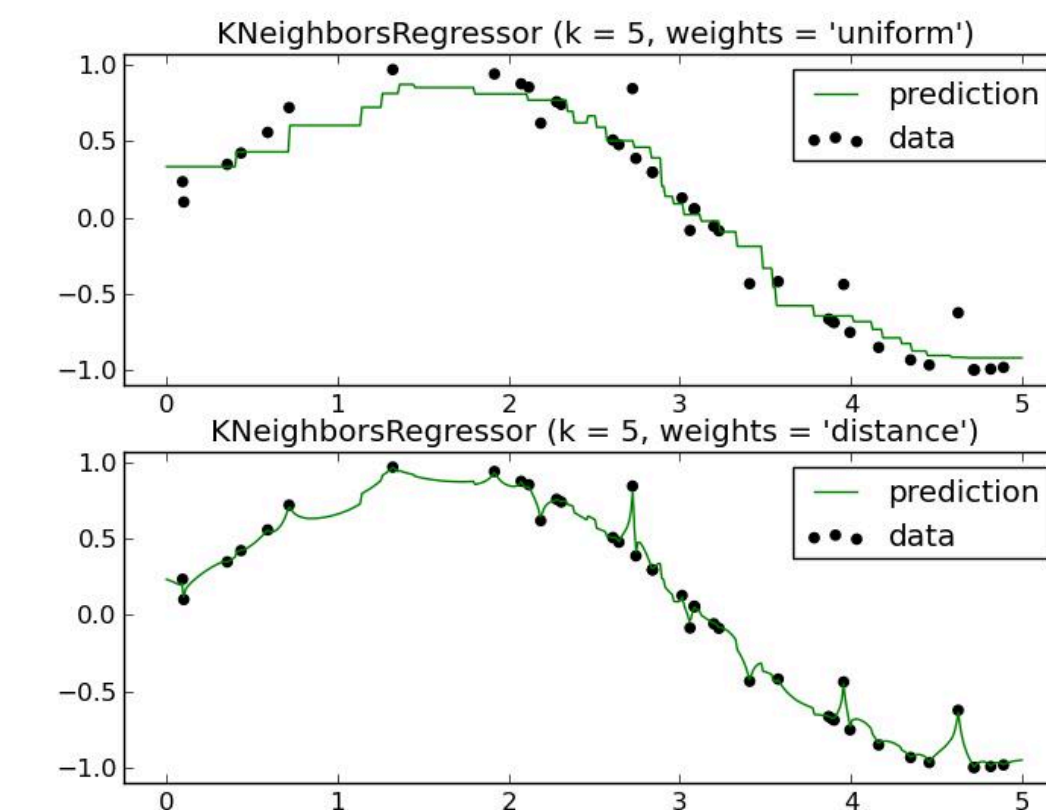
kNN (regression vers.)

Predict value by average of 'k' nearest neighbours

Facebook Prophet

A regression model for time series

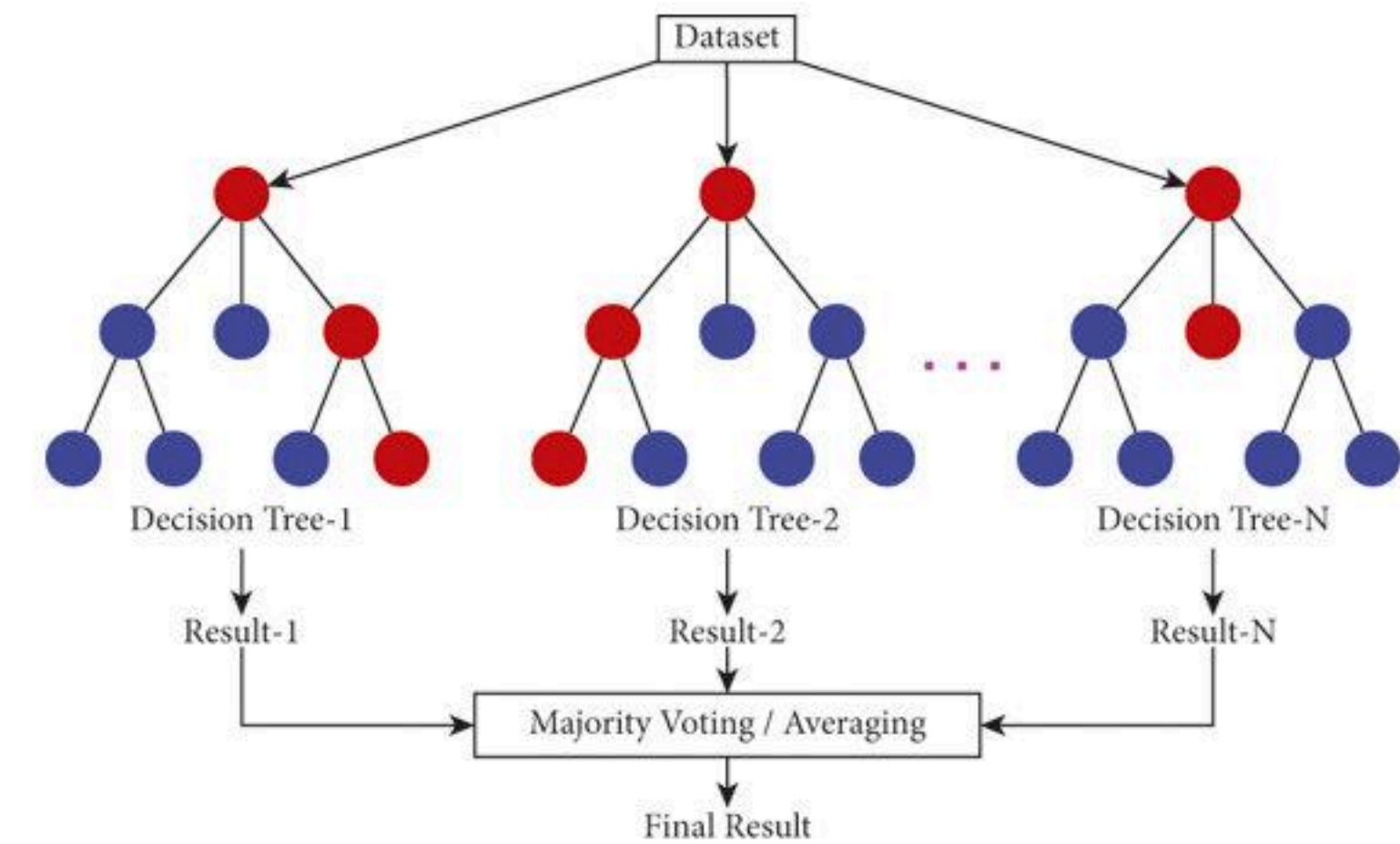
Breaks data into: Trend + seasonality + special events + error



Classification Models

Random Forests

A bunch of decision trees

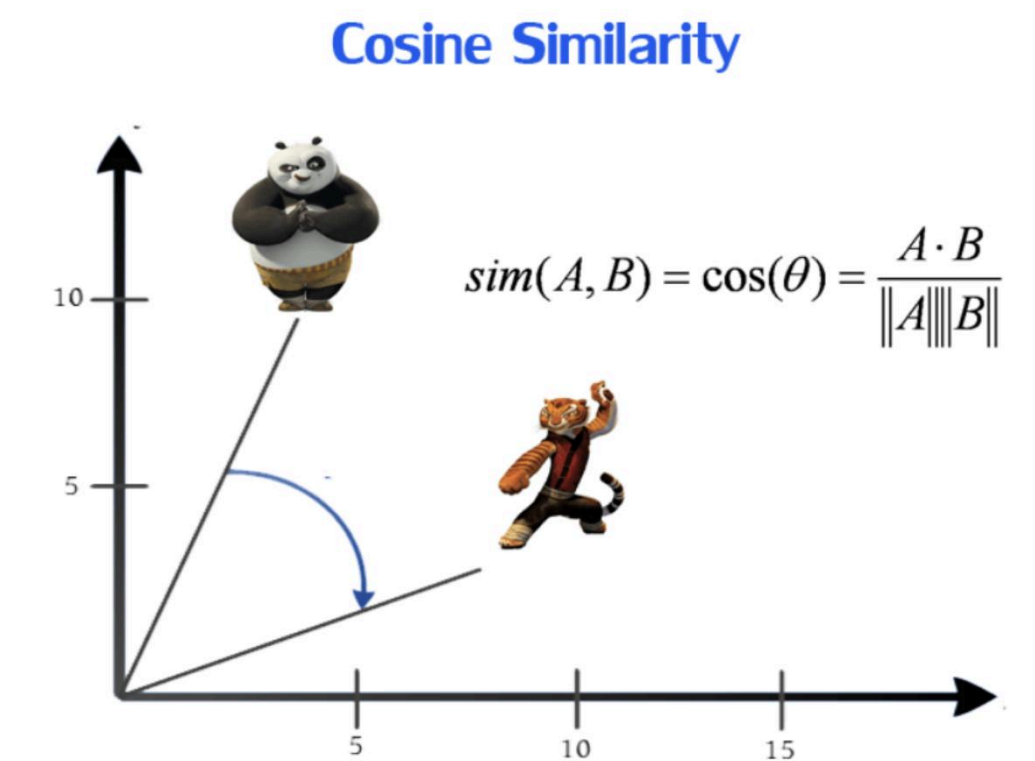
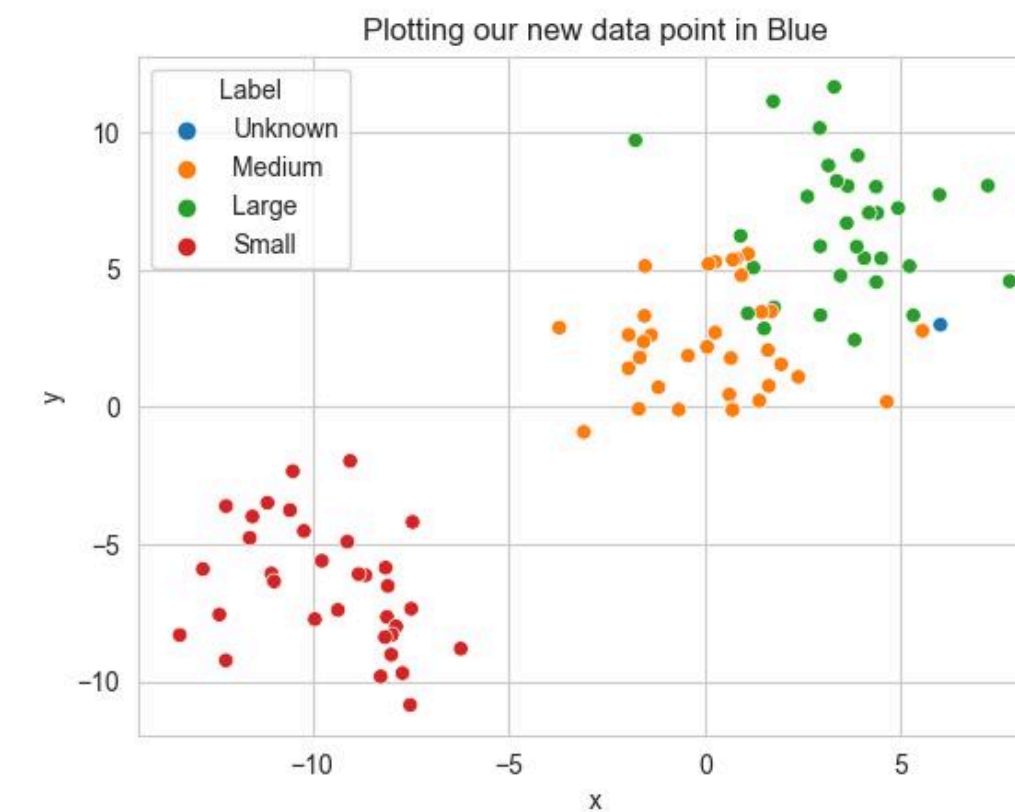


TF-IDF + Naive Bayes (Gaussian)

Classification of word weights

kNN (classification vers.)

Labels all closest neighbors.





Demo Time!