

Demystifying AI - Day 1





Topics

- AI in our everyday lives
- What is AI?
- Data
- Subsets of AI
- Introduction to ML algorithms
- Identifying & Productizing ML Use-Cases
- Conclusion



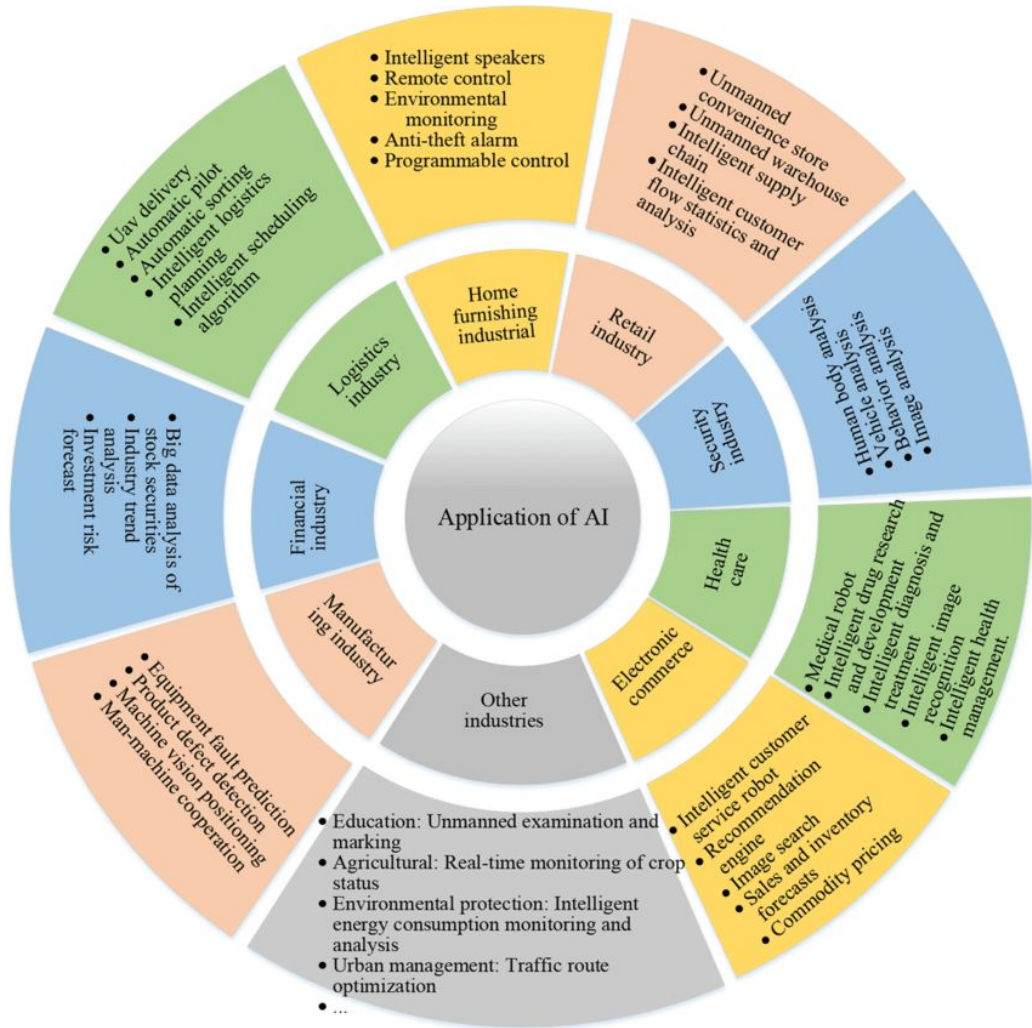
Goals

- Understand and notice AI in our daily lives
- Distinguish between the different types of AI tasks
- Understand data in the context of AI
- Conceptualize what algorithms are used on different types of data
- Be able to recognize good ML use-cases



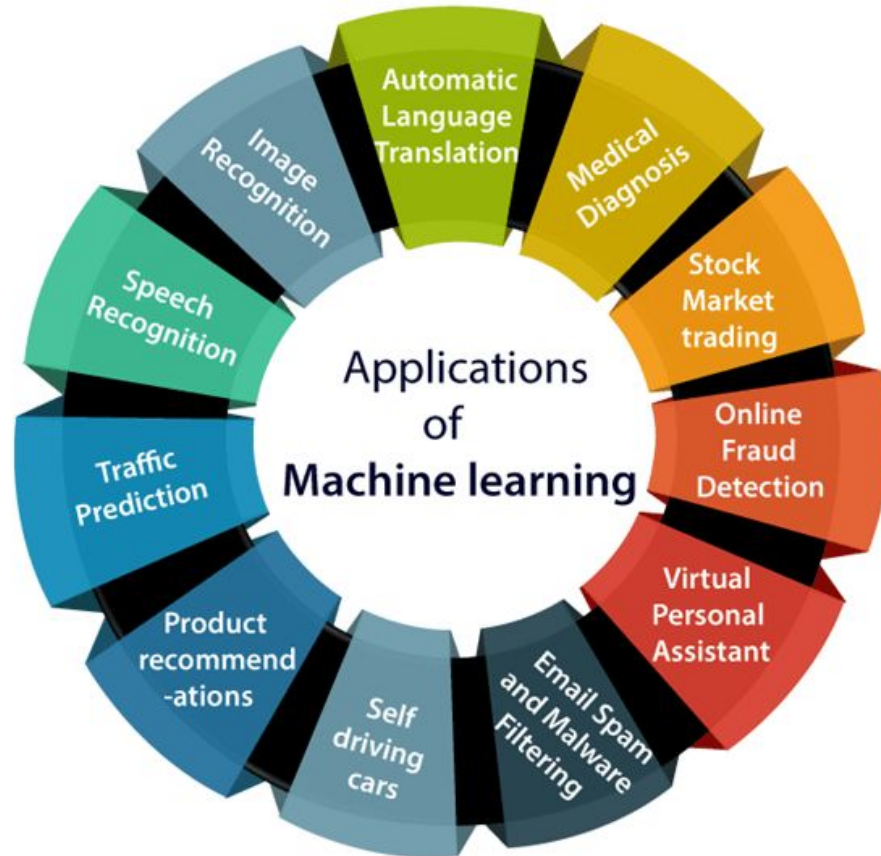
AI in our everyday lives

Overview





Overview





Recommendation Systems









Examples

- Amazon
- YouTube
- Netflix
- Spotify
- LinkedIn
- Pandora

Benefits

- Increased conversion
- Reduced churn
- Overall increased customer satisfaction

Recommended for you, Thomas

 <p>Literature & Fiction 62 ITEMS</p>	 <p>Exercise & Fitness Equipment 8 ITEMS</p>	 <p>Health, Fitness & Dieting Books 37 ITEMS</p>	 <p>Tableware 12 ITEMS</p>
 <p>Prime Video – Unlimited Streaming for Prime Members 12 ITEMS</p>	 <p>Coffee, Tea & Espresso 96 ITEMS</p>	 <p>Biographies & Memoirs 17 ITEMS</p>	 <p>Engineering Books 7 ITEMS</p>

Computer Vision & Natural Language

- Self-driving cars
- Social media filters (instagram/snapchat)
- Siri
- Alexa



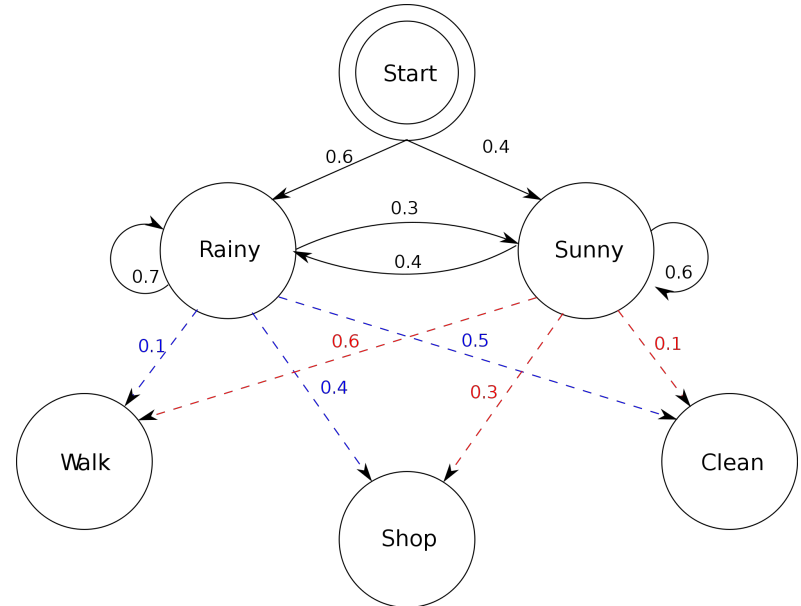
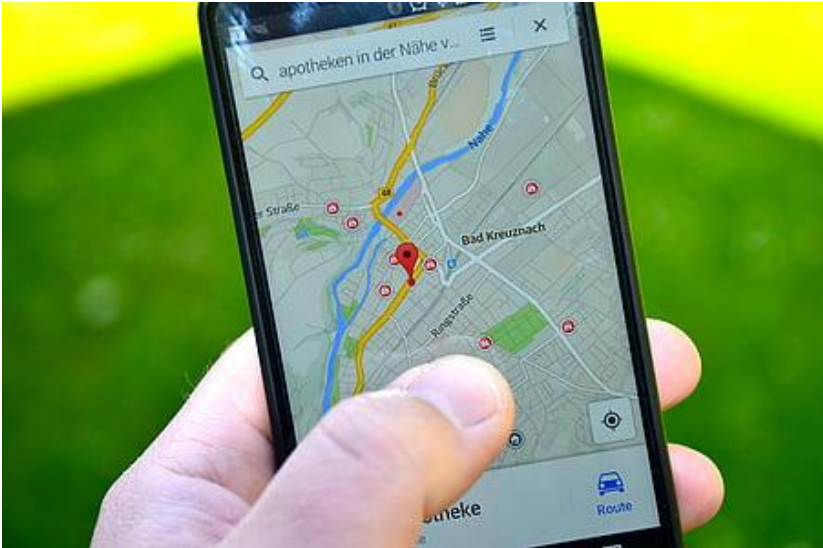
Inventory Optimization

- Grocery Shopping
- Online Shopping
- Supply-chain management



GPS Navigation

- How do our phones always know the fastest way to get from point A to B? AI, of course!
- Modern GPS navigation systems such as Waze and Google Maps are utilizing GNNs (graph-neural-networks).



Fraud Detection

- Fraud detection in financial transactions, loan approvals, insurance rates & more
- When your bank calls and says ‘is this really you’, yes, that is AI!
- When emails get sent to your spam folder, that is AI!





What is AI?

AI = COMPUTERS + MATH + DATA

- Statistics (probabilities)
- Linear Algebra (vectors, matrices)
- Calculus (optimization)

$$A = [x_0, x_1, x_2, \dots, x_n]$$

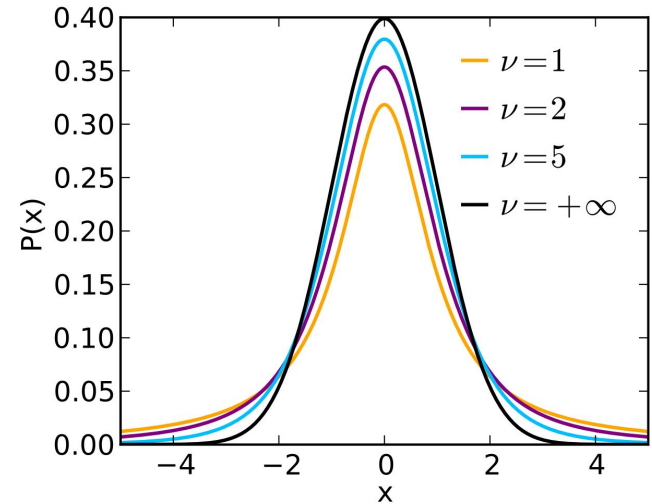
$$B = [y_0, y_1, y_2, \dots, y_n]$$

$$A \odot B = \sum_{i=0}^n x_i \cdot y_i$$

Dot Product

$$H(X) = H(p_1, \dots, p_n) = - \sum_{i=1}^n p_i \log_2 p_i$$

Shannon Entropy Formula



Probability Density Function



How Does It Work?

- Optimization
 - Minimizing or Maximizing a function
- Give algorithm input data, and labeled output data
- Randomly initialize parameters
- Iteratively adjust parameters until the mathematical “distance” between algorithm output and labeled output is small.

$$y = mx + b$$

Simple Regression

$$MSE = \frac{1}{n} \sum \left(\underbrace{y - \hat{y}}_{\substack{\text{The square of the difference} \\ \text{between actual and} \\ \text{predicted}}} \right)^2$$

Measure of “incorrectness”



Trainable Parameters

- In the previous slide, it was mentioned we iteratively adjust our parameters to tighten the gap between the predicted output and the desired output.

- **Example:**

- $X = [1.2, 1.7, 2.9, 4]$
- $Y = [0]$

$$\hat{y} = X \bullet W + b$$

$$\hat{y} = [1.2, 1.7, 2.9, 4] \bullet [-0.1, -0.1, -0.1, -0.1] + [0, 0, 0, 0]$$

$$\hat{y} = -0.98$$

Now we adjust our 'W' and 'b' parameters to minimize the MSE, and run the formula again. This is called training.

$$MSE = 0.9603999999999999 = (0 - (-0.98))^2/1$$



Data



What is data in AI?

- Data is the most vital piece of any machine learning algorithm
- Everything you do everyday, generates data that can be utilized by a ML algorithm
- **Fun Facts:**
 - 2.5 quintillion bytes of data created everyday
 - Humans on avg generate 1.7MB of data per second

Our current love affair with social media certainly fuels data creation. According to Domo's [Data Never Sleeps 5.0 report](#), these are numbers generated **every minute** of the day:

- Snapchat users share 527,760 photos
- More than 120 professionals join LinkedIn
- Users watch 4,146,600 YouTube videos
- 456,000 tweets are sent on Twitter
- Instagram users post 46,740 photos

Source: [How Much Data Do We Create Every Day? The Mind-Blowing Stats Everyone Should Read \(forbes.com\)](#)



Different Types of Data

- Structured Data & Unstructured data
- **Structured Data**
 - Tabular data (any data in table format)
 - Think excel
- **Unstructured Data**
 - Documents (PDFs/invoices/receipts)
 - Text data (articles, books, reviews, comments, tweets, etc)
 - Image data (images/videos)
 - Audio data (mp3/wav)



Data Sources

- YOU are the data source
- Here is an example:
 - You are planning a vacation - you google different hotels, routes, and activities
 - You go on the trip, download a book on audible, order ubers, go shopping, maybe eat some nice meals, post pictures, tweet about it, etc
- Every step above generated tons of data that will be used to sell you things, show you ads, and essentially aim to “program” you.
- Companies of all kinds across all industries will purchase this data, and then develop ML algorithms to increase business capacity or reduce costs



Data Storage

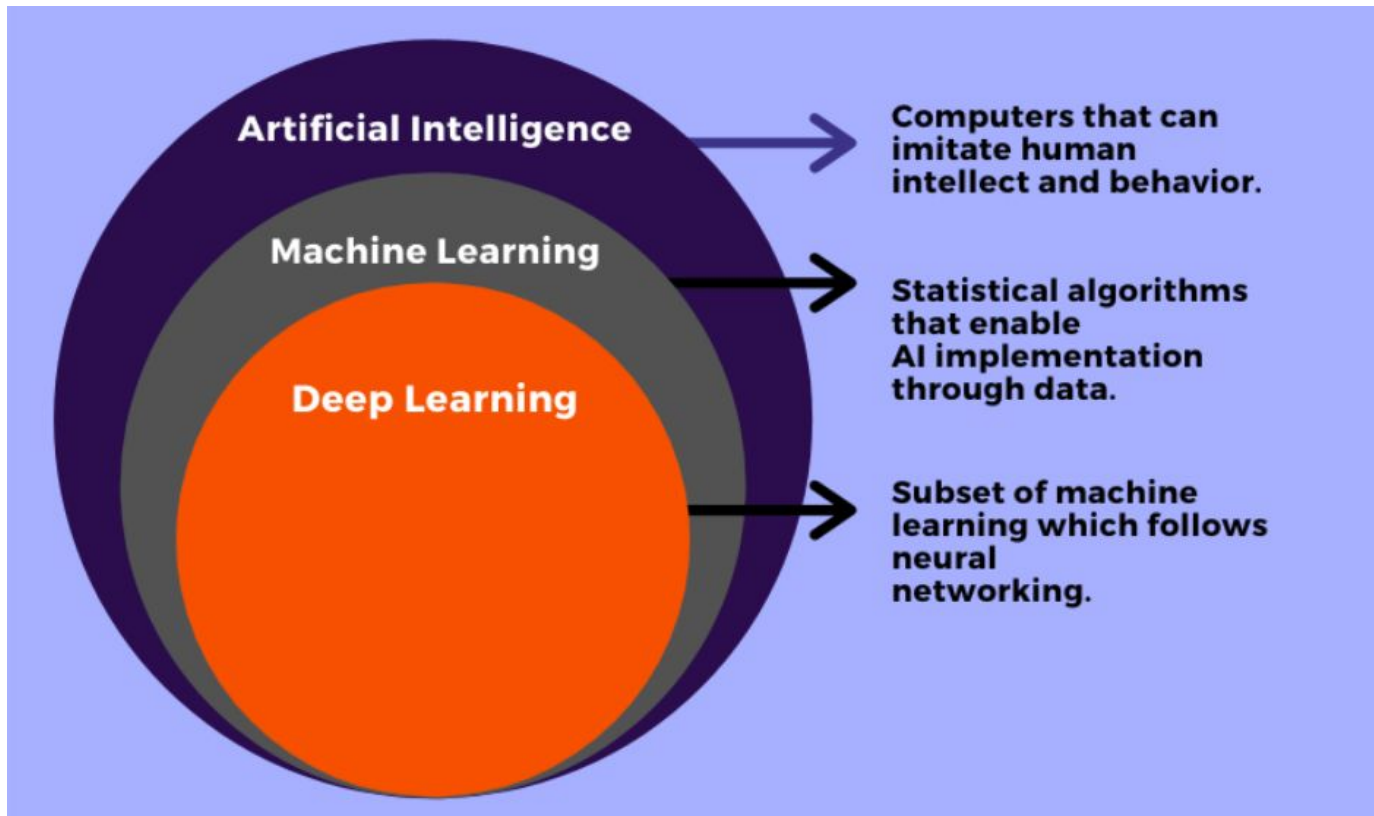
- Generally, we are dealing with large volumes of data, too large to maintain on a local machine
- Companies have different systems
 - CRM (salesforce, SAP)
 - SQL (GraphQL, PostgreSQL)
 - Cloud data-store (Azure, AWS, GCP)
- When building an ML model, we generally need to access multiple types of data, from multiple systems. The more data the better the algorithm.
- **Accessing the data**
 - API
 - Queries
 - & sometimes it's as simple as hitting 'export'! Thank you Salesforce!



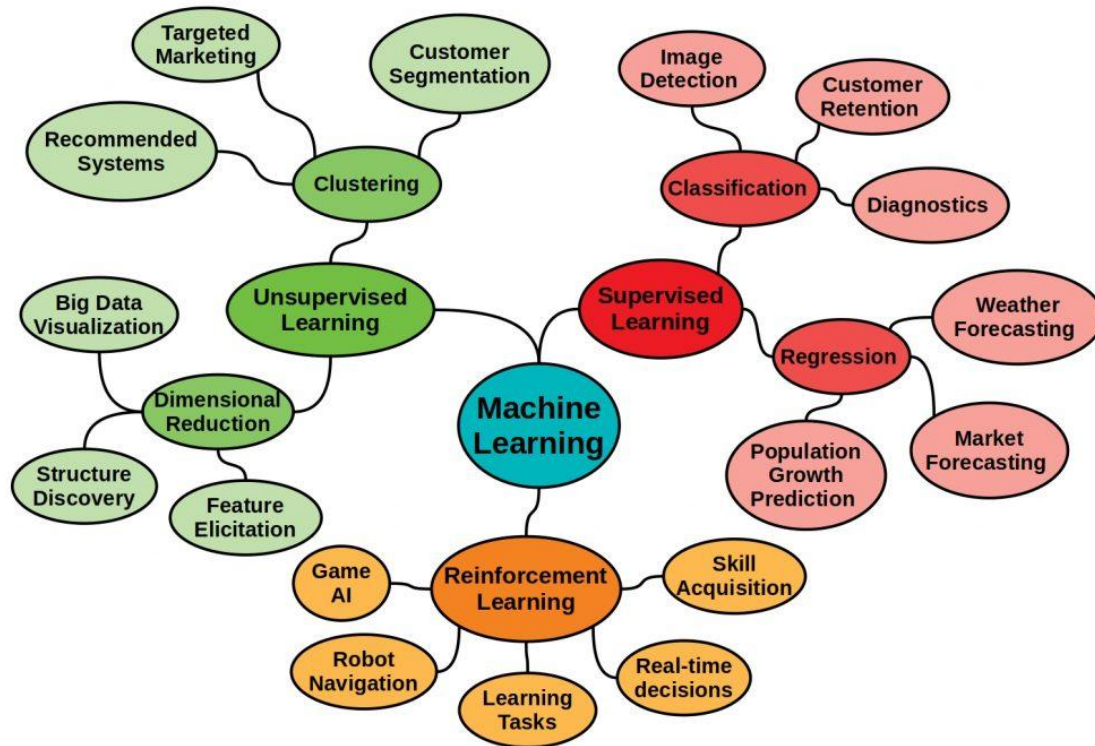
Subsets of AI



Overview



Supervised vs Unsupervised vs RL





Supervised Learning

- Supervised learning means we have labeled data
- **Classification**
 - Seeks to classify inputs
 - Is this an airplane?
 - Will the Yankees win or lose?
 - Will it rain tomorrow?
- **Regression**
 - Seeks to predict a continuous value
 - Close price of \$AAPL tomorrow
 - Quarterly revenue
 - Age



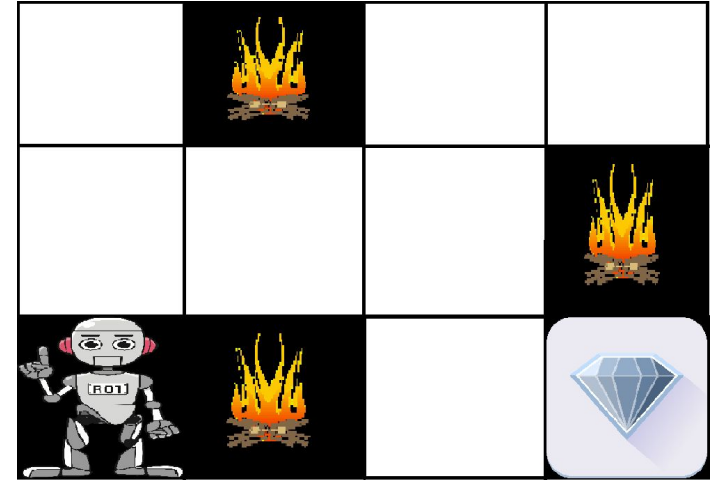
Unsupervised Learning

- Unsupervised learning means the dataset is not labeled
- **Clustering**
 - Discover hidden patterns and data-groupings without the need for human intervention
 - Clustering is a data mining technique which groups unlabeled data based on their similarities or differences.
 - Its ability to discover similarities and differences in information make it the ideal solution for exploratory data analysis, cross-selling strategies, customer segmentation, and image recognition.
- We are essentially creating “classes”, using the raw data.



Reinforcement Learning

- AREA
 - Agent, Reward, Environment, Action
- We are training a model to respond to a changing environment.
- Taking suitable action to maximize reward in a particular situation.
- Differs from supervised learning because we do not have an answer key (labels)
- The agent learns from its experience in the environment



The above image shows the robot, diamond, and fire. The goal of the robot is to get the reward that is the diamond and avoid the hurdles that have fire. The robot learns by trying all possible paths and then choosing the path which gives the reward with the least hurdles. Each right step will give the robot a reward and each wrong step will subtract the reward of the robot. The total reward will be calculated when it reaches the final reward that is the diamond.



Introduction - ML Algorithms

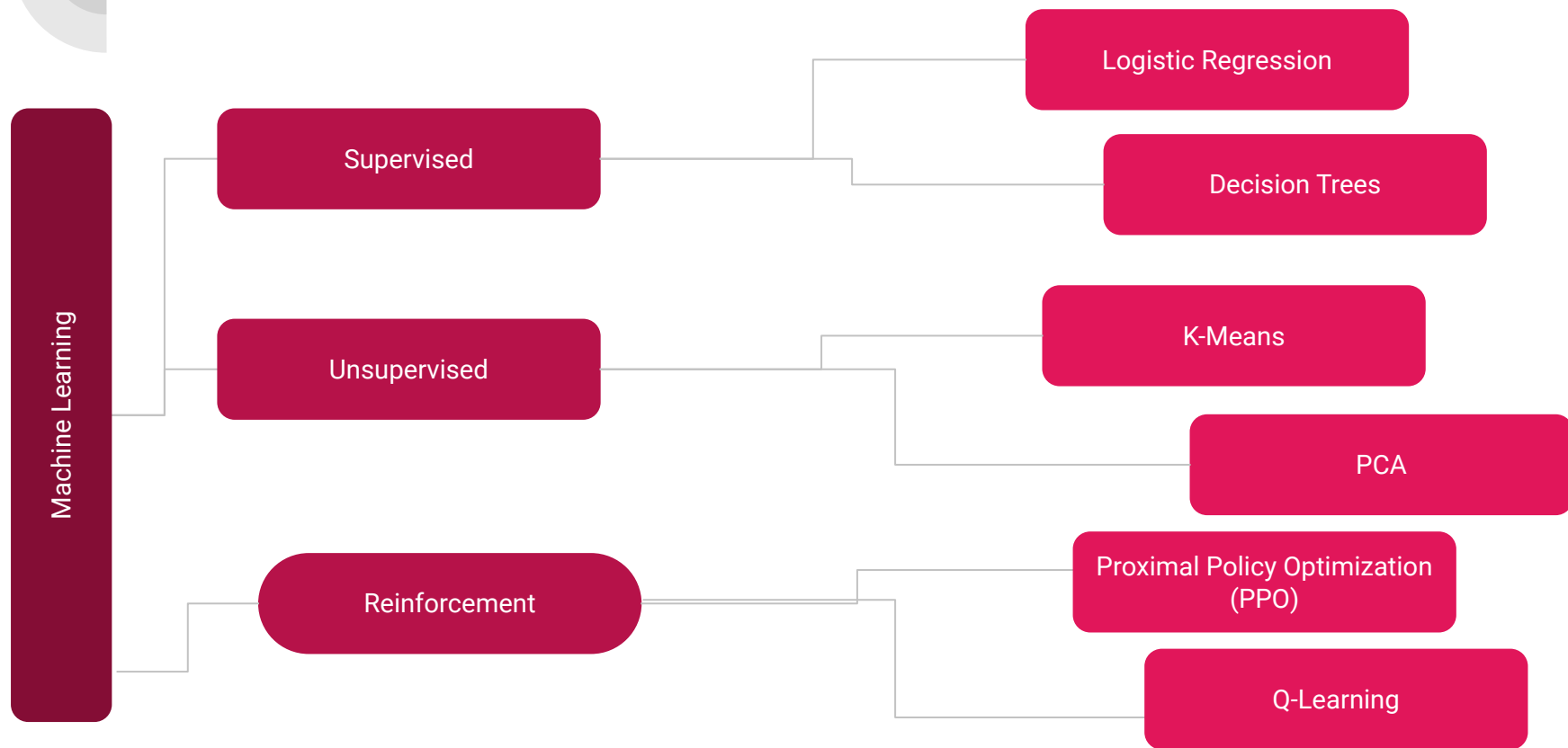




Overview

- Now that we have defined the subsets of ML, we can identify some prevalent algorithms used in each subset.
- **Different tasks require different algorithms**
- We cannot use the same algorithm to classify breeds of dogs and also predict next months revenue.
- Generally, we use the following algorithms
 - Decision Trees/Logistic Regression for tabular data
 - Convolutional-Neural-Net (CNN) for vision tasks
 - Recurrent-Neural-Net (RNN) for language tasks
 - Deep-Q Learning for RL tasks

Some Flavors of ML Algorithms

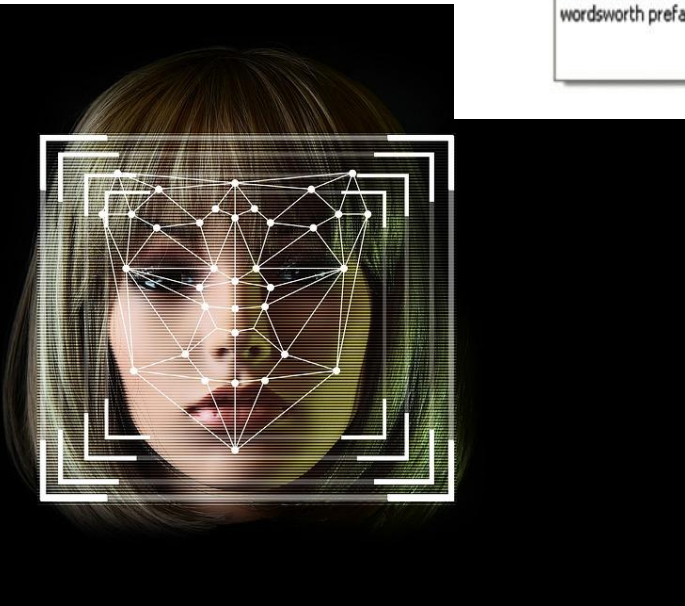


RNN

Examples



CNN



RL



Identifying ML Use-Cases





How To Identify ML Use-Cases

- A process where large, consistent volume of data is being ingested
- A tedious/time-consuming/repetitive process
- A decision making process
- Any process involving predictions & forecasting



Productizing ML Use-Cases

The Machine Learning Process

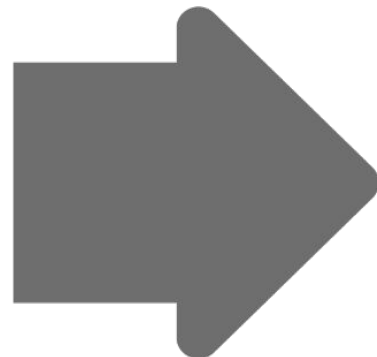
Step 1
Gathering data from
various sources

Step 2
Cleaning data to
have homogeneity

Step 3
Model Building-
Selecting the right ML
algorithm

Step 4
Gaining insights from
the model's results

Step 5
Data Visualization-
Transforming results
into visuals graphs





Types of ML Deployments

- Interactive vs Non-Interactive
- Single Record vs Batch
- Interactive
 - Client facing
 - Client sends data and gets results
- Non-Interactive
 - Not client facing
 - Client only receives results
- Single Record
 - Model is receiving one input at a time and returning one prediction at a time
- Batch
 - Model is receiving large batches of input data and returning batches of predictions



Productizing ML Use-Cases Cont.

- This will be discussed in depth later on in the course.
- The ML process starts with data engineering, and ends with deploying an API.
- **Birds-eye view:**
 - Collect data from data source(s)
 - Develop ML model
 - Build query pipelines for the data
 - Integrate pipeline + model into API
 - Deploy the API



Conclusion





Questions

- I want to forecast my branch's next quarter revenue. What subset of ML is this?
 - a. Unsupervised Learning
 - b. Supervised Learning
 - c. Reinforcement Learning
- I have identified the subset of ML for this project, what type of task is this?
 - a. Regression
 - b. Clustering
 - c. Classification
- I have a dataset on my customer behaviors, what task is best suited to group the customers together and profile them based on behavior?
 - a. Classification
 - b. Reinforcement Learning
 - c. Clustering
- **BONUS:** What algorithm should I use to generate text? Think auto-complete on emails/texts.
 - a. CNN (Convolutional)
 - b. Decision Tree
 - c. RNN (Recurrent)