



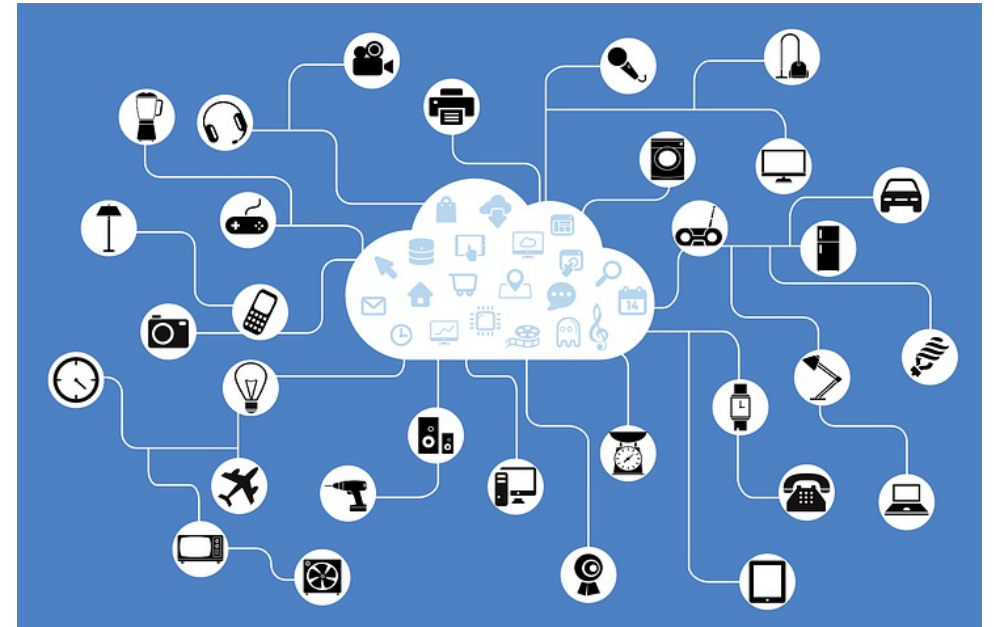
JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

EN.540.635
Software Carpentry

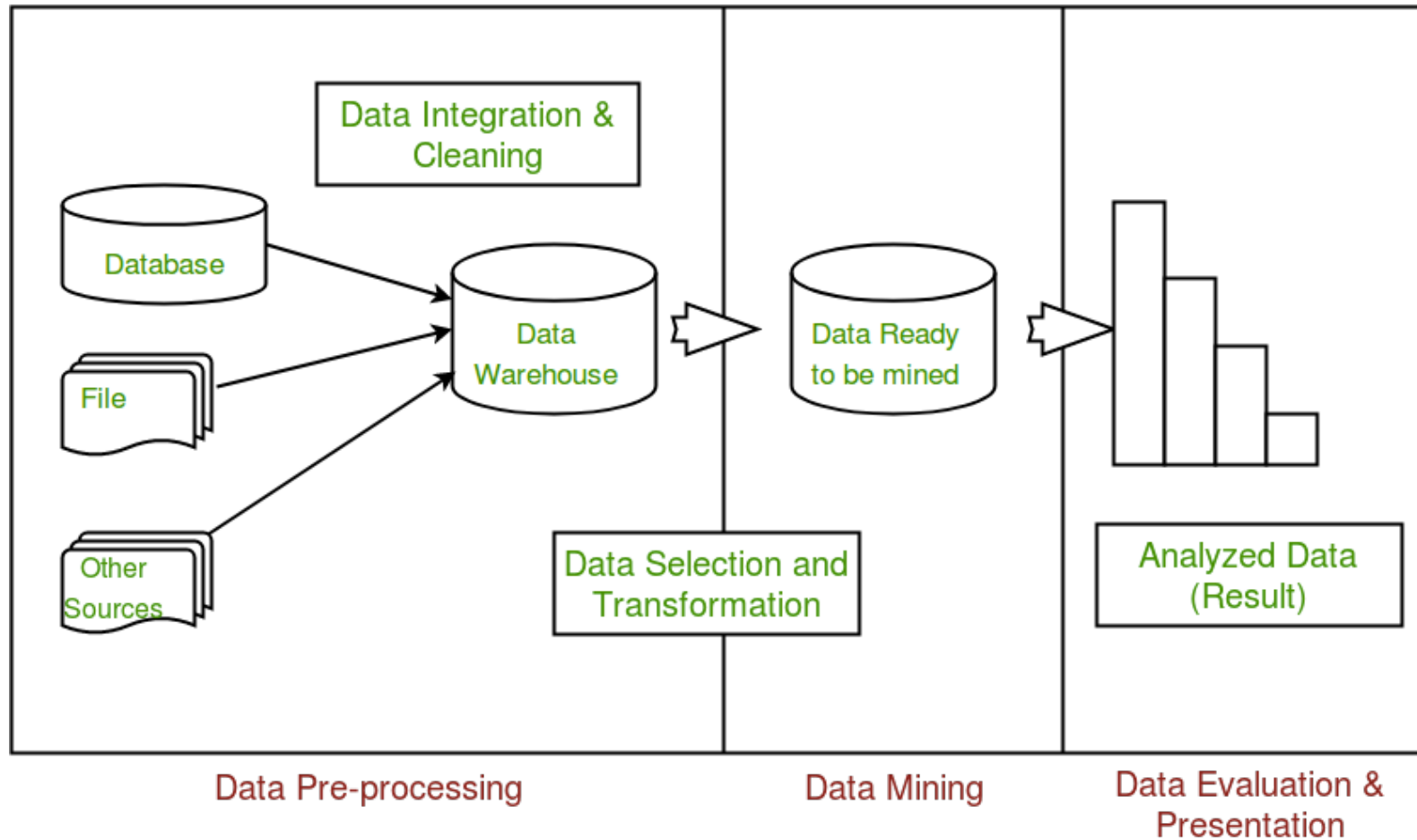
Lecture 12
Data Mining, APIs, and an Introduction to
Machine Learning

What is Data Mining?

- Data mining is defined as the exploration and analysis of large quantities of data (an intersection of statistics, machine learning, and artificial intelligence).
- The goal is to extract meaningful conclusions and find patterns that exist in data.
- Applications:
 - Consumer research and marketing
 - Finance
 - Healthcare
 - Telecommunications
 - Manufacturing and engineering
 - Bioinformatics



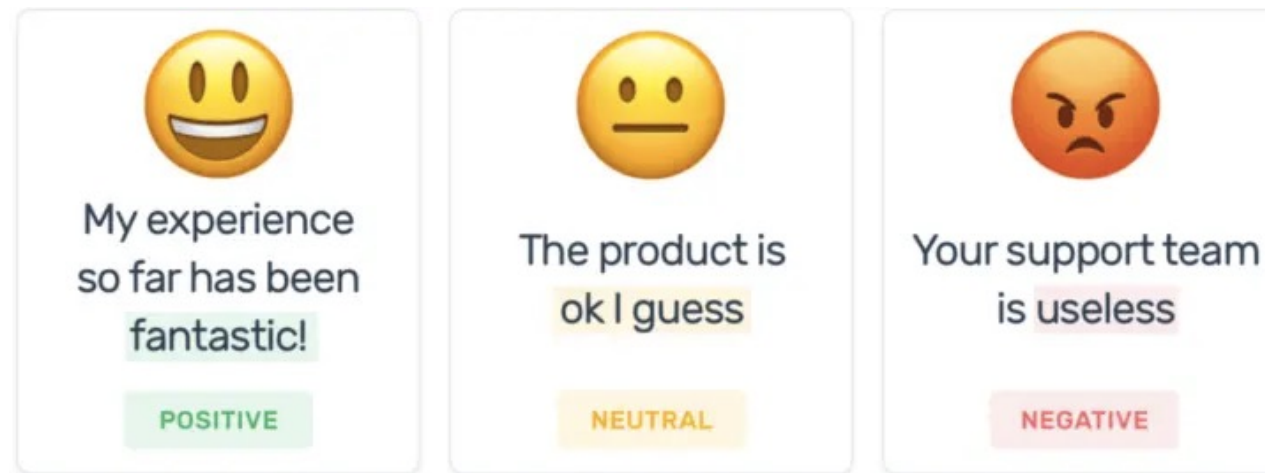
Data Mining Processes



- Anomaly Detection:
 - identification of unusual data
- Association Rule Learning:
 - searching for relationships between variables
- Clustering:
 - discovering groups and structures that are similar
- Classification:
 - applying a classifier to data
- Regression:
 - find a function that can be used to model data
- Summarization:
 - provide a more compact representation of a data set

Example – Sentiment Analysis

- Automated process used to identify positive, negative, and/or neutral opinions from text.



- Examples of text include social media posts, product reviews, survey responses, etc.

Example – Sentiment Analysis

- Here, we propose a data mining scenario:
 - Twitter is a widely-used social network where users give their opinions and thoughts on a broad range of topics.
 - There are millions of tweets posts sent out per day – meaning that there is a large amount of data generated and stored on Twitter's servers.
- Given a certain keyword, we parse through a collection of tweets containing that keyword and then we want to find out which other terms appear most frequently and perform sentiment analysis on the given keyword.



Example – Sentiment Analysis

- Here, we propose a data mining scenario:
 - ~~Twitter~~ X is a widely-used social network where users give their opinions and thoughts on a broad range of topics.
 - There are millions of tweets posts sent out per day – meaning that there is a large amount of data generated and stored on ~~Twitter's~~ X's servers.
- Given a certain keyword, we parse through a collection of tweets containing that keyword and then we want to find out which other terms appear most frequently and perform sentiment analysis on the given keyword.



- Application Programming Interface (API):
 - A set of functions and procedures (protocol) allowing the creation of applications that access the features or data of an operating system, application, or other service.
 - Abstract the underlying implementation so that the client only has access to necessary functions and nothing else.
- Many companies have APIs for interfacing with their own codebase. For example, X has APIs for streaming and searching for posts.

Example – Sentiment Analysis

- Pointwise Mutual Information – measure of how associated two words are:

$$\text{PMI}(\text{word}_1, \text{word}_2) = \log_2 \left[\frac{p(\text{word}_1 \& \text{word}_2)}{p(\text{word}_1) p(\text{word}_2)} \right]$$

- Semantic Orientation – the difference between a word's associations with positive and negative words:

$$SO(w) = \sum_{w' \in V^+} \text{PMI}(w, w') - \sum_{w' \in V^-} \text{PMI}(w, w')$$

- V^+ corresponds to a list of positive words.
- V^- corresponds to a list of negative words.

Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsupervised Classification of Reviews

Peter D. Turney

Institute for Information Technology
National Research Council of Canada
Ottawa, Ontario, Canada, K1A 0R6
peter.turney@nrc.ca

- APIs allow for different software/applications to communicate with one another (programming, internet, operating systems).
- Common examples:
 - Ridesharing (Uber, Lyft)
 - Application development (iOS, Android)
 - Mobile weather applications
 - Web development (embedding applications)
- Internal vs. External APIs



Simple Example of Using APIs



#12 ON TRENDING

This Video Has 3,600,735 Views

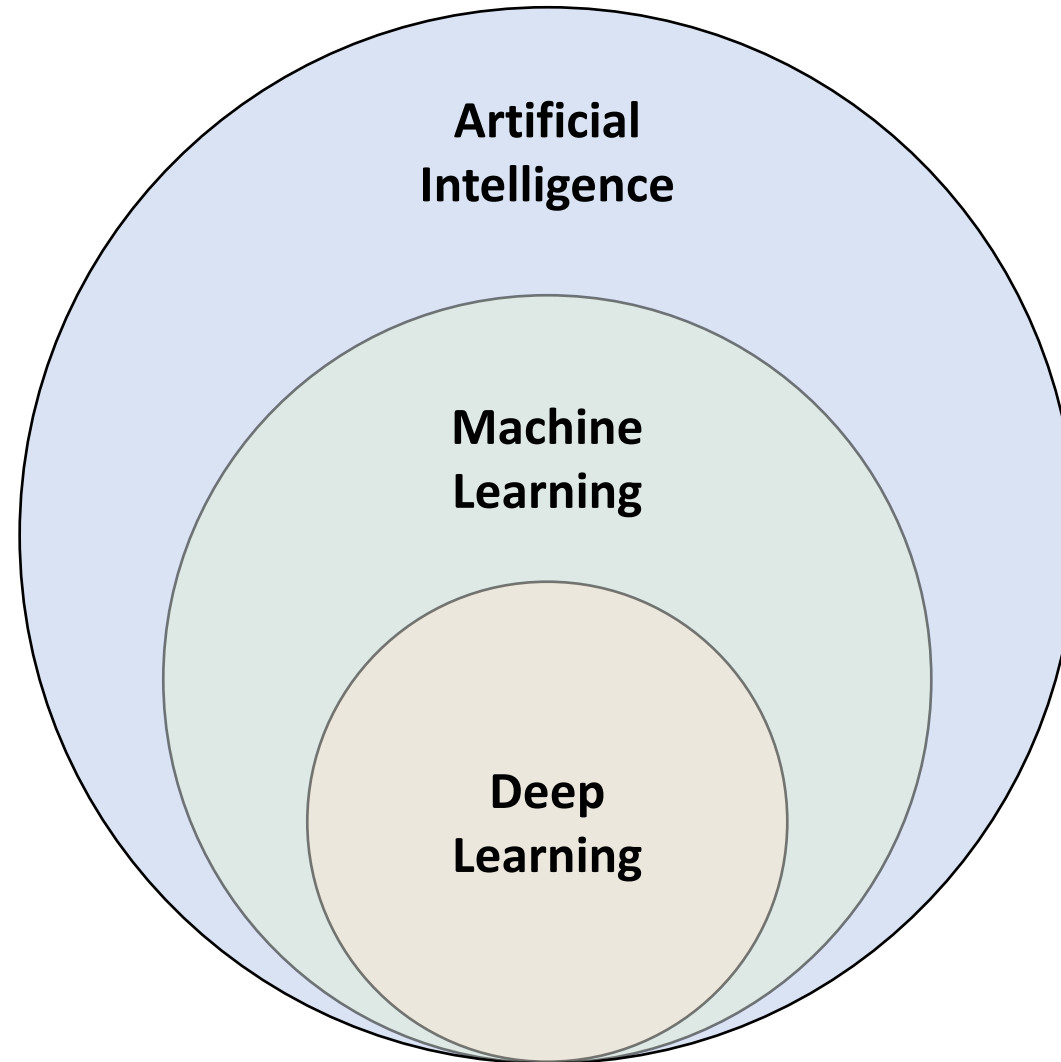
3,600,735 views • Apr 6, 2020

 289K  2.5K  SHARE  SAVE ...

- With data mining, we can learn from large amounts of data and gain insights on patterns that occur in the data.

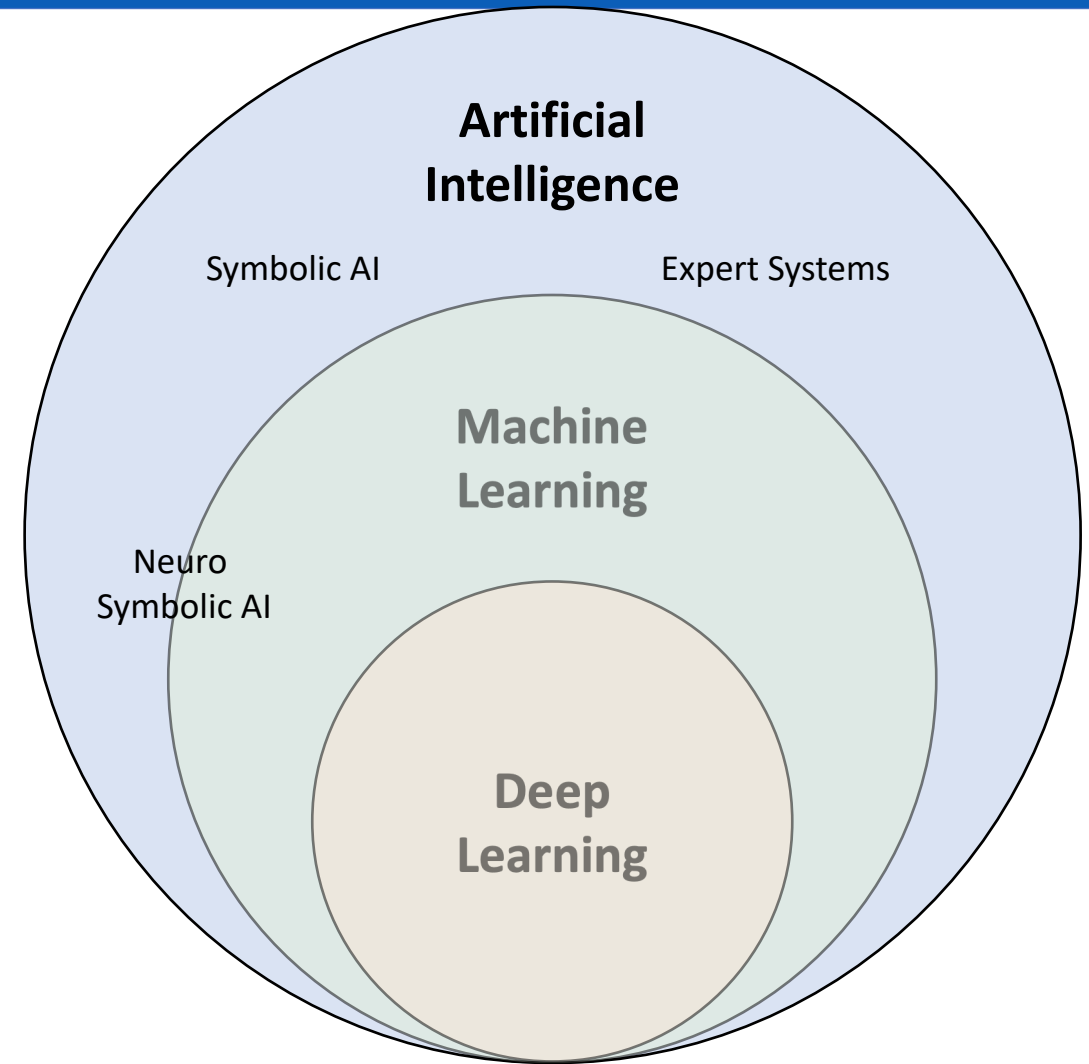
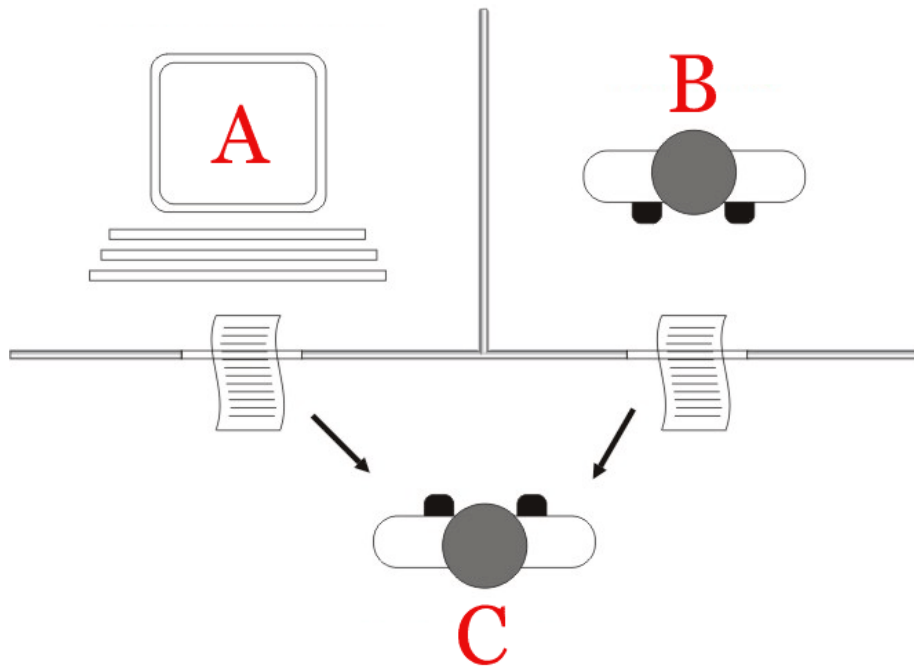


- Can we extend data mining techniques further?



Theory and development of computer systems to perform tasks requiring human intelligence [1]

Can machines think ? [1]



[1] Turing, A. M. I.—COMPUTING MACHINERY AND INTELLIGENCE. *Mind* **LIX**, 433–460 (1950).

[2] Russell, S. J., Norvig, P. & Davis, E. *Artificial intelligence: a modern approach*. (Prentice Hall, 2010).

[3] Garnelo, M. & Shanahan, M. Reconciling deep learning with symbolic artificial intelligence: representing objects and relations. *Current Opinion in Behavioral Sciences* **29**, 17–23 (2019).

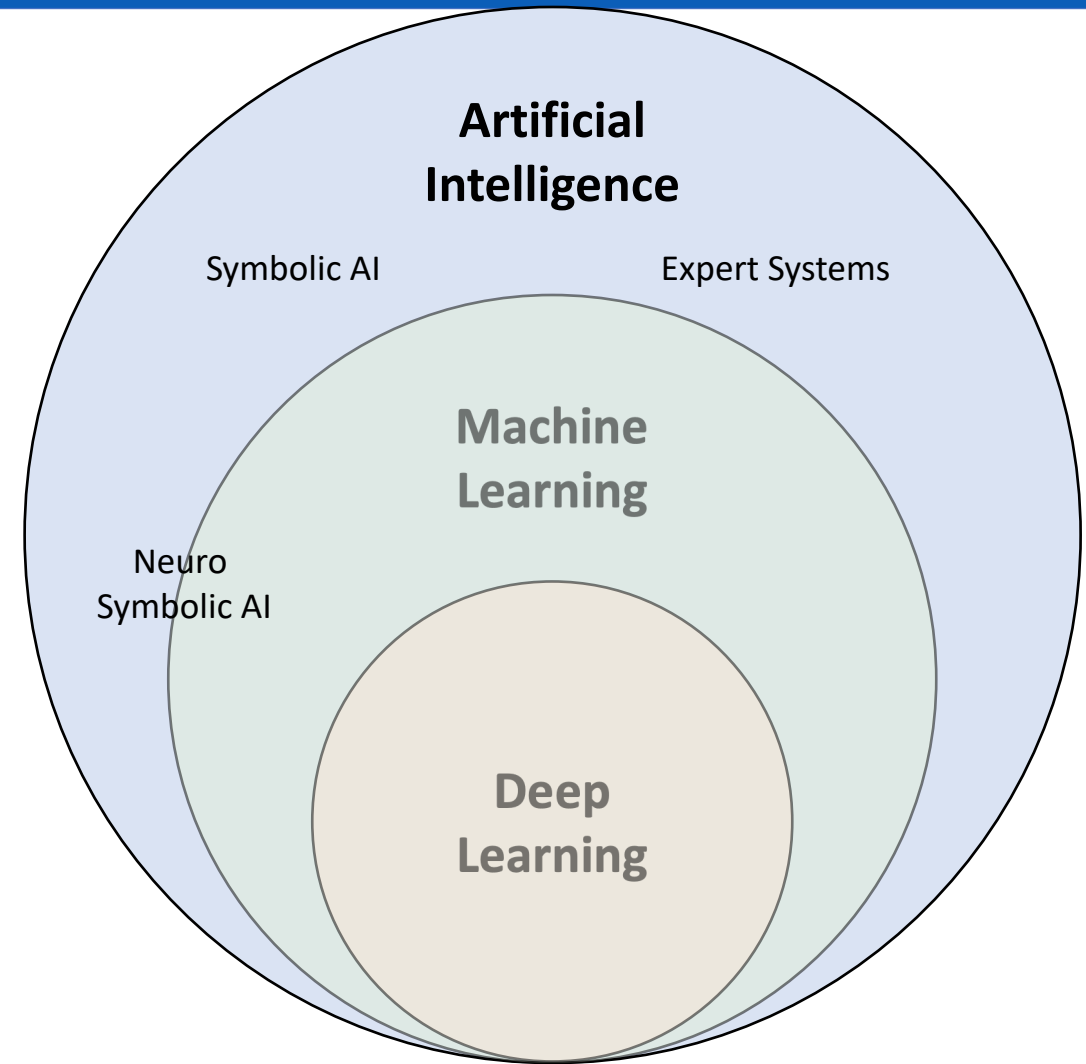
Theory and development of computer systems to perform tasks requiring human intelligence [1]

Can machines think ? [1]

Symbolic AI : Collection of methods that carry out logical reasoning by symbolic representations of problems. [2]

Expert Systems : Computer programs that emulate human decision making through large bodies of conditional statements (Rule based systems)

Neuro-symbolic AI : Combines strength of symbolic AI with deep learning architectures.[3]



[1] Turing, A. M. I.—COMPUTING MACHINERY AND INTELLIGENCE. *Mind* **LIX**, 433–460 (1950).

[2] Russell, S. J., Norvig, P. & Davis, E. *Artificial intelligence: a modern approach*. (Prentice Hall, 2010).

[3] Garnelo, M. & Shanahan, M. Reconciling deep learning with symbolic artificial intelligence: representing objects and relations. *Current Opinion in Behavioral Sciences* **29**, 17–23 (2019).

Use and development of computer systems that are able to learn and adapt without following explicit instructions. (Oxford dictionary)

IMPORTANT ML LINGO

Feature : Measurable property

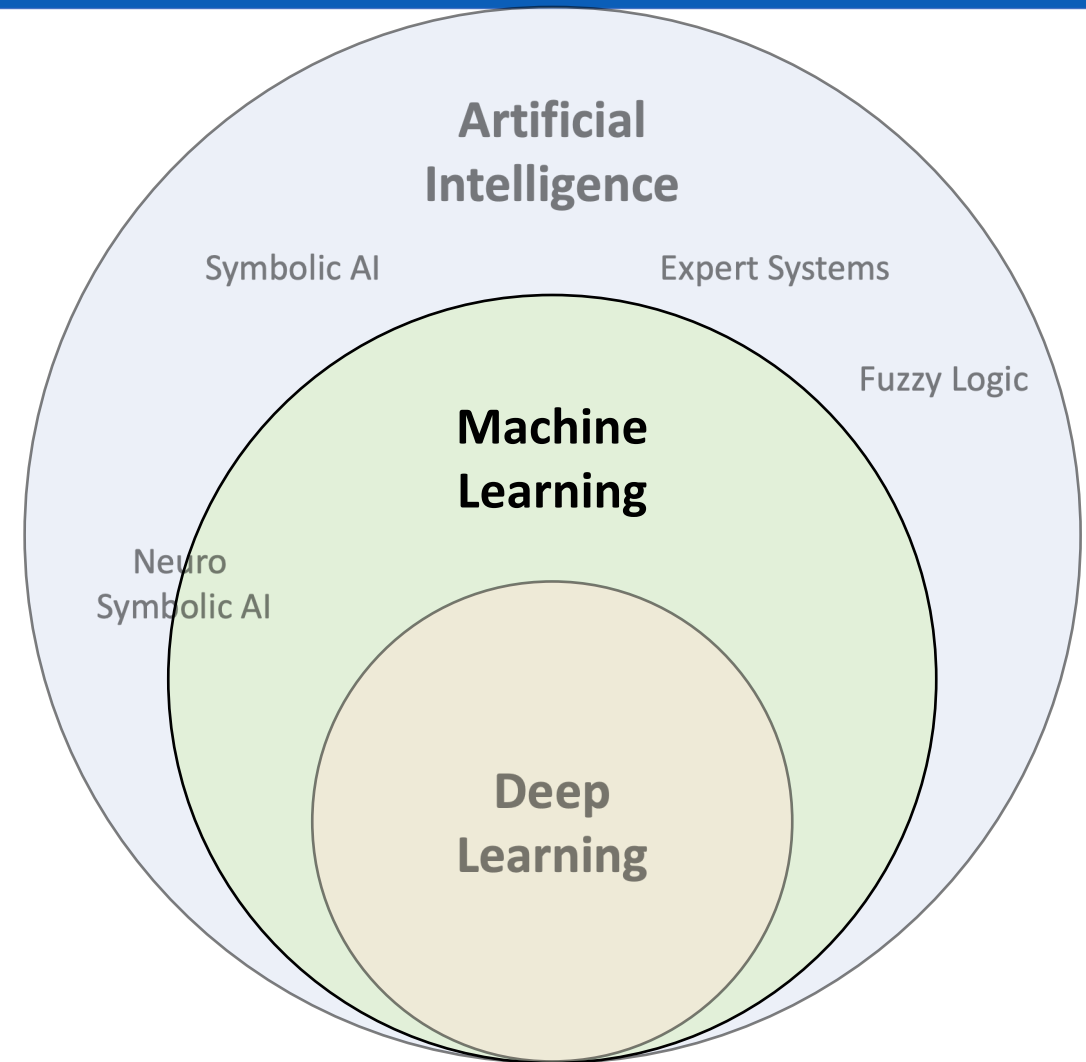
Labels : Value associated with a feature or set of features

Dataset : Set of features and labels.

Training set : Set of features and labels the model is trained on.

Test set : Set of features and labels that we use to evaluate the model predictions.

Hyperparameters : Model parameters

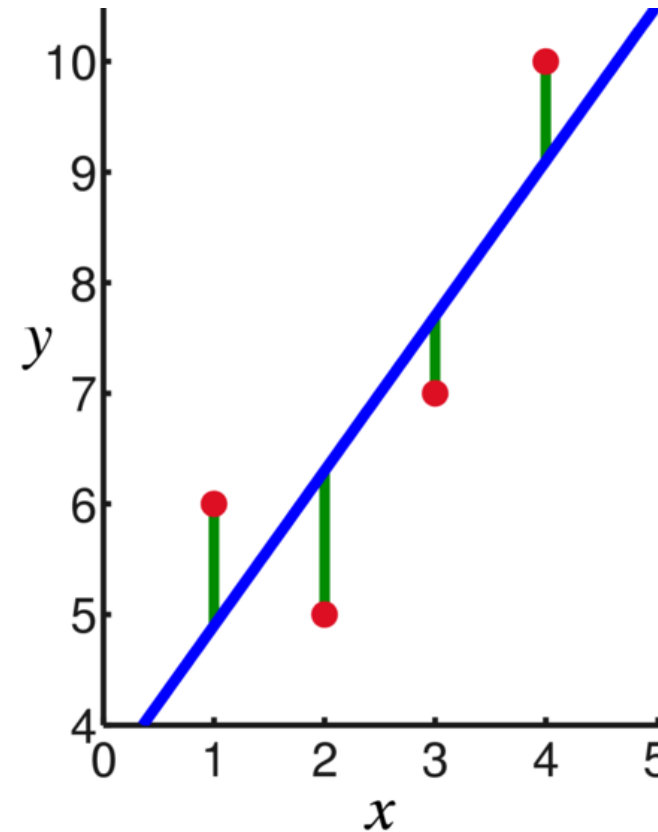


The goal:

- Given sampled points, can we find a best-fit line?

The process:

- Using two parameters (m and b), find the line that minimizes the deviation.



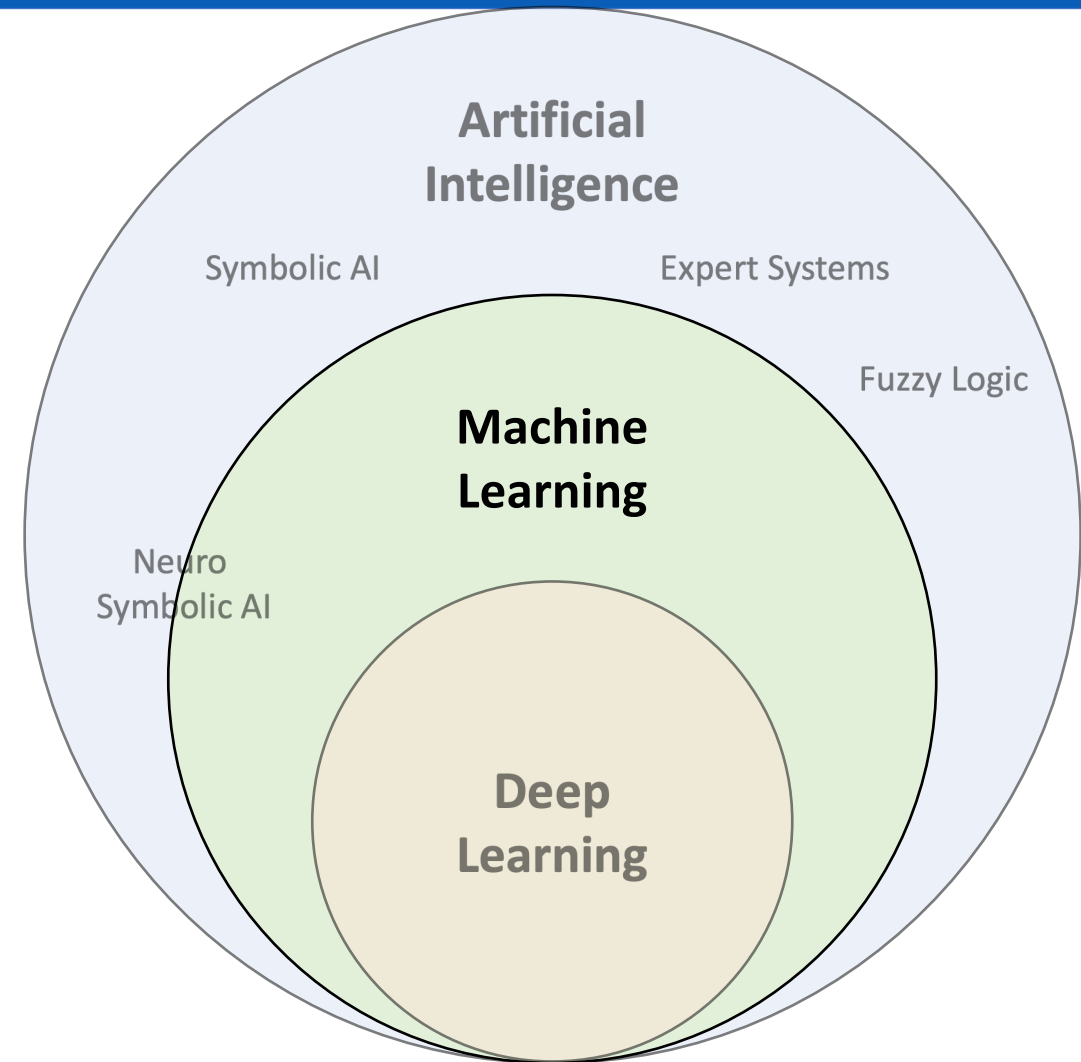
Numerical Data!

Task : Predict the number of fruits sold

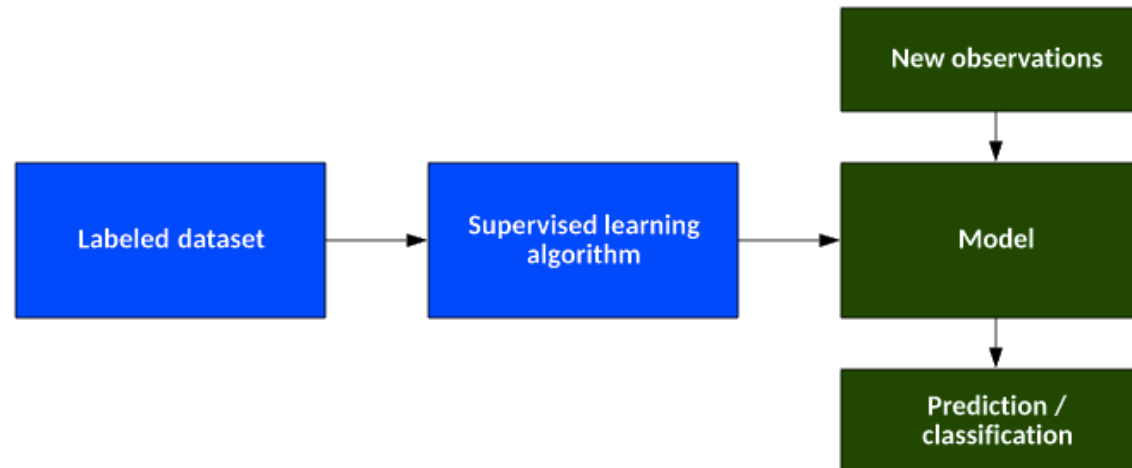
Training Set	Fruit (Feature 1)	Day of week (Feature 2)	Price (in \$) (Feature 3)	Num sold (Label)
	'apple'	'Monday'	1	20
	'banana'	'Monday'	0.30	20
	'apple'	'Tuesday'	1	10
	'pear'	'Wednesday'	1.5	5
	'oranges'	'Thursday'	1.5	10

Note :

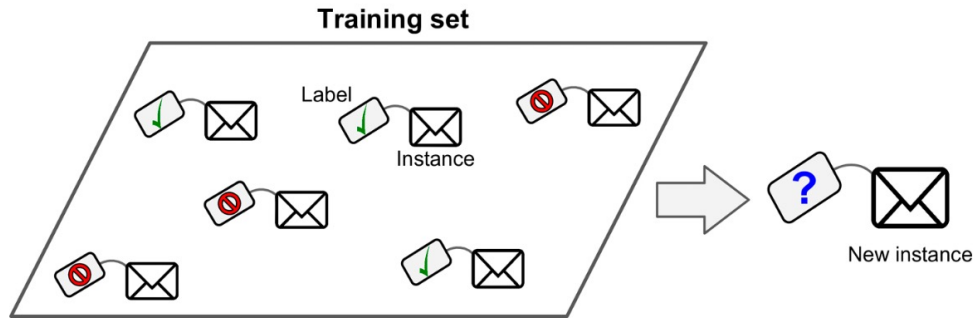
- A machine learning algorithm can only work with numerical data.
- Need to convert all the categories of fruits and days of weeks into numerical vectors (One hot encoding)



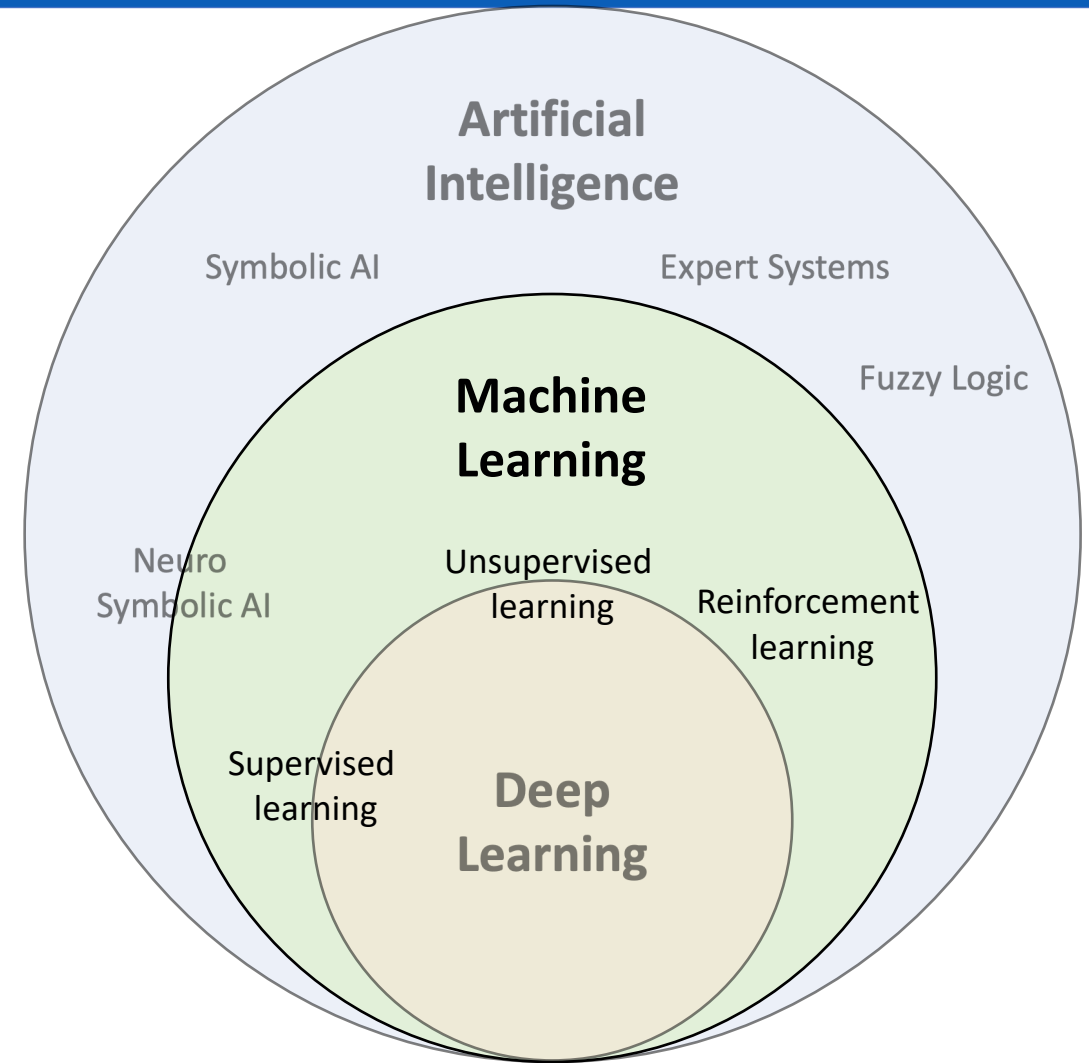
- We have a dataset with inputs and desired outputs.
- We train a model to predict the output from the input (normally using optimization techniques to minimize an error function).
- We use a different data set to test the model and see how accurate it is.



Classification

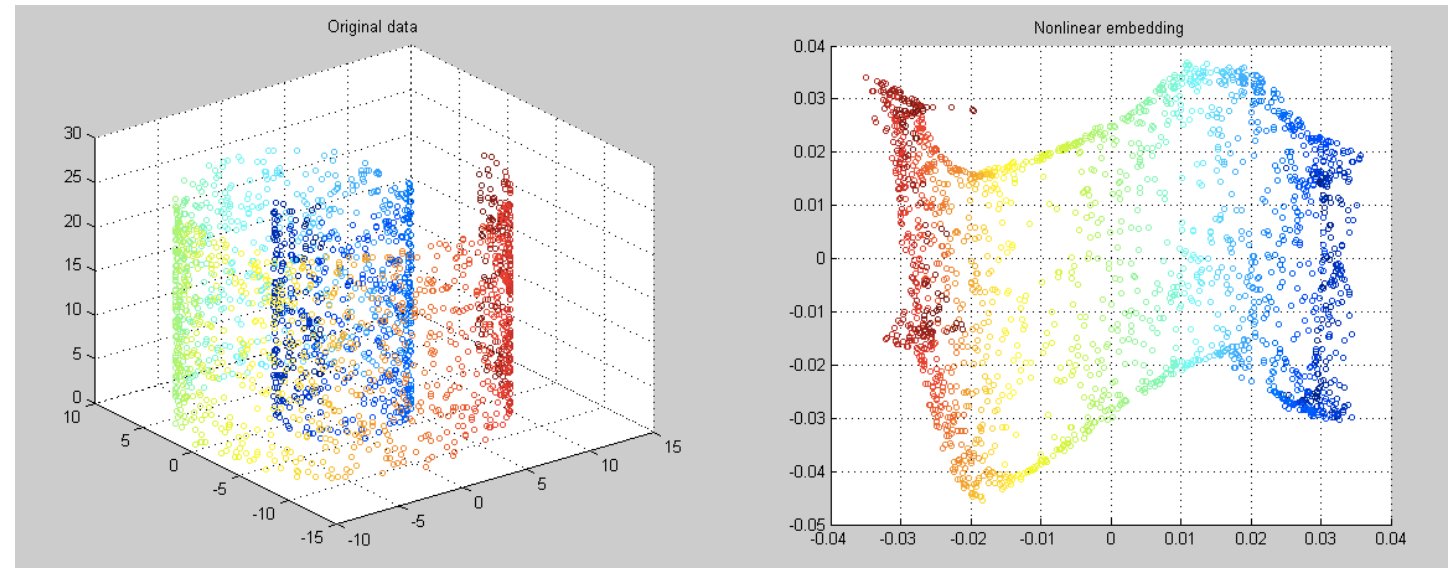


Regression

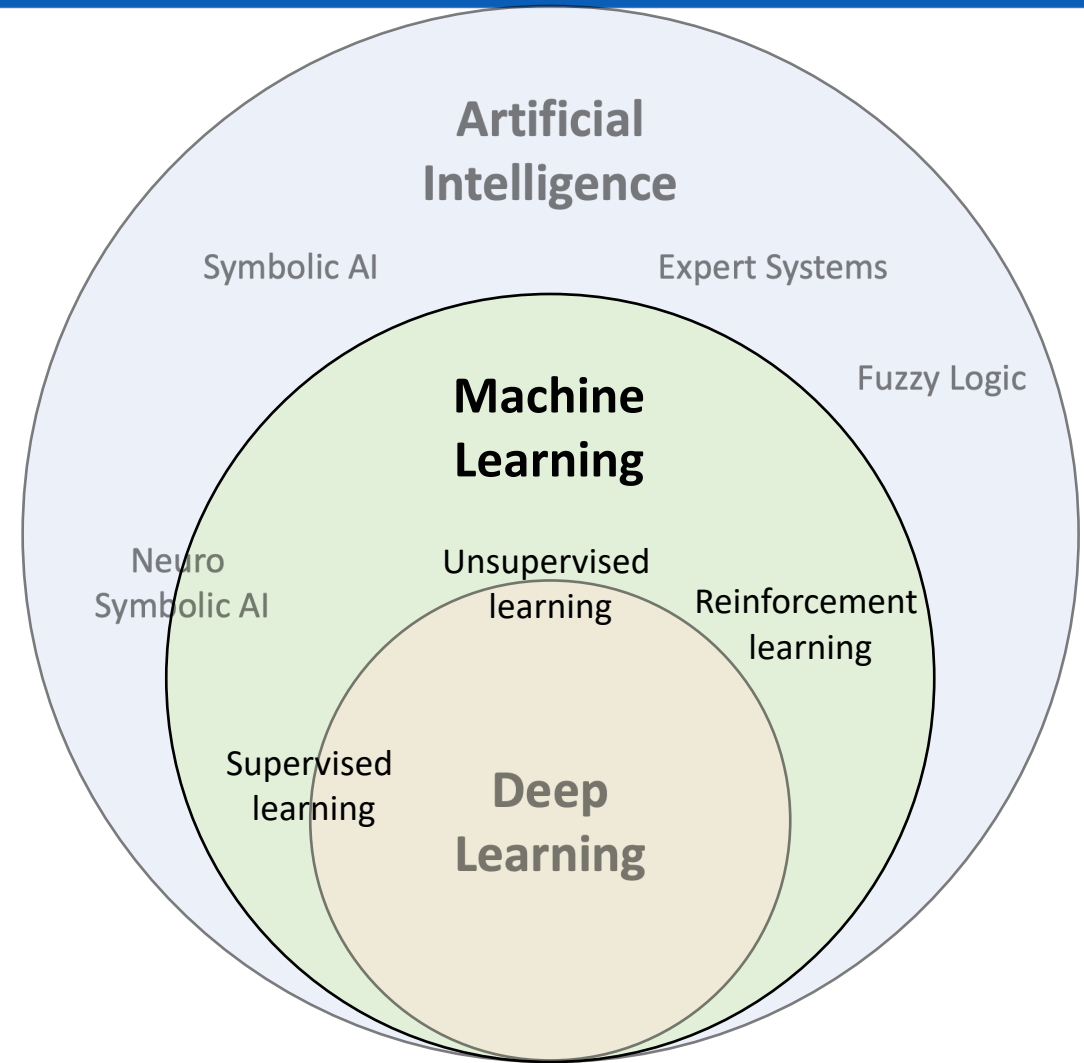
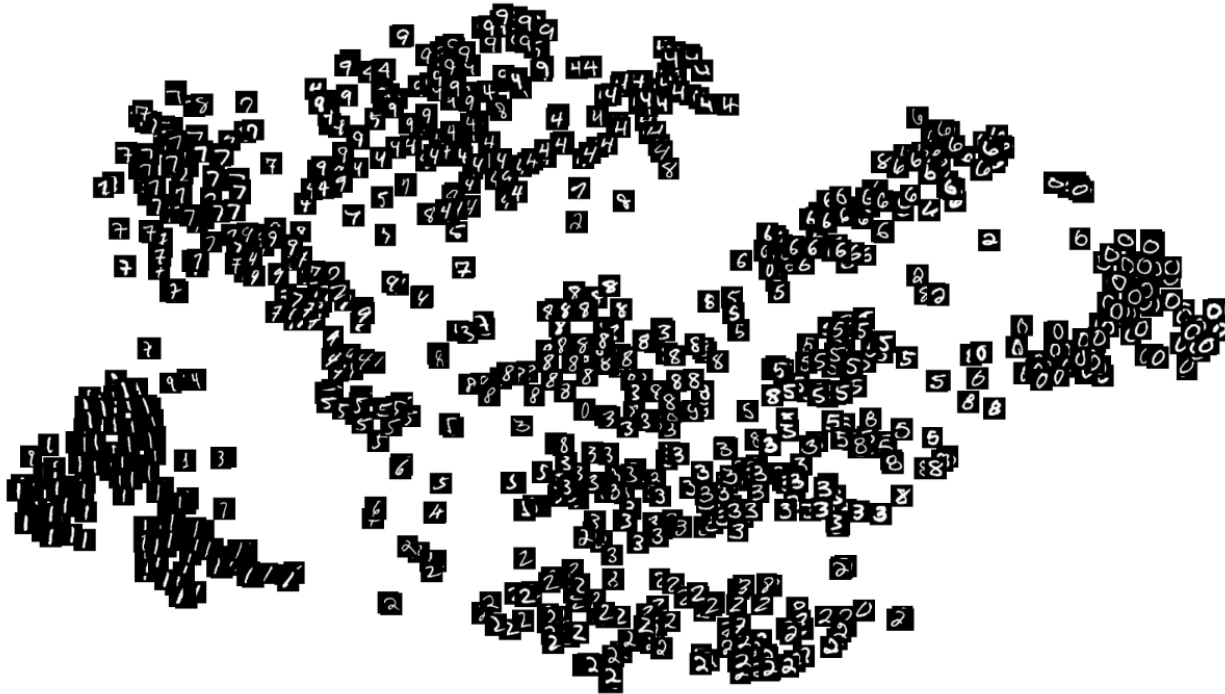


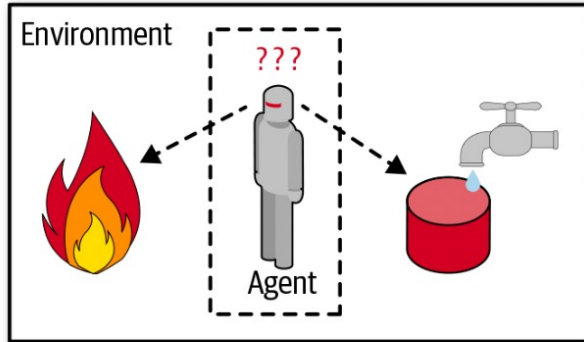
- We have a data set with only inputs.
- The algorithm finds patterns in the data and reacts accordingly.

- Examples:
 - Clustering analysis
 - Principal component analysis

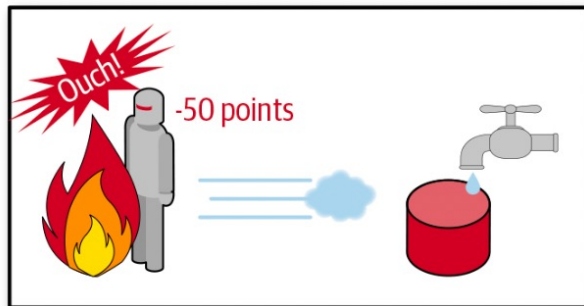


Unsupervised Learning

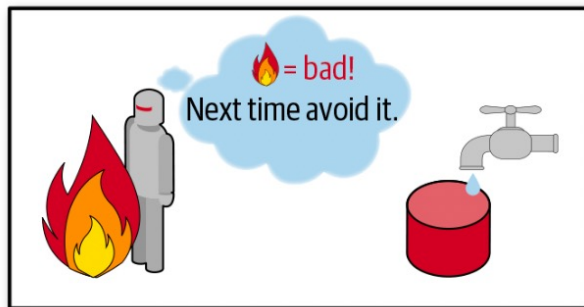




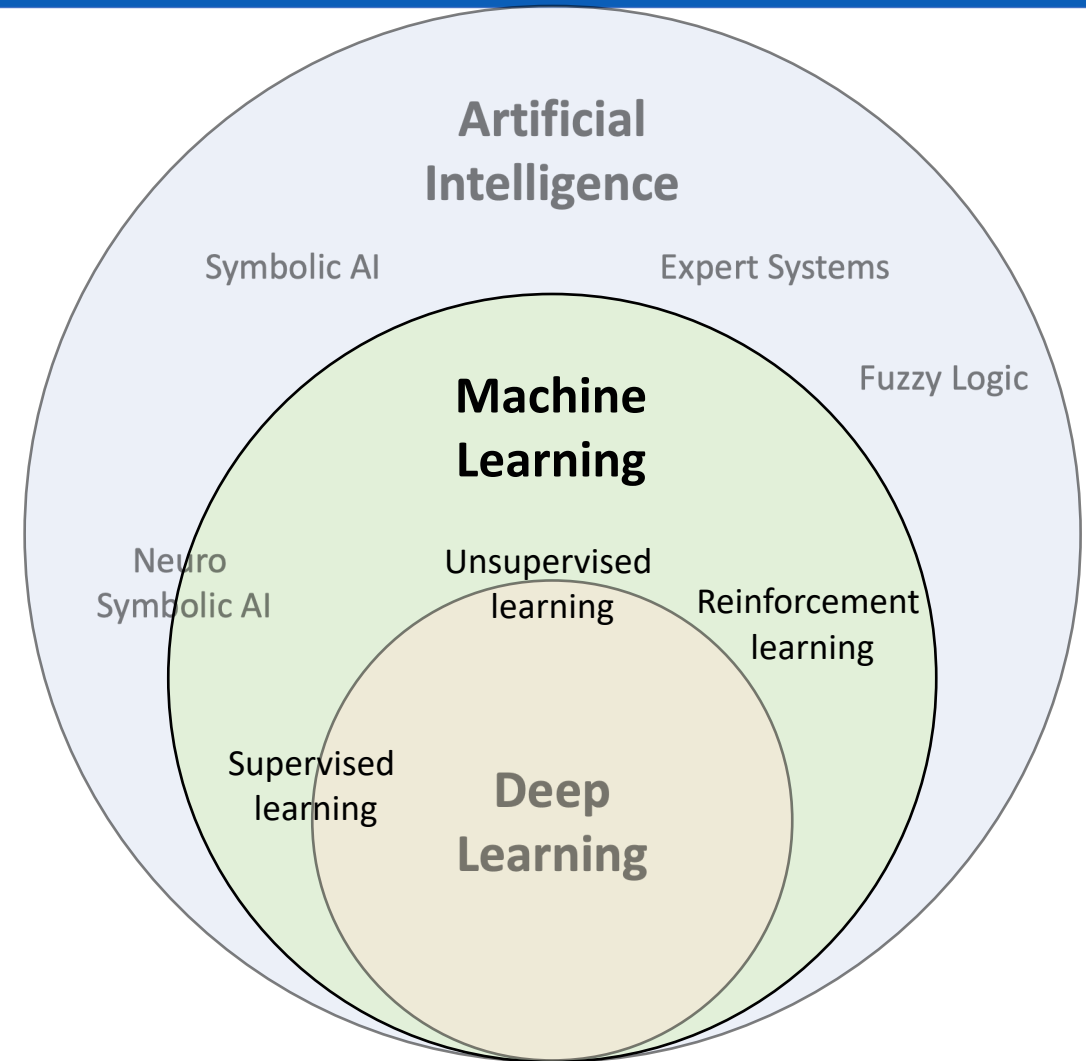
- 1 Observe
- 2 Select action using policy



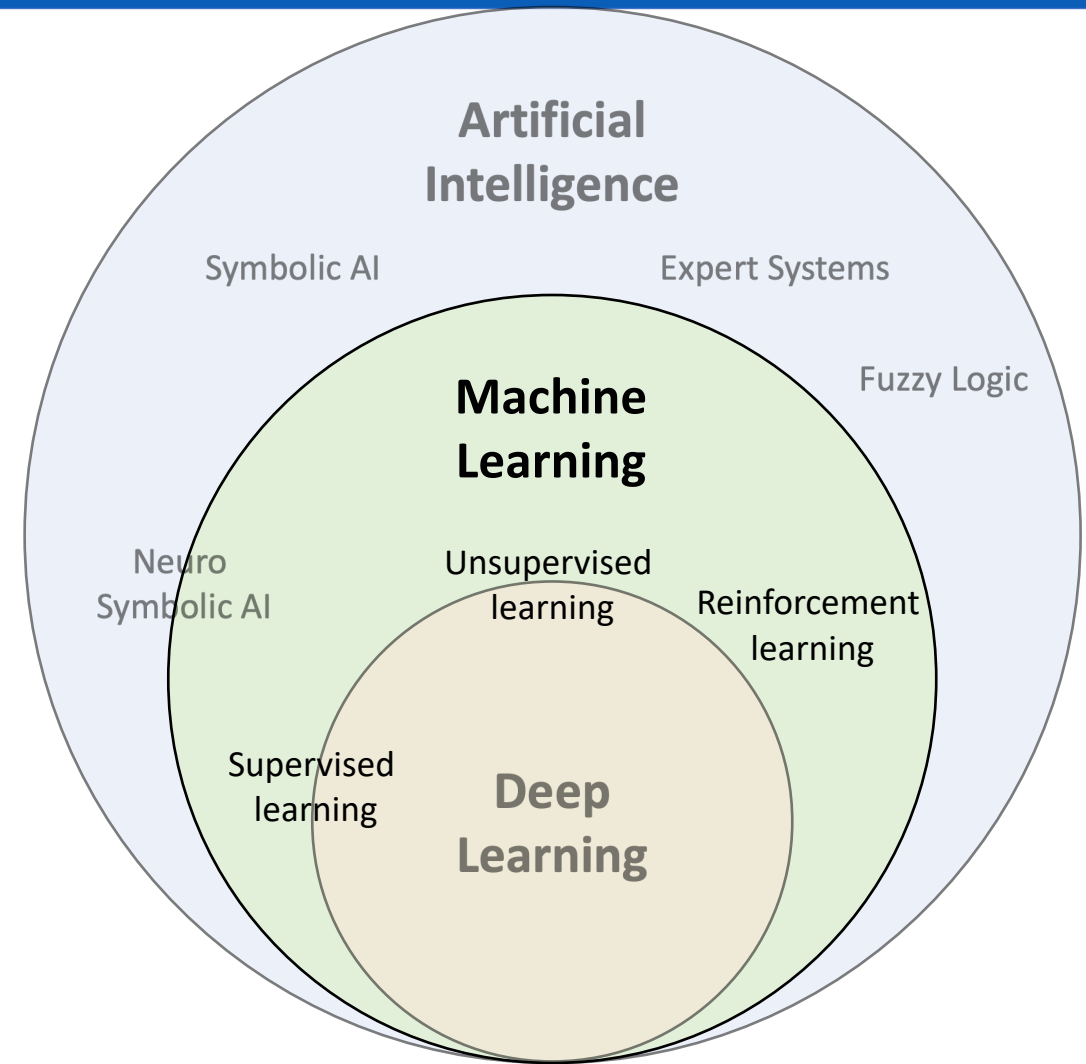
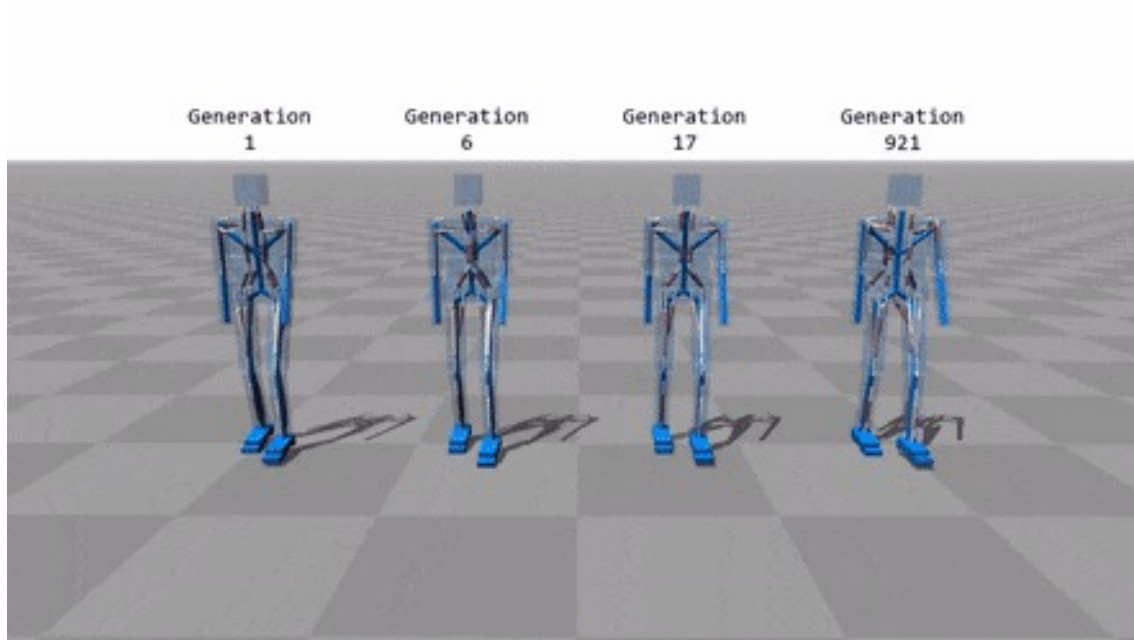
- 3 Action!
- 4 Get reward or penalty

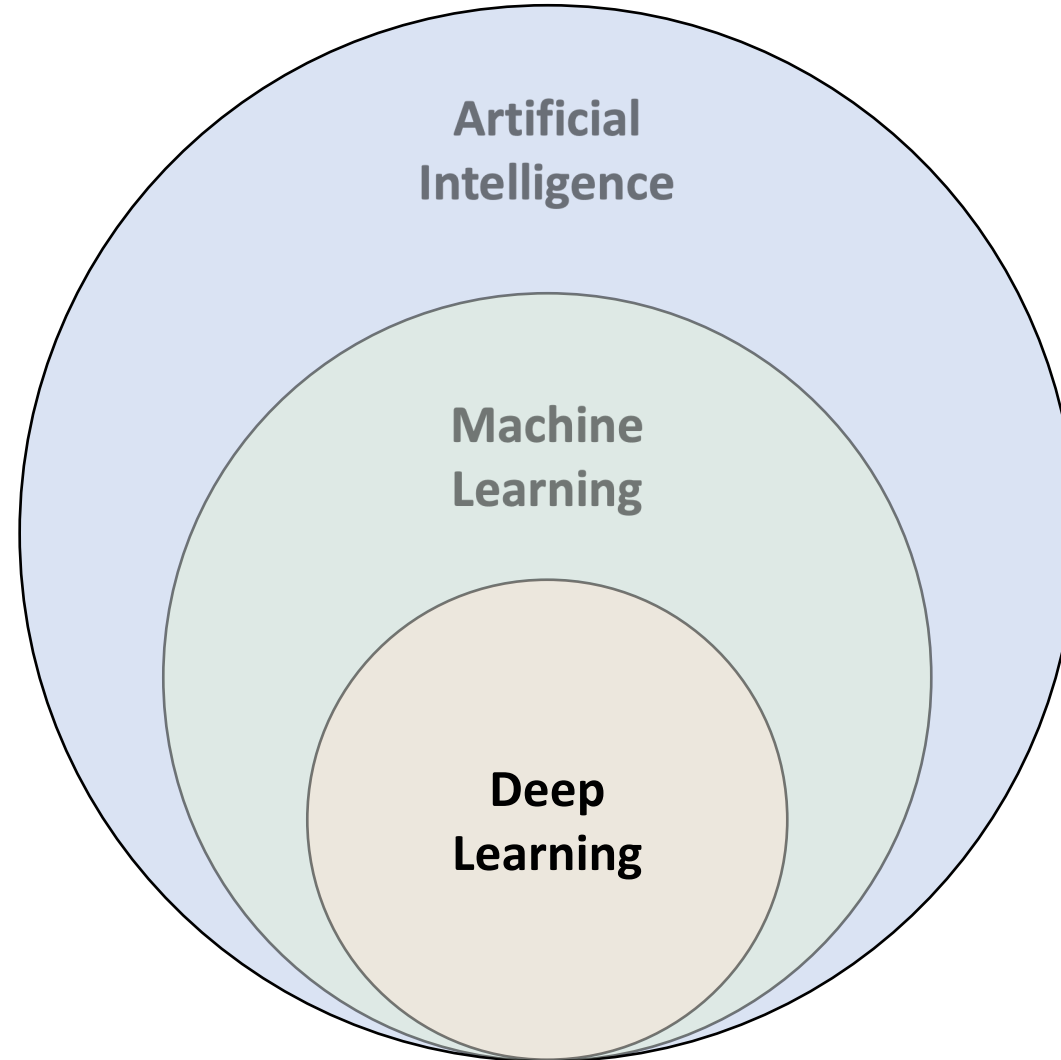


- 5 Update policy (learning step)
- 6 Iterate until an optimal policy is found



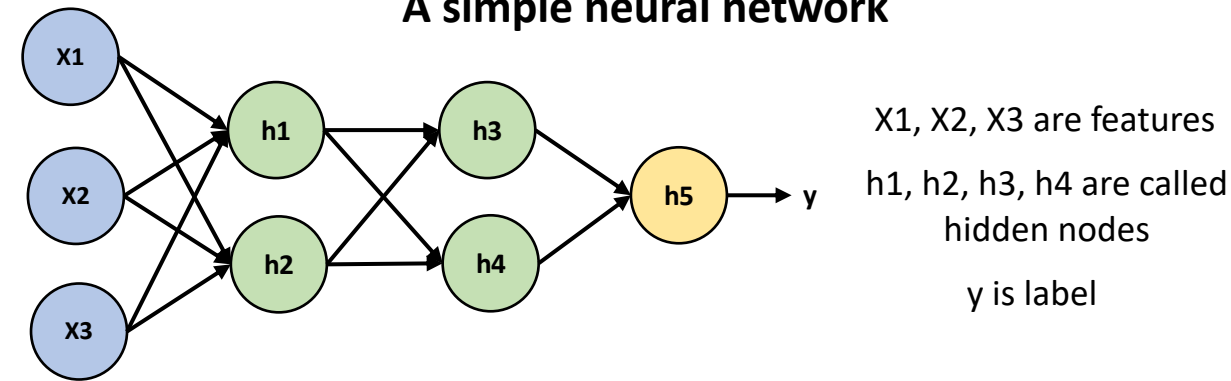
Unsupervised Learning



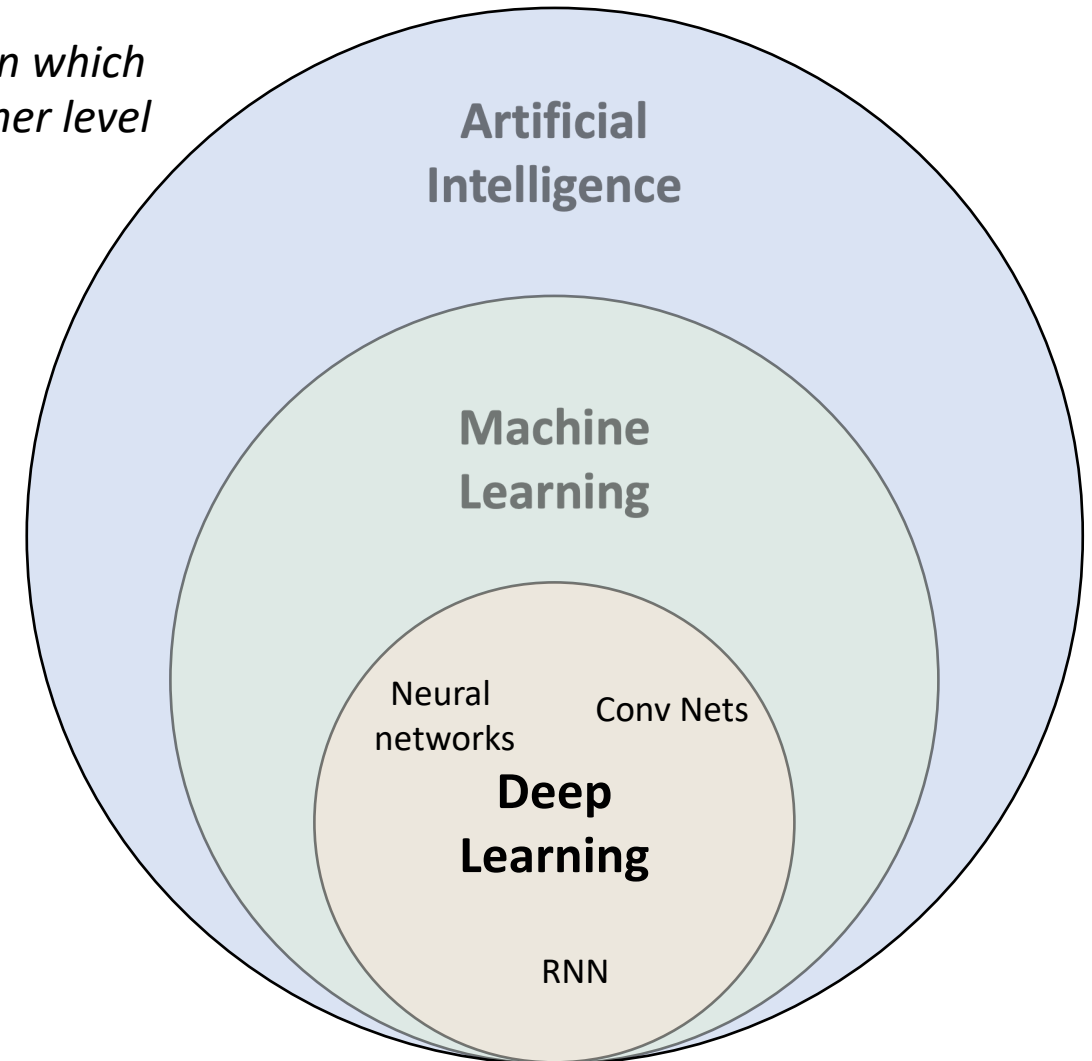
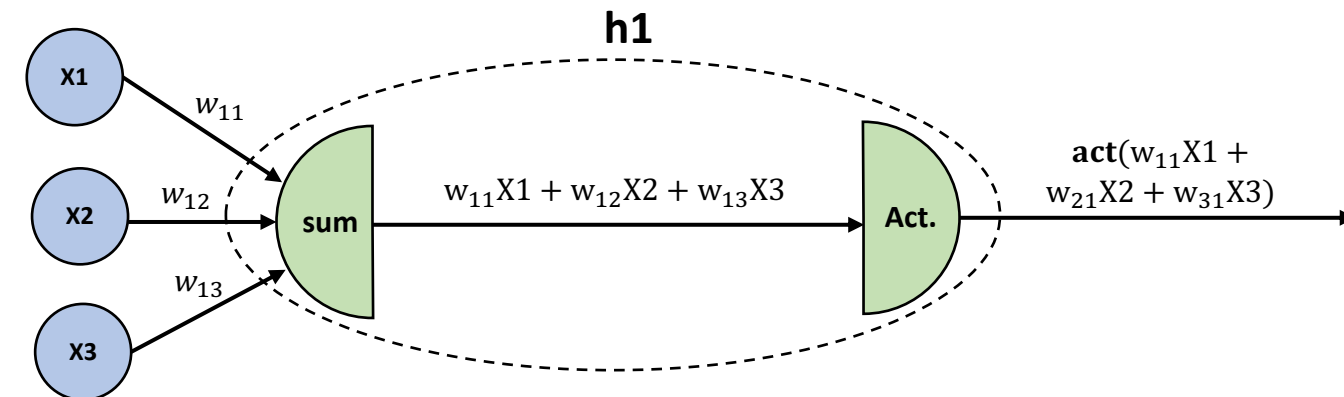


A subset of machine learning based on artificial neural networks in which multiple layers of processing are used to extract progressively higher level of features from data (Oxford dictionary)

A simple neural network

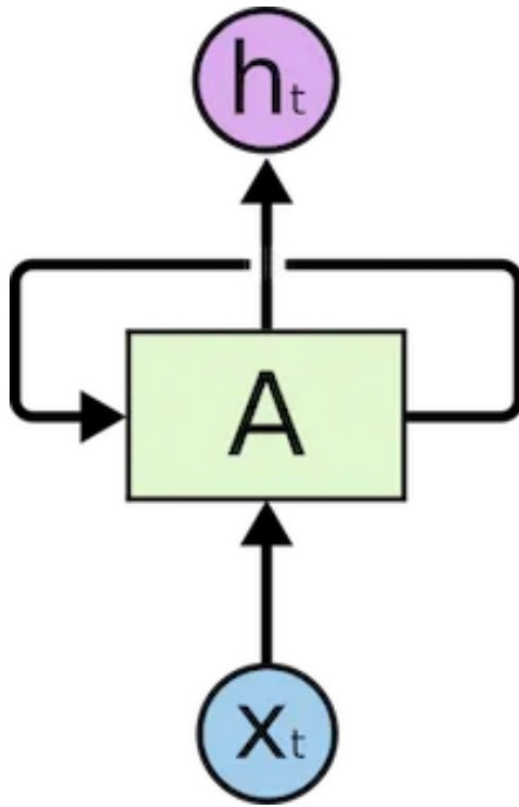


x_1, x_2, x_3 are features
 h_1, h_2, h_3, h_4 are called hidden nodes
 y is label

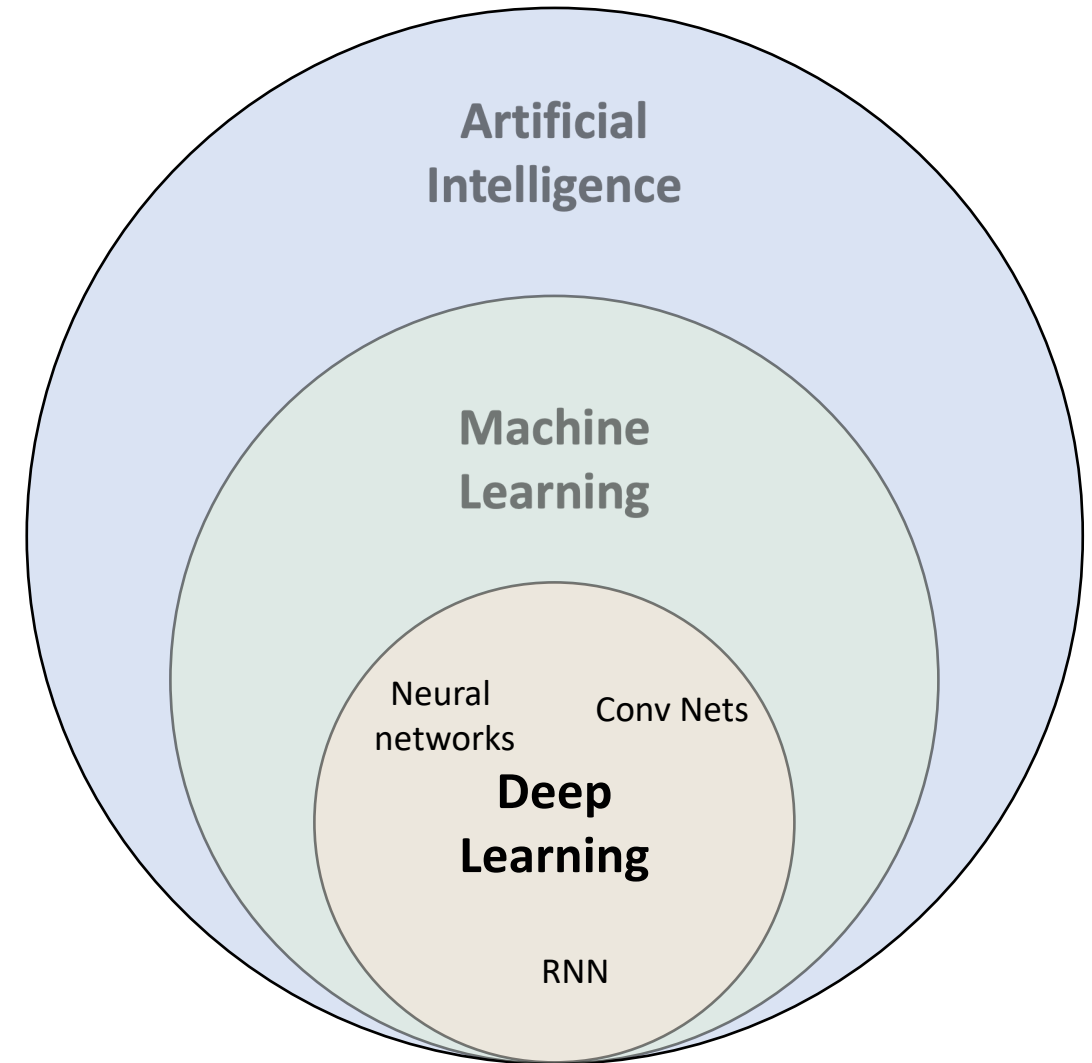


Training a neural network is simply learning the appropriate weights

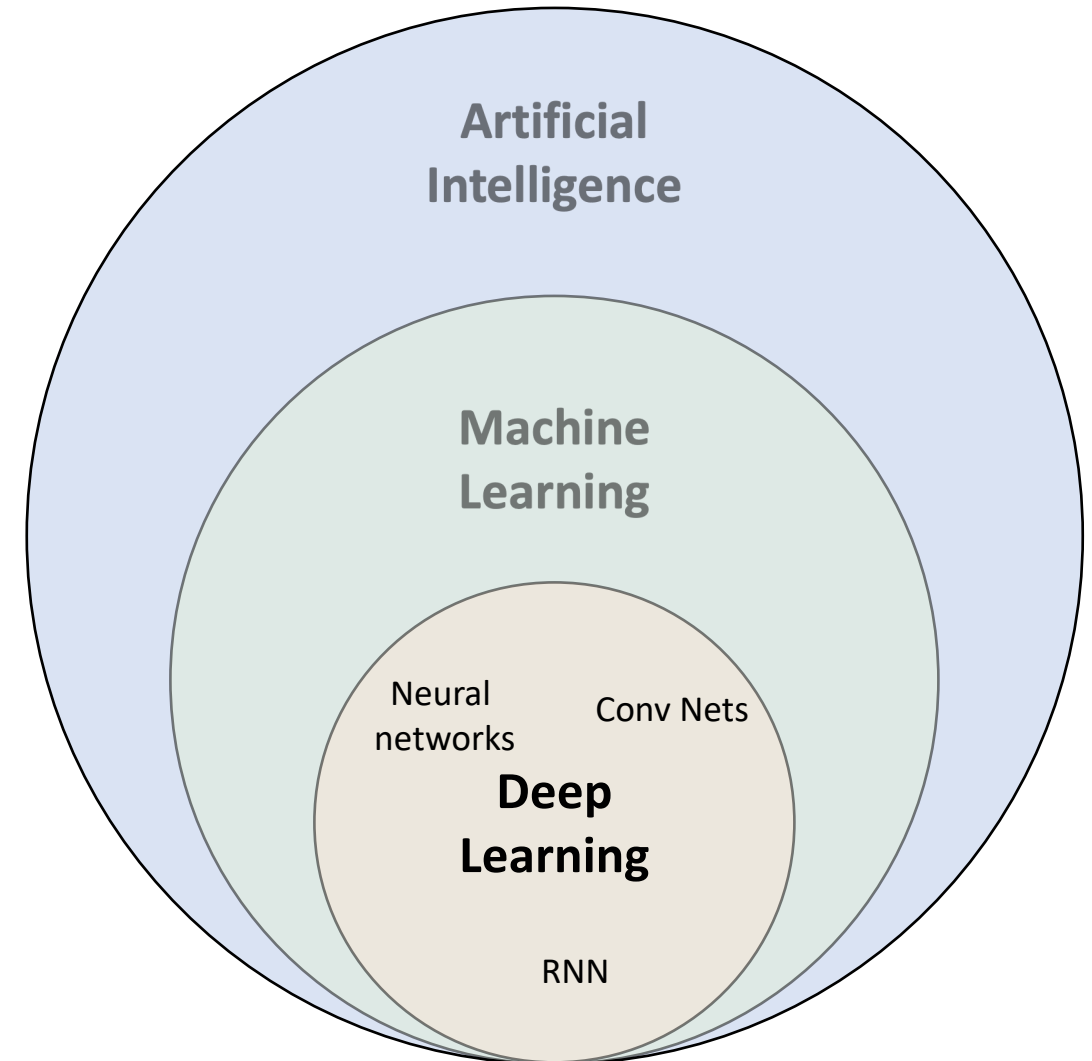
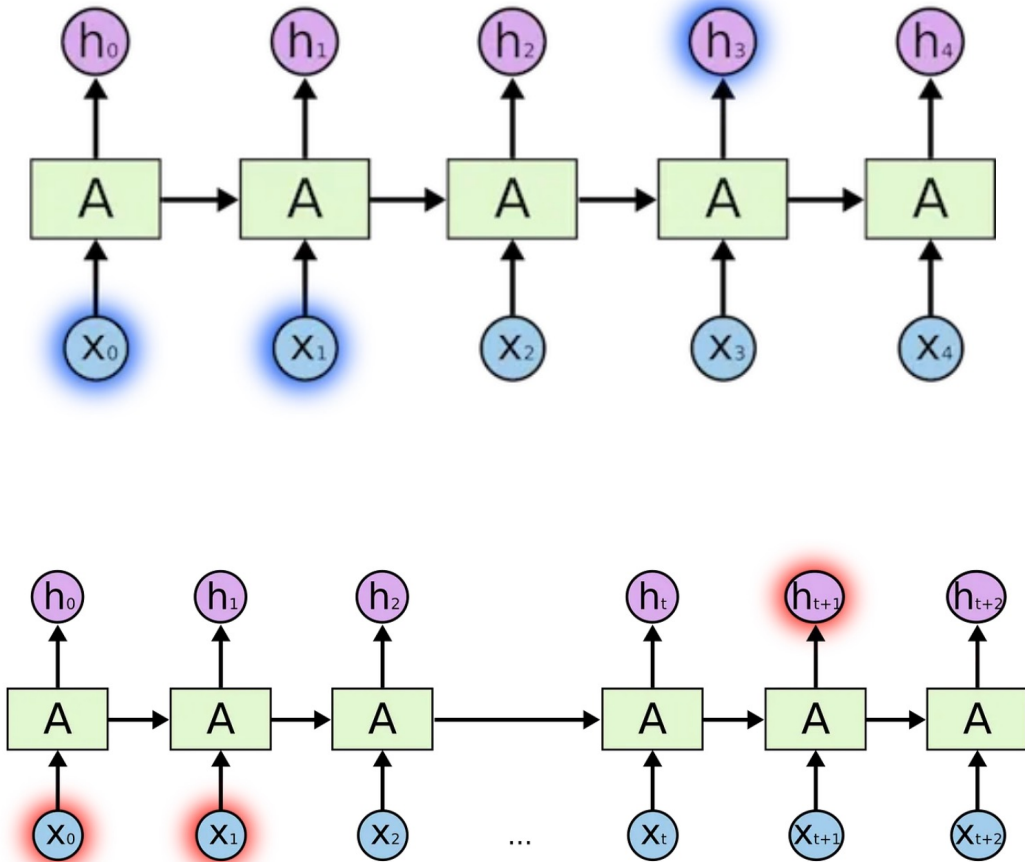
Recurrent Neural Net



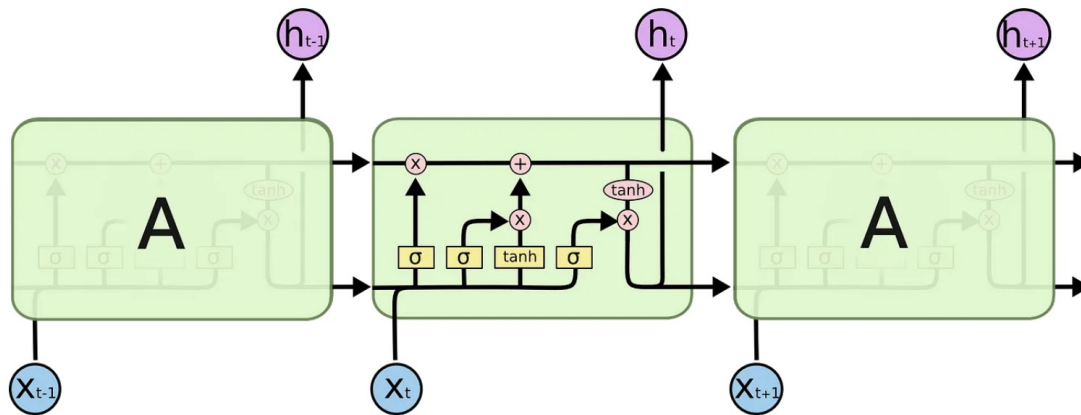
The input is represented as x_t



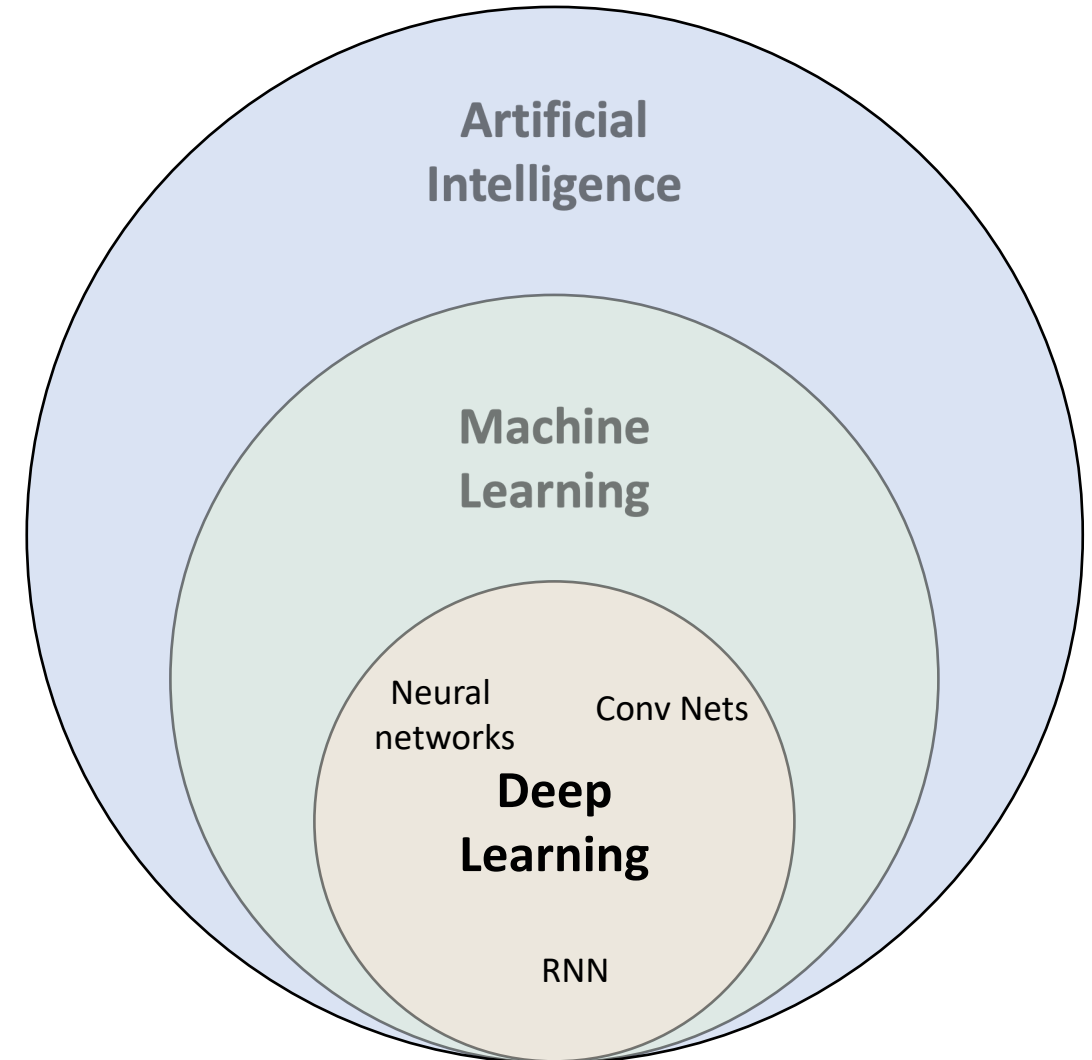
Recurrent Neural Nets ... are Inefficient



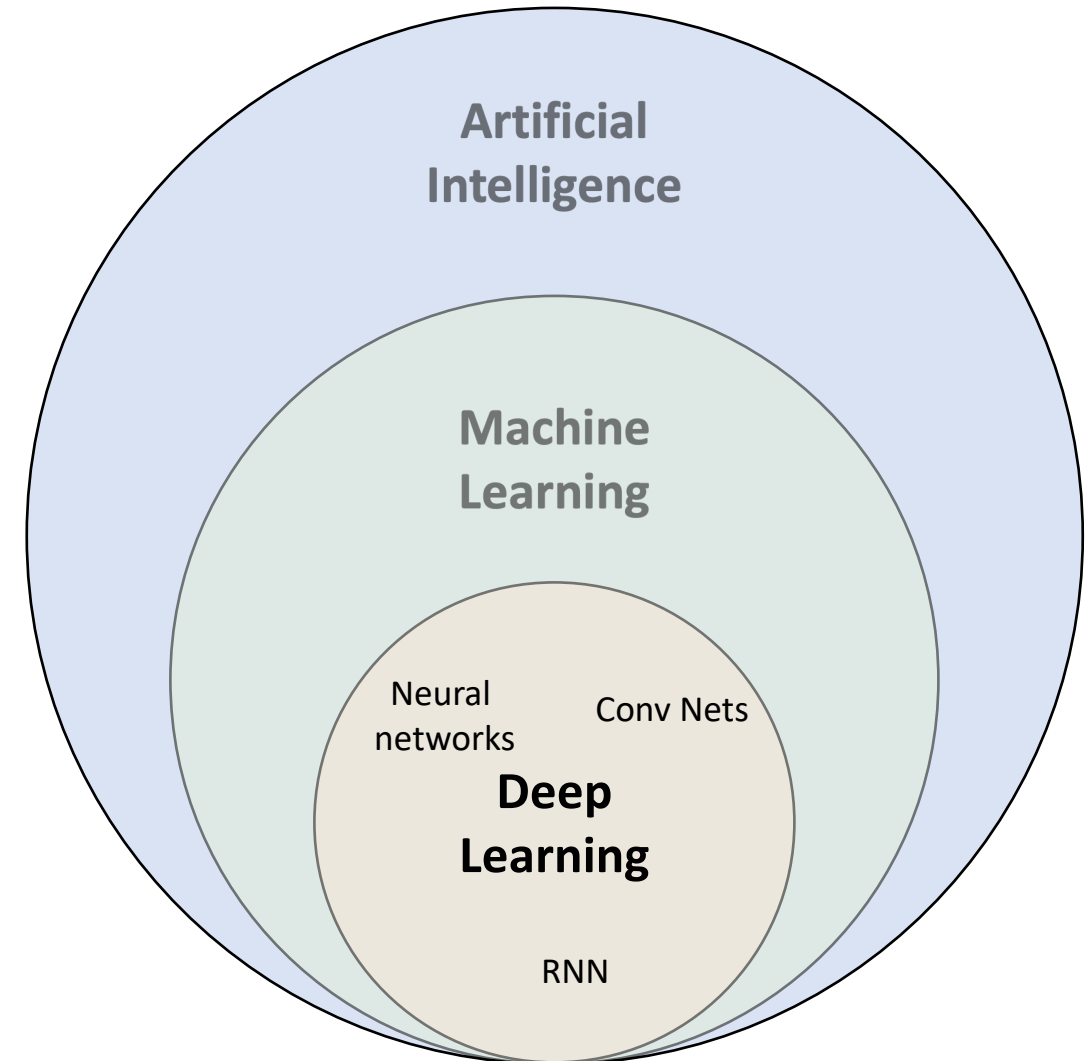
Long-Short Term Memory



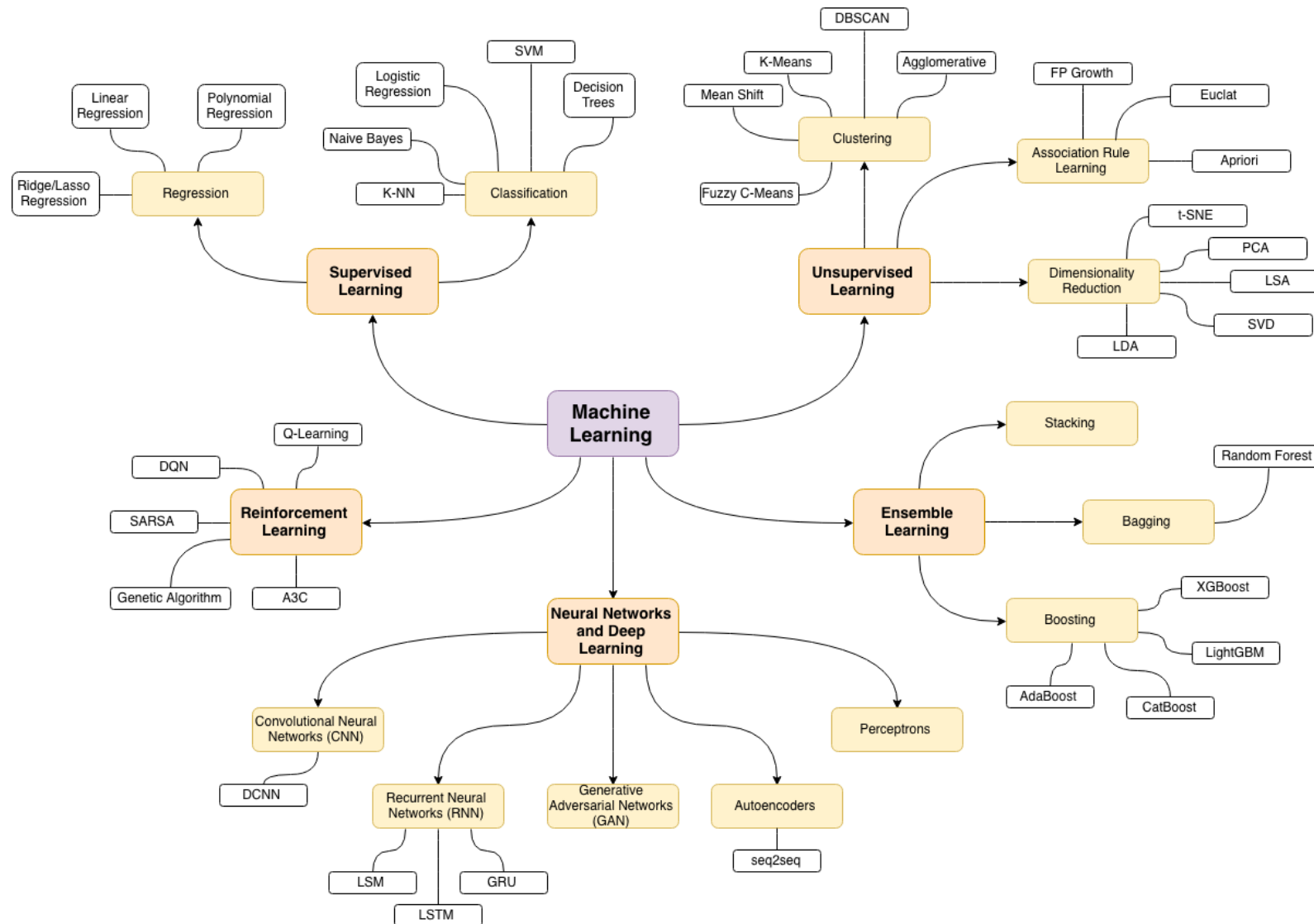
- Cannot parallelize
- No explicit modeling of long & short range dependencies
- Distance between positions is linear



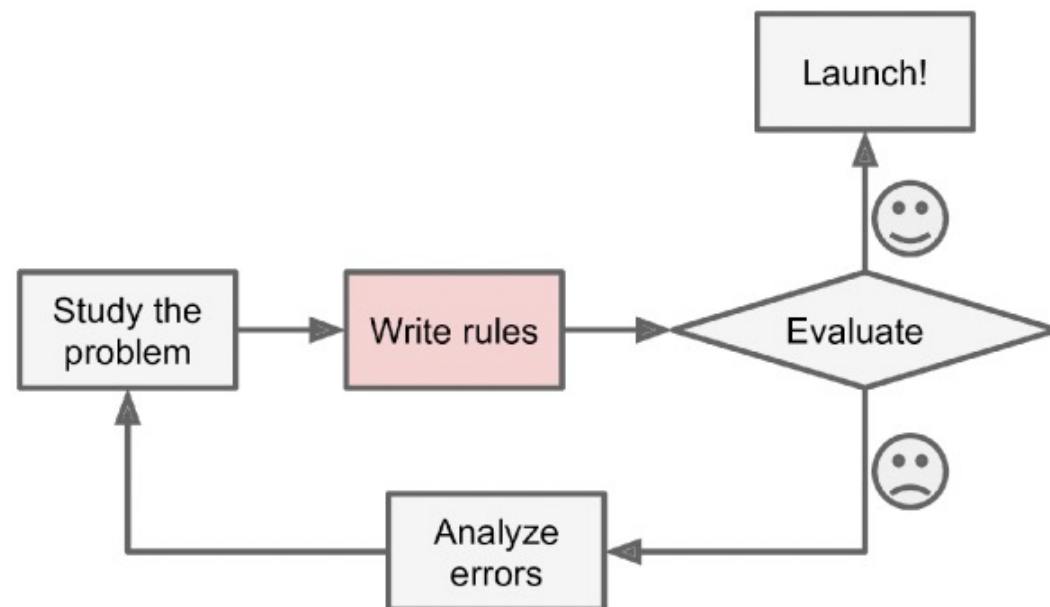
https://www.youtube.com/watch?v=fjJOgb-E41w&ab_channel=GoogleCloudTech



Machine Learning... in a Nut Shell

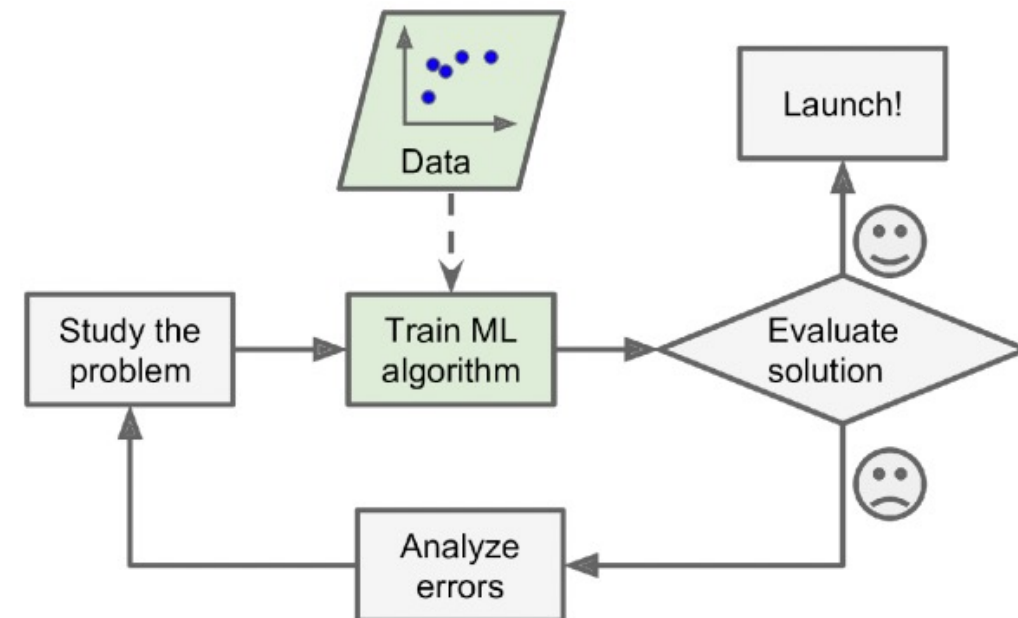


Simplified traditional program workflow

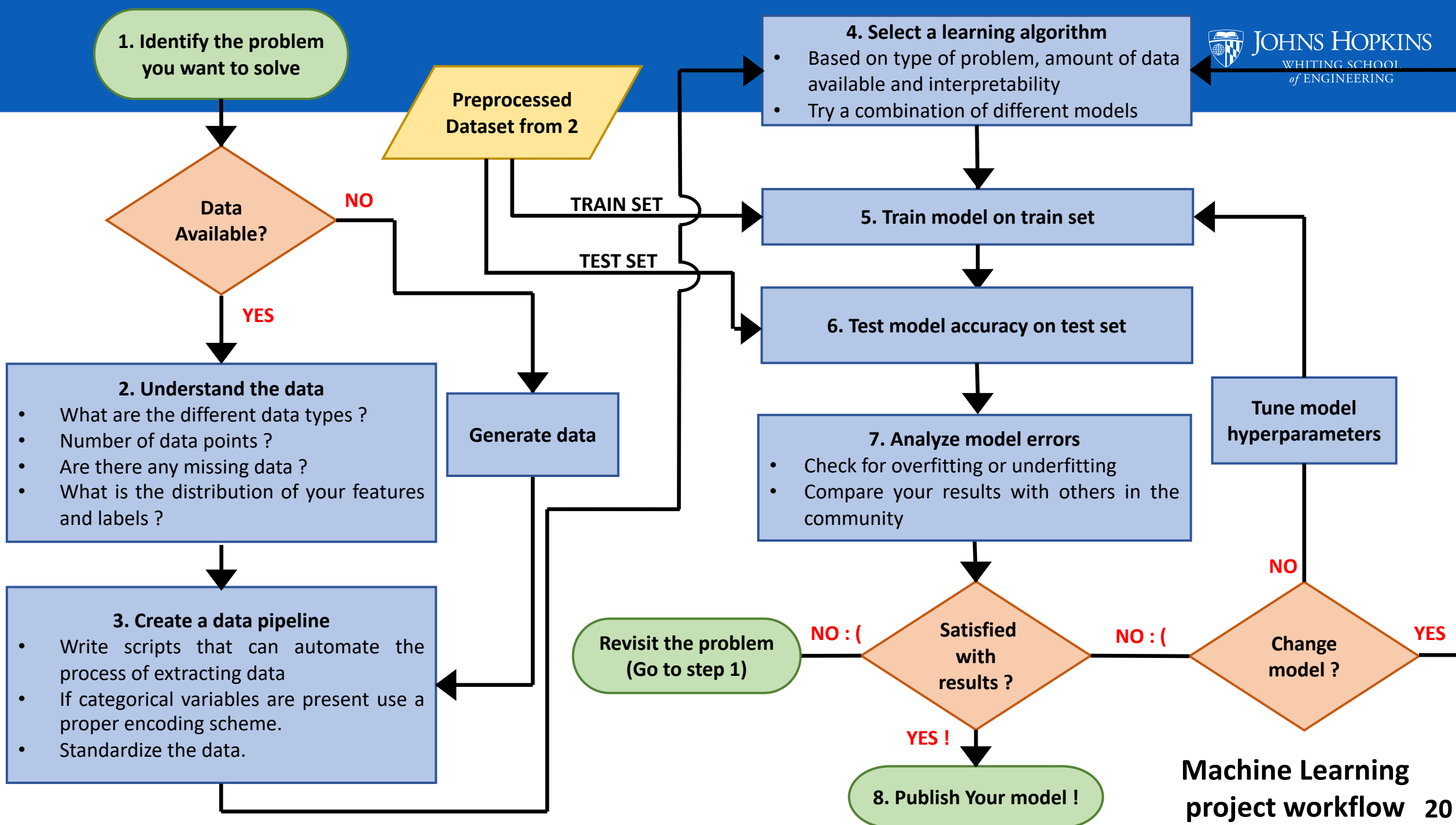


- Requires defining conditional statements
- Building model takes a significant amount of time.
- Must rewrite model code to incorporate new data.
- Cannot solve problems with no analytical or numerical solution.

Simplified machine learning program workflow



- No hard coded rules, model learns patterns from the data.
- Building model is fast due to well documented python libraries
- A well trained model is generalizable to perform a task on new data.
- Can perform complex tasks such as image recognition, language translation and can generate new insights from the data.

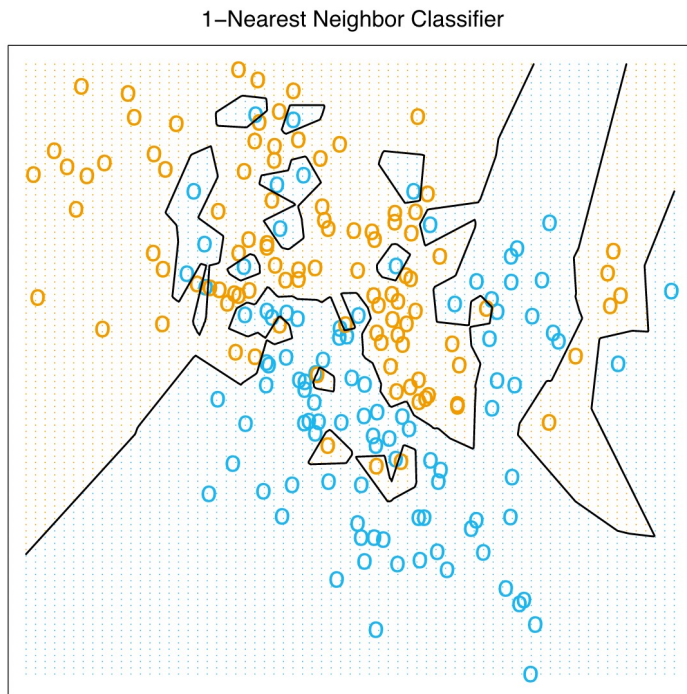


Challenges related to data collection

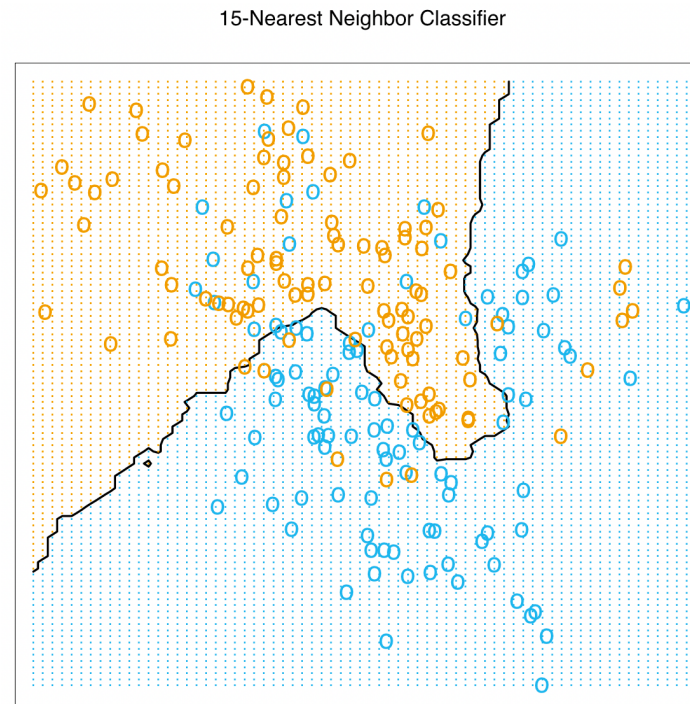
1. Insufficient quantity
2. Non representative
3. Poor quality data
4. Irrelevant features (GIGO)

Challenges related to training machine learning model

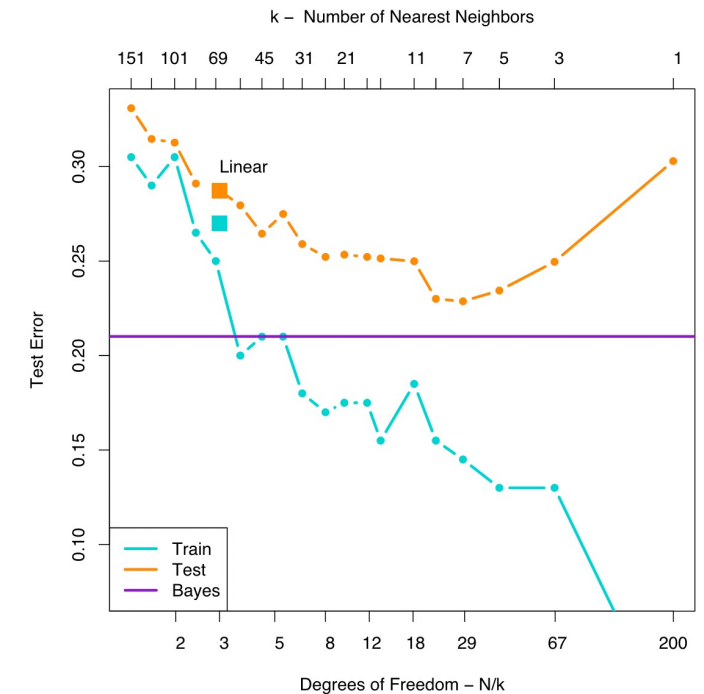
Overfitting on training data



Underfitting on training data



Bias Variance Tradeoff



Common Machine Learning Modules/Toolboxes

- [Keras](#)
 - Neural network library.
 - High-level abstractions, easy to use.
- [scikit-learn](#)
 - Designed to work with NumPy and SciPy.
 - Has basic regression, classification, and clustering algorithms.
- [TensorFlow](#)
 - Google's internal software library for machine learning applications.
 - Similar to Keras, but a bit harder to use.
- [Theano](#)
 - Efficient optimization of mathematical expressions using multi-dimensional arrays.
- [PyTorch](#)
 - Based on Torch library (developed by Facebook).
 - Used for computer vision and natural language processing.
- [Shogun](#)
 - Focus on kernel methods.
 - Originally developed for bioinformatics purposes.

- Linear Regression
- Logistic Regression
- Support Vector Machines
- Decision Trees
- Random Forests
- Naïve Bayes Classification
- Boosting
- Clustering
- Principal Component Analysis
- Singular Value Decomposition
- Independent Component Analysis
- Neural Networks
- Generative Adversarial Networks
- Bayesian Optimization
- Gaussian Process Regression

Textbooks

1. **Machine Learning : A Probabilistic Approach** by **Kevin P. Murphy**
 - Excellent book that dives right into the math of some common machine learning algorithms
2. **Elements of Statistical Learning** by **Trevor Hastie and Robert Tibshirani and Jerome Friedman**
 - A more advanced level introduction to machine learning. Skips a lot of the derivation parts. (Not for faint hearted !)
3. **Hands on Machine Learning with Scikit-Learn, Keras and Tensorflow** by **Aurelien Geron**
 - Great book to start with ! Guided tutorials and explanations to building and training machine learning models without diving too much into the math.
4. **Deep Learning with Python** by **Francois Chollet**
5. <https://www.deeplearningbook.org> by **Ian goodfellow and Yoshua Bengio and Aaron Courville**
 - Excellent book to follow for building and training neural networks !

Neural net visualization

Tensorflow documentation

PyTorch documentation

Youtube video links by 3Blue1Brown on Deep Learning

1. [But what is a neural network ? | Chapter 1, Deep Learning](#)
2. [Gradient Descent, how neural networks learn | Chapter 2, Deep Learning](#)
3. [What is backpropagation really doing ? | Chapter 3, Deep Learning](#)
4. [Backpropagation calculus | Chapter 4, Deep Learning](#)

Courses at Johns Hopkins for Upper Level Undergraduates

- EN.601.475 (01) Machine Learning -> CSE
- EN.601.475 (01) Machine Learning for Medical Applns -> ECE
- EN.540.405 (01) Machine Learning -> ChemBE
- EN.540.405 (01) Machine Learning : Deep Learning -> CSE