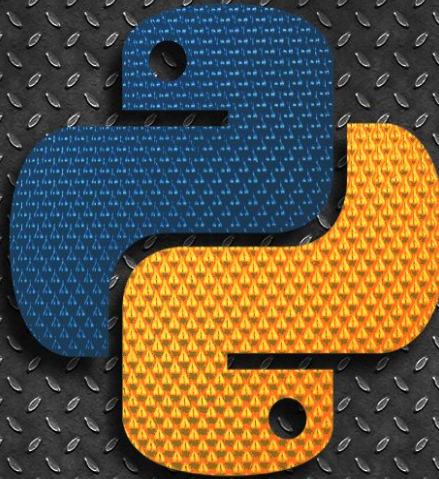


# Machine Learning in Python



# Outline

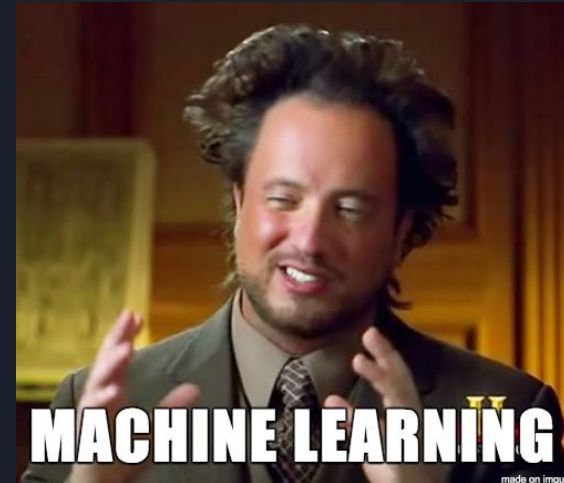
- Foundational Machine Learning Concepts
- Machine Learning!
  - Decision Trees
  - Decision Tree - Exercise
  - Random Forests
- Questions



# What is Machine Learning (ML)?

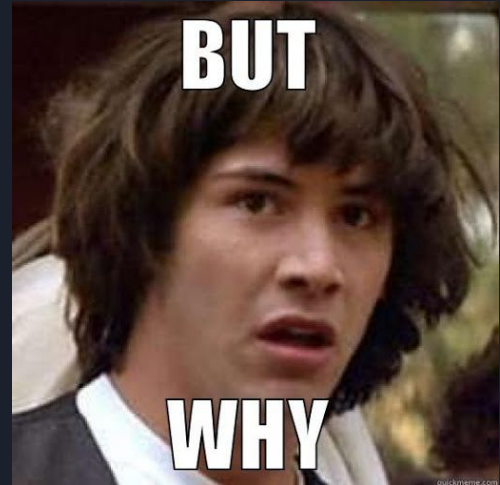
- Computers "learning" from data
- "Machine learning is an application of artificial intelligence (AI) that provides systems the ability to automatically learn and improve from experience without being explicitly programmed."

- <https://expertsystem.com/machine-learning-definition/>



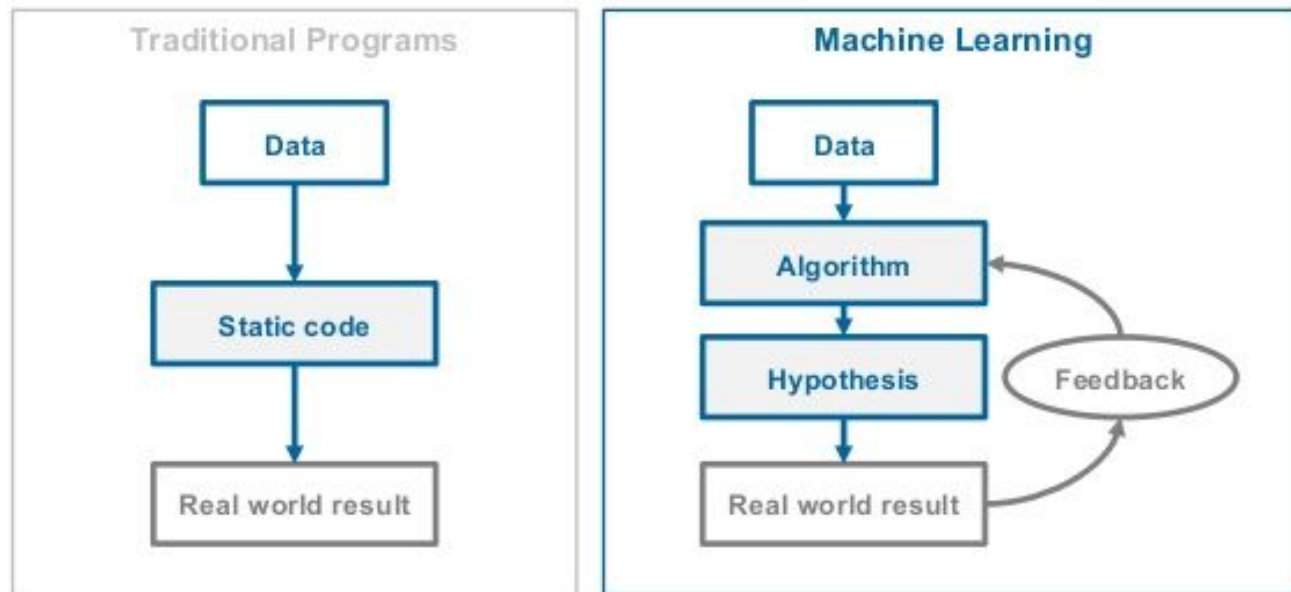
# Why use Machine Learning?

- You don't have to explicitly program something
  - ML is powered by data - not instructions
- When problems are much too complex to be explicitly programmed
- When the environment changes over time





## Traditional Programs vs. Machine Learning



# Applications of Machine Learning



**Sarah Michelle Gellar**  
74 mutual friends



**Stephen King**  
13 mutual friends

# AI vs. ML vs. DL

## Artificial Intelligence

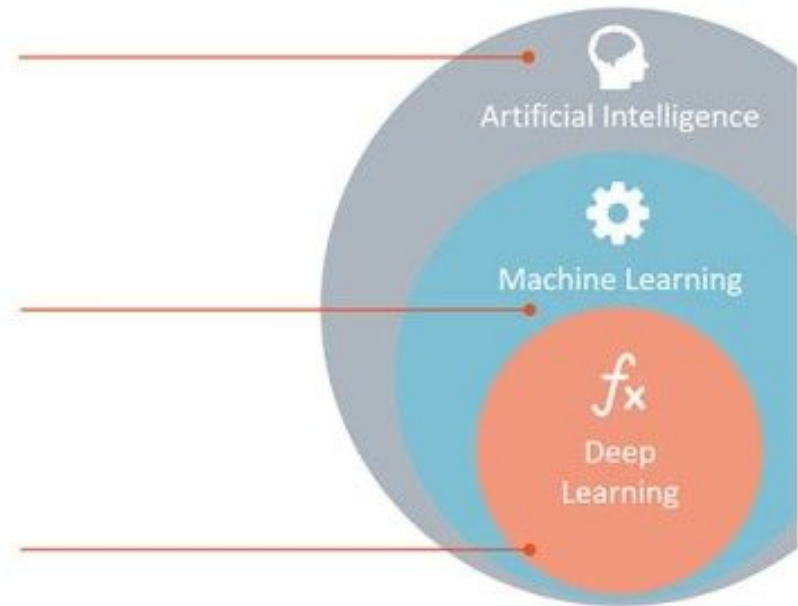
Any technique which enables computers to mimic human behavior.

## Machine Learning

Subset of AI techniques which use statistical methods to enable machines to improve with experiences.

## Deep Learning

Subset of ML which make the computation of multi-layer neural networks feasible.





# ML Terminology

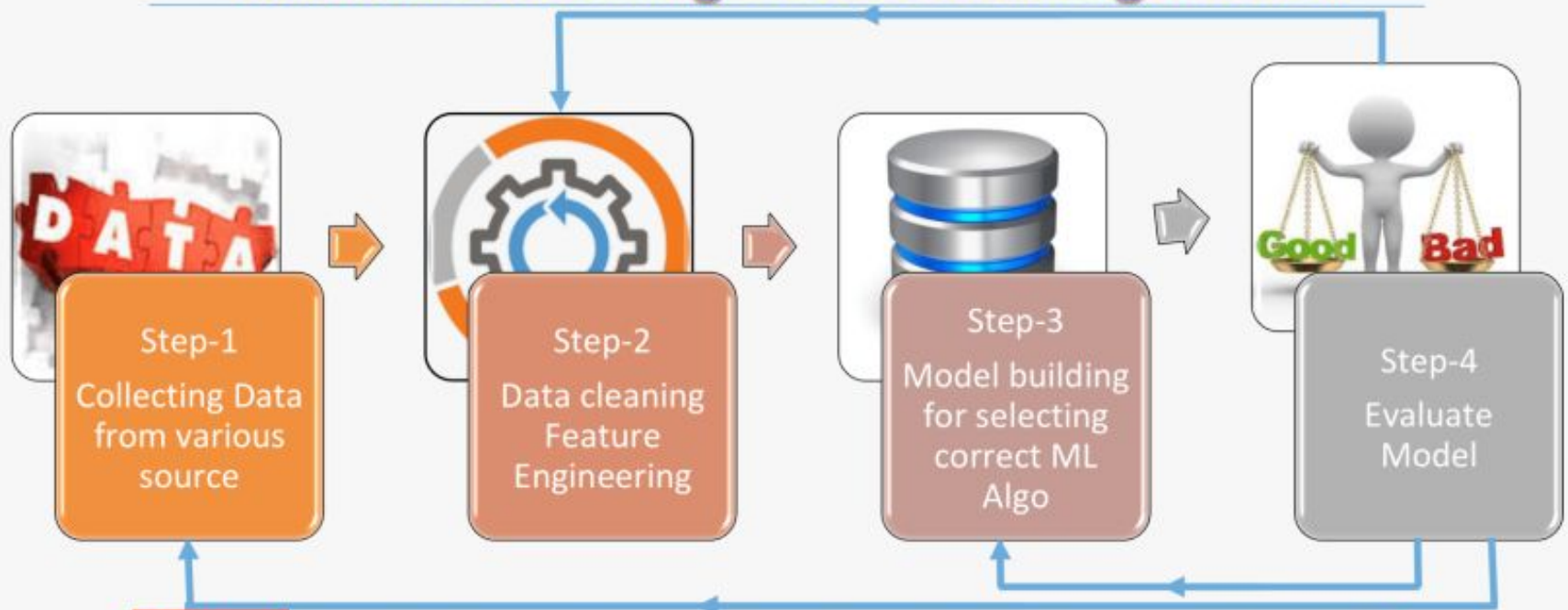
- Training Data: the examples that the system uses to learn
- Features: The input "variables"
- Labels: The value(s) you're trying to predict given the features
- Algorithm: The model that is trained and used to predict the label(s)



# ML Process



# Machine Learning Process at High Level



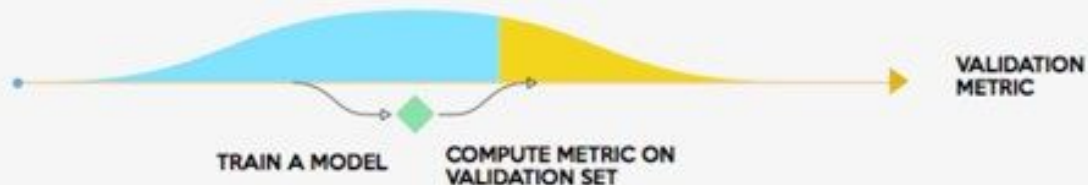
# HOLDOUT STRATEGY

**1** Split your data into train / validation / test



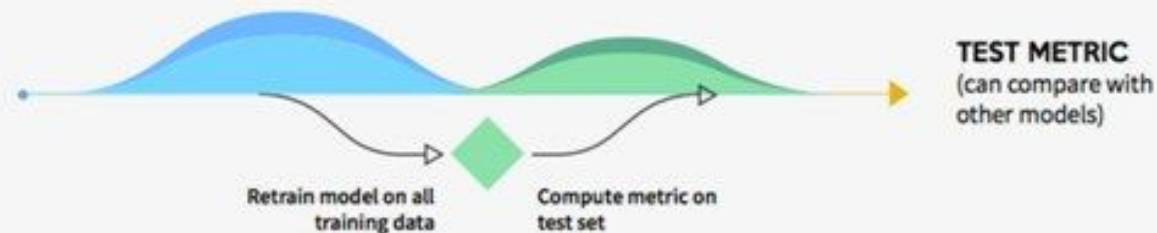
**2** For each parameter combination

Parameter (e.g., depth) A	10 1	10 11	Parameter (e.g., n trees) B
	5	15	
	6	16	
	7	17	



**3** Choose the parameter combination with the best metric

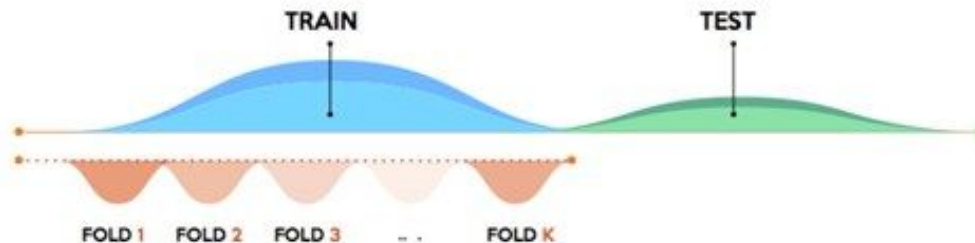
A	6	14	B
---	---	----	---



# K-FOLD STRATEGY

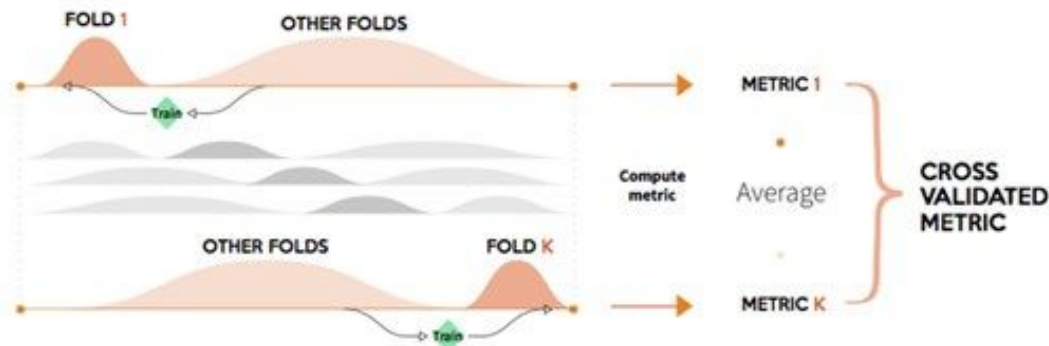
1

Set aside the test set and split the train set into k folds



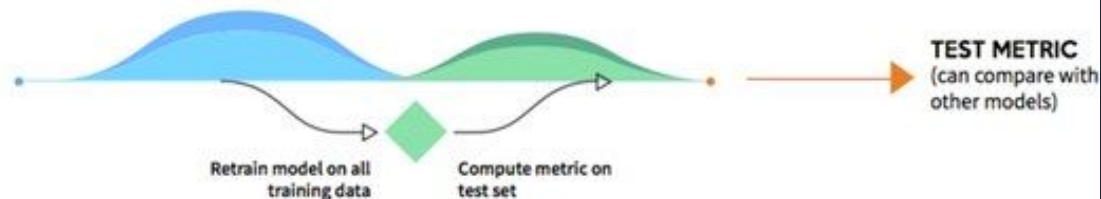
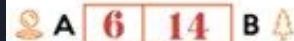
2

For each parameter combination

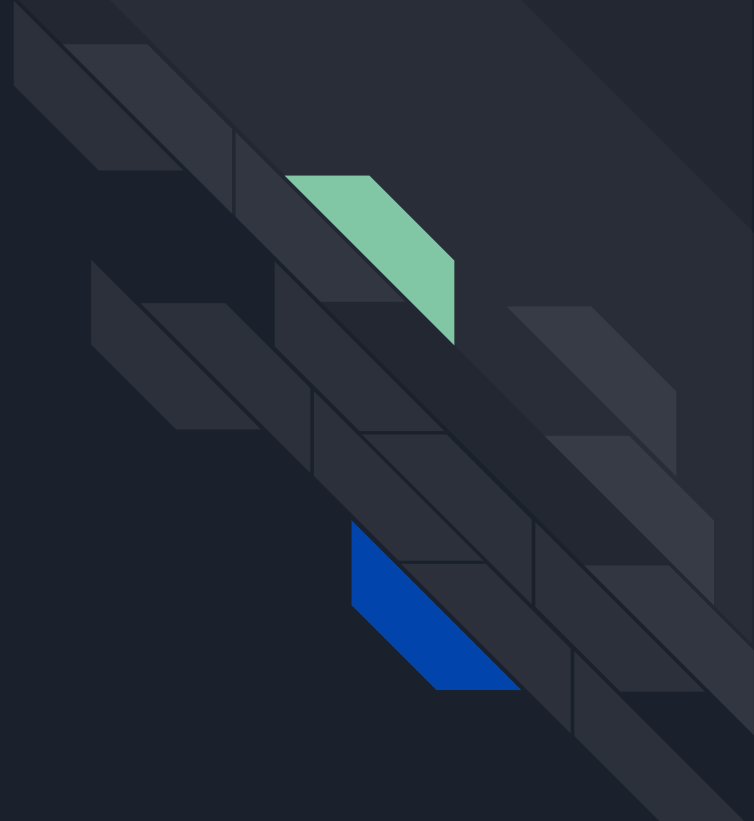


3

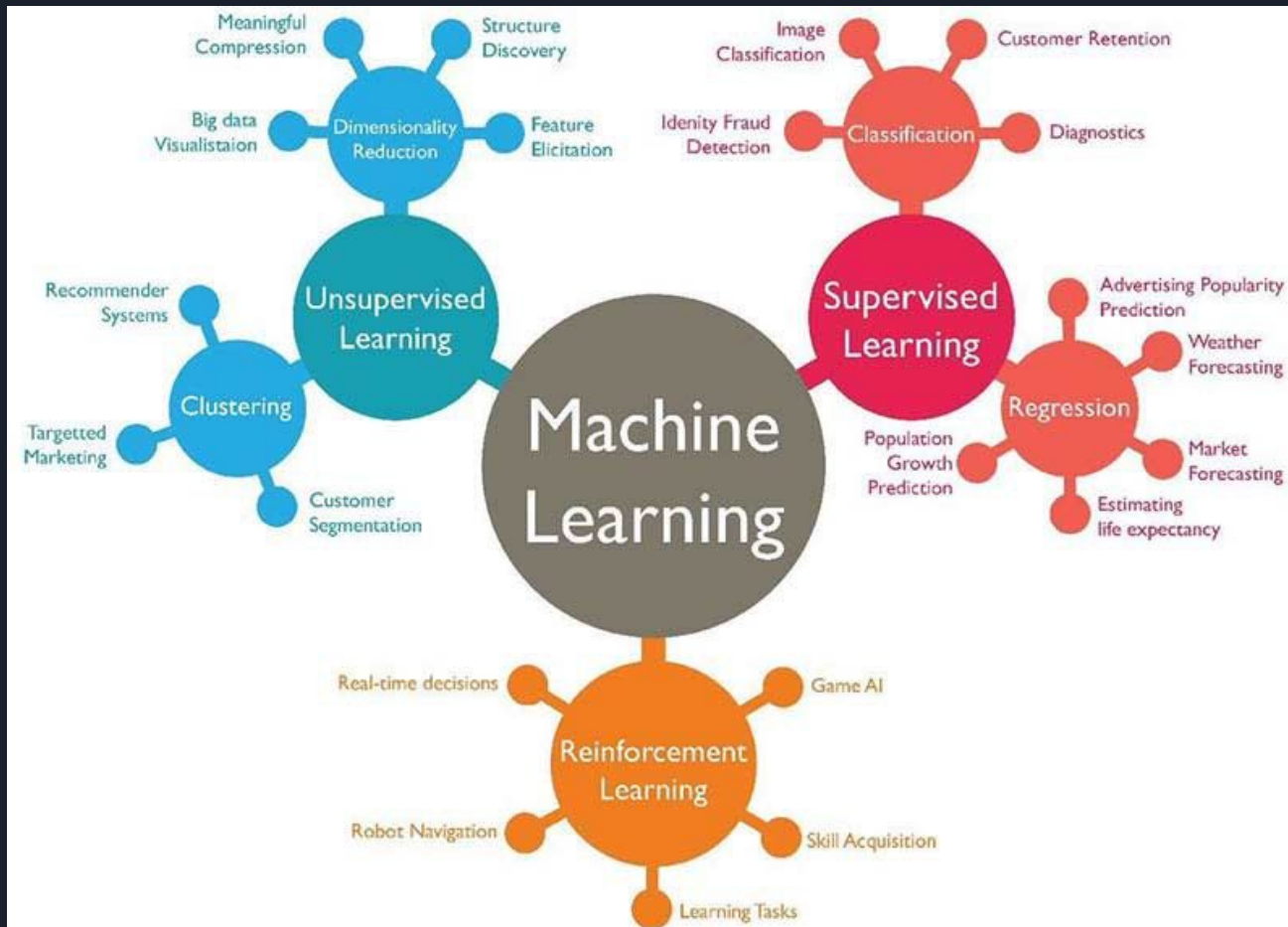
Choose the parameter combination with the best metrics



# Types of ML









# Types of Machine Learning

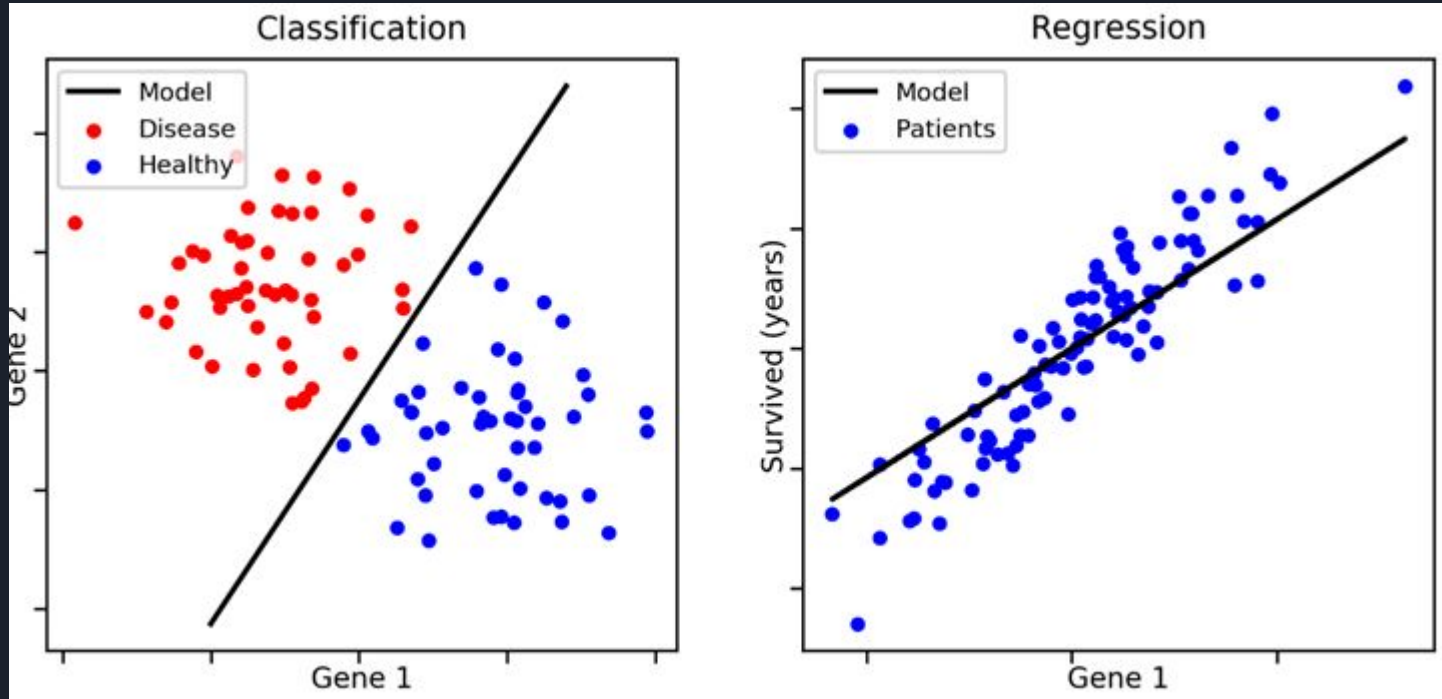
- **Supervised Learning**
  - Labeled Training Data ("right" answer)
  - Predict outcome
  - Ex: Classification, Regression
- **Unsupervised Learning**
  - Unlabeled Training Data (no "right" answer, or "right" answer is unknown)
  - Find structure and patterns in data
  - Ex: Clustering, Anomaly Detection
- **Reinforcement Learning**
  - Machine Agent explores environment
  - Learning is based on rewards and punishments



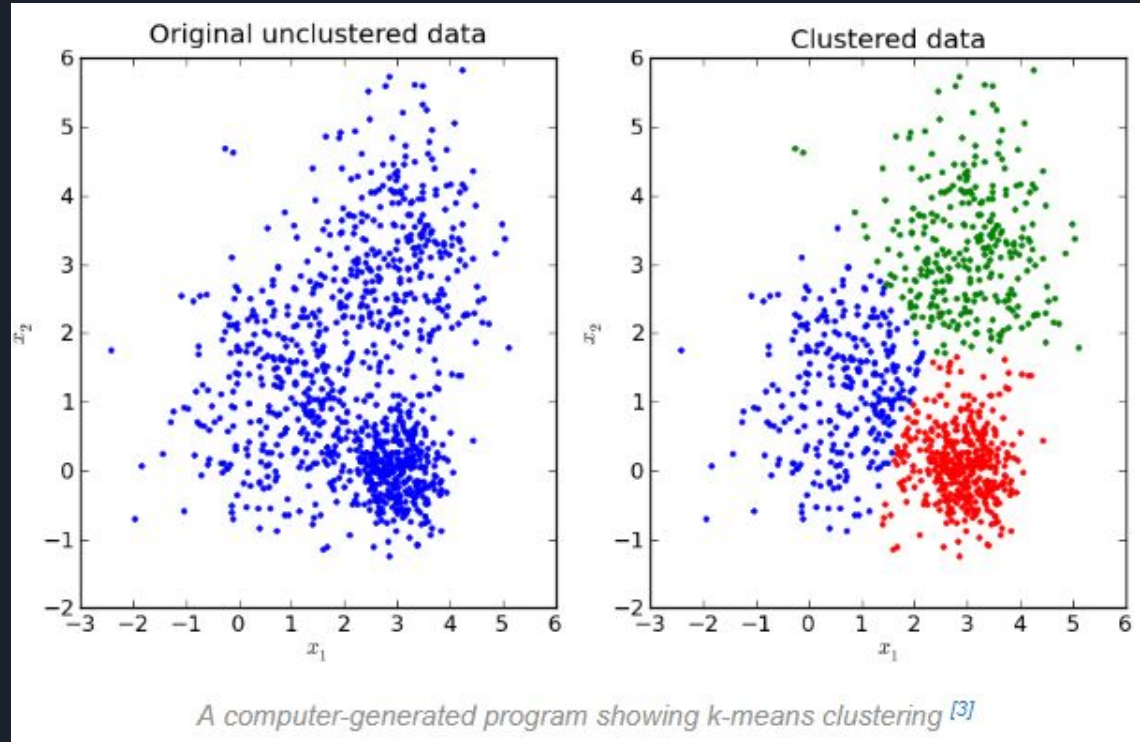
# Common ML Tasks

- Supervised:
  - Classification
  - Regression
- Unsupervised:
  - Clustering
  - Anomaly Detection
  - Association Rule Learning

# Common Supervised ML Tasks: Classification vs. Regression

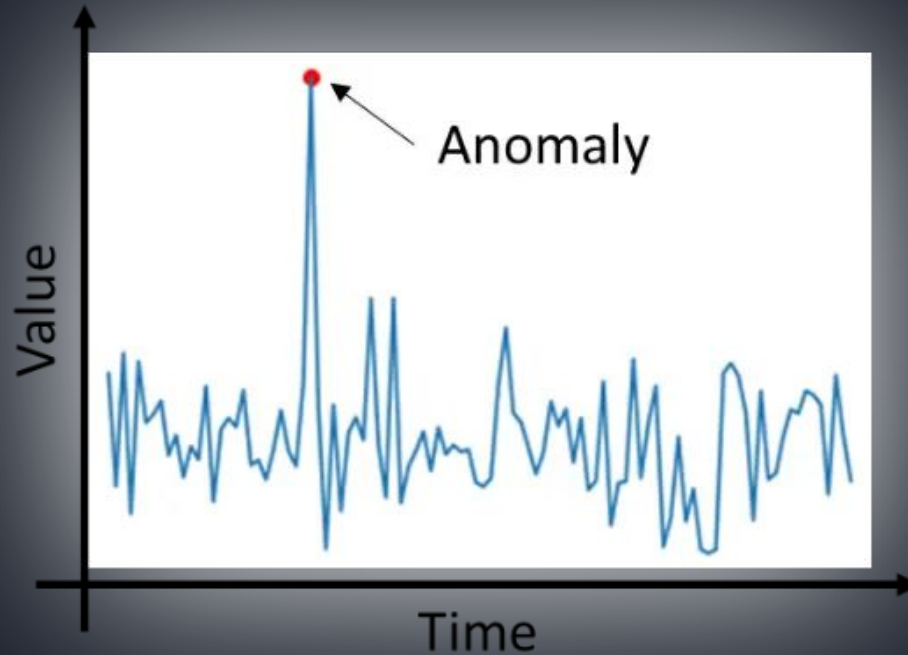


# Common Unsupervised ML Tasks: Clustering





# Common Unsupervised ML Tasks: Anomaly Detection



# Common Unsupervised ML Tasks: Association Rule Mining

ID	Items
1	{Bread, Milk}
2	{Bread, <b>Diapers</b> , <b>Beer</b> , Eggs}
3	{Milk, <b>Diapers</b> , <b>Beer</b> , Cola}
4	{Bread, Milk, <b>Diapers</b> , <b>Beer</b> }
5	{Bread, Milk, Diapers, Cola}
...	...

market  
basket  
transactions

**{Diapers, Beer}**

Example of a frequent itemset

**{Diapers} → {Beer}**

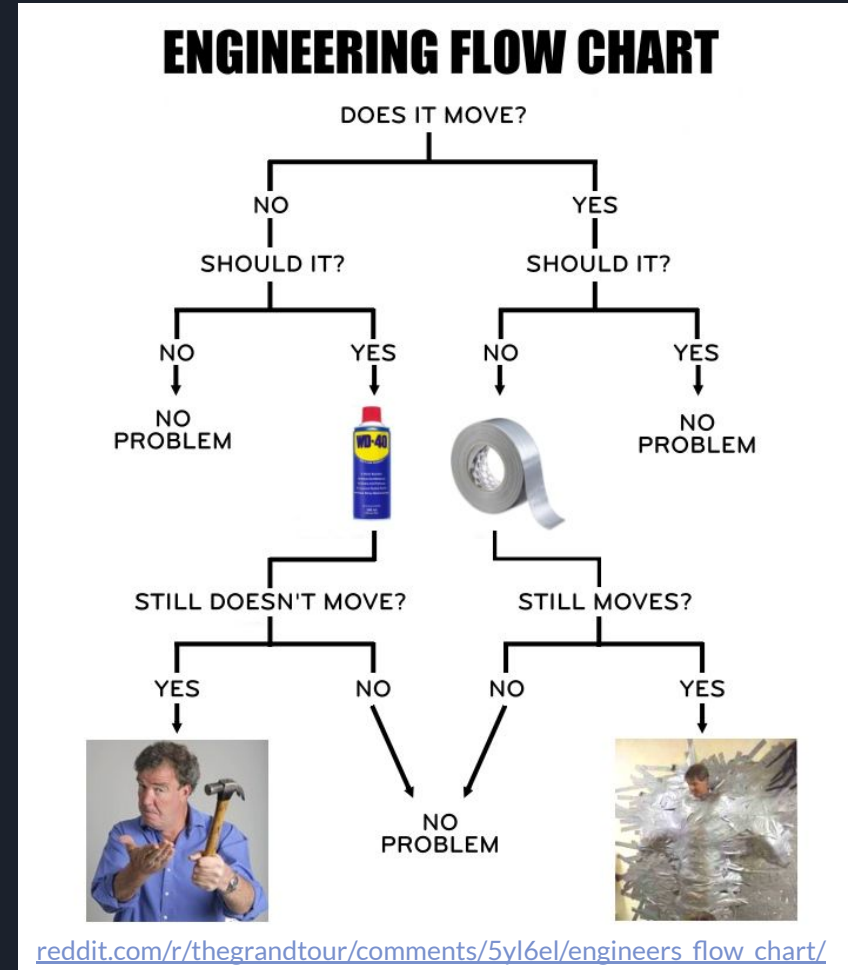
Example of an association rule

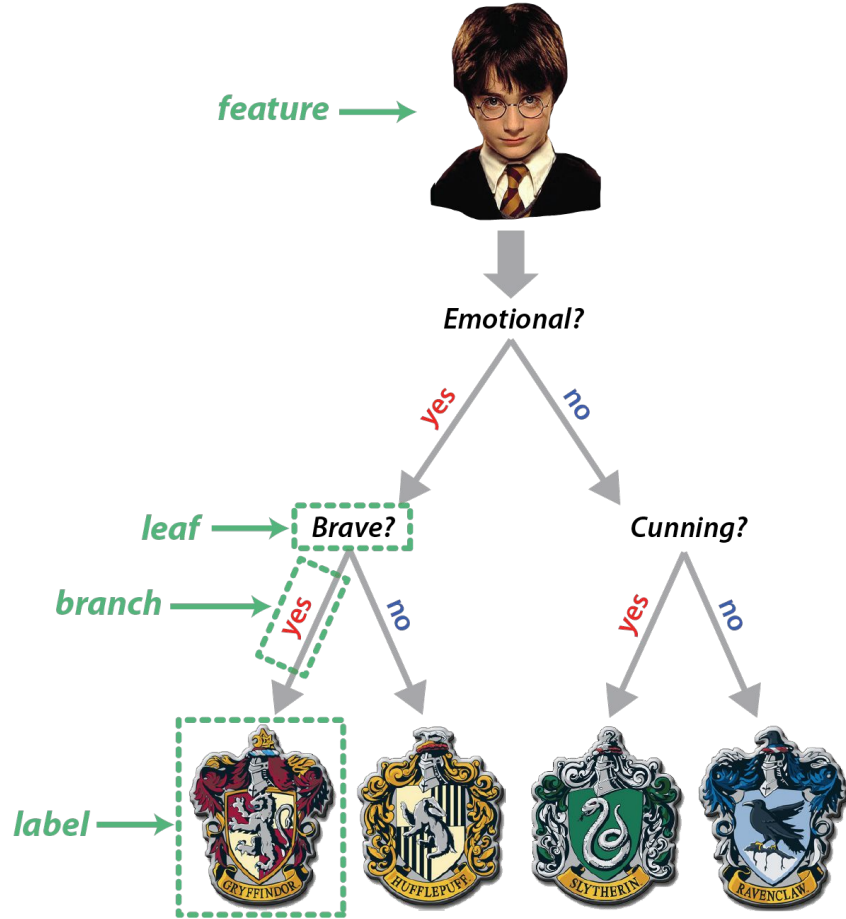
# Decision Trees



# Decision Trees

- Series of Q's
- Versatile: Classification or Regression
- Powerful: Can fix complex datasets
- "White-box"
- Few assumptions about training data
- Backbone of Random Forests







# Ensemble Learning

- A group of predictors
  - AKA training more than one model and combining the predictions of them

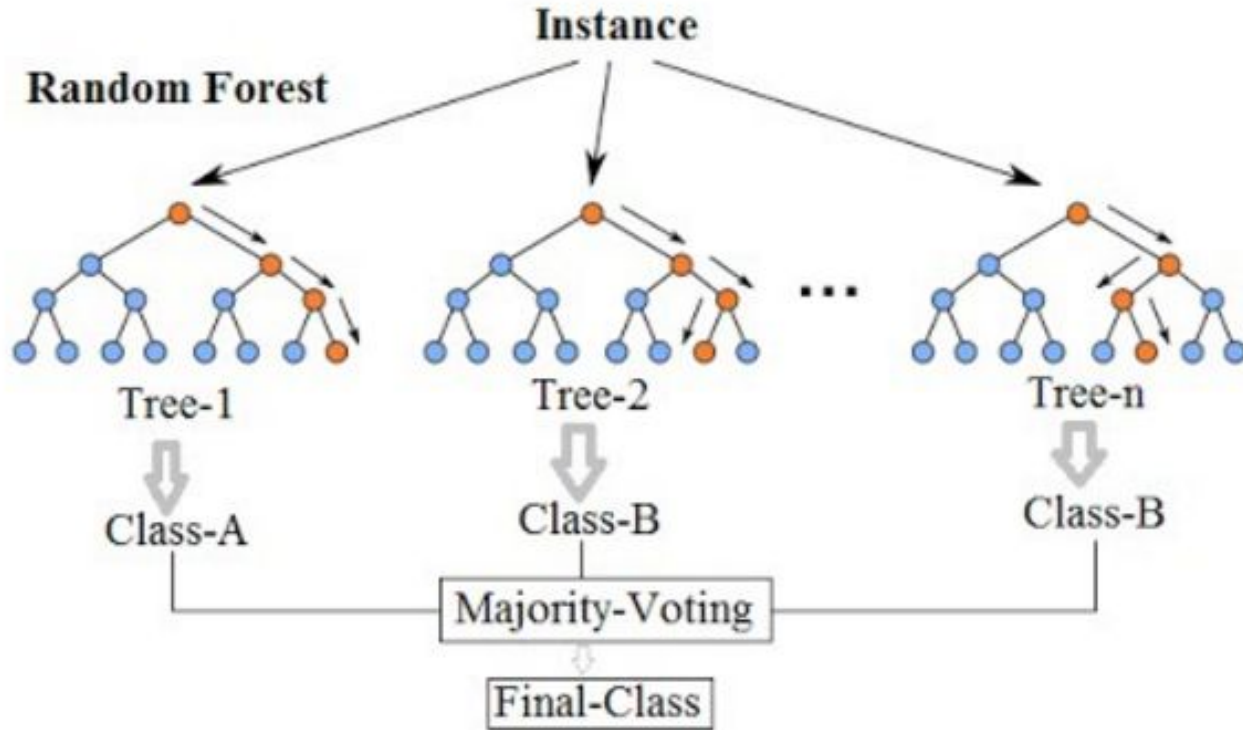




# Random Forests (RF)

- Each decision tree in the RF:
  - has random subset of features
  - uses random set of the training data points
- Final prediction is a majority vote for the predicted class

# Random Forest Simplified





# Further Resources

- Online Classes:
  - [Coursera](#)
- Books:
  - [\*Data Science for Business\* - Foster Provost and Tom Fawcett](#)
  - [\*Hands-On Machine Learning with Scikit-Learn and TensorFlow\* - Aurélien Géro](#)
- Internet:
  - [Youtube: Siraj Raval](#)
- Classes at OSU (have prereqs):
  - CSE 5243: Introduction to Data Mining (3 cr hrs)
  - CSE 3521: Survey of Artificial Intelligence I: Basic Techniques (3 cr hrs)
  - CSE 5523: Machine Learning and Statistical Pattern Recognition (3 cr hrs)
  - CSE 5524: Computer Vision for Human-Computer Interaction (3 cr hrs)
  - CSE 5526: Introduction to Neural Networks (3 cr hrs)
- Projects:
  - Kaggle

# Questions?

- Any Questions?
- My Contact Info:
  - Leo Glowacki
  - Message me on Slack!
  - [www.LeoGlowacki.com](http://www.LeoGlowacki.com)

