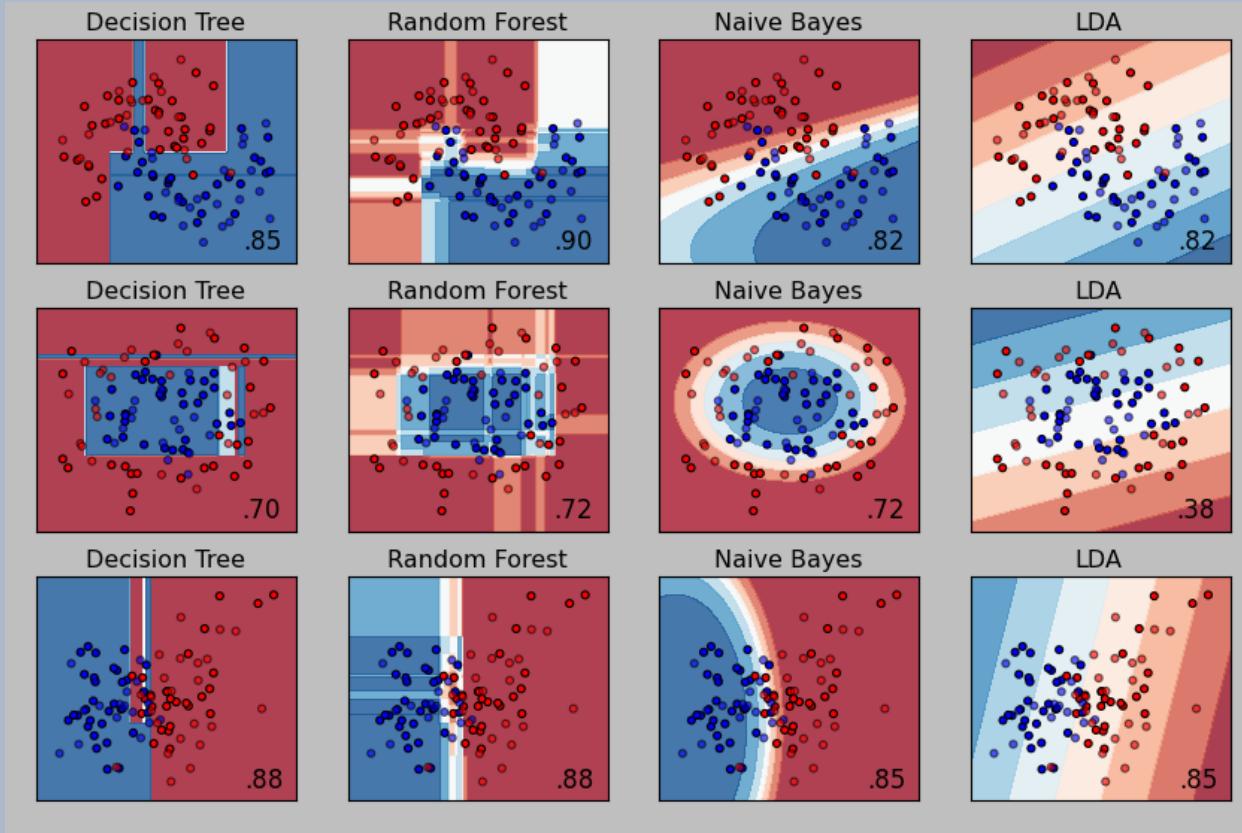


An Introduction to Machine Learning



Ryan Urbanowicz, PhD



Perelman
School of Medicine
UNIVERSITY OF PENNSYLVANIA

PA CURE Machine Learning Workshop: December 17

Overview

- Fundamentals of Machine Learning (ML)
- Focus: Decision Tree
- Choosing an ML algorithm
- Common ML Pitfalls



Terminology and Definitions

- **Instance:** an individual or example in data.
 - E.g. A subject/patient in a drug trial.
- **Feature:** one of the attributes describing an aspect of the instance. E.g. height, weight, age.
- **Outcome:** In supervised learning, this is endpoint value, a.k.a. the dependent variable, or the target being predicted.
 - Label/Class: Terms used for outcome in classification.
 - In regression, the outcome would be real-valued numbers.
- **Model:** A representation or simulation of reality. Typically a simplification based on a number of assumptions.

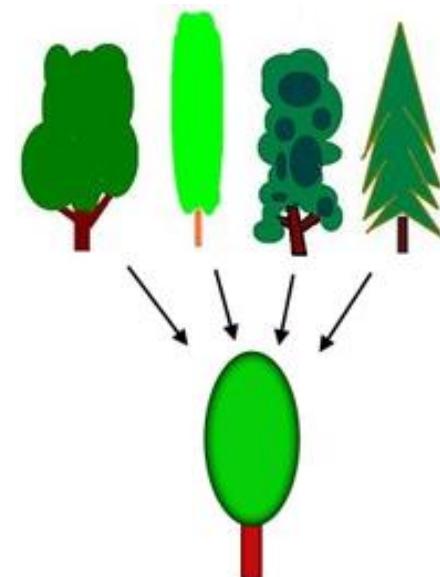


What is Machine Learning (ML)?

- A subset of **artificial intelligence** in the field of **computer science** that often uses **statistical** techniques to give computers the ability to "learn" (i.e., progressively improve performance on a specific task) with data, without being explicitly programmed¹.

¹ Samuel Arthur – 1959 – ML in Checkers

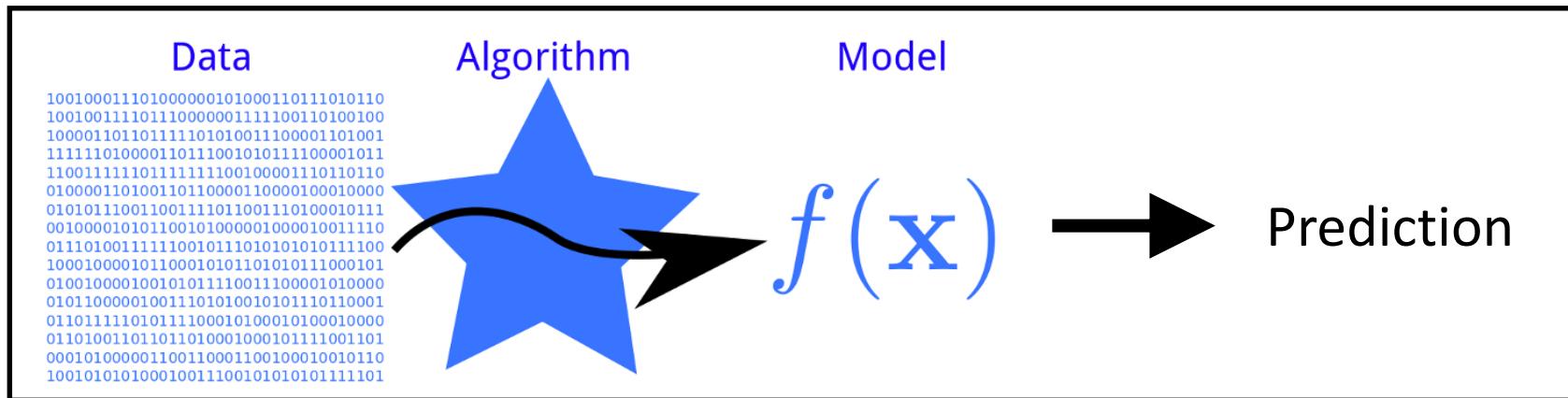
- ML is a general term → many algorithms/methods.
- Big Picture Goal: Learning useful **generalizations**.



An Important Clarification

- Machine Learning is...
 - Finding patterns or **associations** that can be used to make **predictions**.

Example: Predictive Modeling of Outcome



- Mostly **NOT**
 - Designed to demonstrate causality.
 - At best: associations are candidates for causality.

<http://phdp.github.io/posts/2013-07-05-dtl.html>



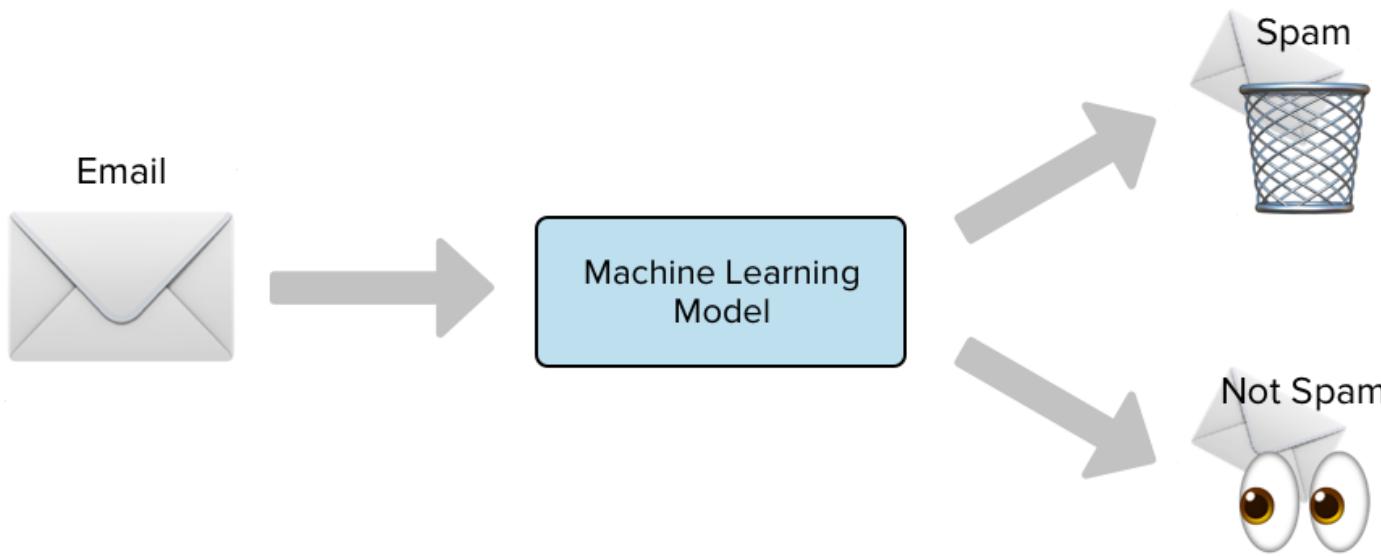
Example: Email Spam Detection

From: cheapsales@buystufffromme.com
To: ang@cs.stanford.edu
Subject: Buy now!

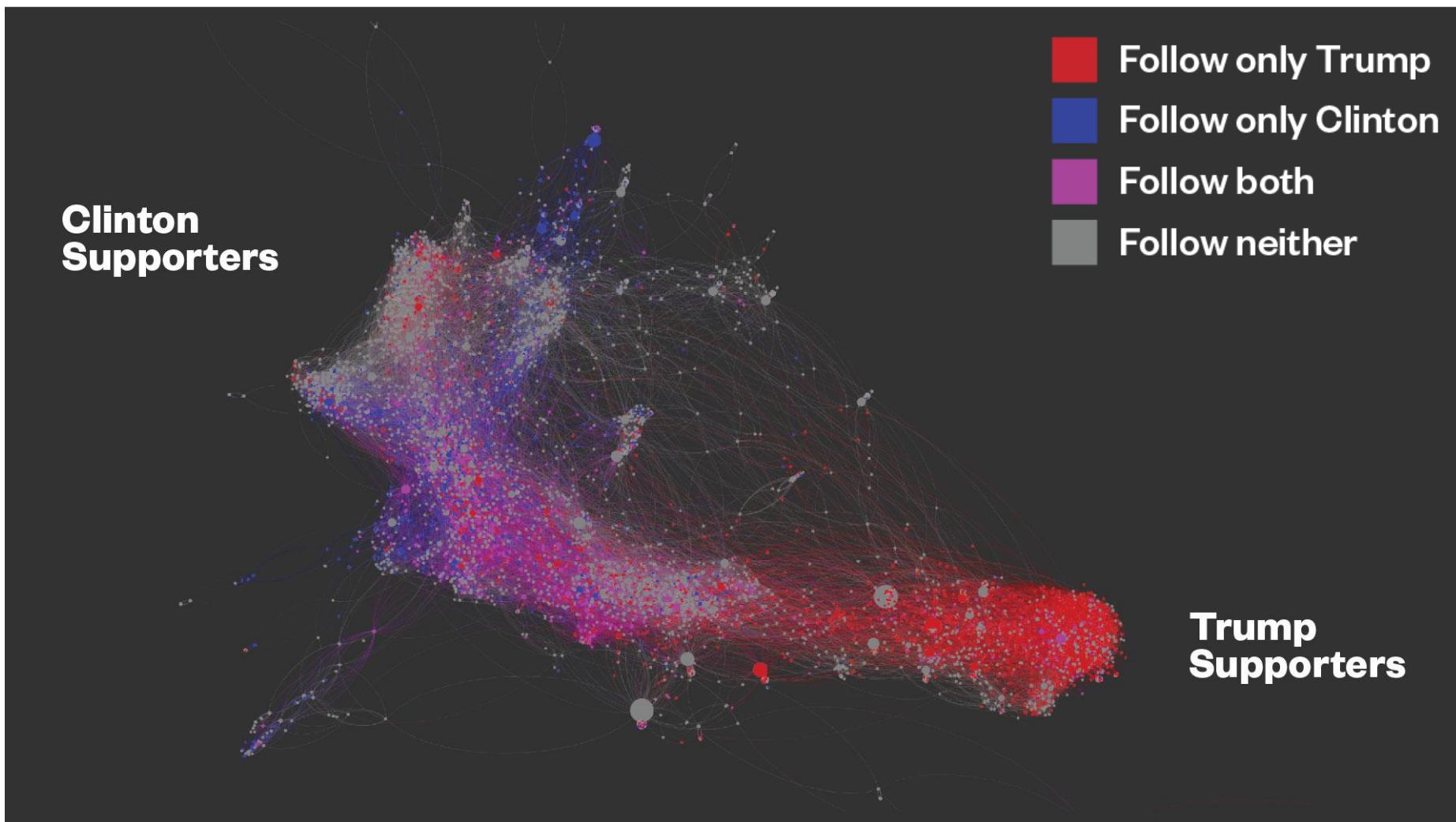
Deal of the week! Buy now!
Rolex w4tchs - \$100
Med1cine (any kind) - \$50
Also low cost M0rgages available.

From: Alfred Ng
To: ang@cs.stanford.edu
Subject: Christmas dates?

Hey Andrew,
Was talking to Mom about plans for Xmas. When do you get off work. Meet Dec 22?
Alf



Example: Community Detection



https://news.vice.com/en_us/article/d3xamx/journalists-and-trump-voters-live-in-separate-online-bubbles-mit-analysis-shows

Example: Association Mining

- Given a set of transactions, find rules that will predict purchase associations among items.

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

market basket transactions

{Diapers, Beer}

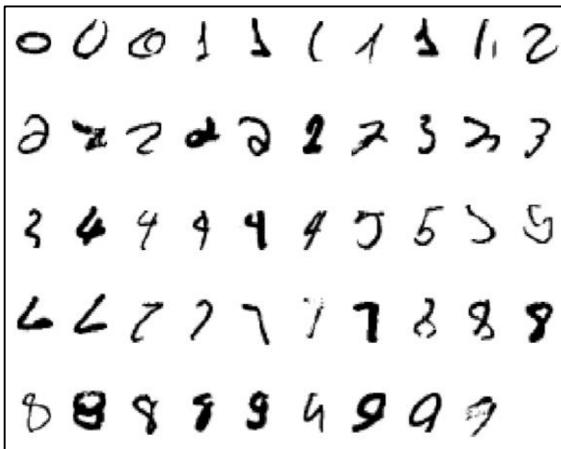
{Diapers} → {Beer}

<https://www.datacamp.com/community/tutorials/market-basket-analysis-r>

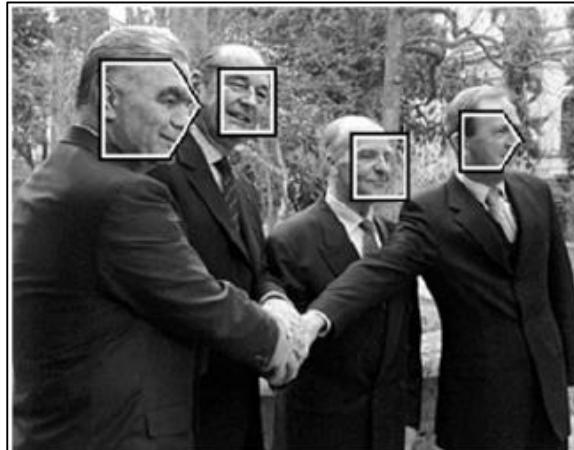


Other Examples of Applied ML

Image Classification



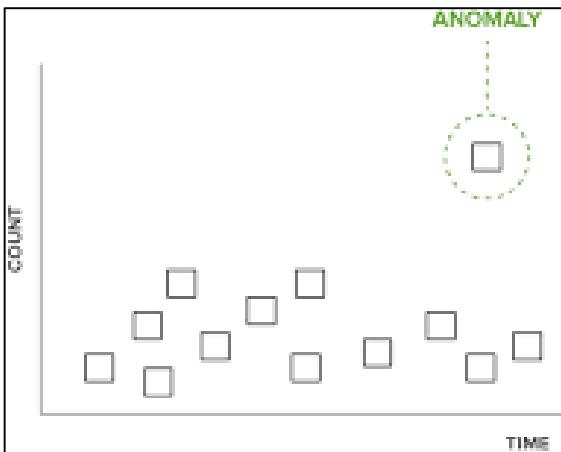
Face Detection



Stock Prediction



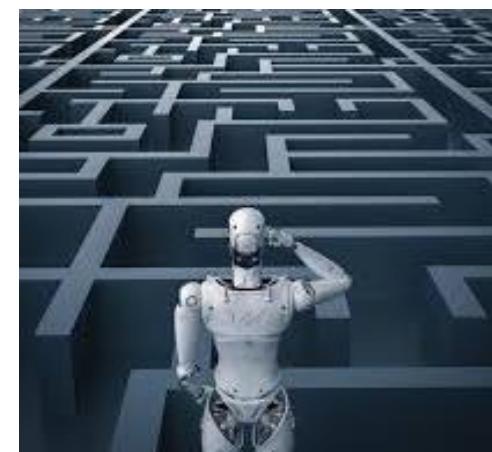
Fraud Detection



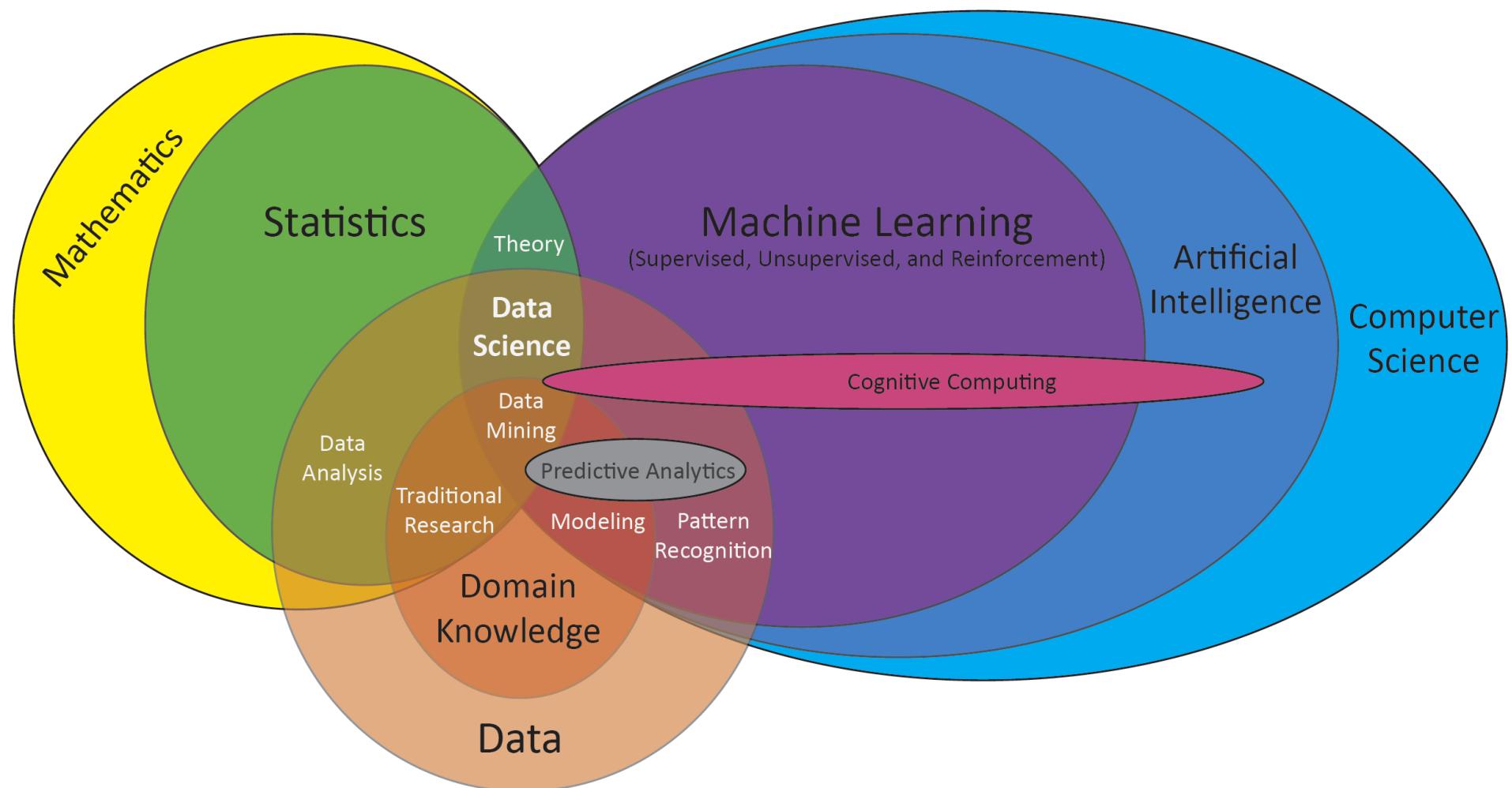
Risk Analysis

		Probability				
		Very High	High	Medium	Low	Very Low
EXAMPLE RISK		Very High	Very High	Very High	High	High
Conse- quence	Very High	Very High	Very High	Very High	High	High
	High	Very High	High	High	Medium	Medium
	Medium	High	High	Medium	Medium	Low
	Low	High	Medium	Medium	Low	Very Low
	Very Low	Medium	Low	Low	Very Low	Very Low

Navigation



Fields & Terms Related to Machine Learning

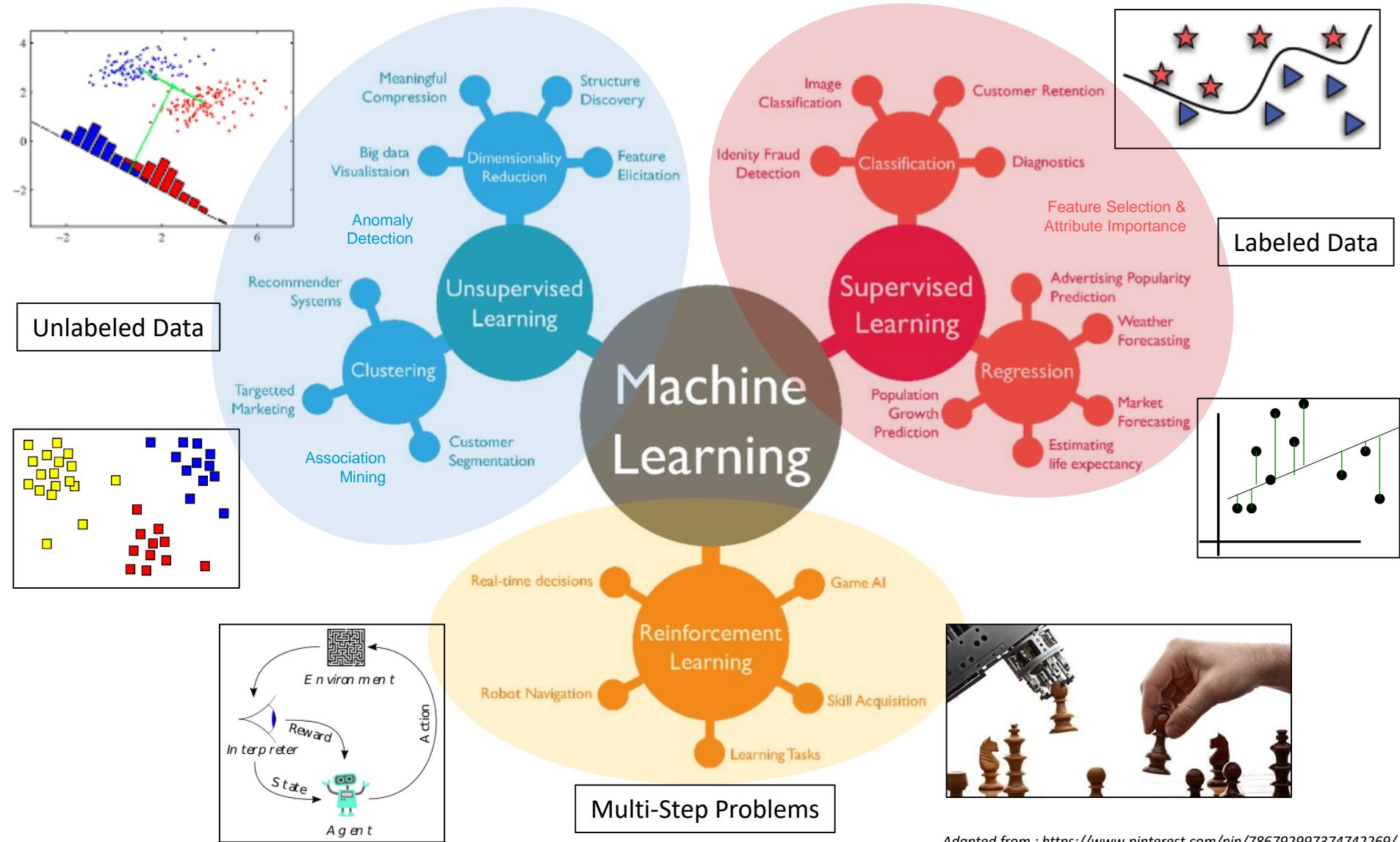


Statistics vs. Machine Learning

- Largely overlapping fields:
 - Both concerned with **learning from data**
 - Philosophical difference on ‘focus’ and ‘approach’.
- Statistics:
 - Founded in mathematics
 - Drawing **valid conclusions** based on analyzing **existing data**.
 - Making inference about a ‘population’ based on a ‘sample’
 - Tends to focus on fewer variables at once.
 - Precision and uncertainty are measures of model goodness.
- Machine Learning:
 - Founded in computer science
 - Focused on **making predictions** or **seeking patterns** (generalization).
 - Often considers a large number of variables at once.
 - Prediction accuracy to measure model goodness.

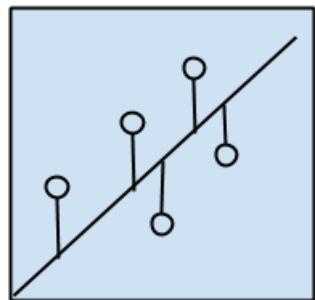


Types of Machine Learning

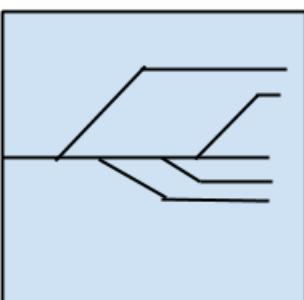


Adapted from : <https://www.pinterest.com/pin/7867929973742269/>

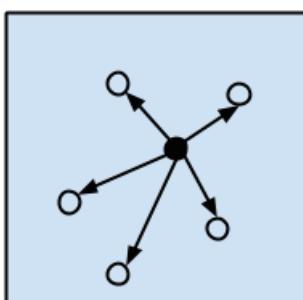
Machine Learning Algorithm Families



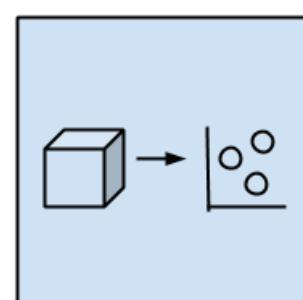
Regression Algorithms



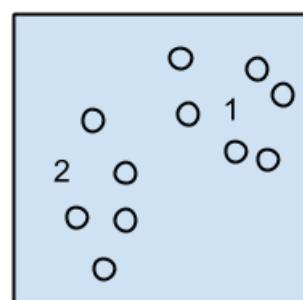
Regularization
Algorithms



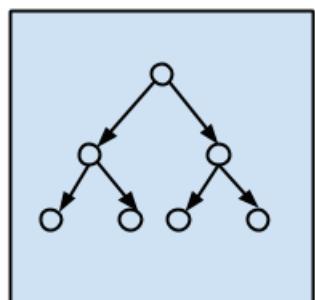
Instance-based
Algorithms



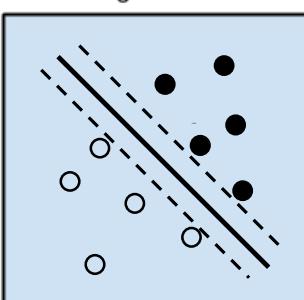
Dimensional Reduction
Algorithms



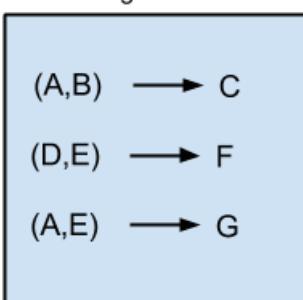
Clustering Algorithms



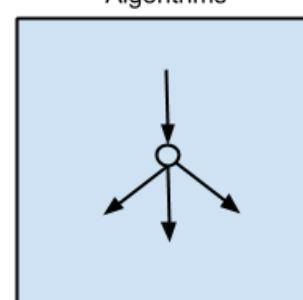
Decision Tree
Algorithms



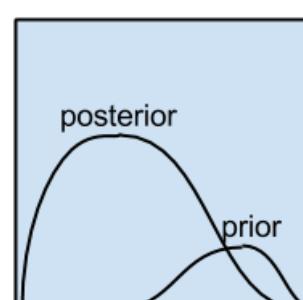
Support Vector
Machines



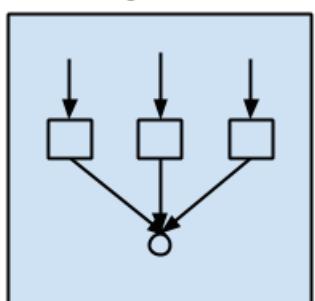
Association Rule
Learning Algorithms



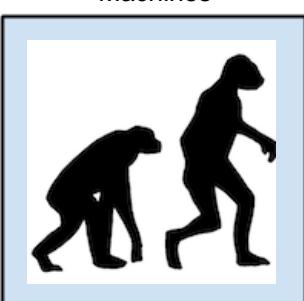
Artificial Neural Network
Algorithms



Bayesian Algorithms

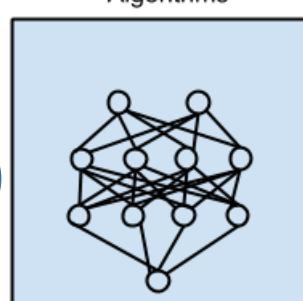


Ensemble Algorithms

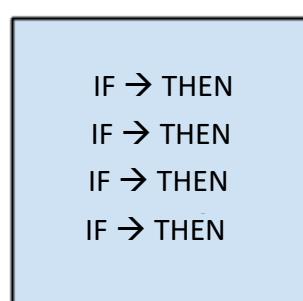


Evolutionary
Algorithms

Non-exhaustive
list of ML families



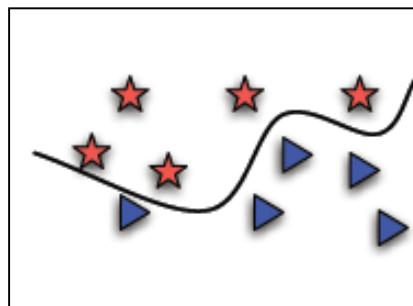
Deep Learning
Algorithms



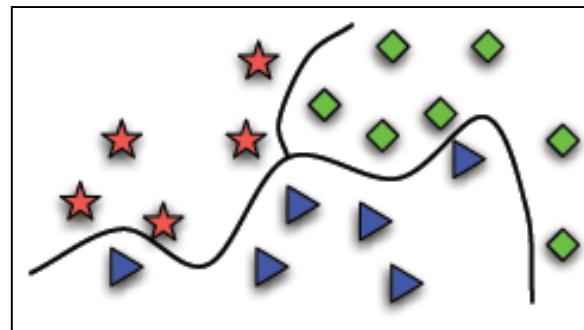
Learning Classifier
Systems

Supervised Learning: Prediction

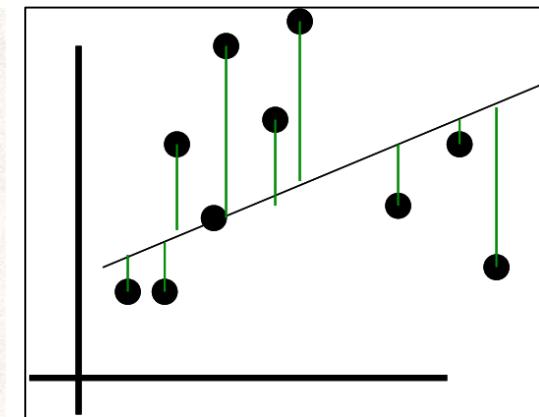
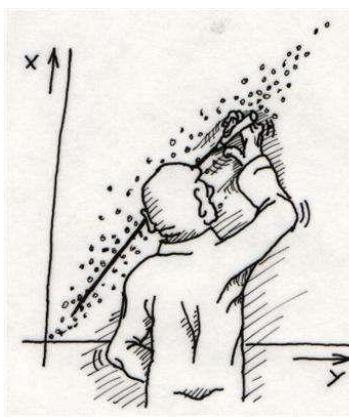
- Binary classification
 - Discriminate between two discrete classes/labels



- Multiclass classification
 - Allows for more than 2 discrete classes.
 - E.g. Cancer classes may be healthy, early state, late stage.

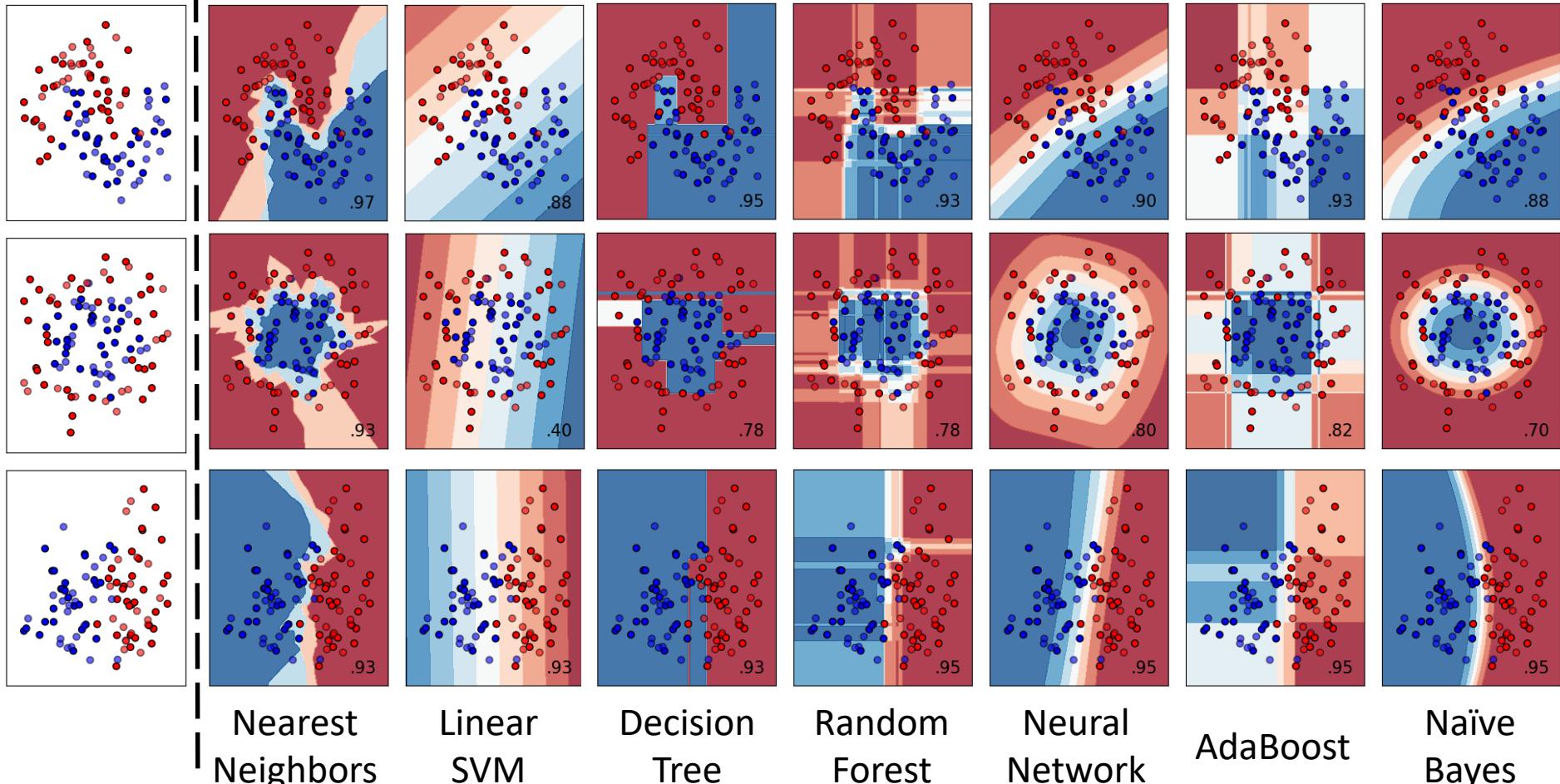


- Regression
 - Estimate a real-valued output variable



Modeling with Machine Learning

Input
Data



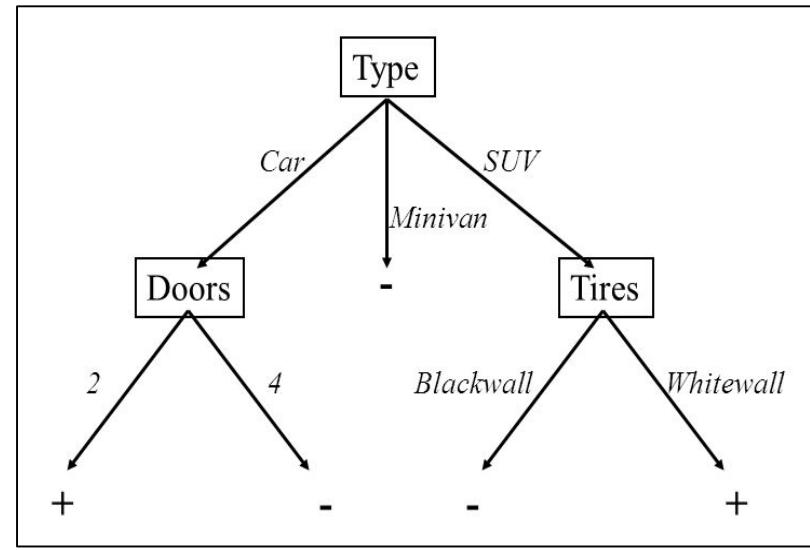
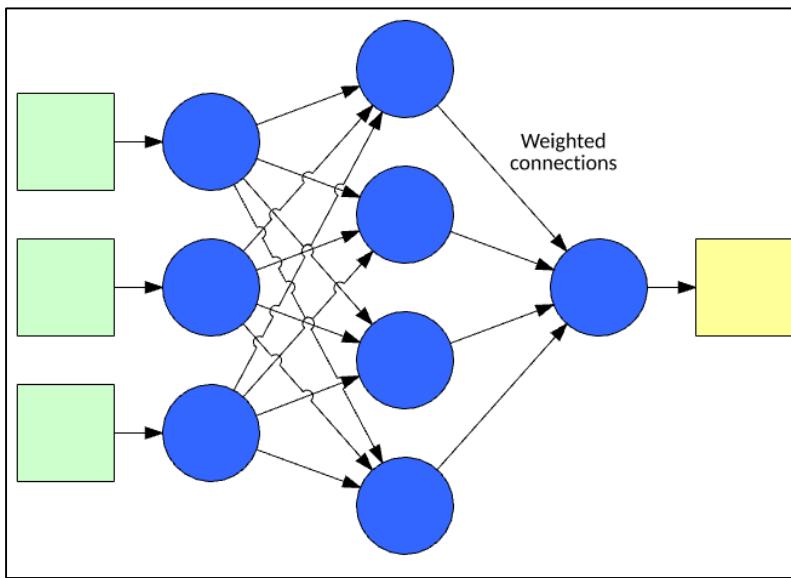
Models/ML: Representation

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

Annotations for the regression equation:

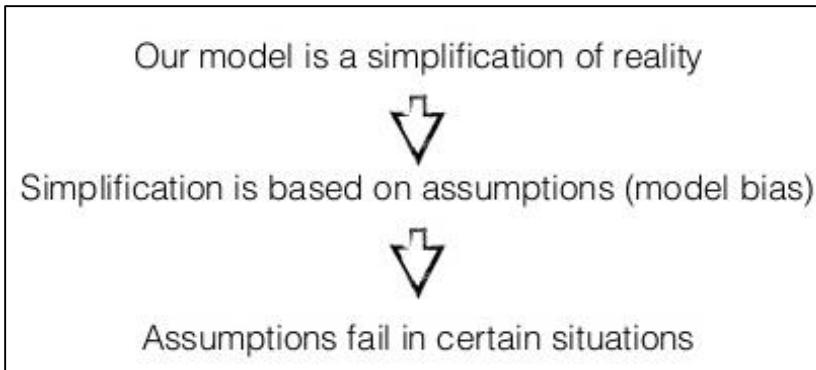
- Dependent Variable
- Population Y intercept
- Population Slope Coefficient
- Independent Variable
- Random Error term
- Linear component
- Random Error

- R1: IF THEN the animal has hair it is a mammal
- R2: IF THEN the animal gives milk it is a mammal
- R3: IF THEN the animal has feathers it is a bird
- R4: IF THEN the animal flies the animal lays eggs it is a bird
- R5: IF THEN the animal is a mammal the animal eats meat it is a carnivore



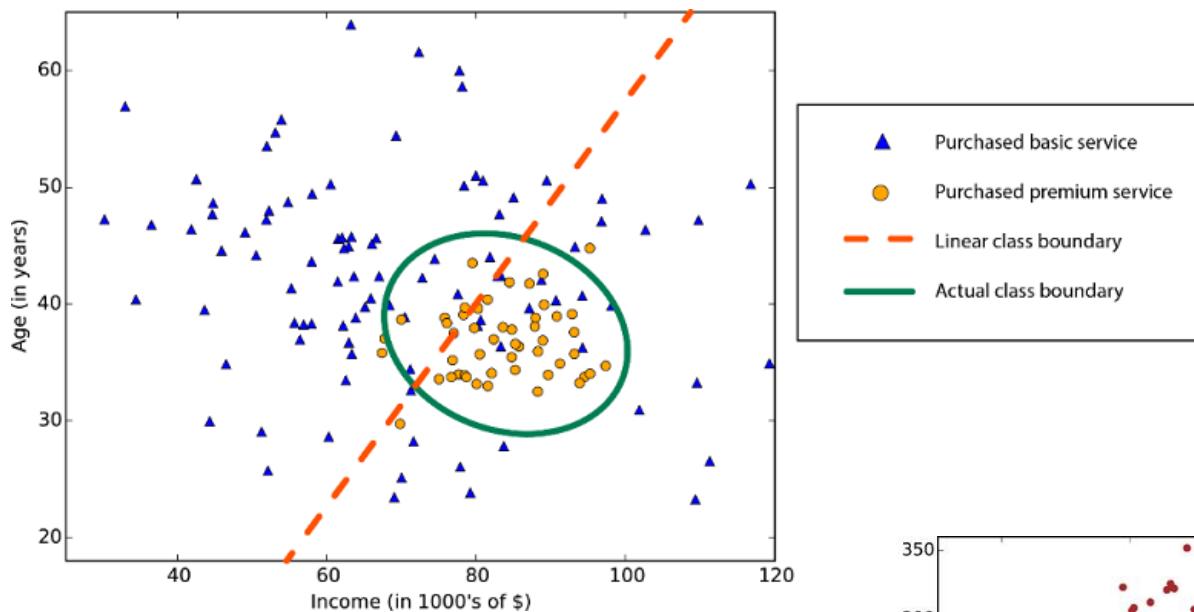
Models and the NFL

“All models are wrong, but some models are useful” – George Box



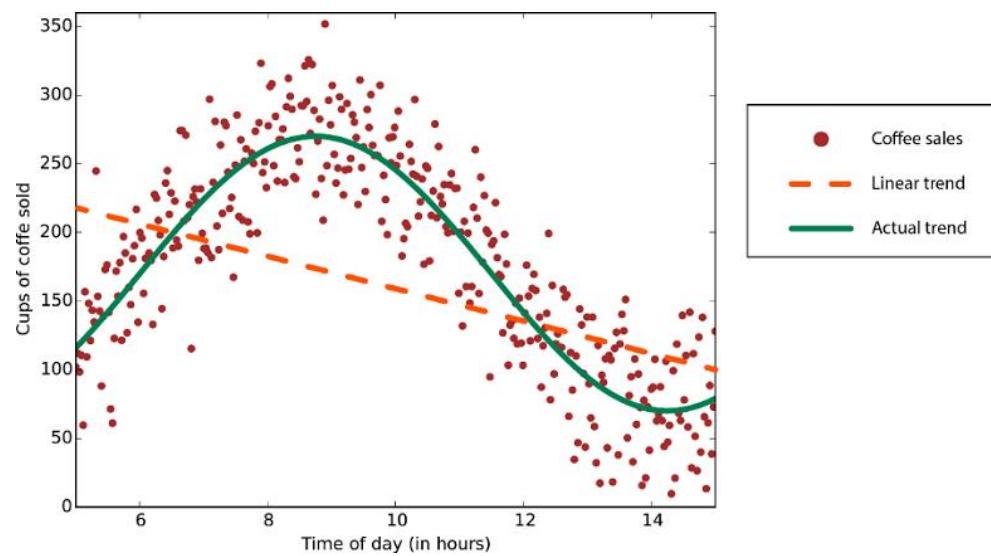
- Assumptions that work well in one domain may fail in another.
- **No Free Lunch Theorem (NFL):**
 - No single algorithm/model can perform optimally across all problems.
- Try:
 - More than one modeling approach
 - Different run parameters
 - “The knobs a data scientist gets to turn when setting up an algorithm to run”
 - Ensemble methods.

Non-Linear Class Boundaries



Linear classification
algorithm (e.g. SVM)

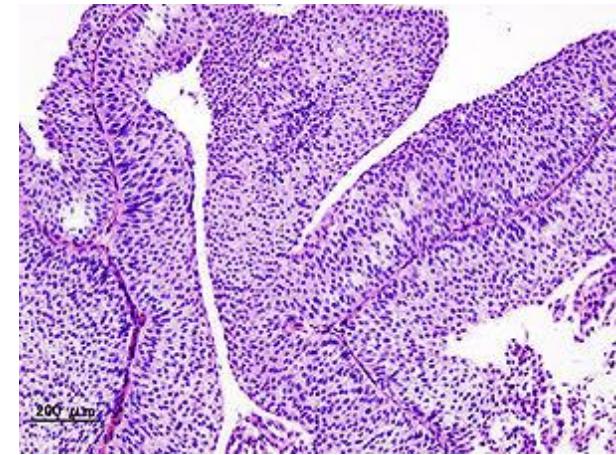
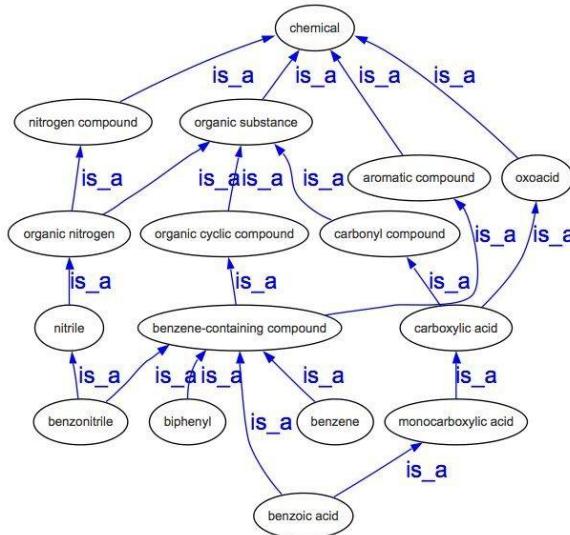
Linear regression
algorithm



Data: Types

[0, 1, 1, 1, 2, 1, 0, 0]

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$



data practice ehr vendor
health record health system hospital ehr incentive report
meaningful years
health care system information technology technology healthcare
medicare and medicaid health information clinical physicians system
ehr health information exchange medical electronic new
providers health information technology medical center emr
electronic medical records records information work
ehr incentive program management
accountable care organizations
clinical decision support
patient health care
time services improve
department of health
privacy and security
patient care care
doctors



Feature Extraction/Engineering

Example:

Email Spam Detection

From unstructured text...

...To meaningful features
for ML to interrogate.

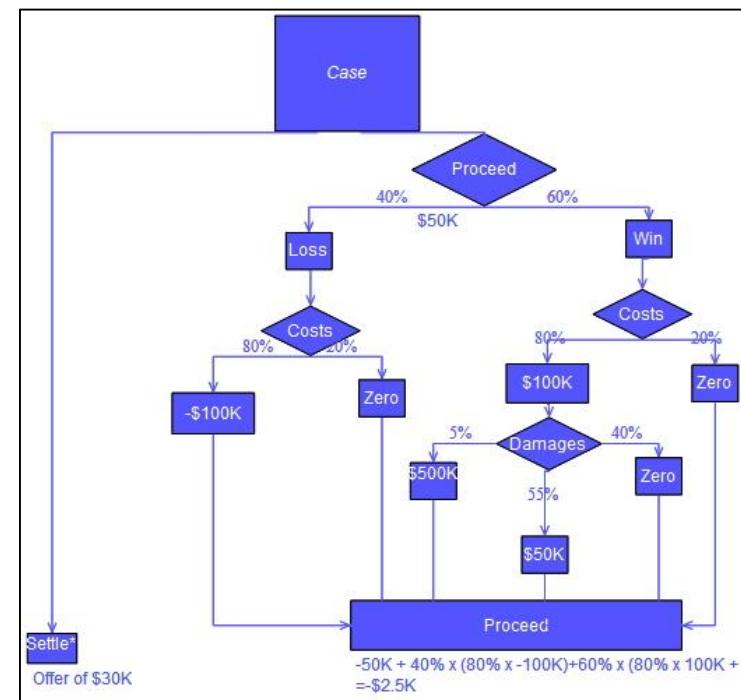
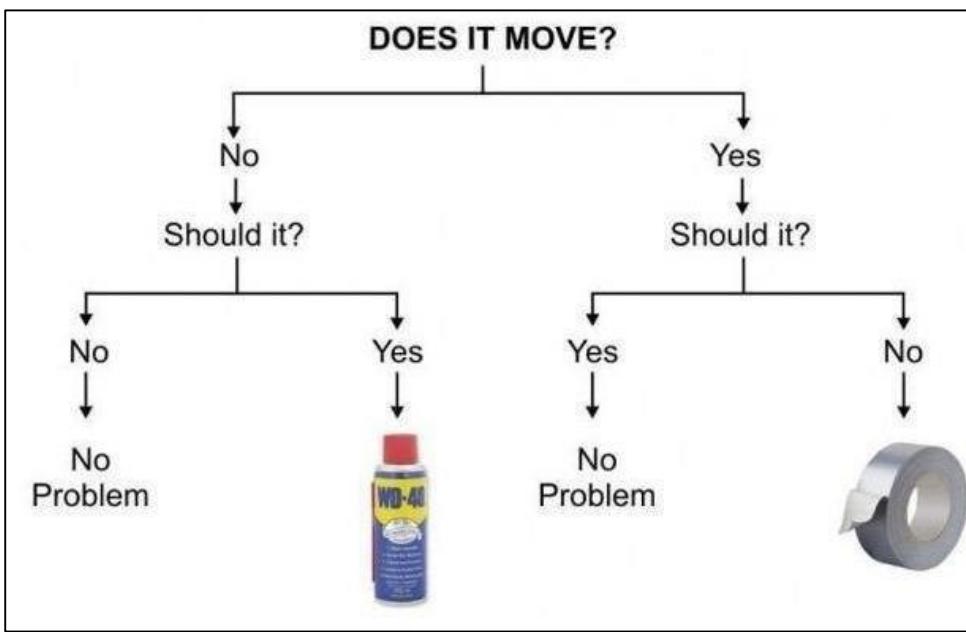
From: cheapsales@buystufffromme.com To: ang@cs.stanford.edu Subject: Buy now! Deal of the week! Buy now! Rolex w4tchs - \$100 Medlcine (any kind) - \$50 Also low cost M0rgages available.	From: Alfred Ng To: ang@cs.stanford.edu Subject: Christmas dates? Hey Andrew, Was talking to Mom about plans for Xmas. When do you get off work. Meet Dec 22? Alf
---	--

	“money”	“pills”	“Mr.”	bad spelling	known-sender	spam?
	Y	N	Y	Y	N	Y
	N	N	N	Y	Y	N
	N	Y	N	N	N	Y
example	Y	N	N	N	Y	N
	N	N	Y	N	Y	N
	Y	N	N	Y	N	Y
	N	N	Y	N	N	N



Decision Tree: What is it?

- A decision support tool: way to present information for decision making and evaluate their consequences (e.g. cost)

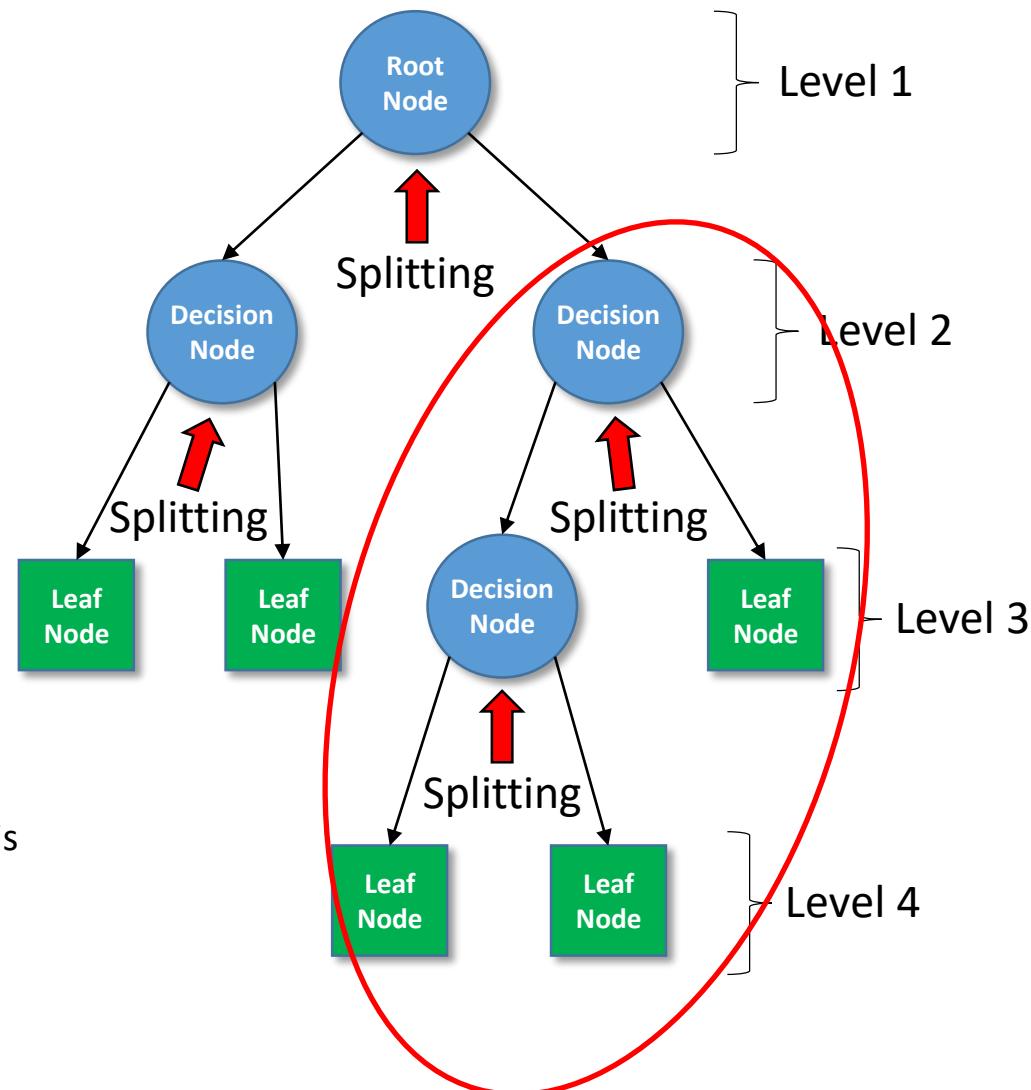


- A supervised, machine learning algorithm to model and predict outcomes

Decision Tree: Terminology

- **Nodes:**

- **Root:** It represents entire population or sample. Will get divided into two or more homogeneous sets.
- **Decision:** When a sub-node splits into further sub-nodes, then it is called decision node.
 - (AKA: Sub, internal, split, or chance node)
- **Leaf:** Nodes that don't split. Gives class or average value.
 - (AKA: Terminal, or outcome node)
- **Parent and Child:** Parent node splits into offspring nodes.



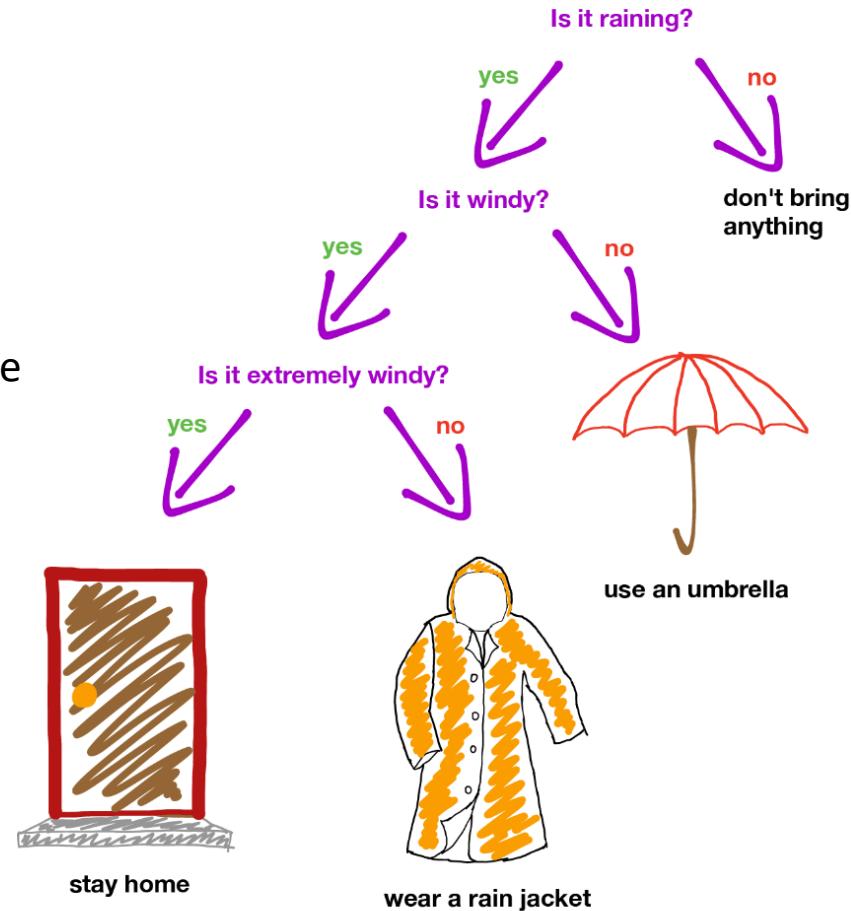
- **Splitting:** It is a process of dividing a node into two or more sub-nodes.

- **Branch / Sub-Tree:** A sub section of entire tree is called branch or sub-tree.

- **Levels/Depth:** The number of splits through a given path down the tree.

Decision Rules: Tree Interpretation

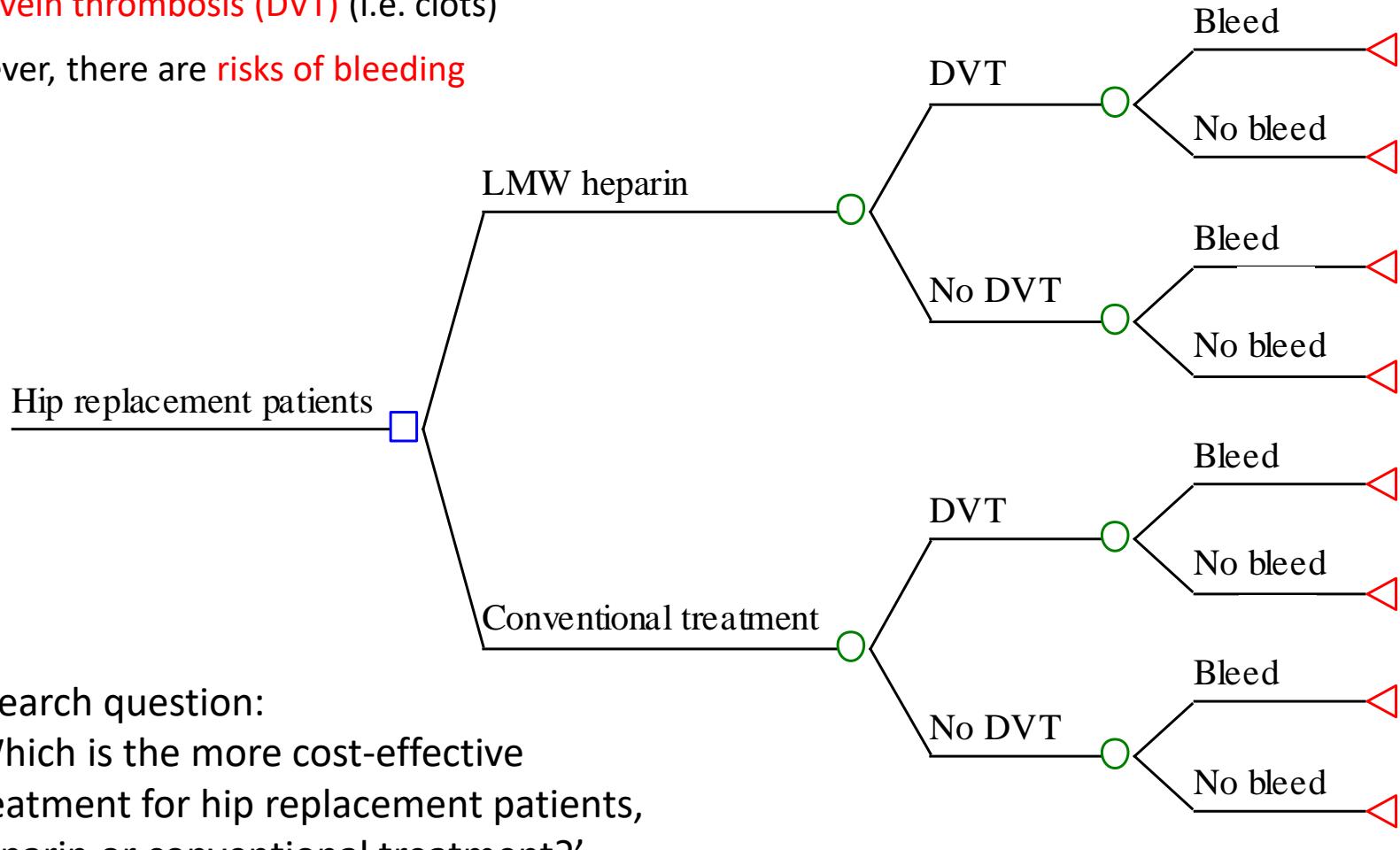
- Decision tree can be ‘linearized’ into **decision rules**.
 - One rule per path from root to leaf.
 - Rule outcome = Leaf node
- Rule:
 - If [condition1] and [condition2] Then: outcome
- Examples:
 - If [not raining] Then: Don’t bring anything
 - If [is raining] and [not windy] Then: use an umbrella



© Machine Learning @ Berkeley

Decision Tree for Heparin

- Heparin (anticoagulant) injection for the prevention of deep vein thrombosis (DVT) (i.e. clots)
- However, there are risks of bleeding



The research question:

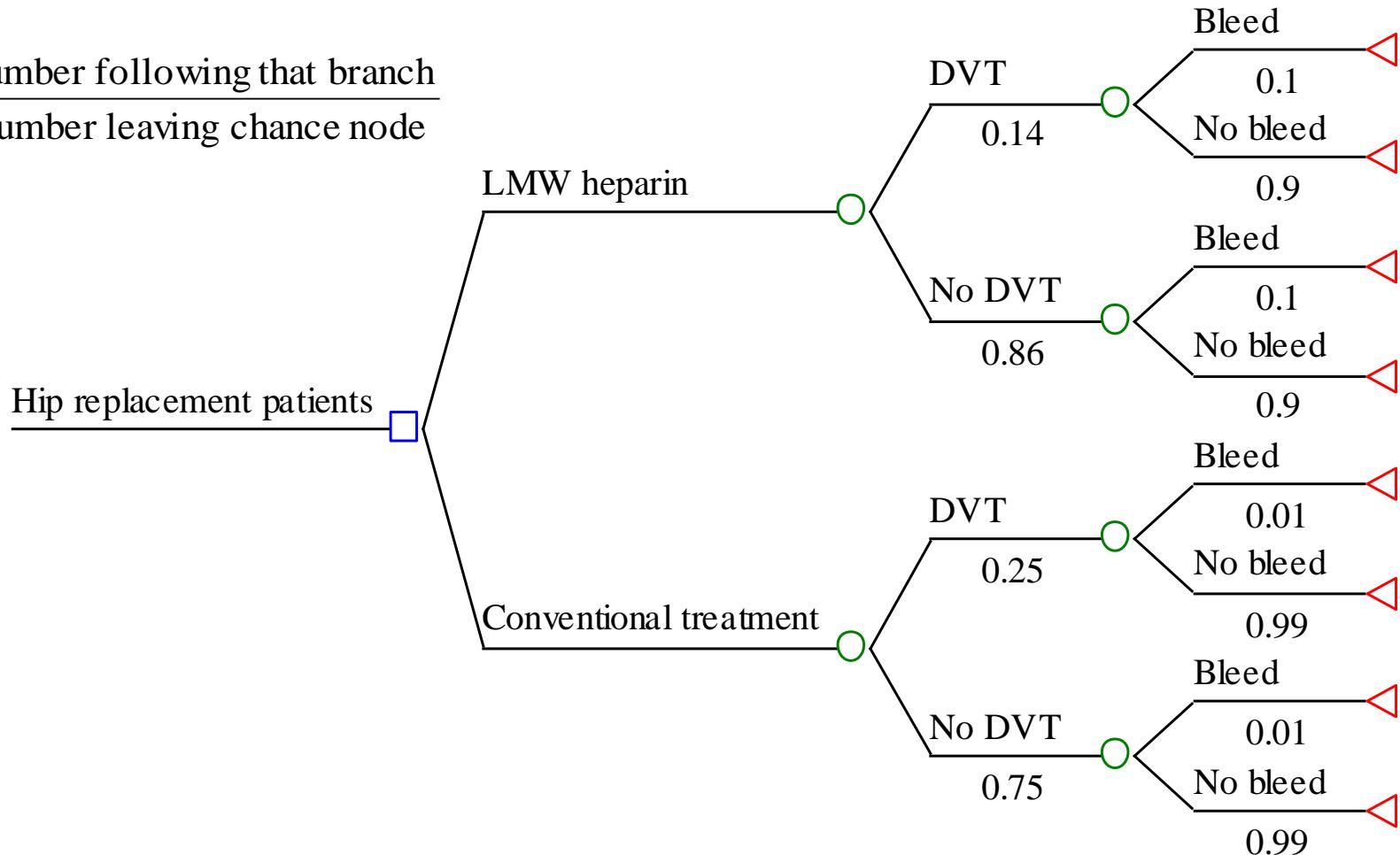
'Which is the more cost-effective treatment for hip replacement patients, heparin or conventional treatment?'



Decision Tree for Heparin

- Entering probabilities

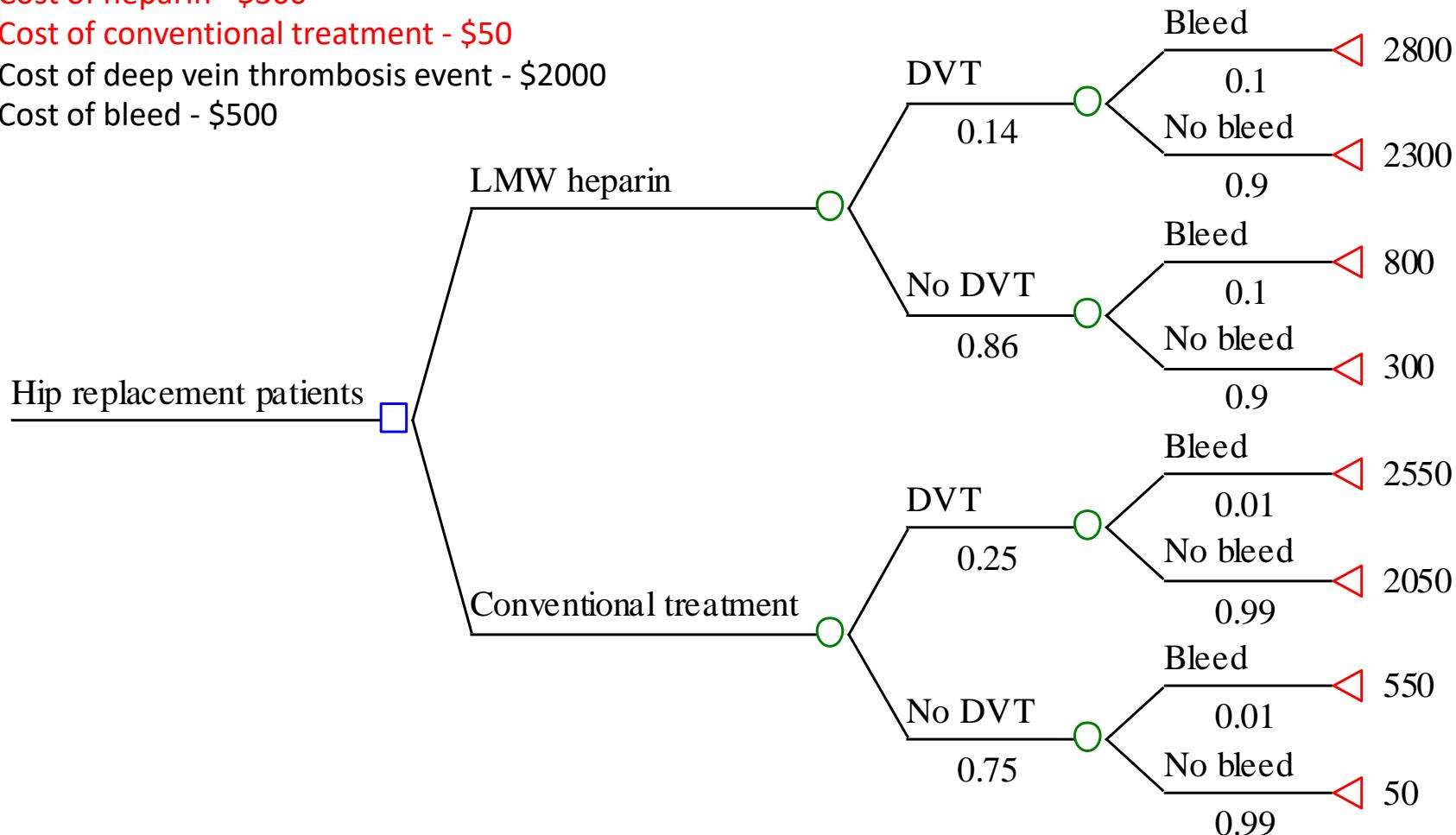
$$P = \frac{\text{Number following that branch}}{\text{Number leaving chance node}}$$



Evaluating Outcome Costs

- Costs assumed

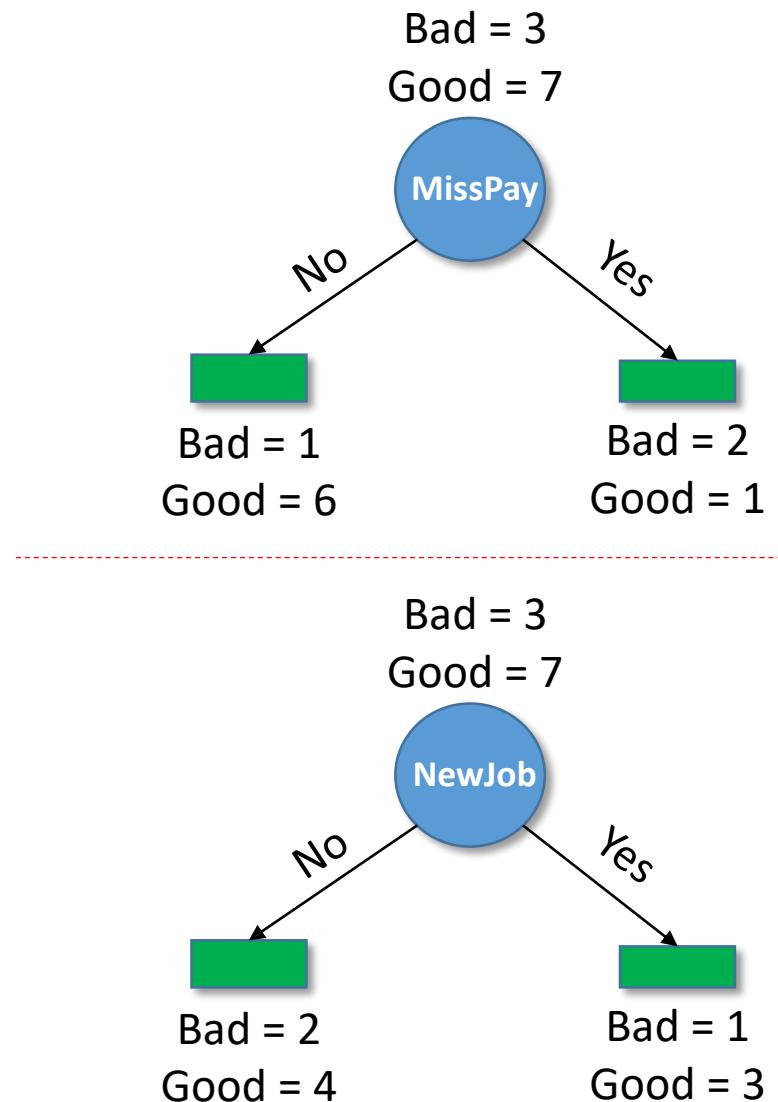
- Cost of heparin - \$300
- Cost of conventional treatment - \$50
- Cost of deep vein thrombosis event - \$2000
- Cost of bleed - \$500



Decision Tree: Choosing a Split

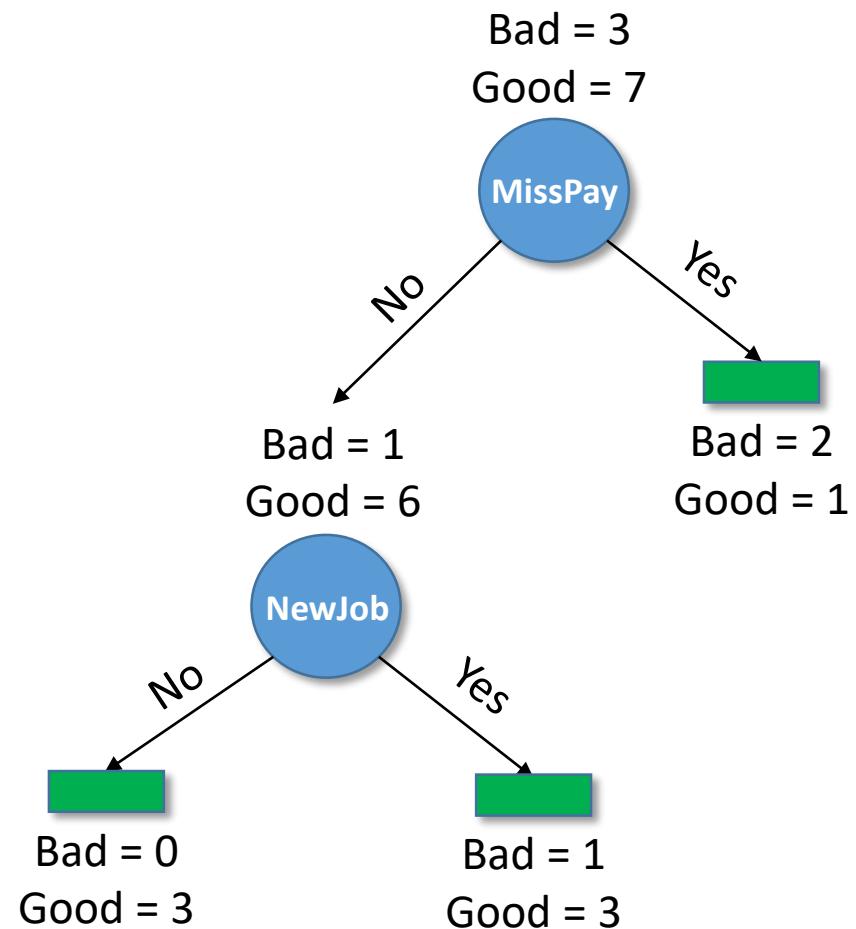
- E.g. Predicting Credit Risk
- What feature to split on?
- Want correct classification in fewest number of tests/branches.

	< 2 years at current job	Missed payments?	Credit
S1	N	N	Good
S2	Y	N	Bad
S3	N	N	Good
S4	N	N	Good
S5	N	Y	Bad
S6	Y	N	Good
S7	N	Y	Good
S8	N	Y	Bad
S9	Y	N	Good
S10	Y	N	Good



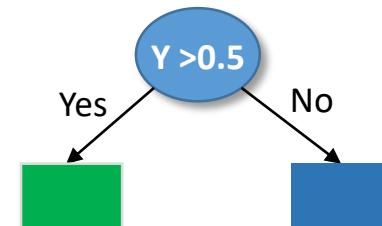
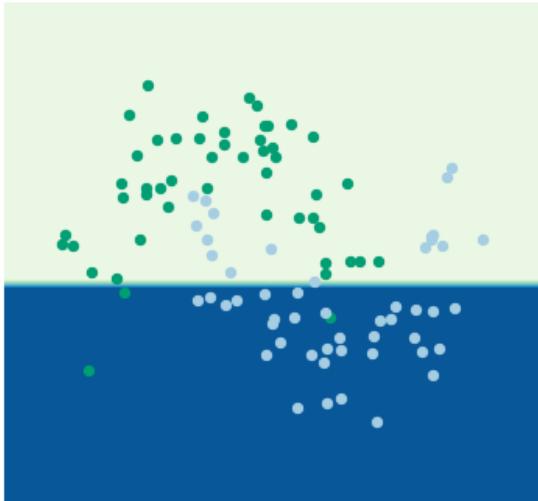
Decision Tree: Choosing a Split

	< 2 years at current job	Missed payments?	Credit
S1	N	N	Good
S2	Y	N	Bad
S3	N	N	Good
S4	N	N	Good
S5	N	Y	Bad
S6	Y	N	Good
S7	N	Y	Good
S8	N	Y	Bad
S9	Y	N	Good
S10	Y	N	Good



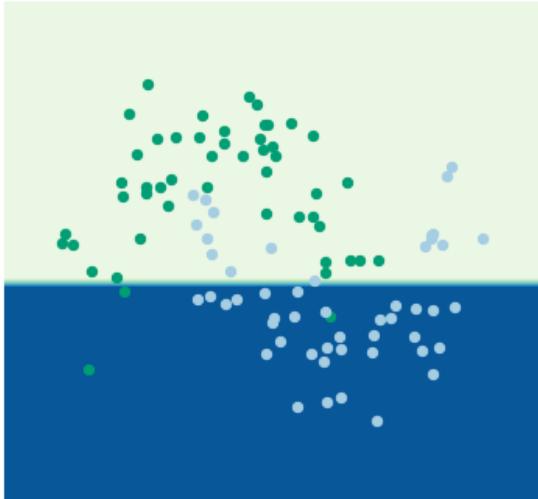
Decision Tree: Fitting with Splits

Max Depth: 1

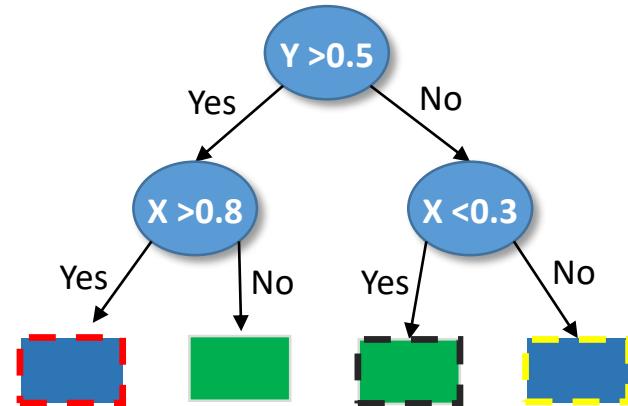
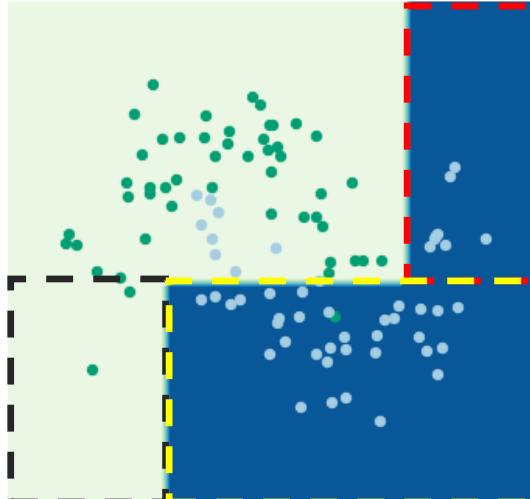


Decision Tree: Fitting with Splits

Max Depth: 1

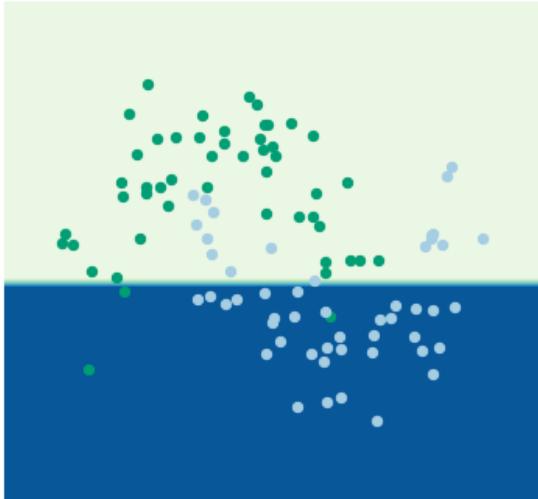


Max Depth: 2

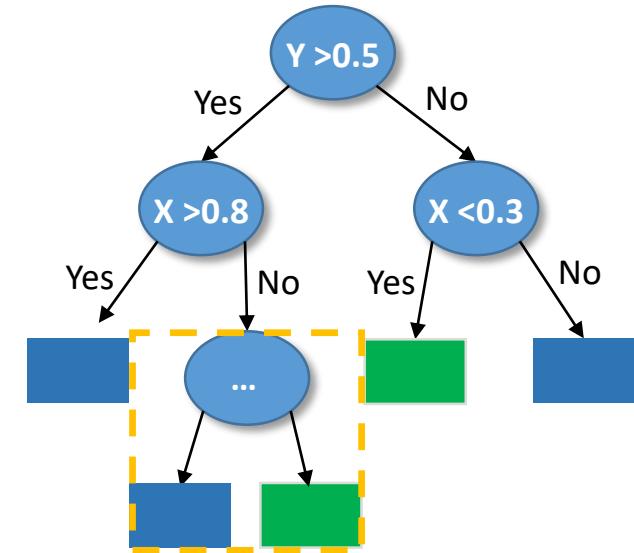
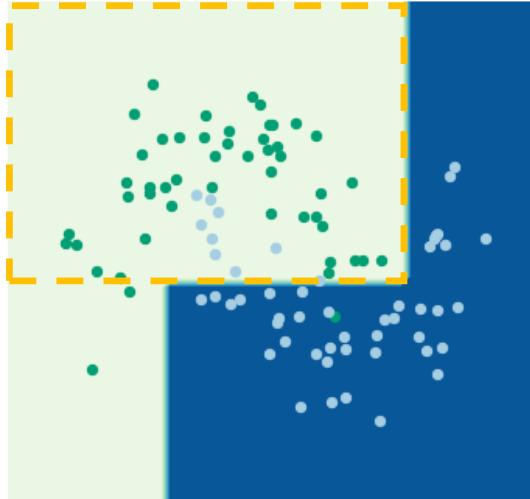


Decision Tree: Fitting with Splits

Max Depth: 1

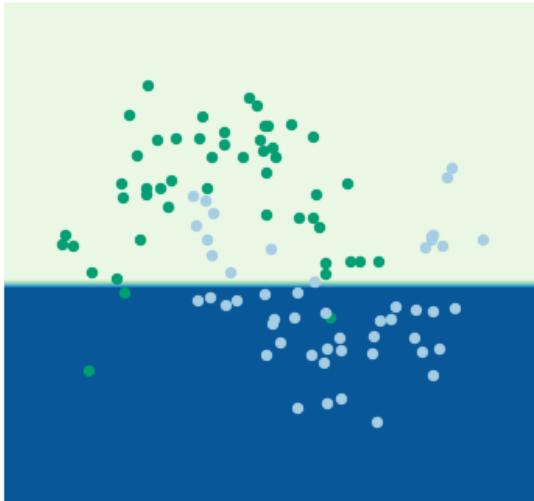


Max Depth: 2

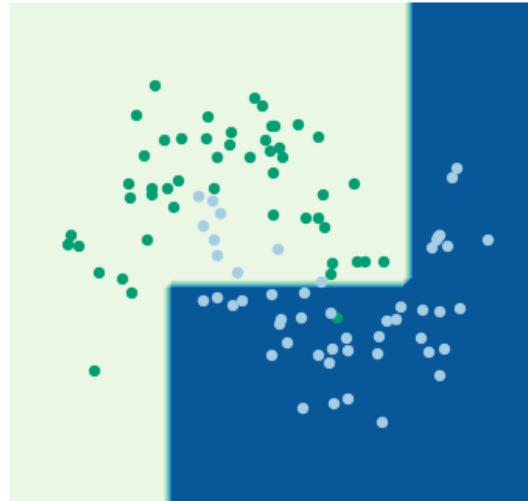


Decision Tree: Fitting with Splits

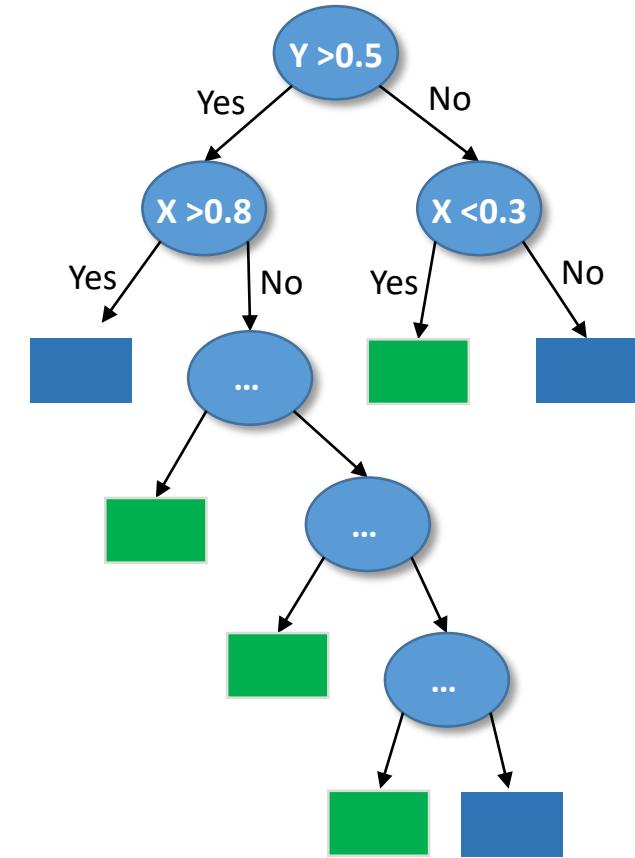
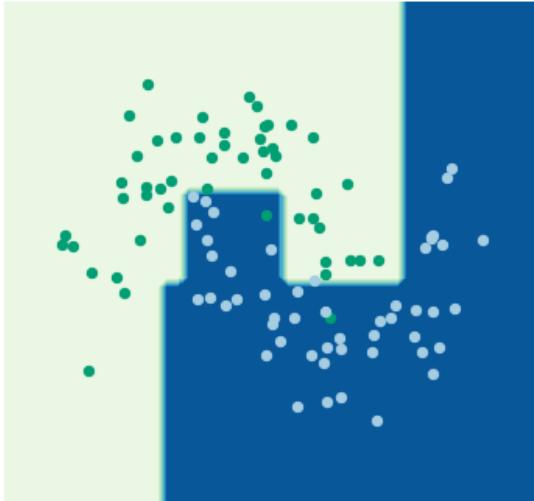
Max Depth: 1



Max Depth: 2

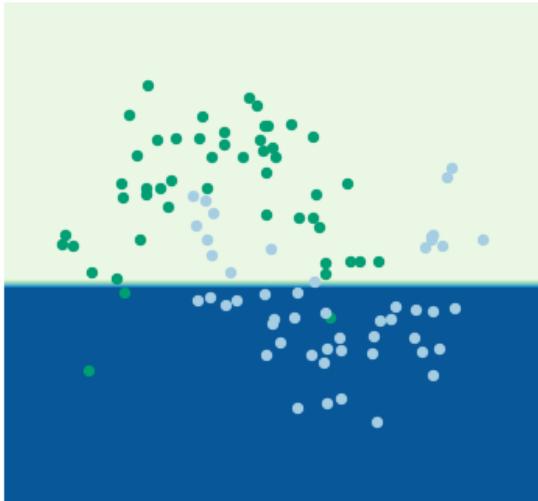


Max Depth: 5

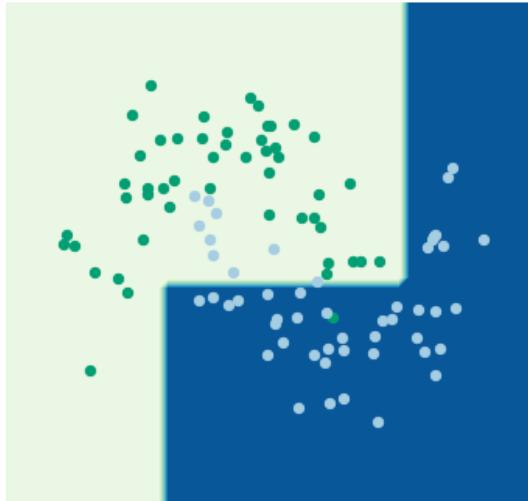


Decision Tree: Fitting with Splits

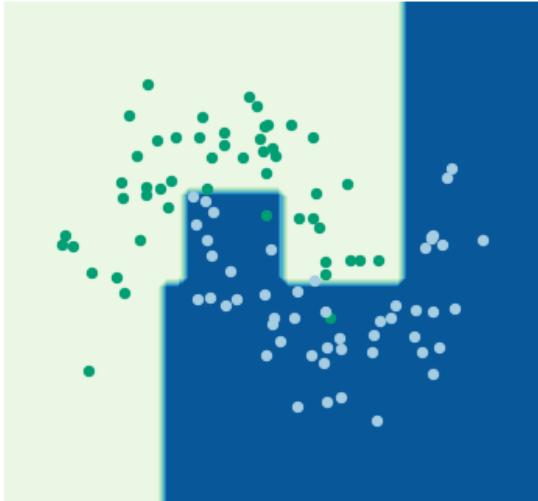
Max Depth: 1



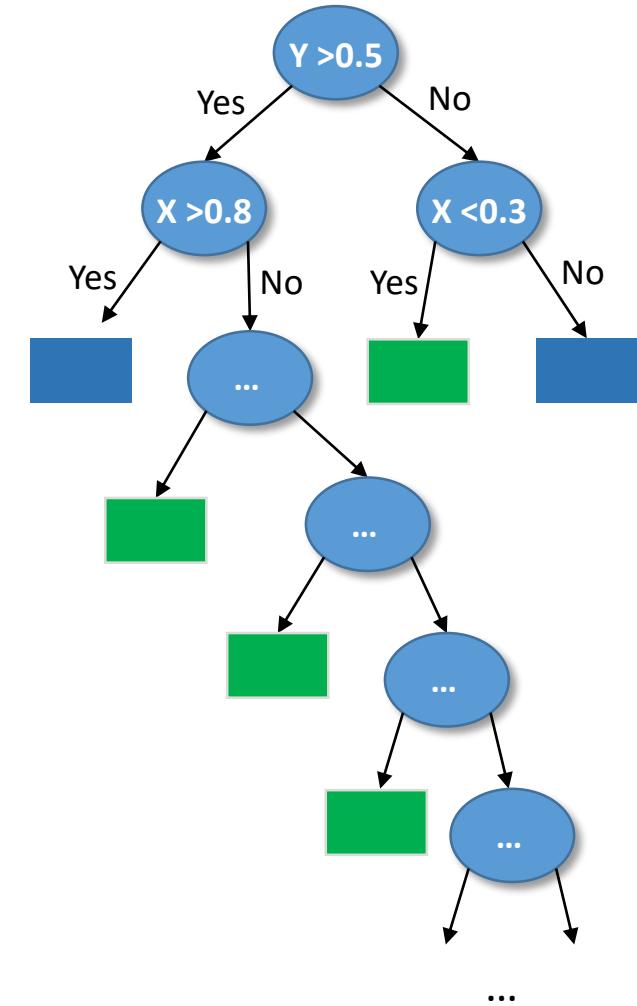
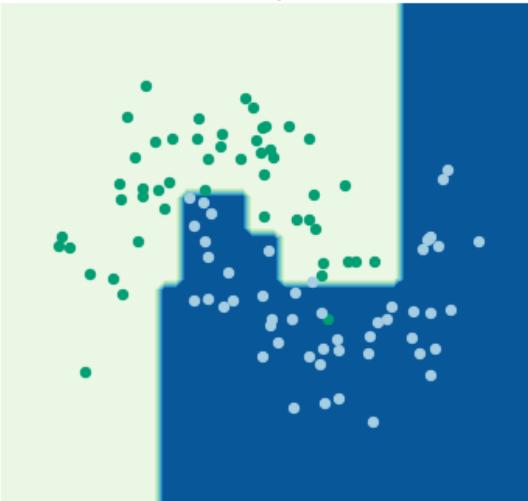
Max Depth: 2



Max Depth: 5



Max Depth: 10



Decision Tree Challenges

- How do we decide best feature or value to split on?
- When should we stop splitting?
- What do we do if we can't achieve perfect classification?
- What if the tree is too large? Can we approximate a smaller one?



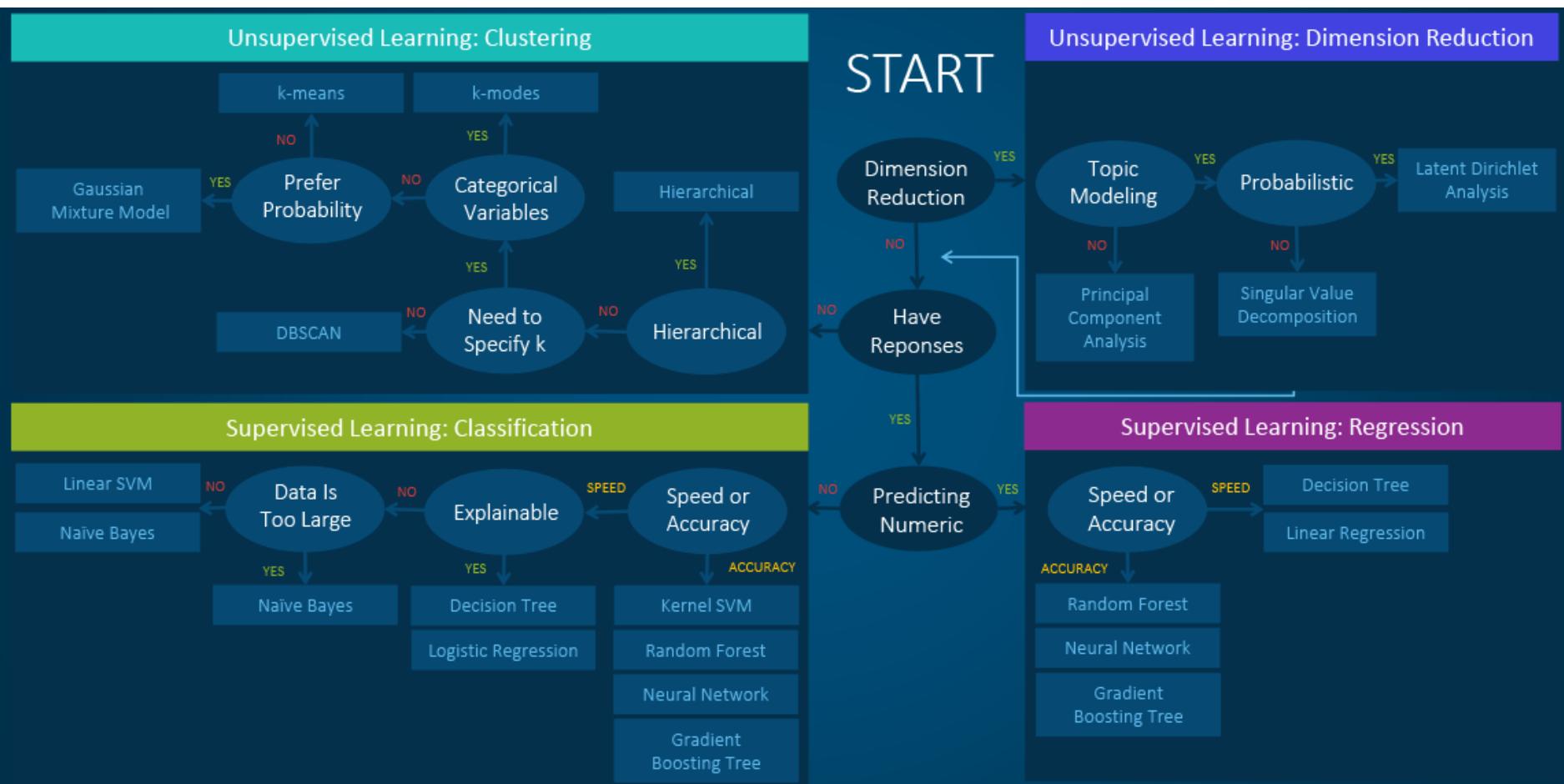
Where to start in selecting a method?

- If there is a strong, simple relationship among variables, most methods will find it.
- Generally start with simpler methods if you know nothing about the problem.
- When possible, **limit the search space with knowledge/assumptions** about the problem.
 - E.g. If we want to know if there are linear patterns, use linear regression.
- **Incorrect assumptions will limit or invalidate what can be found.**



Considerations When Choosing an ML Algorithm

- Data – Labeled?, Endpoint?
- Training Time / Run Speed
- Number and Importance of Parameters
- Data Size – Features, Instances
- Interpretability
- Assumptions



ML Performance Evaluation

$$Accuracy = \frac{\text{Number of Correct predictions}}{\text{Total number of predictions made}}$$

$$\text{LogarithmicLoss} = \frac{-1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} * \log(p_{ij})$$

$$\text{TruePositiveRate} = \frac{\text{TruePositive}}{\text{FalseNegative} + \text{TruePositive}}$$

$$\text{FalsePositiveRate} = \frac{\text{FalsePositive}}{\text{FalsePositive} + \text{TrueNegative}}$$

$$\text{Precision} = \frac{\text{TruePositives}}{\text{TruePositives} + \text{FalsePositives}}$$

$$\text{Recall} = \frac{\text{TruePositives}}{\text{TruePositives} + \text{FalseNegatives}}$$

$$F1 = 2 * \frac{1}{\frac{1}{precision} + \frac{1}{recall}}$$

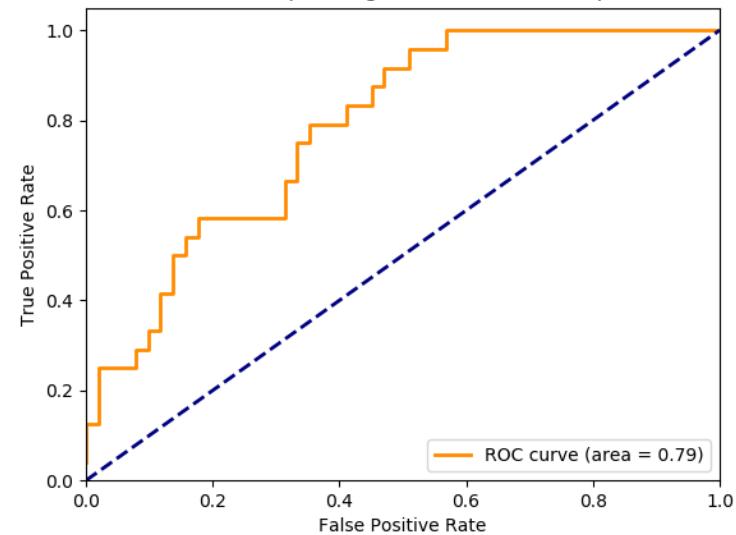
$$\text{MeanAbsoluteError} = \frac{1}{N} \sum_{j=1}^N |y_j - \hat{y}_j|$$

$$\text{MeanSquaredError} = \frac{1}{N} \sum_{j=1}^N (y_j - \hat{y}_j)^2$$

Confusion Matrix

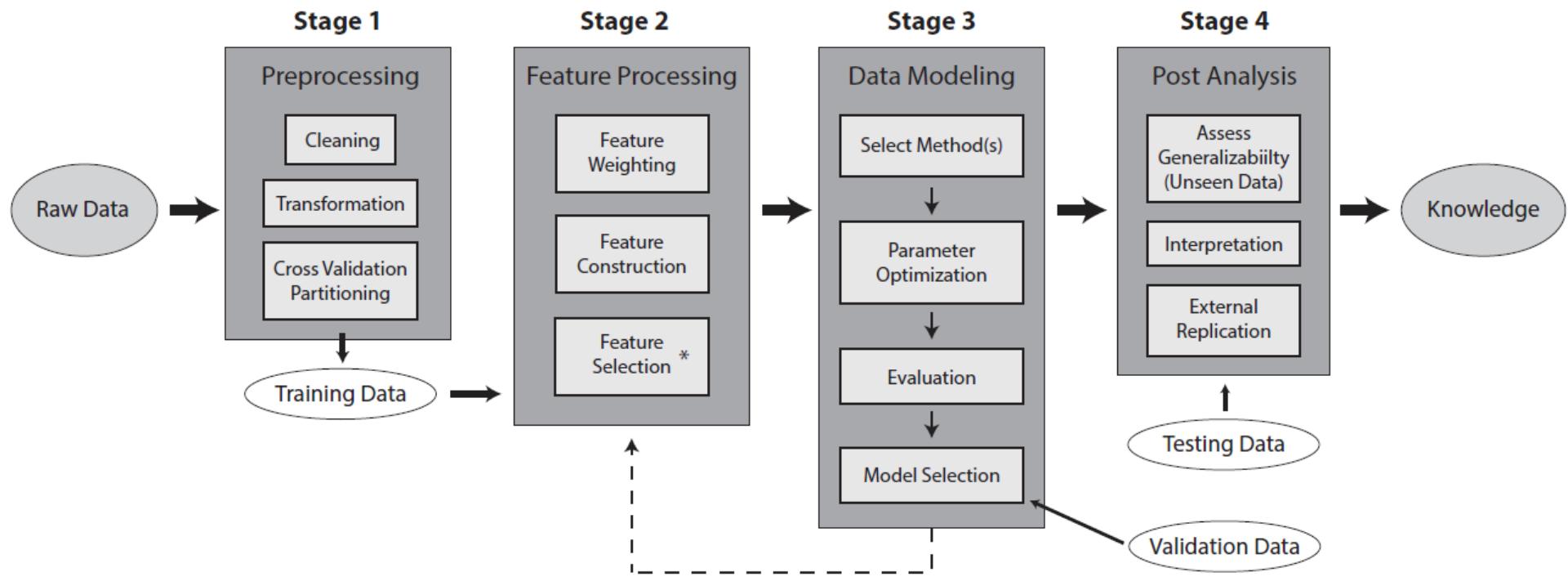
		Predicted: NO	Predicted: YES
n=165	Actual: NO	50	10
	Actual: YES	5	100

Receiver operating characteristic example



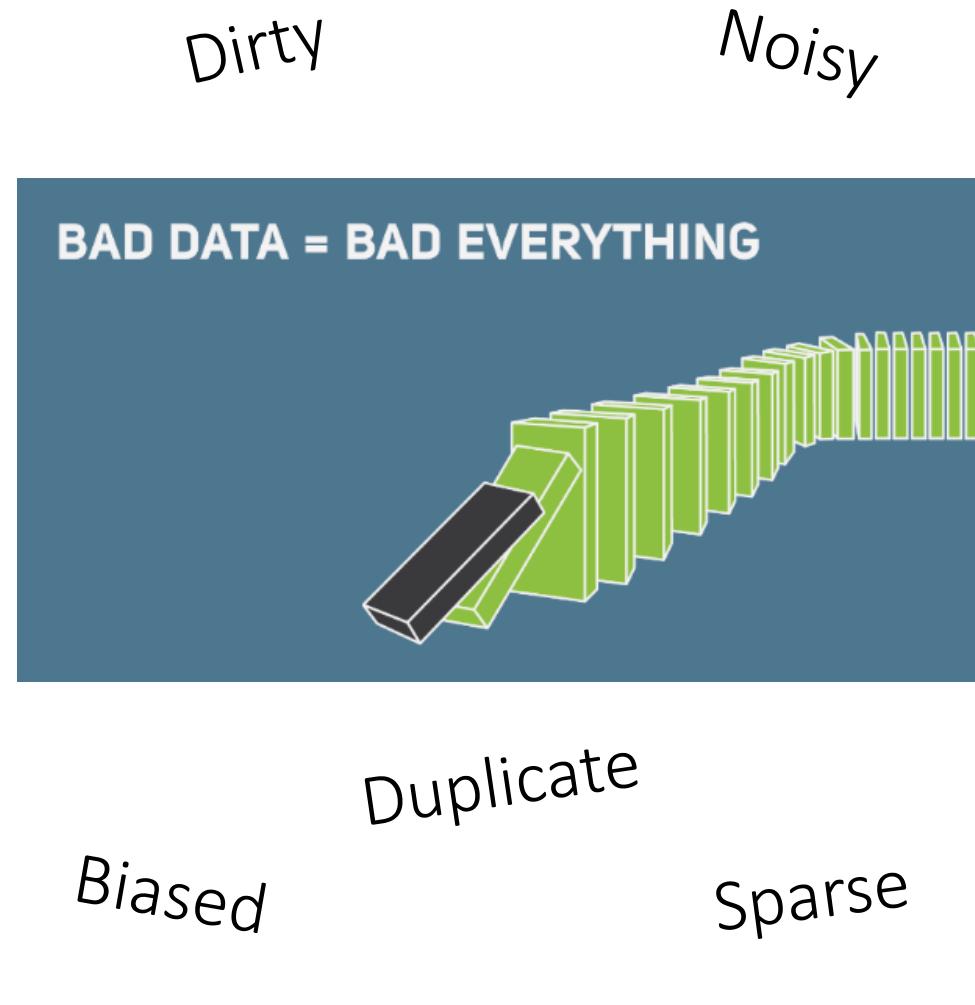
<https://towardsdatascience.com/metrics-to-evaluate-your-machine-learning-algorithm-f10ba6e38234>

Data Mining Pipeline



Common Machine Learning Pitfalls

- **Working with bad data**
- Data leakage
- Not understanding the target problem
- Ignoring exploratory analysis
- Handling missing data
- Ignoring assumptions
- Representable does not imply learnable
- Sampling bias
- Overfitting
- Simplicity does not imply better generalizability
- Using the default parameters
- Failing to use an appropriate evaluation metric
- Data dredging
- Mistaking correlation for causation
- Failing to consider confounding variables



Common Machine Learning Pitfalls

- Working with bad data
- **Data leakage**
- Not understanding the target problem
- Ignoring exploratory analysis
- Handling missing data
- Ignoring assumptions
- Representable does not imply learnable
- Sampling bias
- Overfitting
- Simplicity does not imply better generalizability
- Using the default parameters
- Failing to use an appropriate evaluation metric
- Data dredging
- Mistaking correlation for causation
- Failing to consider confounding variables



Common Machine Learning Pitfalls

- Working with bad data
- Data leakage
- **Not defining the target problem/goals**
- Ignoring exploratory analysis
- Handling missing data
- Ignoring assumptions
- Representable does not imply learnable
- Sampling bias
- Overfitting
- Simplicity does not imply better generalizability
- Using the default parameters
- Failing to use an appropriate evaluation metric
- Data dredging
- Mistaking correlation for causation
- Failing to consider confounding variables

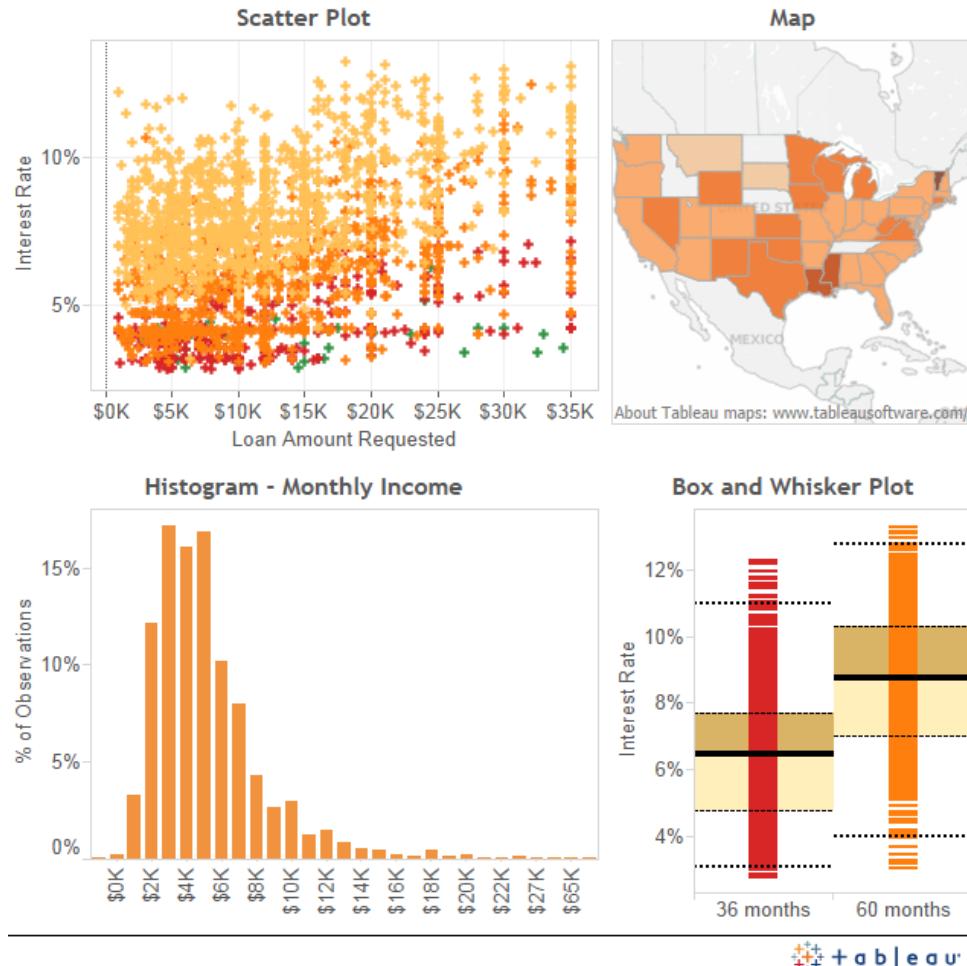
„A problem well stated
is a problem half solved.“

Charles Kettering (1876-1958)



Common Machine Learning Pitfalls

- Working with bad data
- Data leakage
- Not defining the target problem/goals
- **Ignoring exploratory analysis**
- Handling missing data
- Ignoring assumptions
- Representable does not imply learnable
- Sampling bias
- Overfitting
- Simplicity does not imply better generalizability
- Using the default parameters
- Failing to use an appropriate evaluation metric
- Data dredging
- Mistaking correlation for causation
- Failing to consider confounding variables



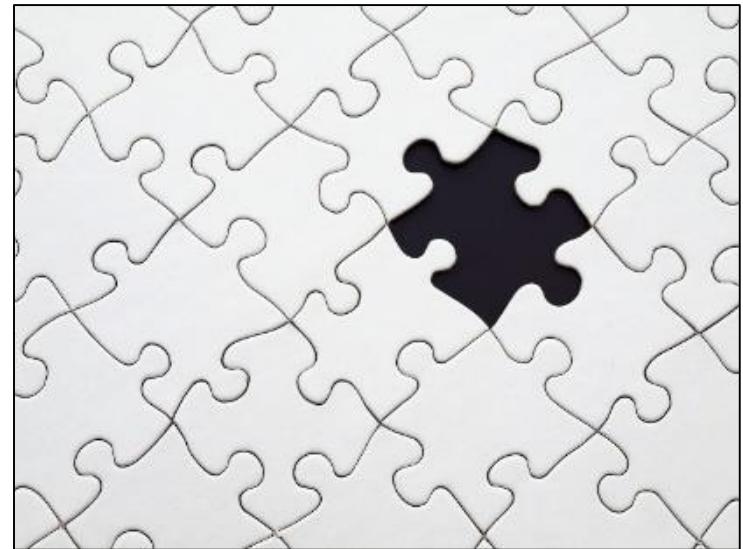
tableau



Common Machine Learning Pitfalls

- Working with bad data
- Data leakage
- Not defining the target problem/goals
- Ignoring exploratory analysis
- **Handling missing data**
- Ignoring assumptions
- Representable does not imply learnable
- Sampling bias
- Overfitting
- Simplicity does not imply better generalizability
- Using the default parameters
- Failing to use an appropriate evaluation metric
- Data dredging
- Mistaking correlation for causation
- Failing to consider confounding variables

- Different types of ‘missingness’



- Handling:
 - Removal
 - Imputation
 - Encoding as Features

Common Machine Learning Pitfalls

- Working with bad data
- Data leakage
- Not defining the target problem/goals
- Ignoring exploratory analysis
- Handling missing data
- **Ignoring assumptions**
- Representable does not imply learnable
- Sampling bias
- Overfitting
- Simplicity does not imply better generalizability
- Using the default parameters
- Failing to use an appropriate evaluation metric
- Data dredging
- Mistaking correlation for causation
- Failing to consider confounding variables



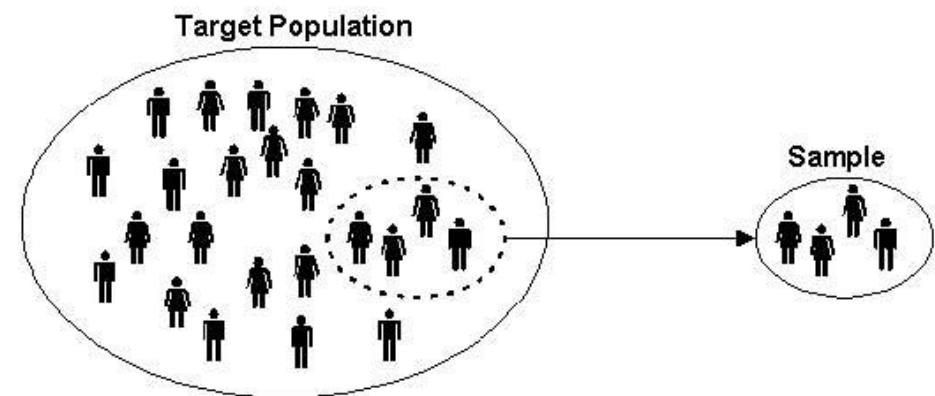
Common Machine Learning Pitfalls

- Working with bad data
- Data leakage
- Not defining the target problem/goals
- Ignoring exploratory analysis
- Handling missing data
- Ignoring assumptions
- **Representable does not imply learnable**
- Sampling bias
- Overfitting
- Simplicity does not imply better generalizability
- Using the default parameters
- Failing to use an appropriate evaluation metric
- Data dredging
- Mistaking correlation for causation
- Failing to consider confounding variables



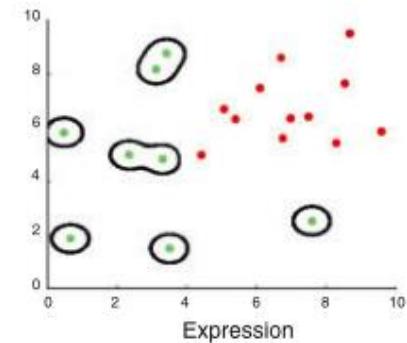
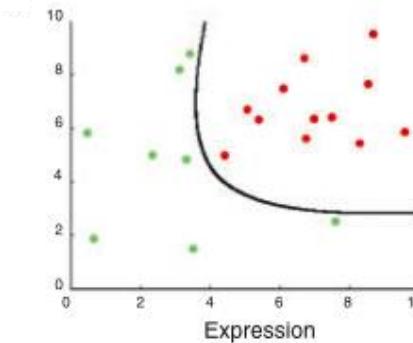
Common Machine Learning Pitfalls

- Working with bad data
- Data leakage
- Not defining the target problem/goals
- Ignoring exploratory analysis
- Handling missing data
- Ignoring assumptions
- Representable does not imply learnable
- **Sampling bias**
- Overfitting
- Simplicity does not imply better generalizability
- Using the default parameters
- Failing to use an appropriate evaluation metric
- Data dredging
- Mistaking correlation for causation
- Failing to consider confounding variables



Common Machine Learning Pitfalls

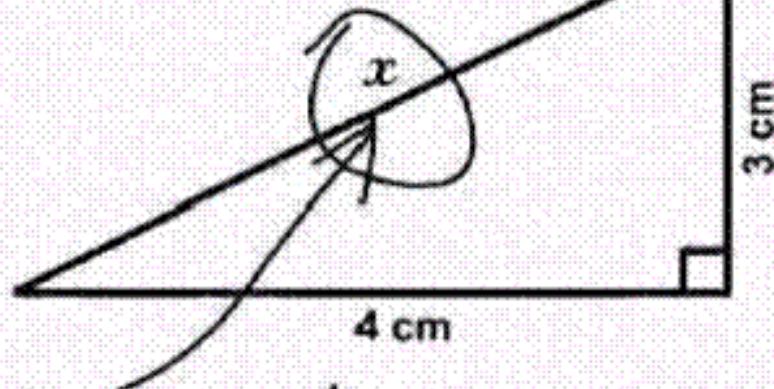
- Working with bad data
- Data leakage
- Not defining the target problem/goals
- Ignoring exploratory analysis
- Handling missing data
- Ignoring assumptions
- Representable does not imply learnable
- Sampling bias
- **Overfitting**
- Simplicity does not imply better generalizability
- Using the default parameters
- Failing to use an appropriate evaluation metric
- Data dredging
- Mistaking correlation for causation
- Failing to consider confounding variables



Common Machine Learning Pitfalls

- Working with bad data
- Data leakage
- Not defining the target problem/goals
- Ignoring exploratory analysis
- Handling missing data
- Ignoring assumptions
- Representable does not imply learnable
- Sampling bias
- Overfitting
- **Simplicity does not imply better generalizability**
- Using the default parameters
- Failing to use an appropriate evaluation metric
- Data dredging
- Mistaking correlation for causation
- Failing to consider confounding variables

3. Find x .



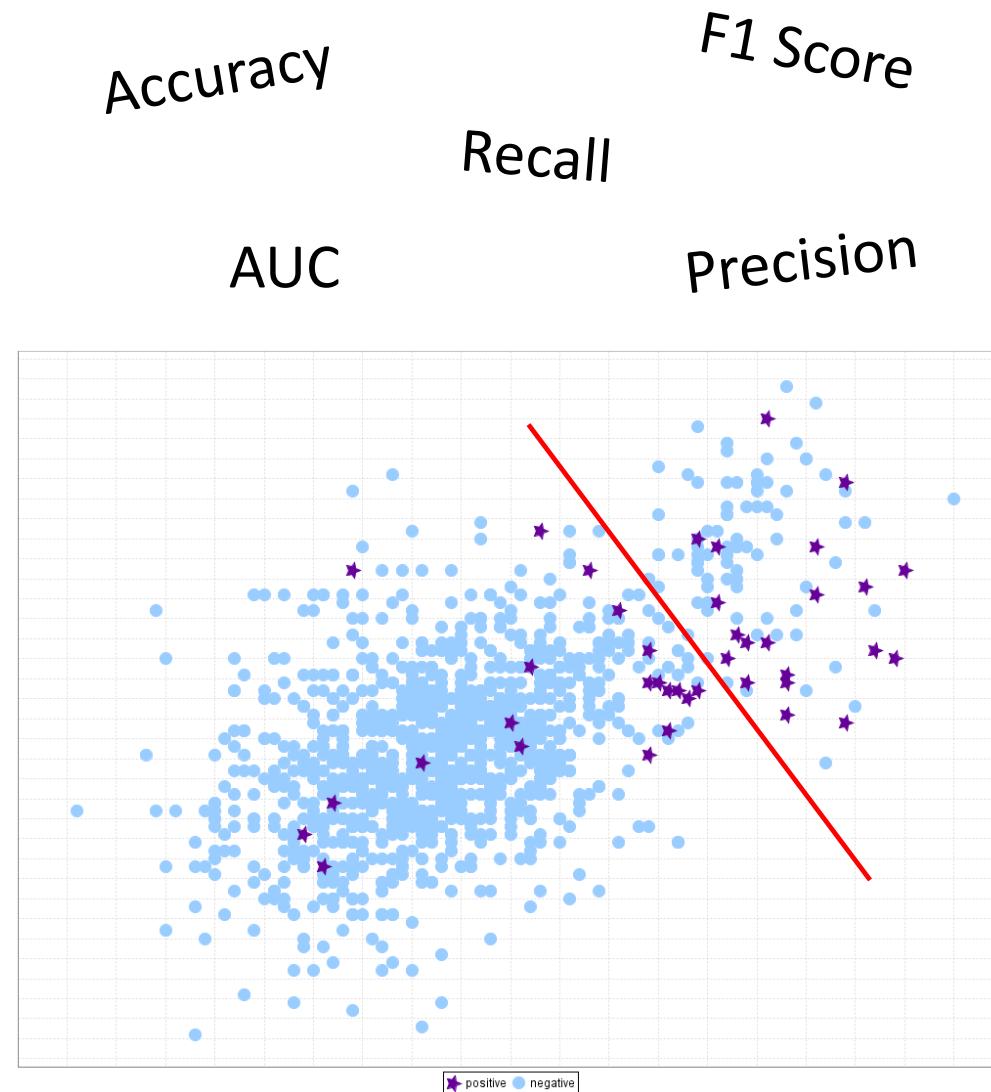
Common Machine Learning Pitfalls

- Working with bad data
- Data leakage
- Not defining the target problem/goals
- Ignoring exploratory analysis
- Handling missing data
- Ignoring assumptions
- Representable does not imply learnable
- Sampling bias
- Overfitting
- Simplicity does not imply better generalizability
- **Using the default parameters**
- Failing to use an appropriate evaluation metric
- Data dredging
- Mistaking correlation for causation
- Failing to consider confounding variables



Common Machine Learning Pitfalls

- Working with bad data
- Data leakage
- Not defining the target problem/goals
- Ignoring exploratory analysis
- Handling missing data
- Ignoring assumptions
- Representable does not imply learnable
- Sampling bias
- Overfitting
- Simplicity does not imply better generalizability
- Using the default parameters
- **Failing to use an appropriate evaluation metric**
- Data dredging
- Mistaking correlation for causation
- Failing to consider confounding variables



Common Machine Learning Pitfalls

- Working with bad data
- Data leakage
- Not defining the target problem/goals
- Ignoring exploratory analysis
- Handling missing data
- Ignoring assumptions
- Representable does not imply learnable
- Sampling bias
- Overfitting
- Simplicity does not imply better generalizability
- Using the default parameters
- Failing to use an appropriate evaluation metric
- **Data dredging**
- Mistaking correlation for causation
- Failing to consider confounding variables

If you torture the data long enough,
it will confess.

— Ronald Coase —

Data Fishing

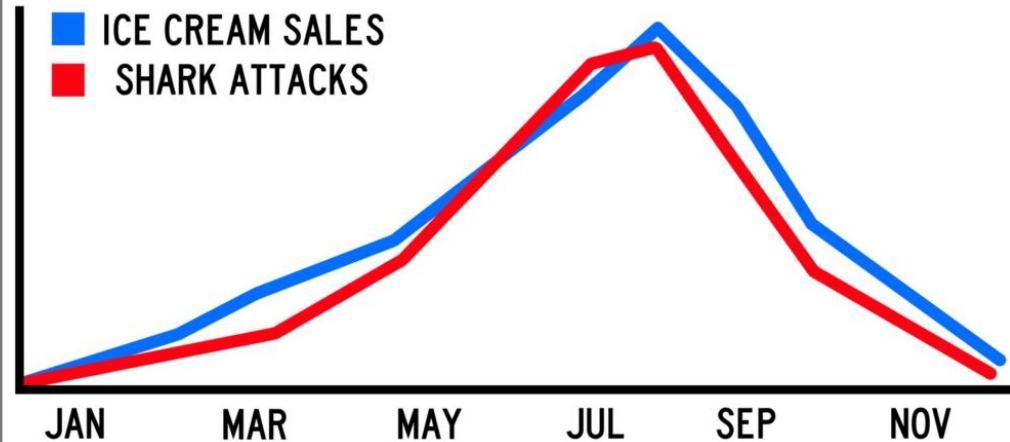
Data Snooping

P-hacking



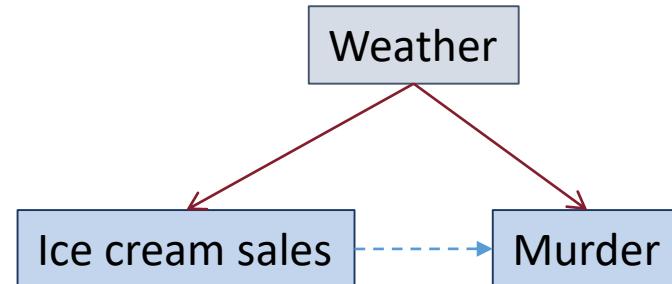
Common Machine Learning Pitfalls

- Working with bad data
- Data leakage
- Not defining the target problem/goals
- Ignoring exploratory analysis
- Handling missing data
- Ignoring assumptions
- Representable does not imply learnable
- Sampling bias
- Overfitting
- Simplicity does not imply better generalizability
- Using the default parameters
- Failing to use an appropriate evaluation metric
- Data dredging
- **Mistaking correlation for causation**
- Failing to consider confounding variables



Common Machine Learning Pitfalls

- Working with bad data
- Data leakage
- Not defining the target problem/goals
- Ignoring exploratory analysis
- Handling missing data
- Ignoring assumptions
- Representable does not imply learnable
- Sampling bias
- Overfitting
- Simplicity does not imply better generalizability
- Using the default parameters
- Failing to use an appropriate evaluation metric
- Data dredging
- Mistaking correlation for causation
- **Failing to consider confounding variables**

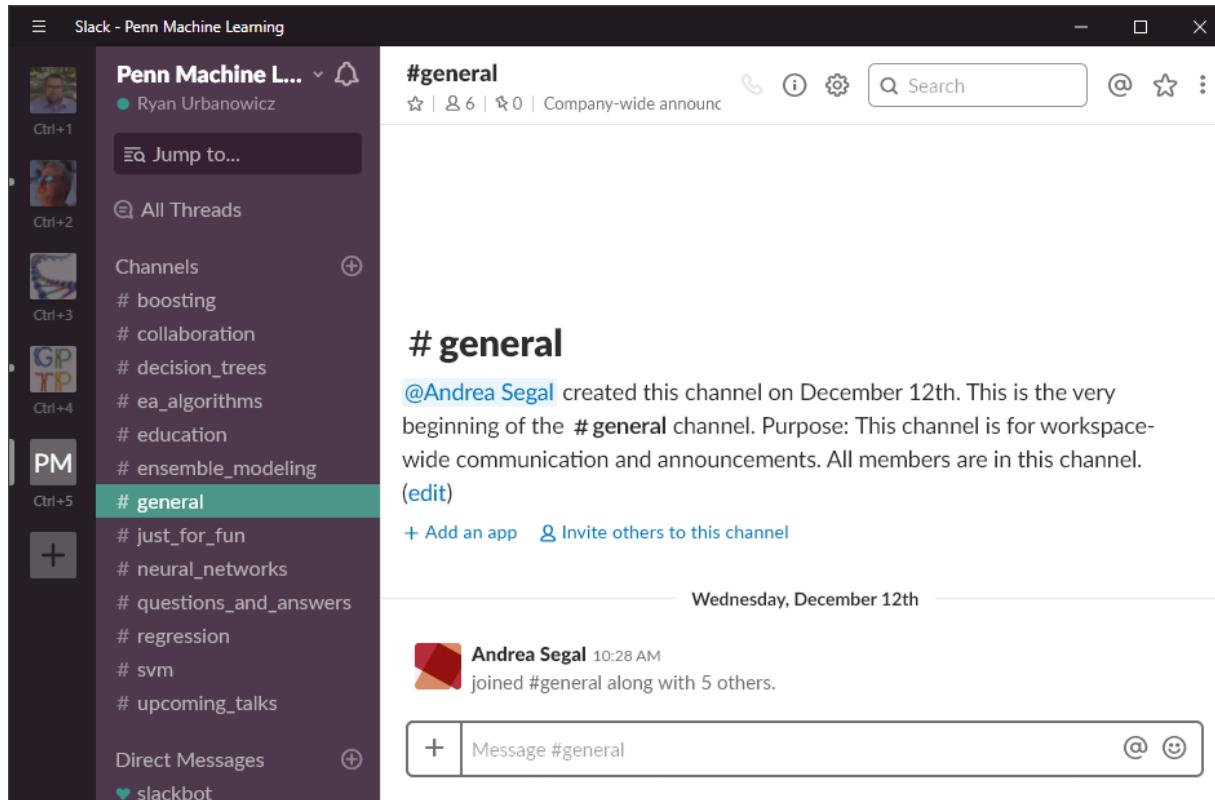


Where do we go from here?

- Data preparation
- How do different ML methods work?
- Feature selection
- Selecting run parameters
- Software/code to run ML
- Evaluation and statistical analysis
- Ensemble learning
- Model interpretation

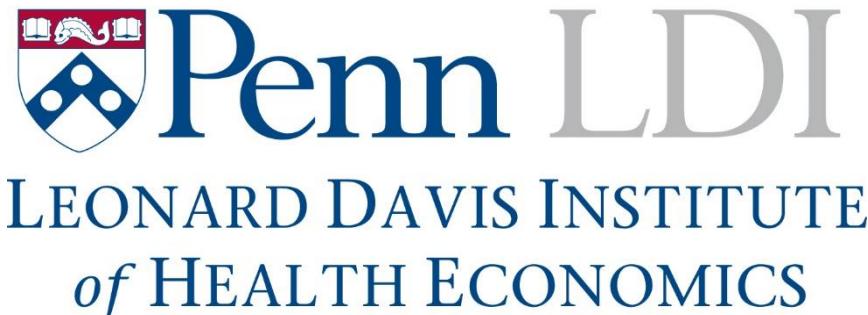


- Penn Machine Learning – Slack Workspace
- pennmachinelearning.slack.com



Acknowledgements and Funding

- Pennsylvania Commonwealth Universal Research Enhancement Program (CURE)



DEPARTMENT of
BIOSTATISTICS
EPIDEMILOGY &
INFORMATICS

