# 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
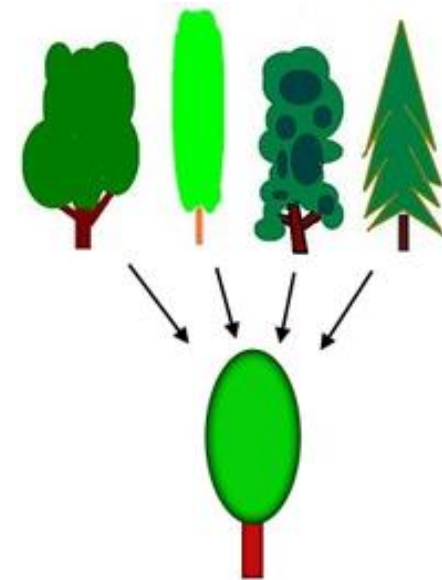
- Common ML Pitfalls

# 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 programmed[1].
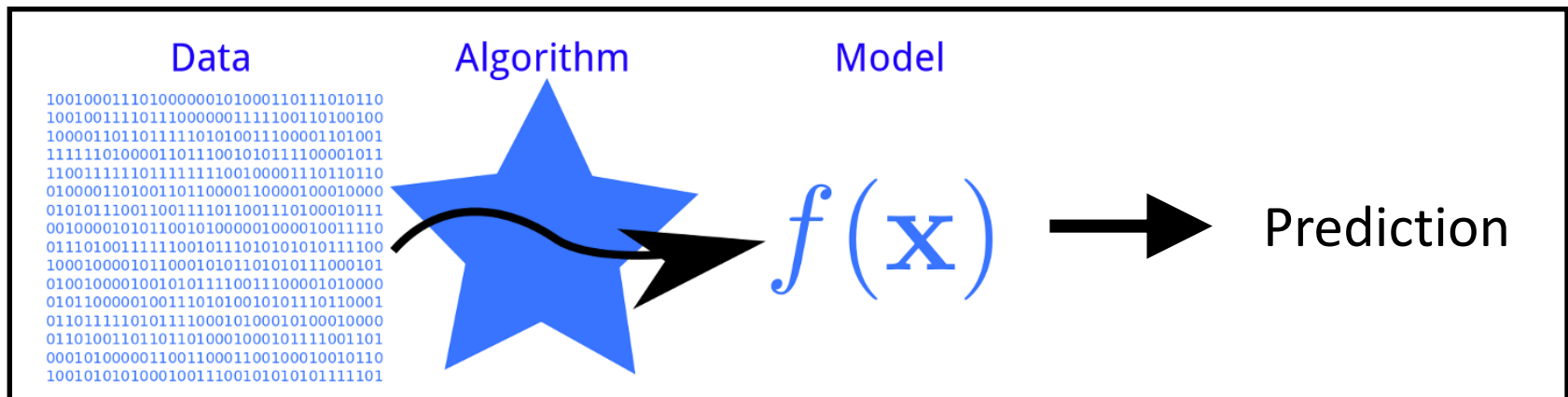
  [1] Samuel Arthur – 1959 – ML in Checkers

- 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



| Data | Algorithm | Model | |
|---|---|---|---|
| 1001000111010000010100011011101011000100111101110000011111001101001001000011010101111010100111000011010011111011010000110111001010111100001011110011111011011111111001000001110110110010000110100110110000110000010001000001010111001100110011110110011101000101110010000101001010010100010001001111010011111110010111101010101011110010001000100101100010101101010101100010100100001001010111100110000011000000100101011110011100011100001000001011000001001110101010101111011000101101111110101110001010001010001000001101001101101101000010001011110011010001010000011001100110010010010110010101010100010011100101010101111101 | f(x) | → | Prediction |

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

# Example: Email Spam Detection

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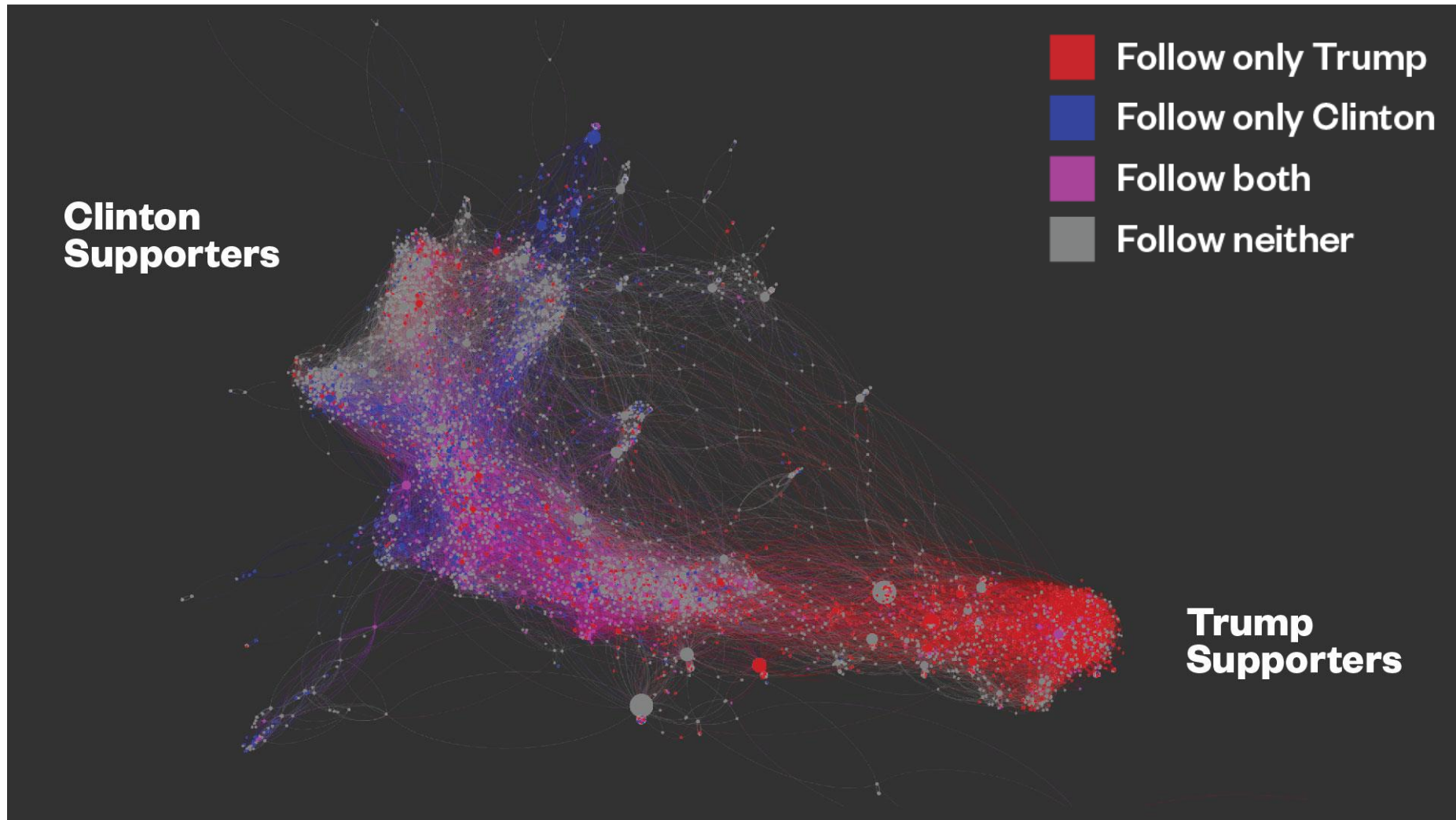To: ang@cs.stanford.edu
Subject: Buy now!

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

Email → Machine Learning Model → Spam / Not Spam

# Example: Community Detection



Follow only Trump
Follow only Clinton
Follow both
Follow neither

Clinton Supporters

Trump Supporters

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}

# Other Examples of Applied ML

**Image Classification**



**Face Detection**



**Stock Prediction**



**Fraud Detection**



**Risk Analysis**



**Navigation**

# Fields & Terms Related to Machine Learning

# Statistics vs. Machine Learning
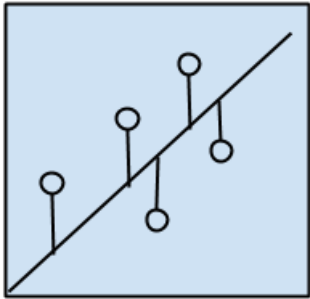
- Largely overlapping fields:
  - Both concerned with learning from data
  - Philosophical difference on 'focus' and 'approach'.

- 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.
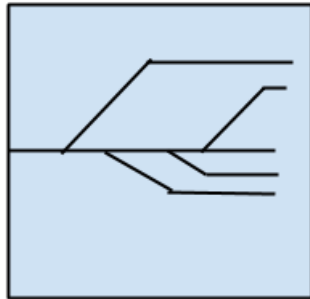
# Types of Machine Learning



Unlabeled Data

Labeled Data

Multi-Step Problems

*Adapted from : https://www.pinterest.com/pin/786792997374742269/*
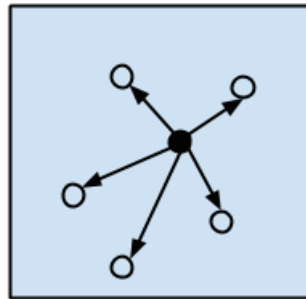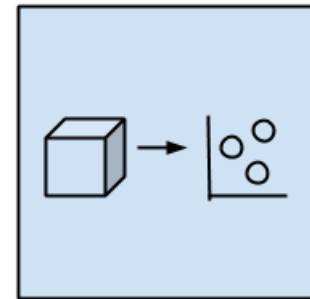
# Machine Learning Algorithm Families
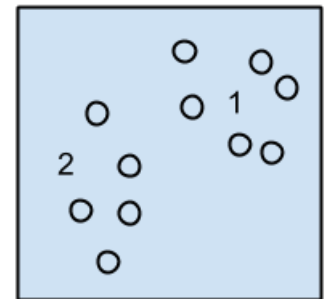


Regression Algorithms
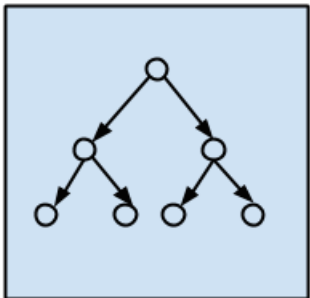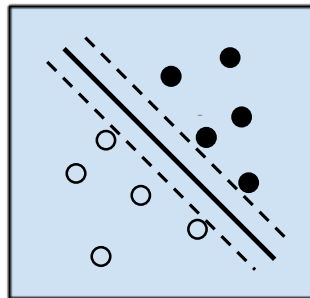
Regularization Algorithms

Instance-based Algorithms

Dimensional Reduction Algorithms
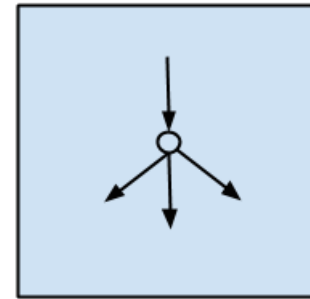
Clustering Algorithms
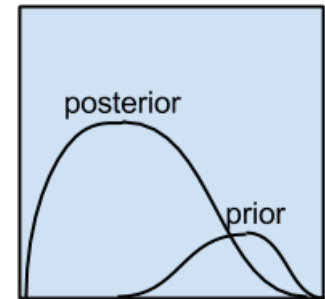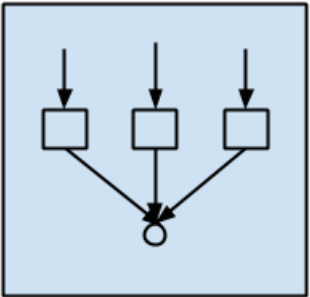
Decision Tree Algorithms

Support Vector Machines

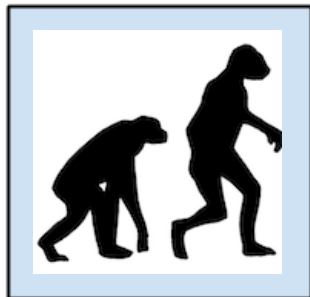Association Rule Learning Algorithms

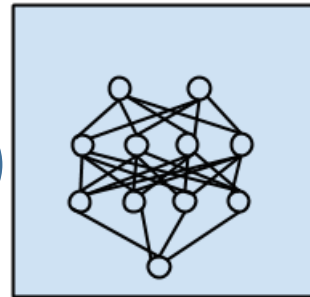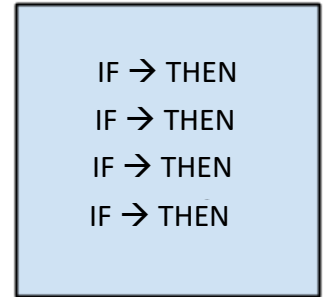Artificial Neural Network Algorithms

Bayesian Algorithms

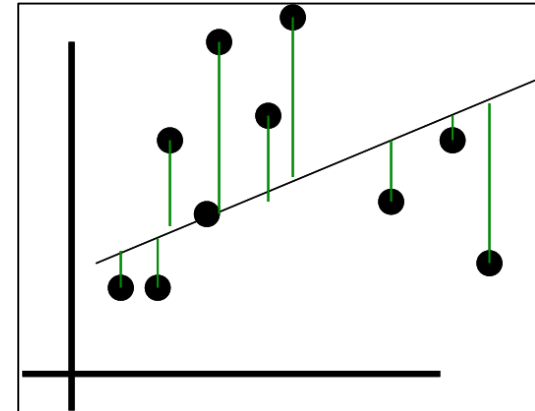Ensemble Algorithms

Evolutionary Algorithms

Non-exhaustive list of ML families

Deep Learning Algorithms

Learning Classifier Systems
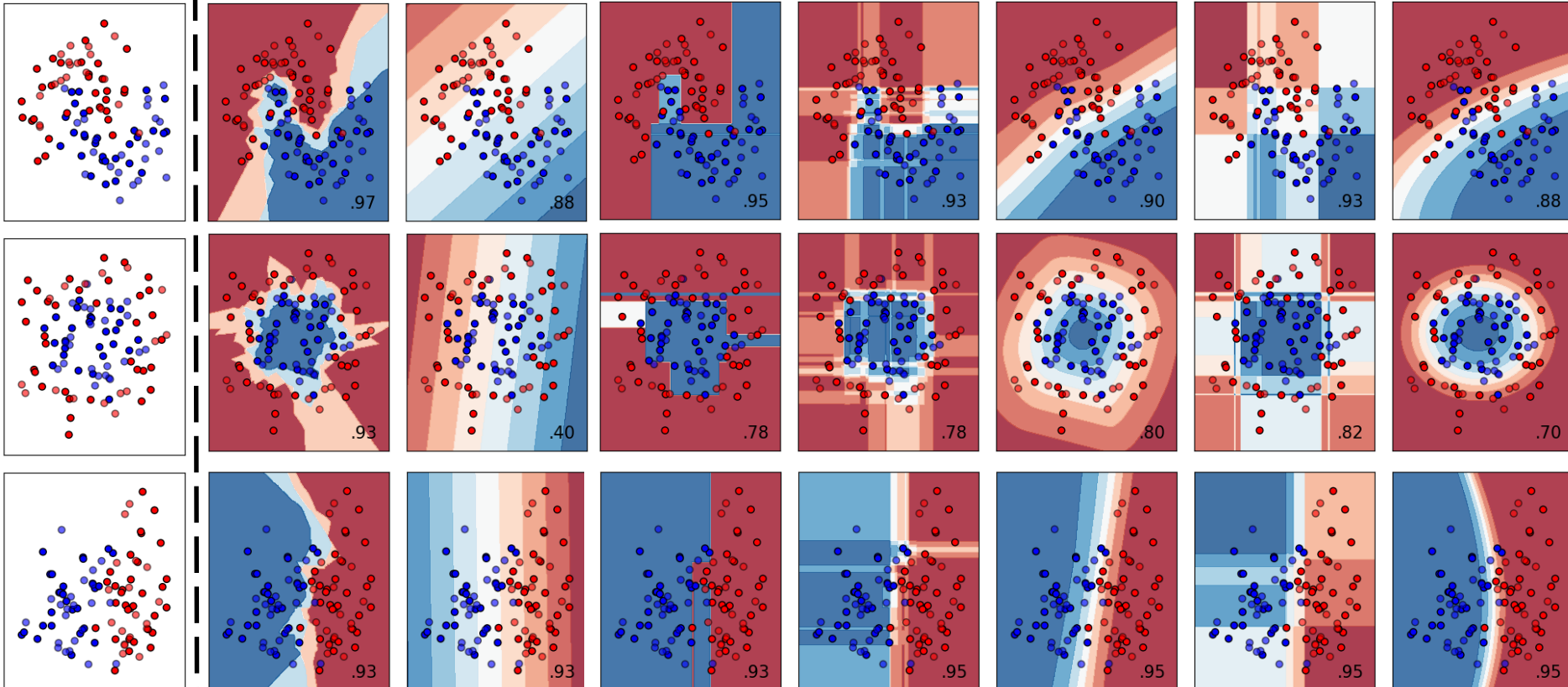
# Supervised Learning: Prediction

- **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

# Models/ML: Representation

# 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]

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

# 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
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
```

| "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 |
| Y | N | N | N | Y | N |
| N | N | Y | N | Y | N |
| Y | N | N | Y | N | Y |
| N | N | Y | N | N | N |

example → (row 4)    label → (column "spam?")

# 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 three.

# 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



**Is it raining?**
yes / no

**Is it windy?** / don't bring anything
yes / no

**Is it extremely windy?** / use an umbrella
yes / no

stay home / wear a rain jacket

© 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

Hip replacement patients

LMW heparin

DVT
- Bleed
- No bleed

No DVT
- Bleed
- No bleed

Conventional treatment

DVT
- Bleed
- No bleed

No DVT
- Bleed
- No bleed

The research question:

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

# Decision Tree for Heparin

- Entering probabilities

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

Hip replacement patients

LMW heparin

DVT
0.14

Bleed
0.1

No bleed
0.9

No DVT
0.86

Bleed
0.1

No bleed
0.9

Conventional treatment

DVT
0.25

Bleed
0.01

No bleed
0.99

No DVT
0.75

Bleed
0.01

No bleed
0.99

Perelman
School of Medicine
University of Pennsylvania

# 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 |

Bad = 3
Good = 7

**MissPay**

No — Yes

Bad = 1
Good = 6

Bad = 2
Good = 1

Bad = 3
Good = 7

**NewJob**

No — Yes

Bad = 2
Good = 4

Bad = 1
Good = 3

# 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 |

Bad = 3
Good = 7

**MissPay**

No          Yes

Bad = 1
Good = 6

Bad = 2
Good = 1

**NewJob**

No          Yes

Bad = 0
Good = 3

Bad = 1
Good = 3

# Decision Tree: Fitting with Splits

Max Depth: 1

# Decision Tree: Fitting with Splits

# Decision Tree: Fitting with Splits

# Decision Tree: Fitting with Splits

# Decision Tree: Fitting with Splits

# 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?

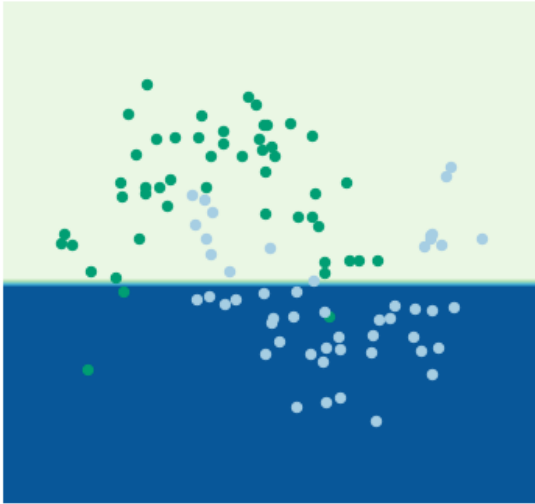- If there is a strong, simple relationship among variables, most methods will find it.

- 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

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

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

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

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

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

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

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

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

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

### Confusion Matrix

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

Receiver operating characteristic example



ROC curve (area = 0.79)

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

Dirty    Noisy

**BAD DATA = BAD EVERYTHING**

Duplicate

Biased    Sparse

# 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
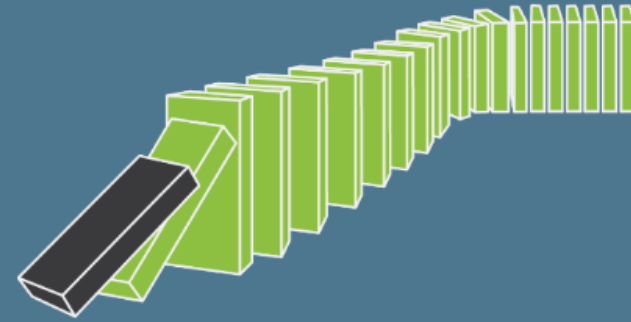- Failing to consider confounding variables

# 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

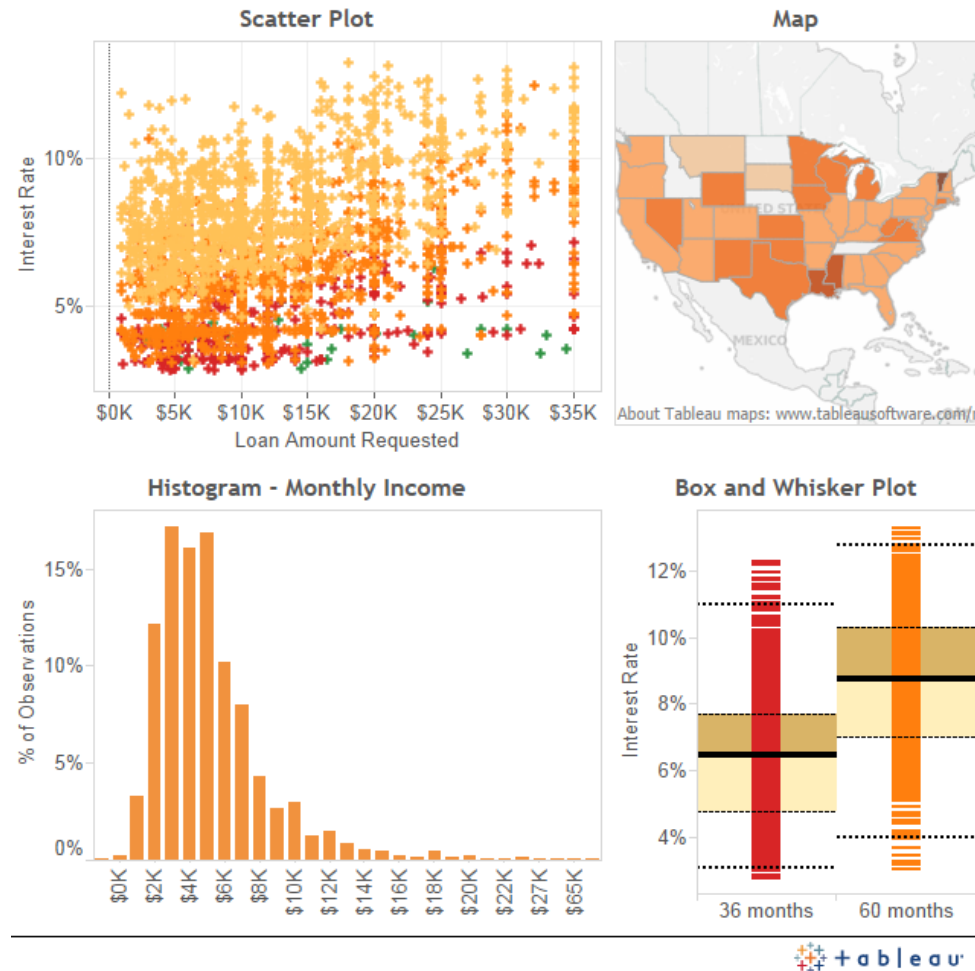# Common Machine Learning Pitfalls

- 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
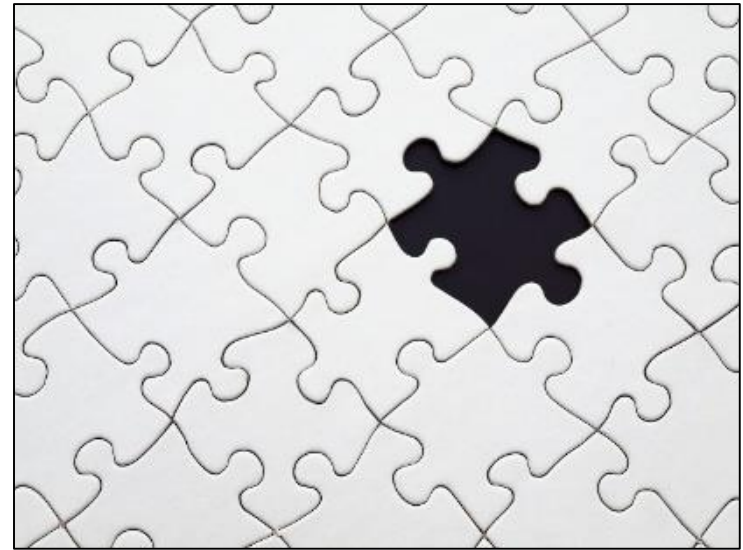- Failing to consider confounding variables



Everything is possible; just not too probable.

# 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
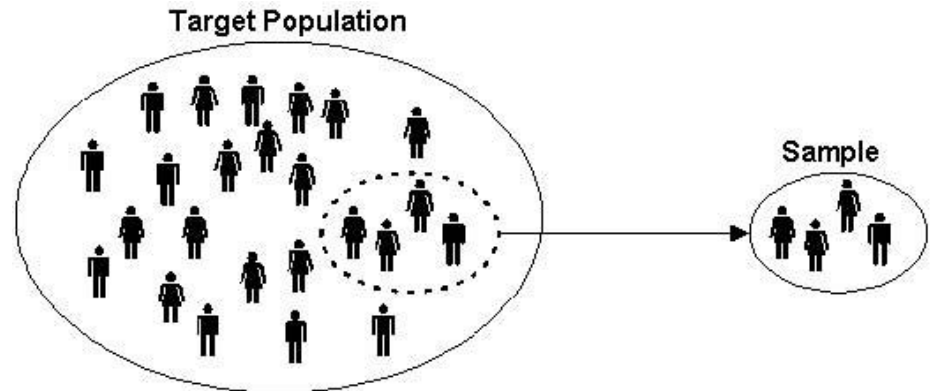- Failing to consider confounding variables

# 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
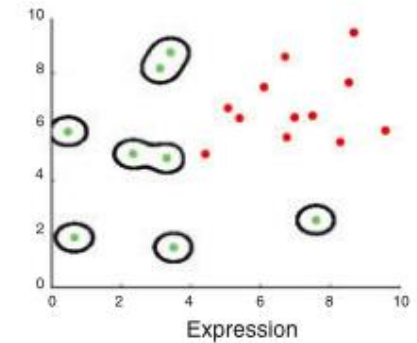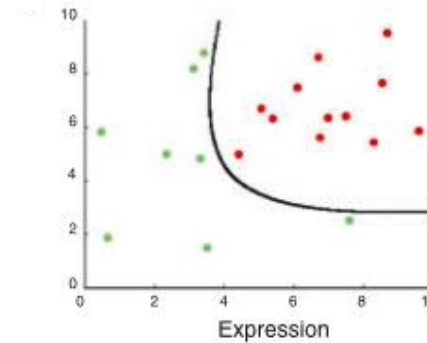- Failing to consider confounding variables

# 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
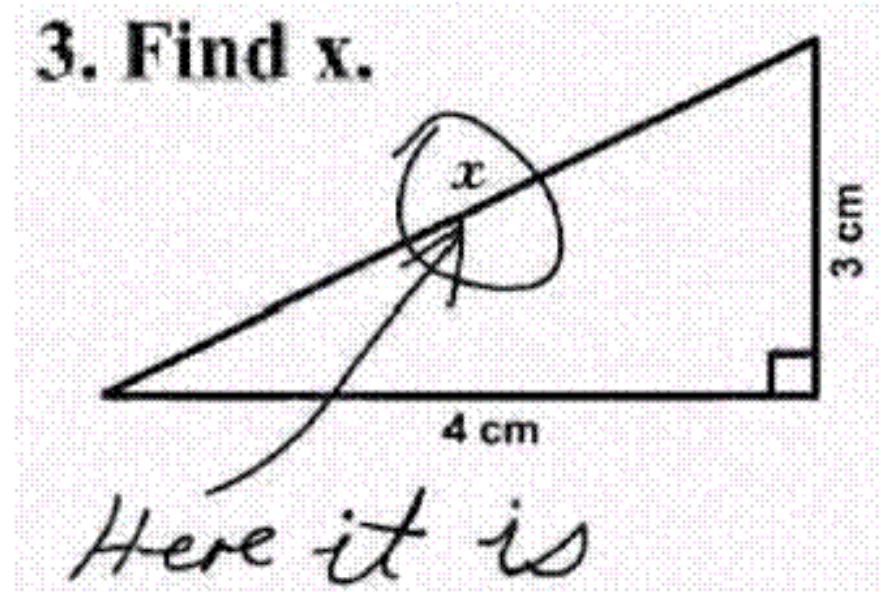- Failing to consider confounding variables

# 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

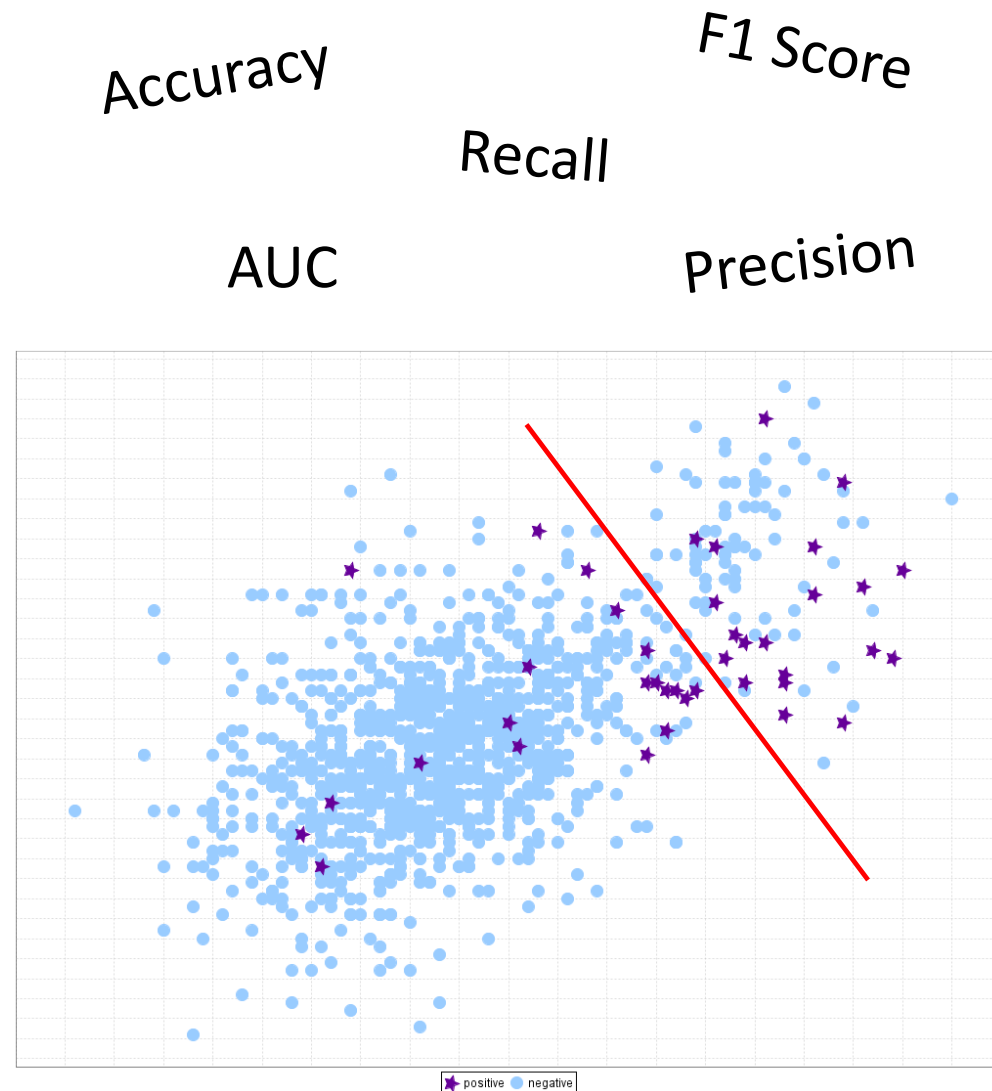Accuracy

F1 Score

Recall

AUC

Precision

# 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

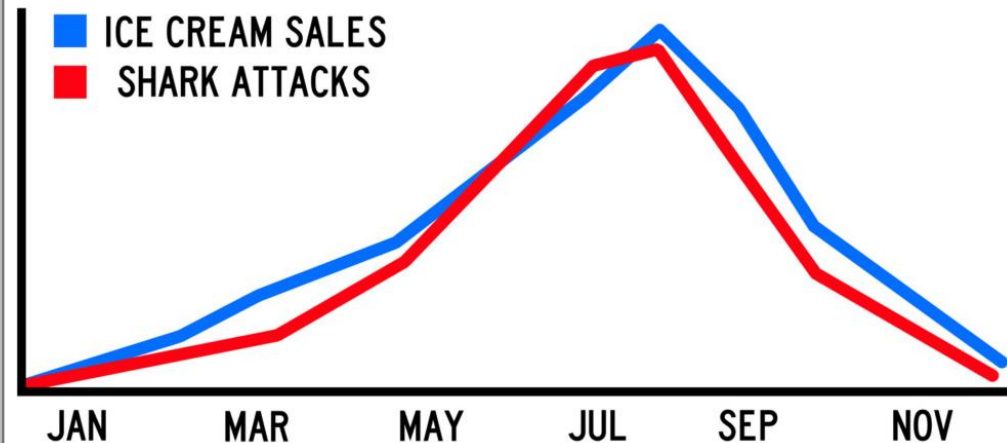If you torture the data long enough,
it will confess.

— *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



ICE CREAM SALES
SHARK ATTACKS
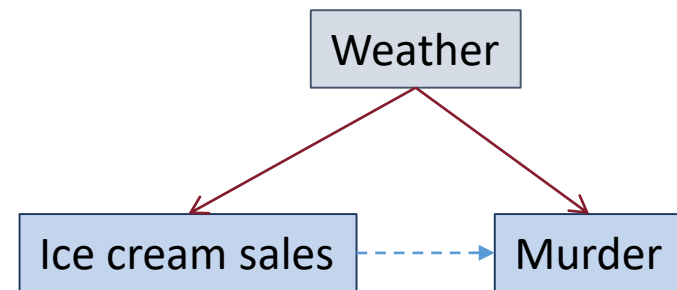
JAN   MAR   MAY   JUL   SEP   NOV

# 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
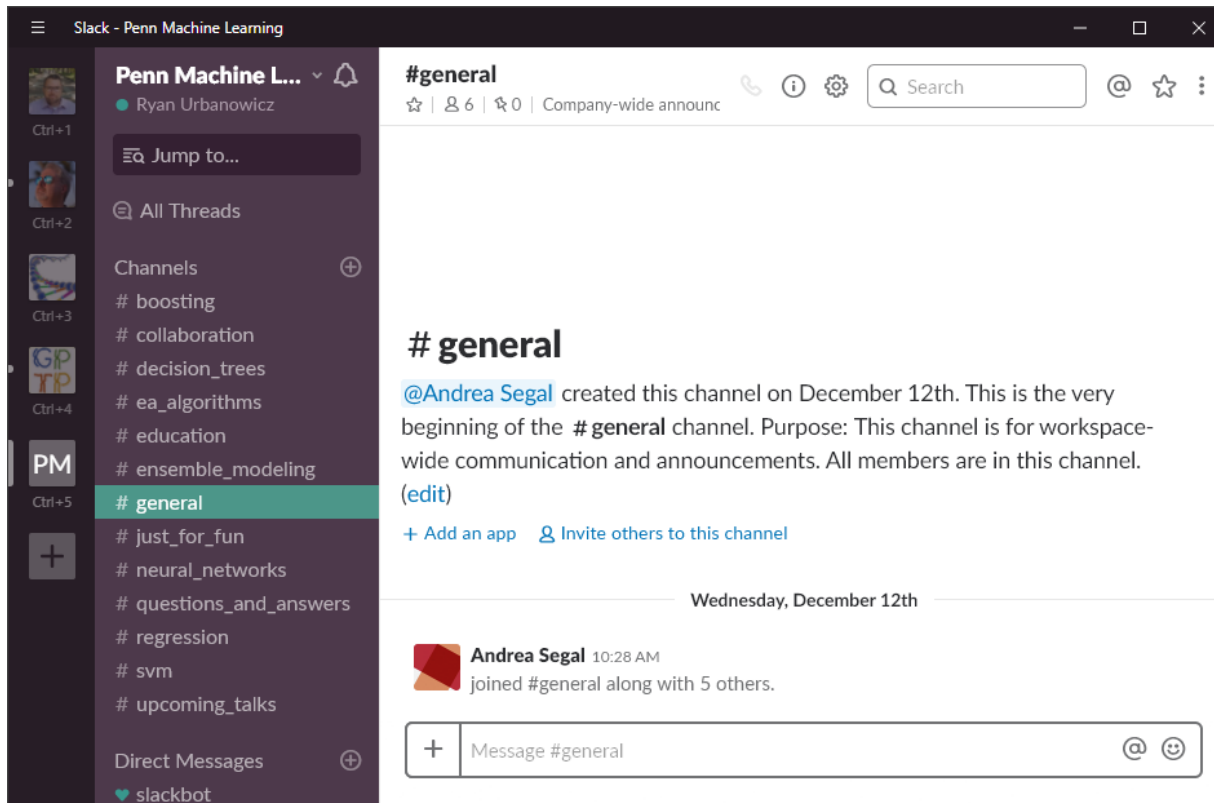- **Failing to consider confounding variables**

# 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

# UPenn ML  slack

- Penn Machine Learning – Slack Workspace

- pennmachinelearning.slack.com

# Acknowledgements and Funding

- Pennsylvania Commonwealth Universal Research Enhancement Program (CURE)