

Plan:

1. Define predictive analysis
2. Explain the relationship between ML and AI
3. Walk through the four general steps to predictive analysis

# Machine Learning: Basics

Shannon E. Ellis, Ph.D  
UC San Diego



Department of Cognitive Science  
[sellis@ucsd.edu](mailto:sellis@ucsd.edu)

Did they summarize the data? **Yes** → Did they report the summaries without interpretation? **No** → Did they quantify whether the discoveries are likely to hold in a new sample? **Yes** → Are they trying to figure out how changing the average of one measurement affects another?

**Predictive:** apply machine learning techniques to data you have currently to generate a model that will be able to make a prediction on future data

Classic Statistics  
(parametric & nonparametric)

Text Analysis

Are they trying to predict measurement(s) for individuals?

**Causal**

**STOP!**  
Not a data analysis

**No** Are the data a corpus of text?

**No** Are the observations spatially related?

**Inferential**

**Predictive**

Supervised Machine Learning

Geospatial Statistics

Unsupervised Machine Learning

Did the computer decide the groups/labels from your data?

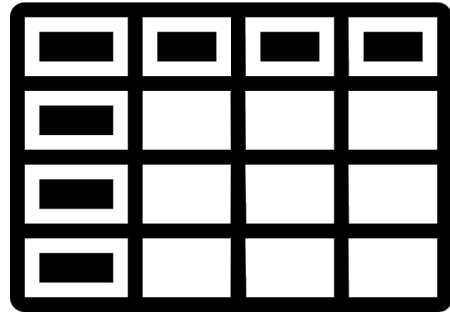
**No**

- **Problem:** Detecting whether credit card charges are fraudulent.
- **Data science question:** Can we use the time of the charge, the location of the charge, and the price of the charge to predict whether that charge is fraudulent or not?
- **Type of analysis:** Predictive analysis



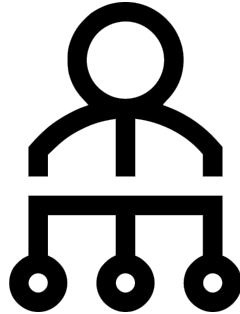
**predictive analysis** uses data  
you have now to make  
predictions in the future

**machine learning**  
approaches are used for  
predictive analysis!



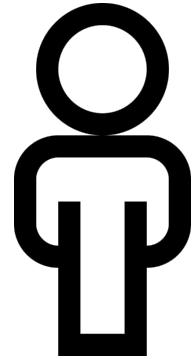
data

train →



model

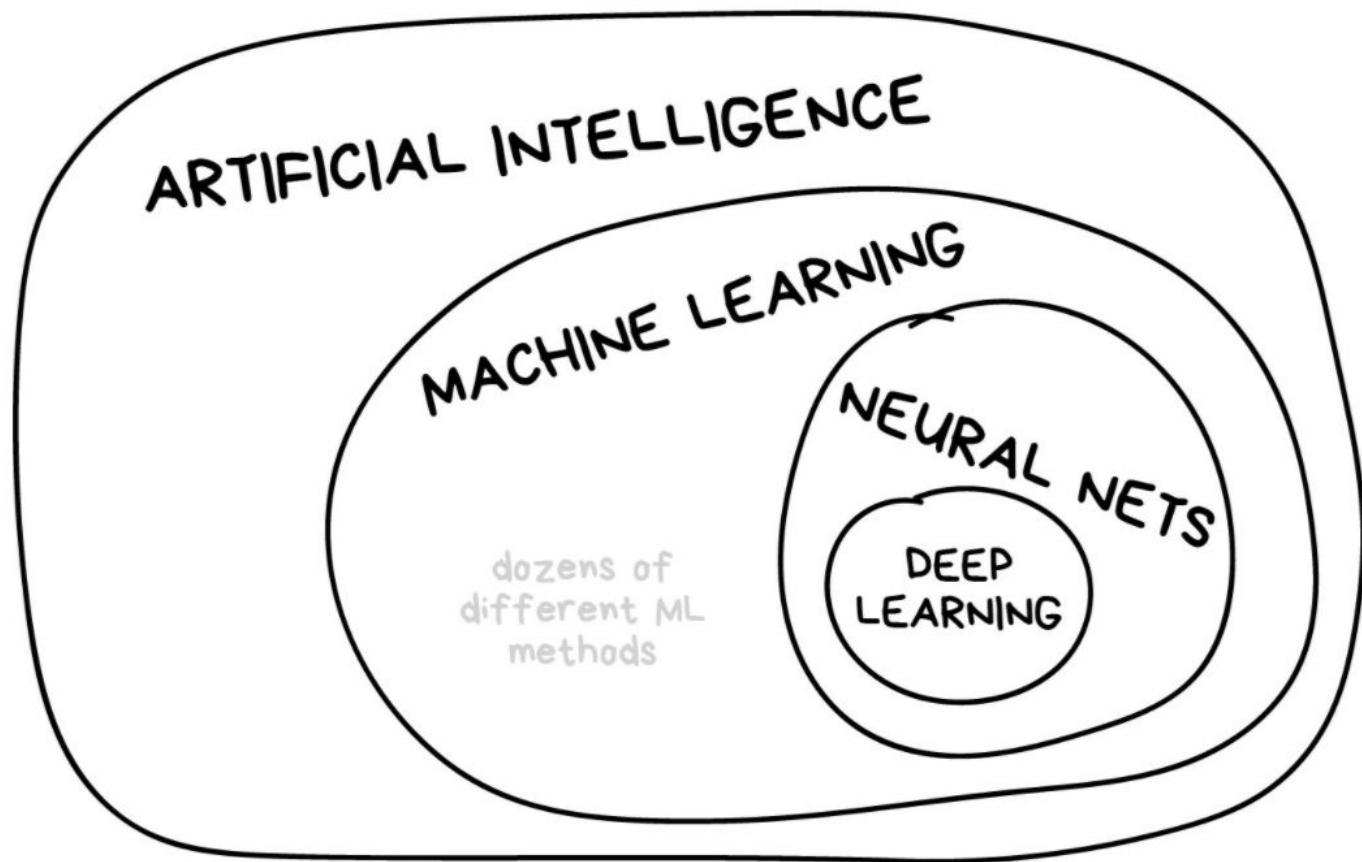
→ predict



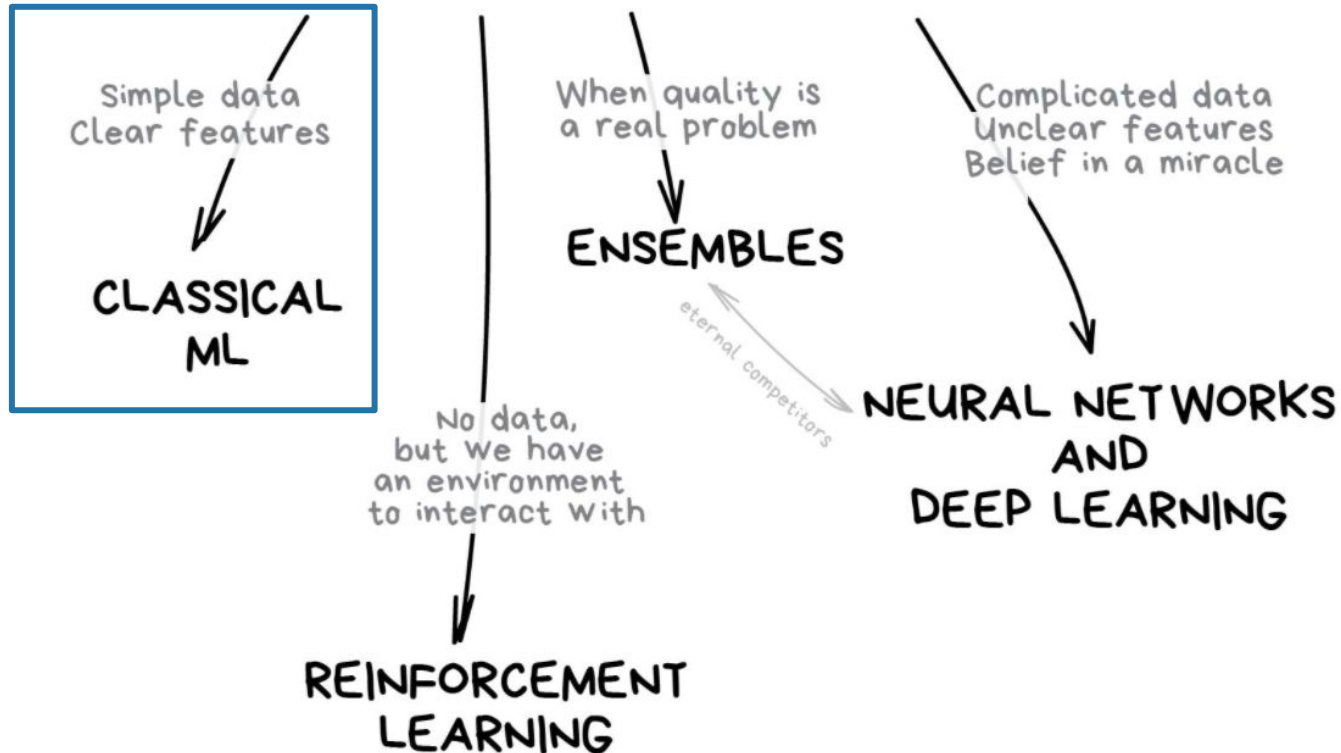
# What is machine learning?

“Machine learning is the science of getting computers to act without being explicitly programmed”

- Andrew Ng, Stanford, ex-Google, chief scientist at Baidu, Coursera founder, Stanford Adjunct Faculty



# THE MAIN TYPES OF MACHINE LEARNING

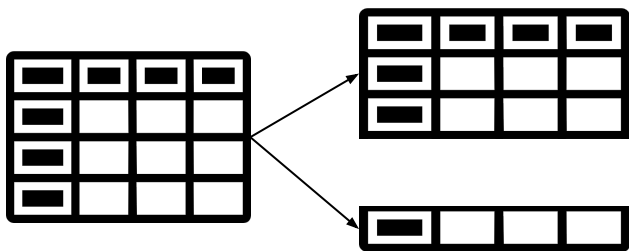


# Machine Learning Generalizations

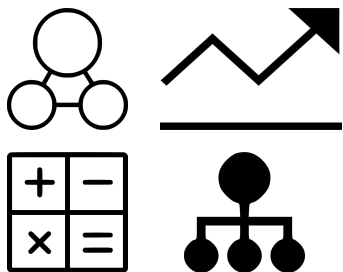
---



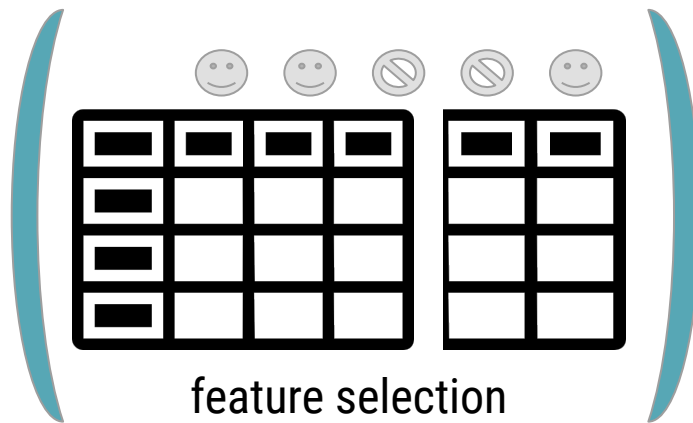
# Basic Steps to Prediction



data  
partitioning



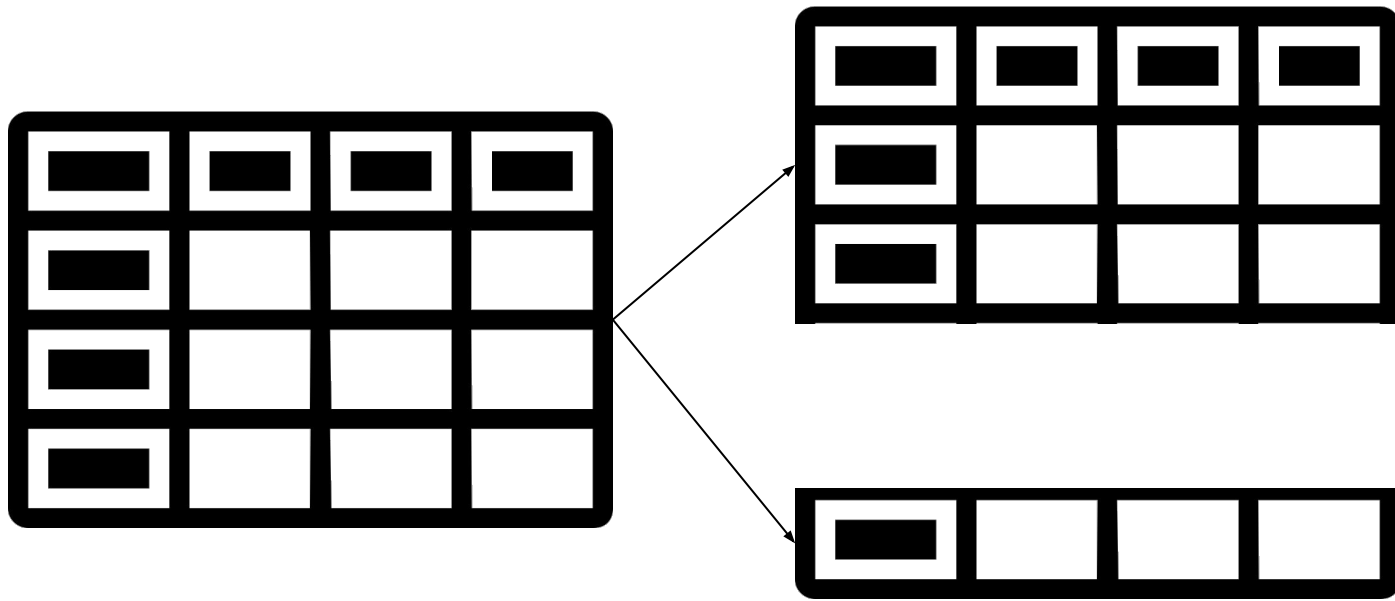
model selection



feature selection



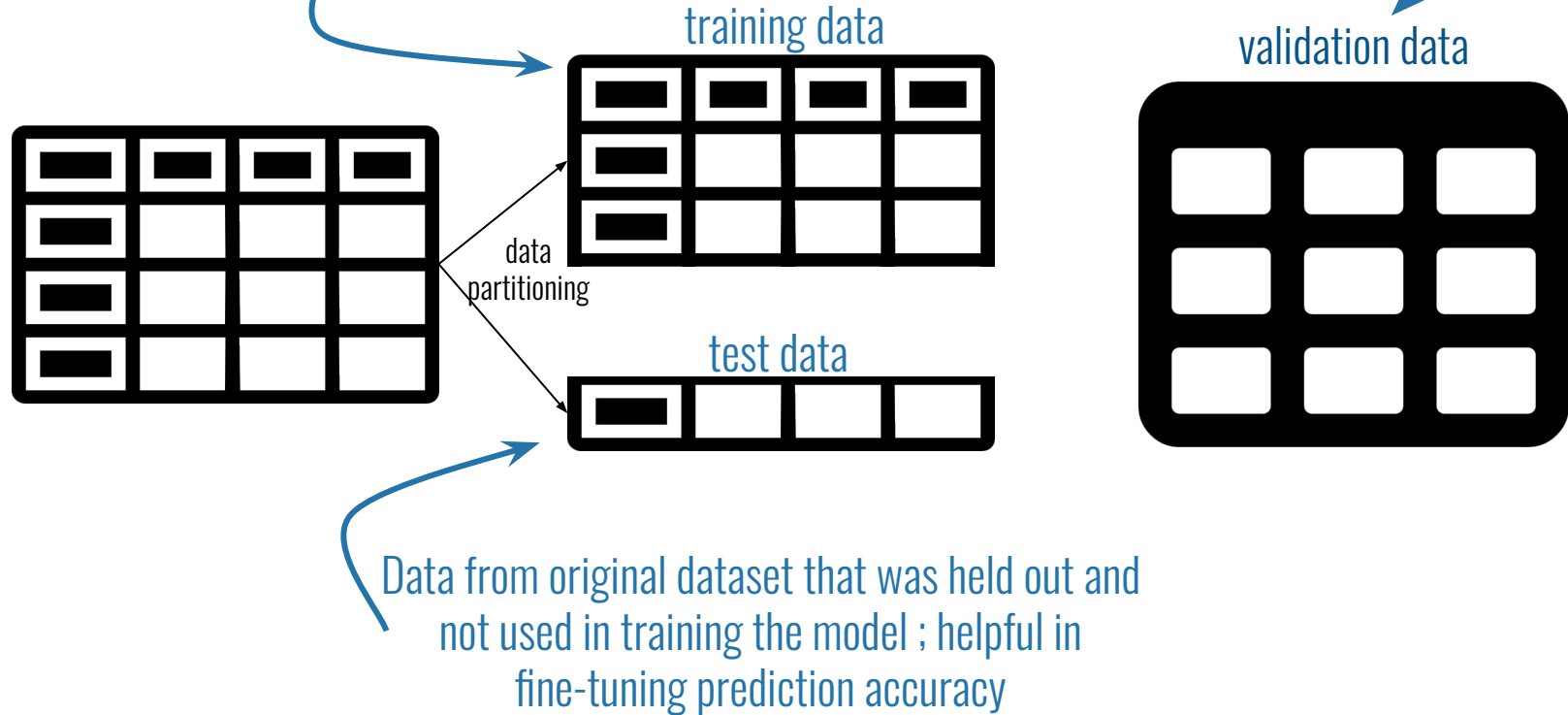
model assessment

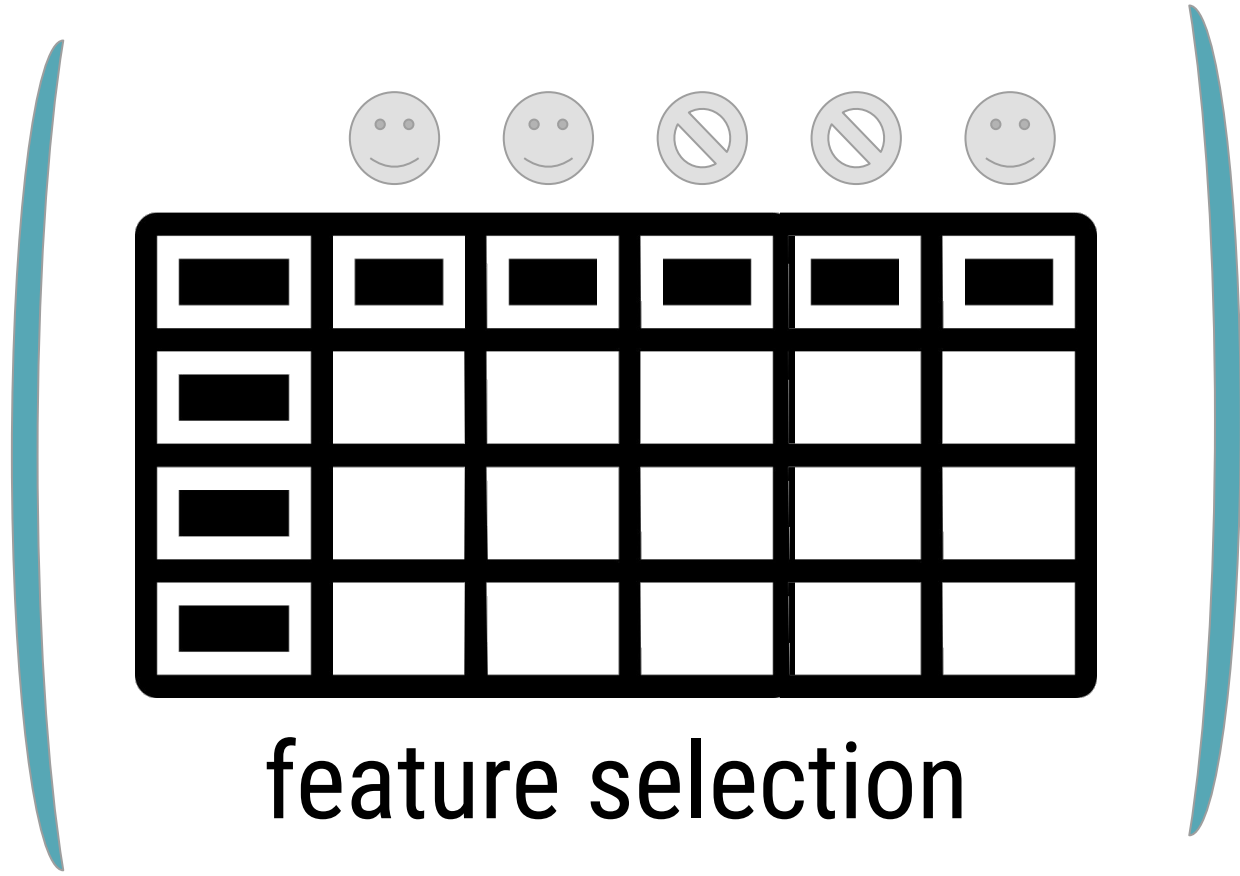


data partitioning

the data used to build  
your predictive model

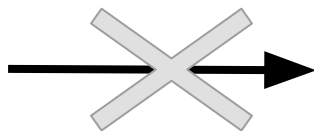
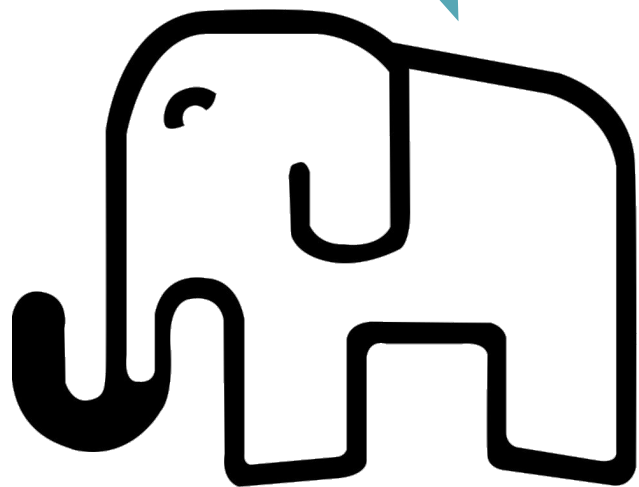
new and independent data set  
used to assess if prediction model  
is generalizable

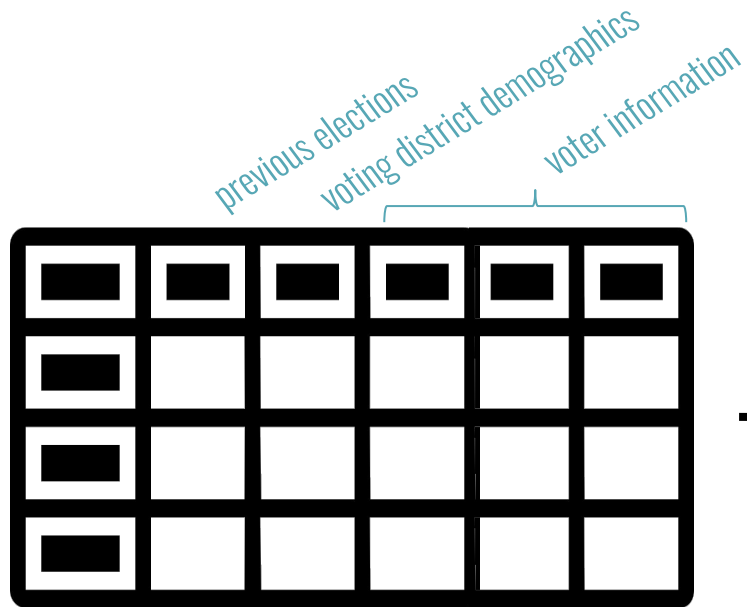




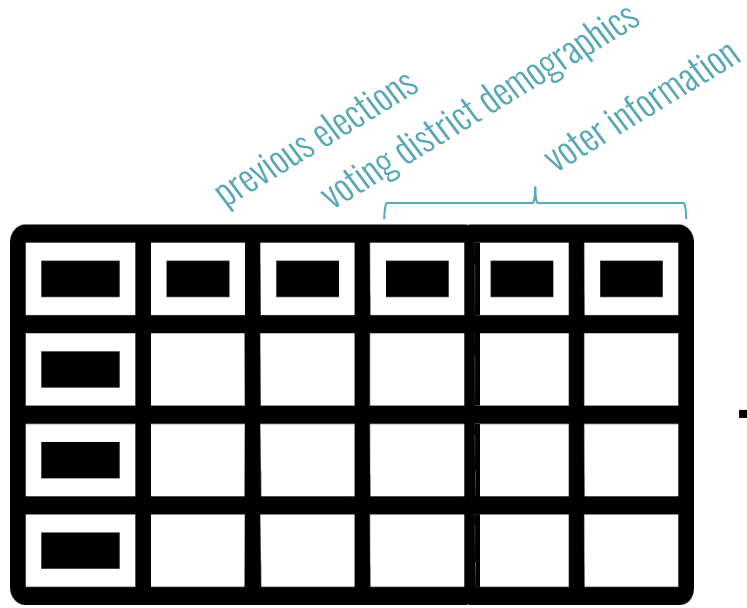
feature selection

elephant height data are likely  
not predictive of US elections





these data are likely  
predictive of US election  
outcomes



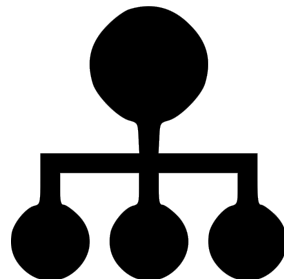
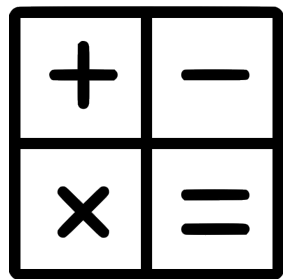
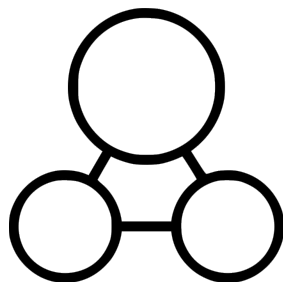
**feature selection** determines which variables are most predictive and includes them in the model

■	■	■	■	■	■
■					
■					
■					

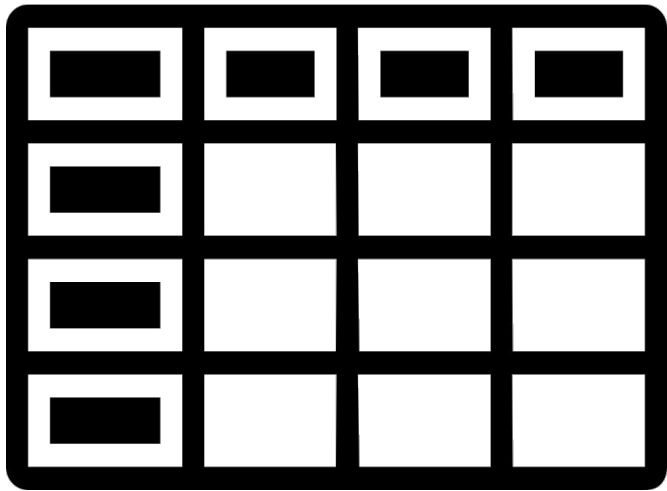


variables that can be used for accurate prediction exploit the relationship between the variables but do NOT mean that one causes the other

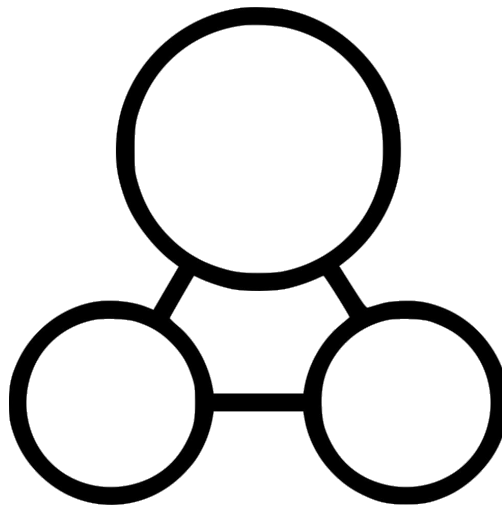




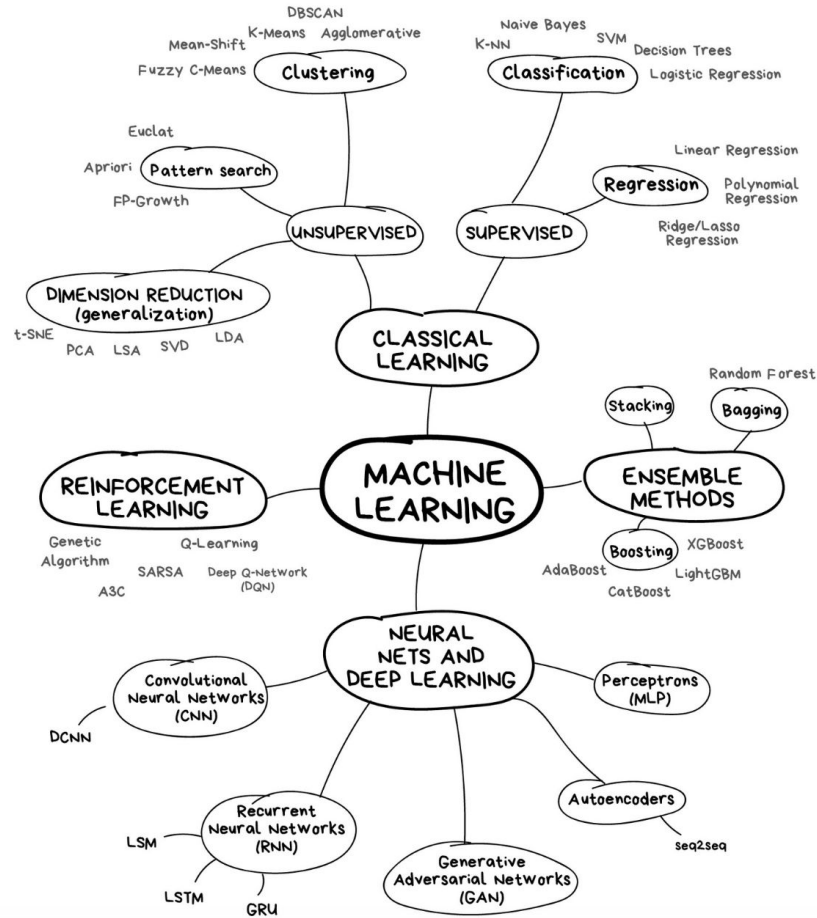
model selection



big  
datasets

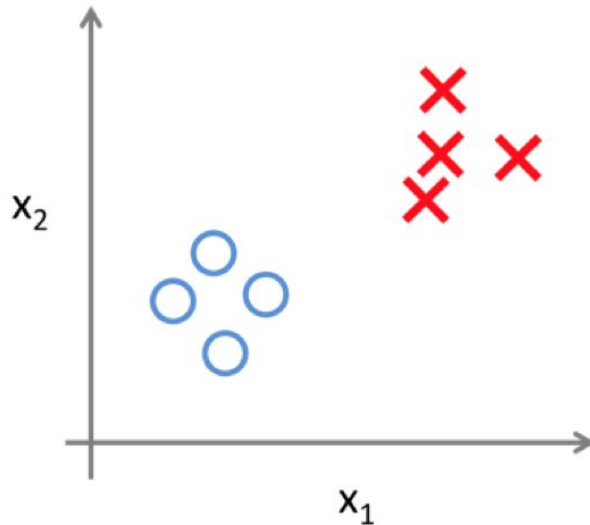


simple  
models



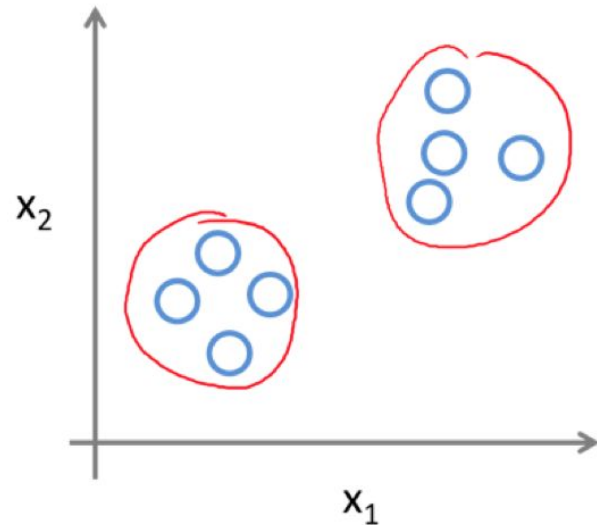
# To modes of machine learning

## Supervised Learning



You tell the computer how to classify the observations

## Unsupervised Learning

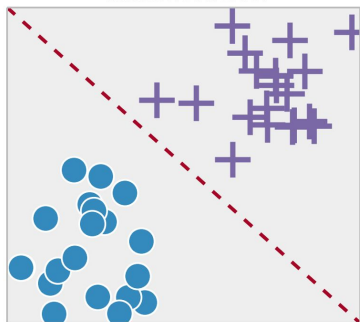


The computer determines how to classify based on properties within the data

# Approaches to machine learning

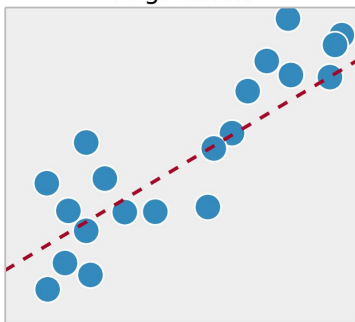
## Supervised Learning

Classification



categorical variables

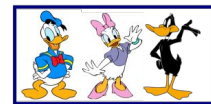
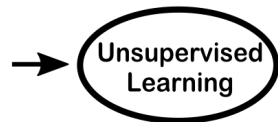
Regression



continuous variables

Prediction accuracy  
dependent on  
training data

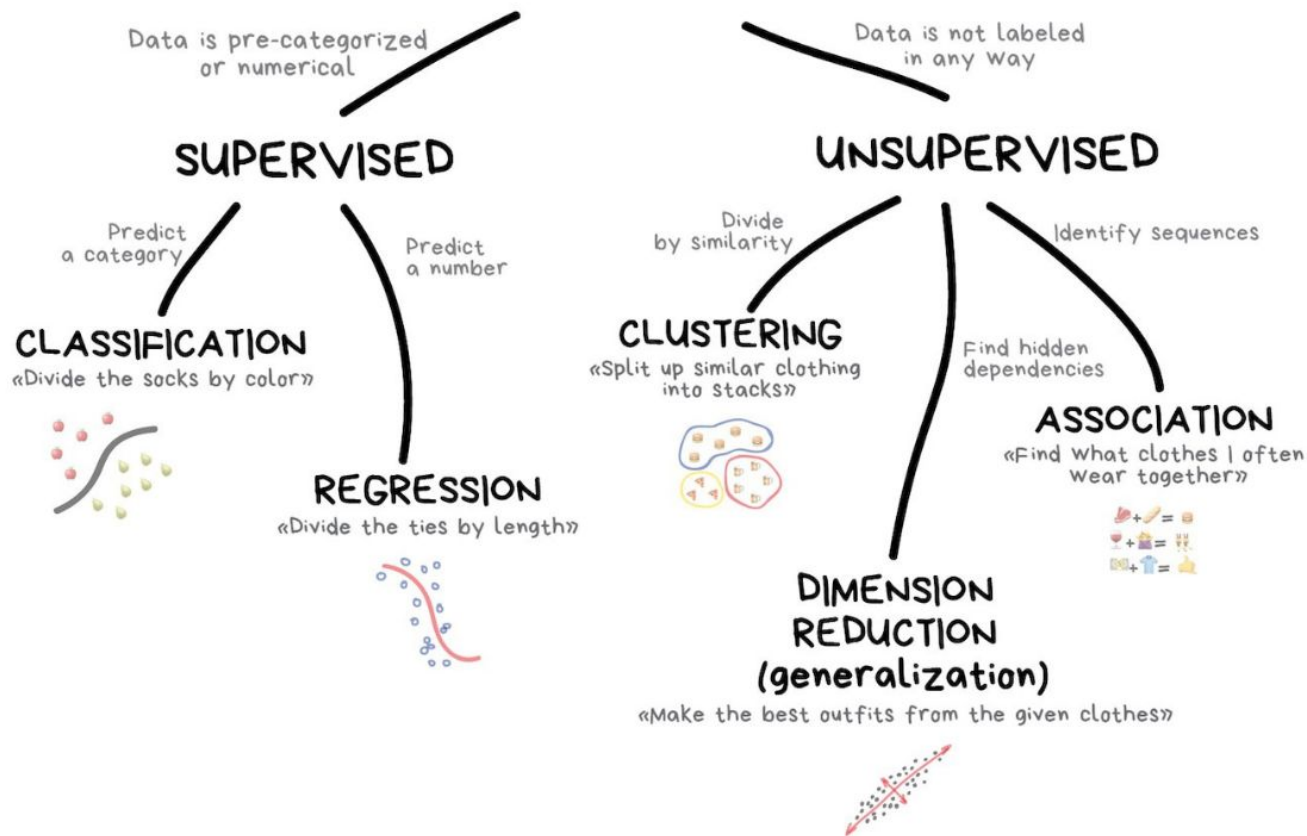
## Unsupervised Learning



Clustering (categorical)  
& dimensionality reduction (continuous)

can automatically  
identify structure in  
data

# CLASSICAL MACHINE LEARNING



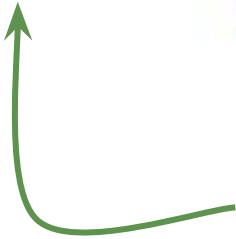


model assessment

---

# Root Mean Squared Error (RMSE)

$$RMSE = \sqrt{\frac{\sum_{i=1}^N (Predicted_i - Actual_i)^2}{N}}$$



A few outliers can lead to a big increase in RMSE, even if all the other predictions are pretty good

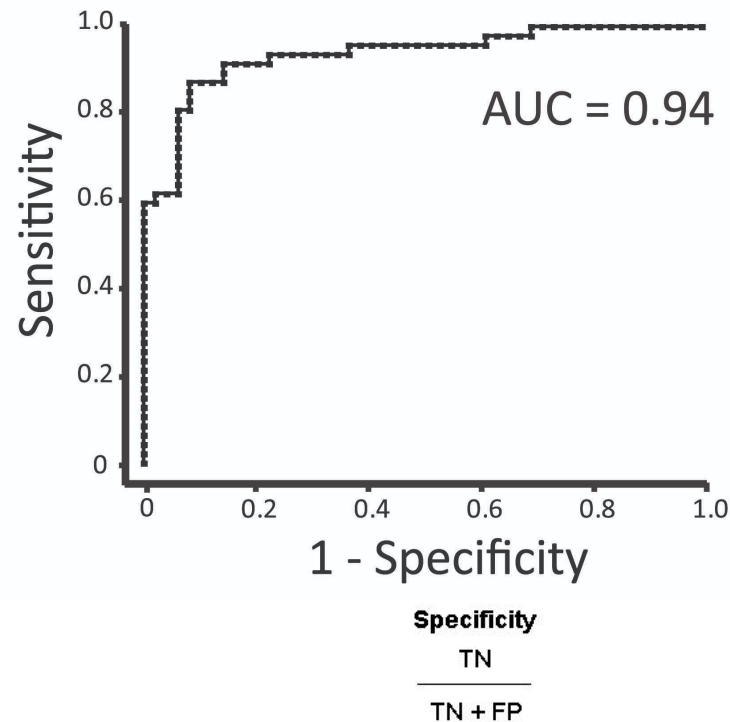


$$\text{Accuracy} = \frac{\# \text{ of samples predicted correctly}}{\# \text{ of samples predicted}} * 100$$

		Actual	
		Positive	Negative
Predicted	Positive	True Positive (TP)	False Positive (FP)
	Negative	False Negative (FN)	True Negative (TN)

A 2x2 table is a type of confusion matrix

$$\text{Sensitivity} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$



<b>Accuracy</b>	What % were predicted correctly?
<b>Sensitivity</b>	Of those that <i>were</i> <b>positives</b> , what % were predicted to be positive?
<b>Specificity</b>	Of those that were <b>negatives</b> , what % were predicted to be negative?