# Machine learning

**Jason G. Fleischer, Ph.D.**
**Asst. Teaching Professor**
**Department of Cognitive Science, UC San Diego**
jfleischer@ucsd.edu

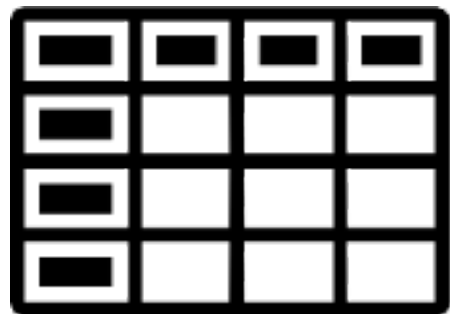@jasongfleischer

https://jgfleischer.com

- **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
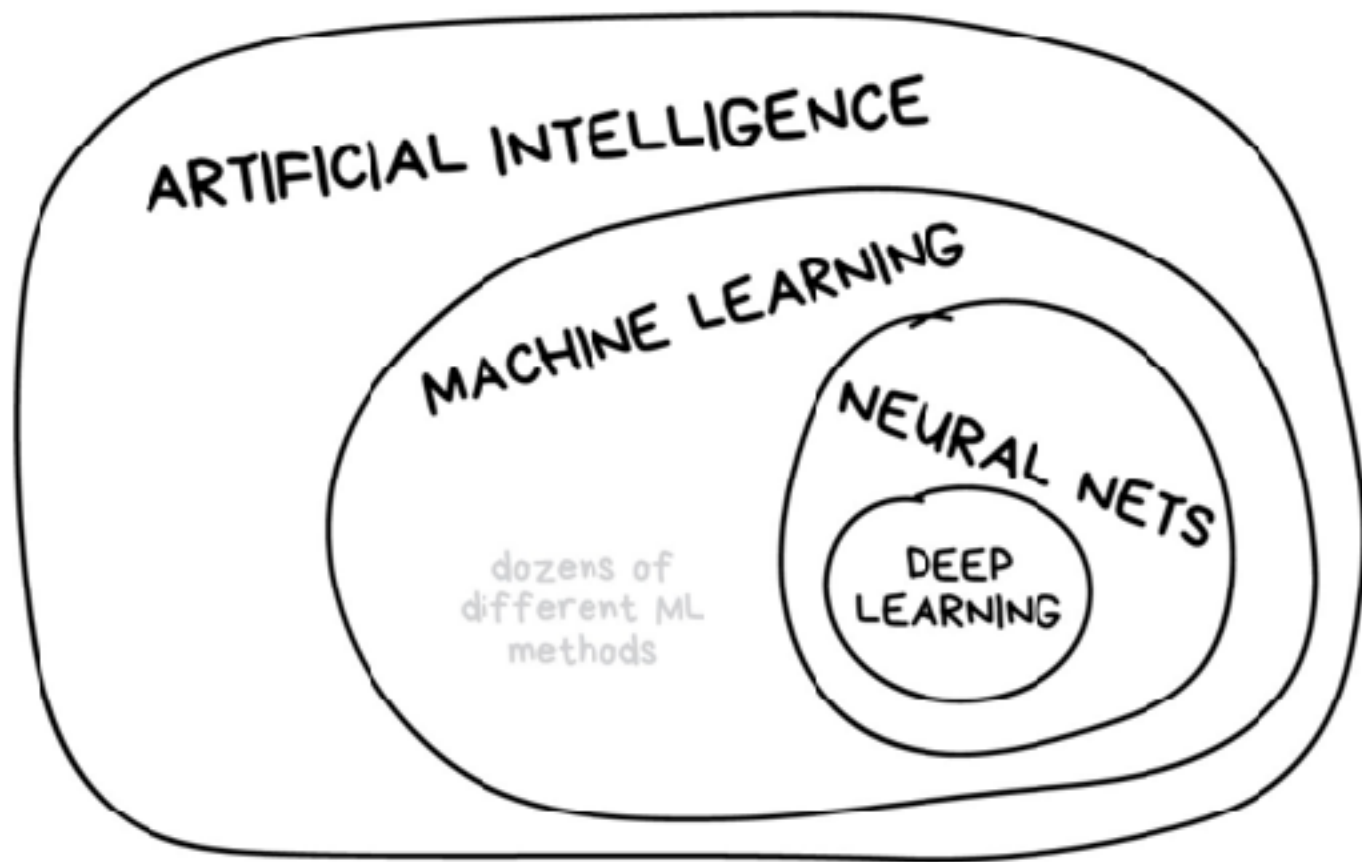predictive analysis!

train

predict

data

model

In contrast to statistical approaches which care more about the model
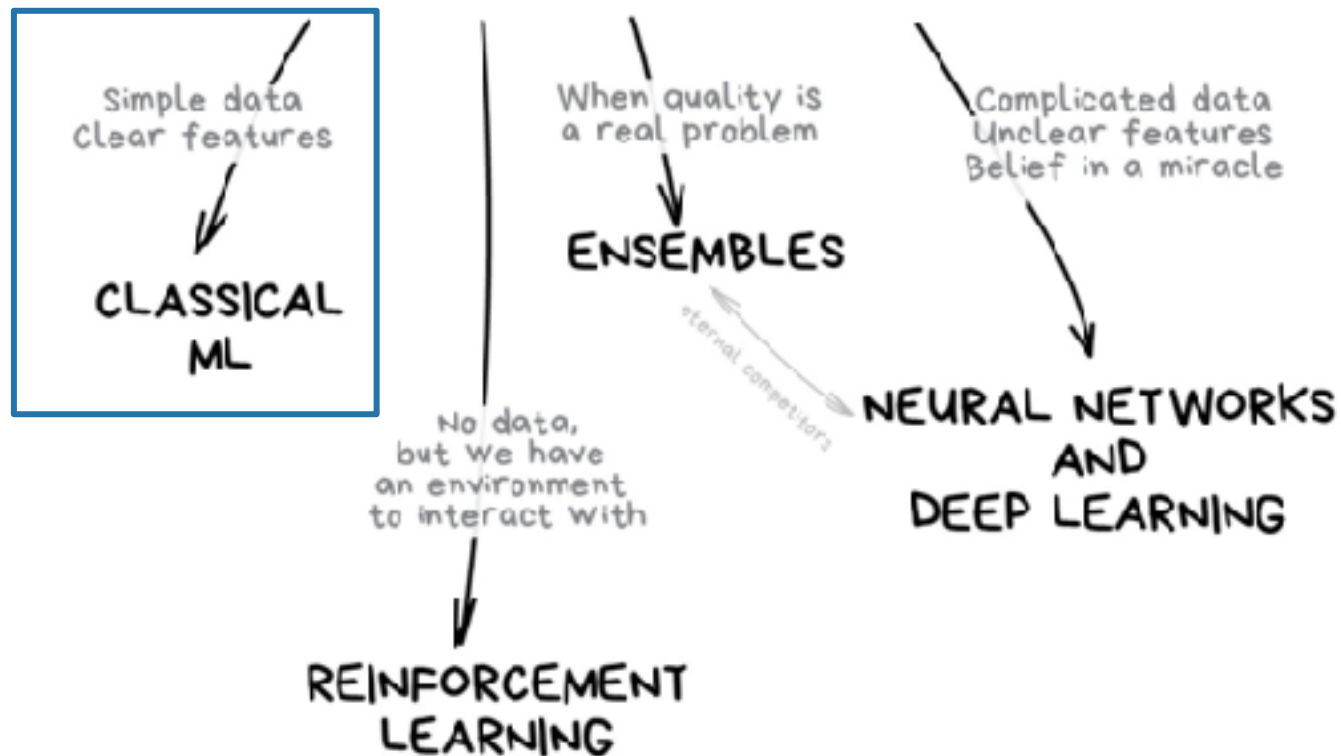accurately reflecting the process than nailing the predictions

# 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

Simple data
Clear features

When quality is
a real problem

Complicated data
Unclear features
Belief in a miracle

CLASSICAL
ML

ENSEMBLES

eternal competitors

NEURAL NETWORKS
AND
DEEP LEARNING

No data,
but we have
an environment
to interact with

REINFORCEMENT
LEARNING

# Prediction Questions

Which of these questions is most appropriate for machine learning?

**A** How common is watching Sesame Street in the US?

**B** What is the effect of watching Sesame Street on children's brains?

**C** What is the relationship between early childhood educational programming and success in elementary school?

**D** Can we use information about one's early childhood to predict their success in elementary school?
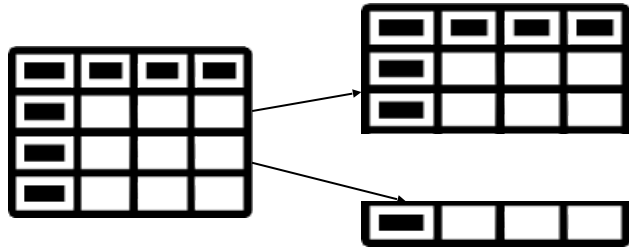
**E** How does Sesame Street cause an increase in educational attainment?

# Machine Learning Generalizations
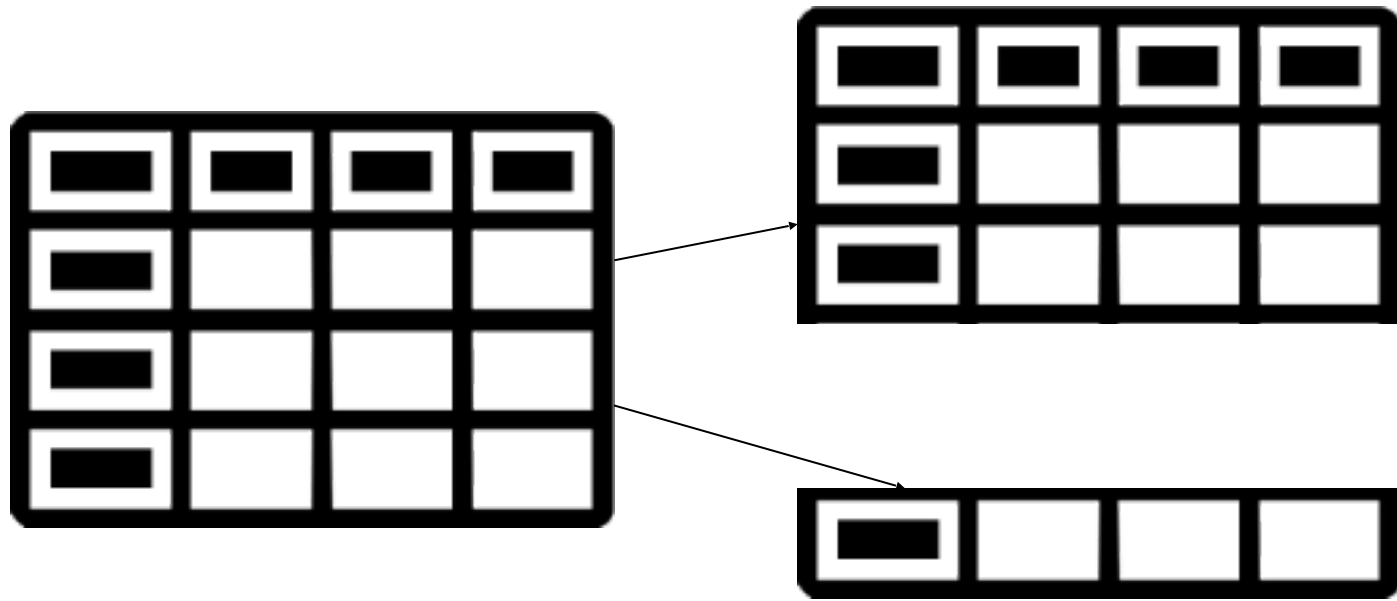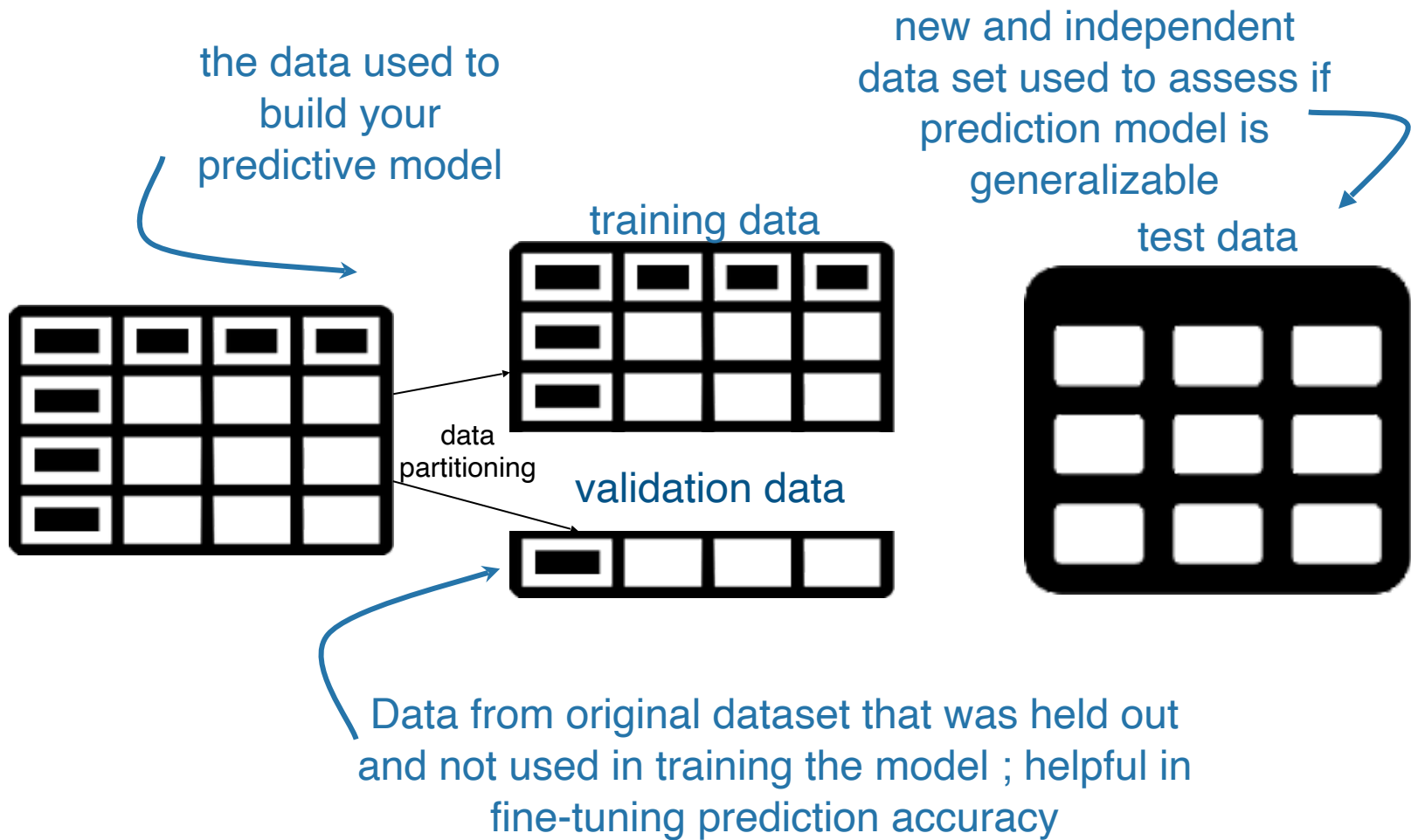
# Basic Steps to Prediction



data partitioning

feature selection

model selection
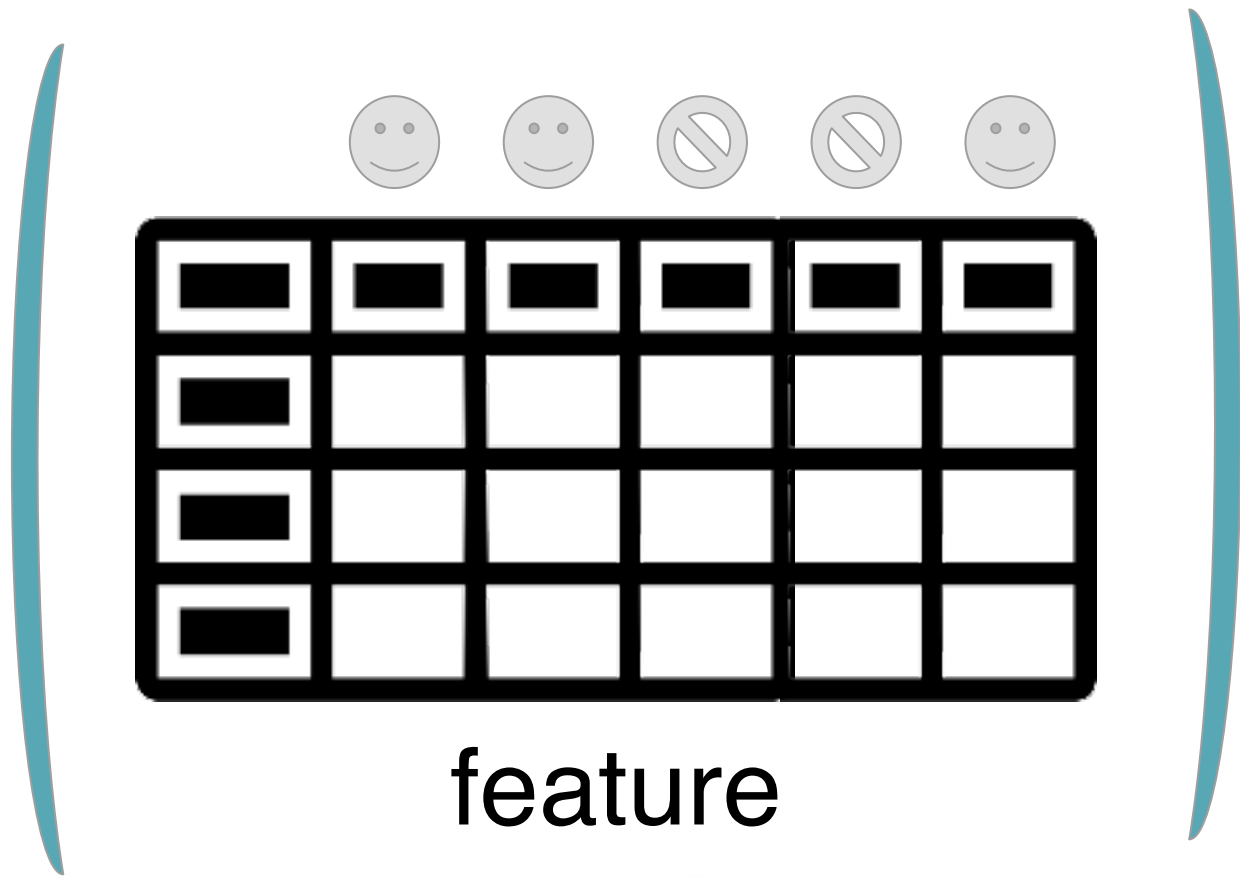
model assessment

data partitioning

the data used to build your predictive model

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

training data

test data

data partitioning

validation data

Data from original dataset that was held out and not used in training the model ; helpful in fine-tuning prediction accuracy

# Data Partitioning

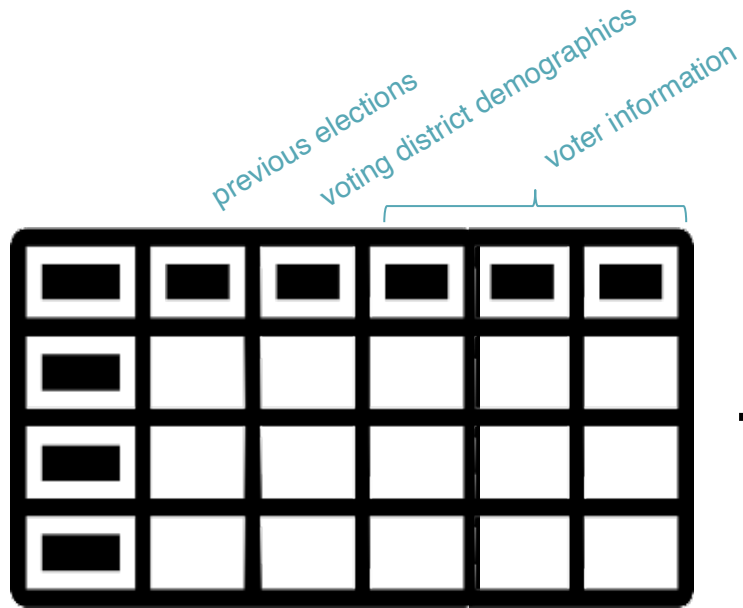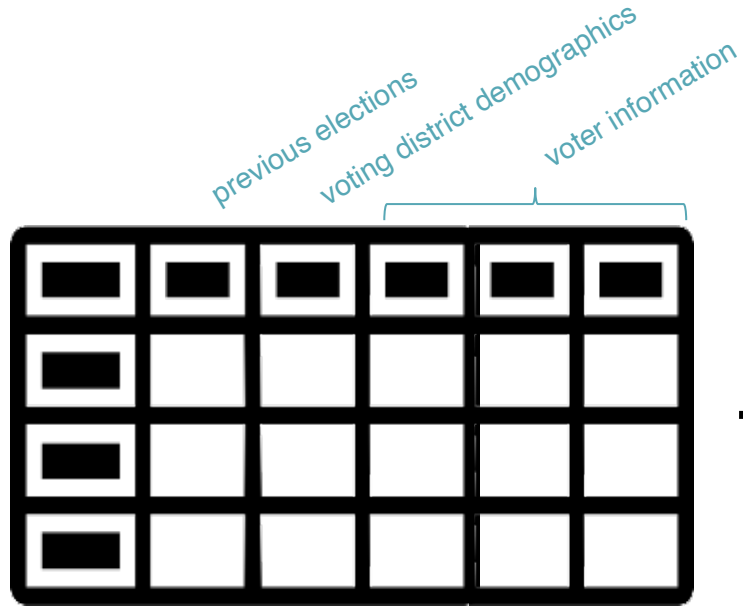## What portion of the data are typically used for generating the model?

A
The entire dataset

B
The training data

C
The testing data

D
The validation data

feature
selection

elephant height data
are likely not predictive
of US elections

previous elections
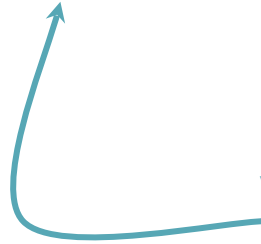
voting district demographics

voter information

these data are likely predictive of US election outcomes

**feature selection** determines which variables are most predictive and includes them in the model
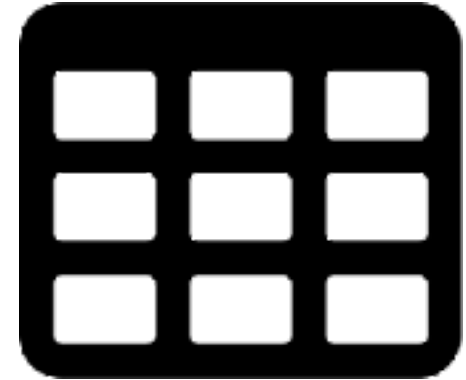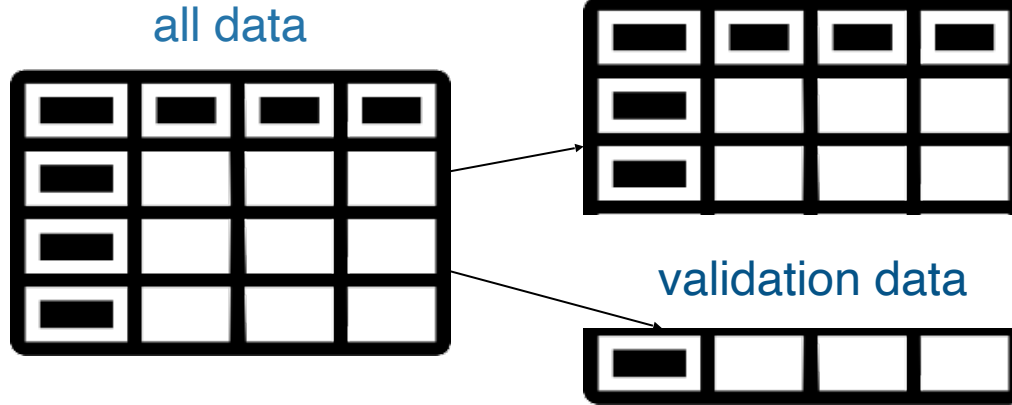
variables that can be used for accurate prediction exploit the relationship between the variables but do NOT mean that one causes the other
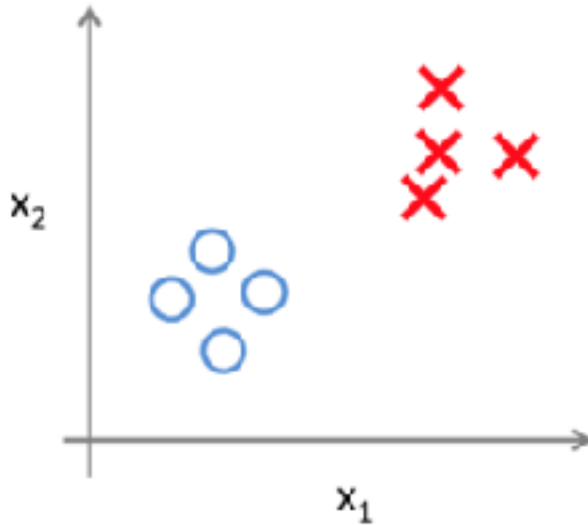
Try different feature sets

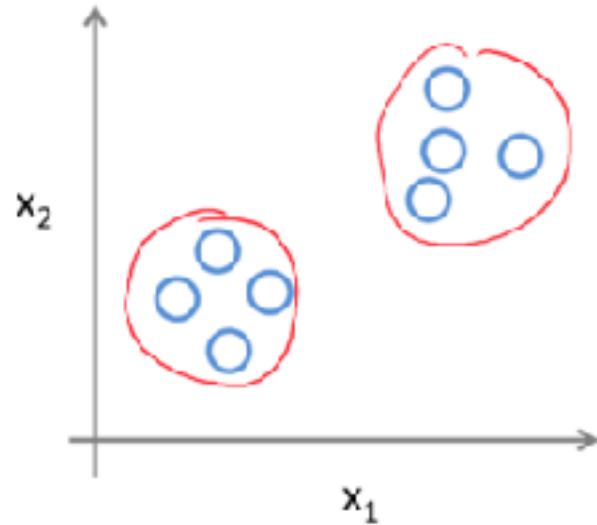Evaluate best feature set on test data

training data

all data

validation data

Use validation set to select the features!

# Two modes of machine learning



Supervised Learning

Unsupervised Learning

You tell the computer what features to use to classify the observations
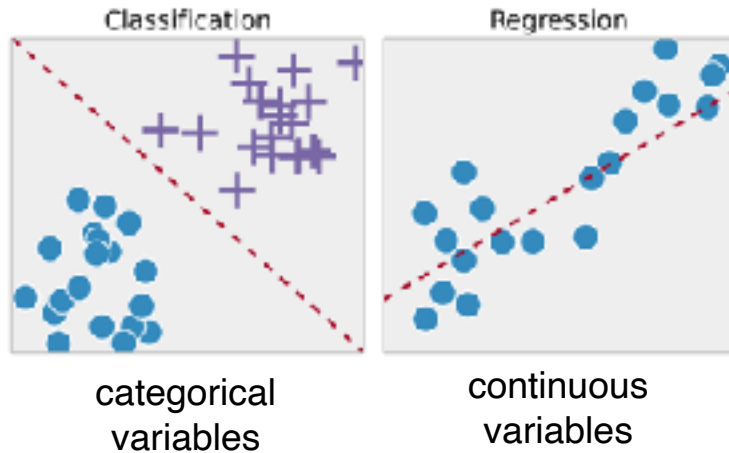
The computer determines how to classify based on properties within the data

# Approaches to machine learning



Supervised Learning

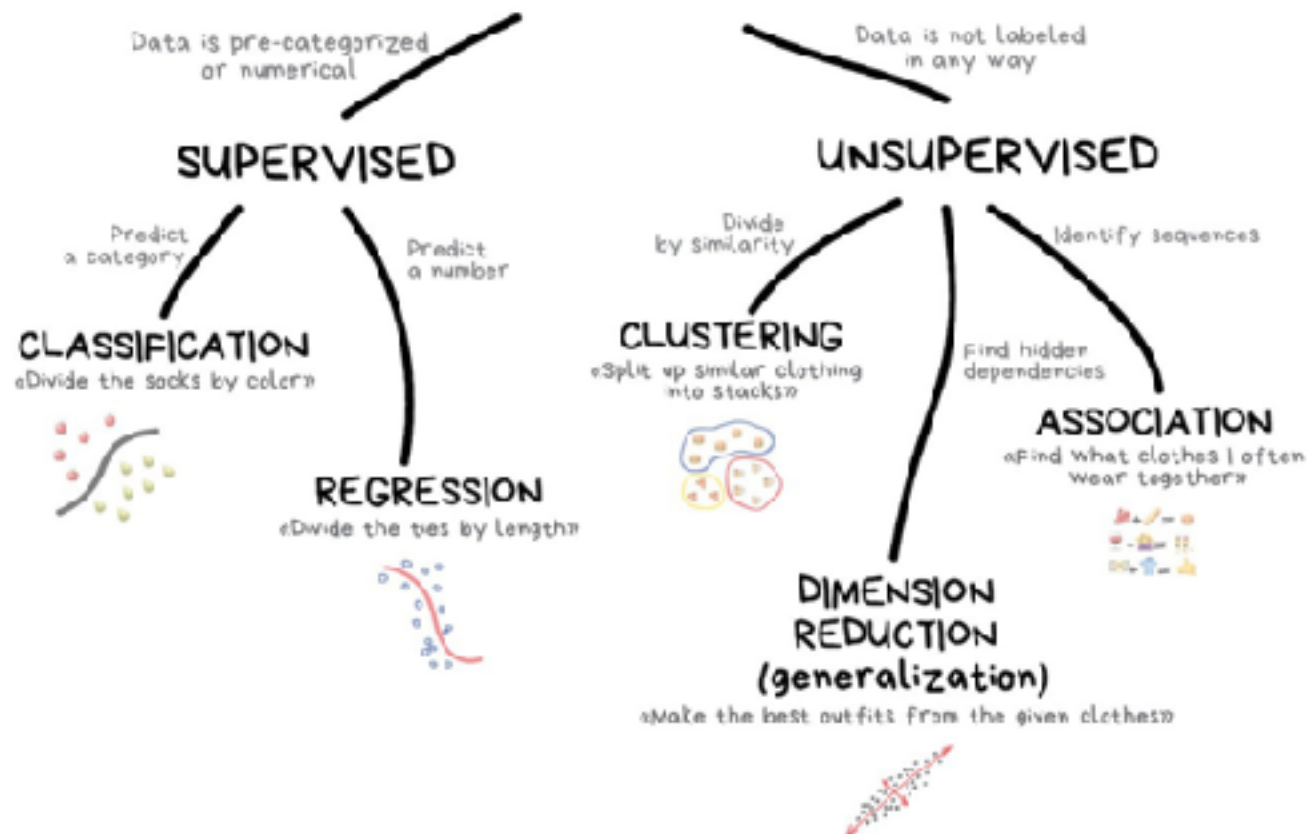Classification — categorical variables

Regression — continuous variables

Unsupervised Learning
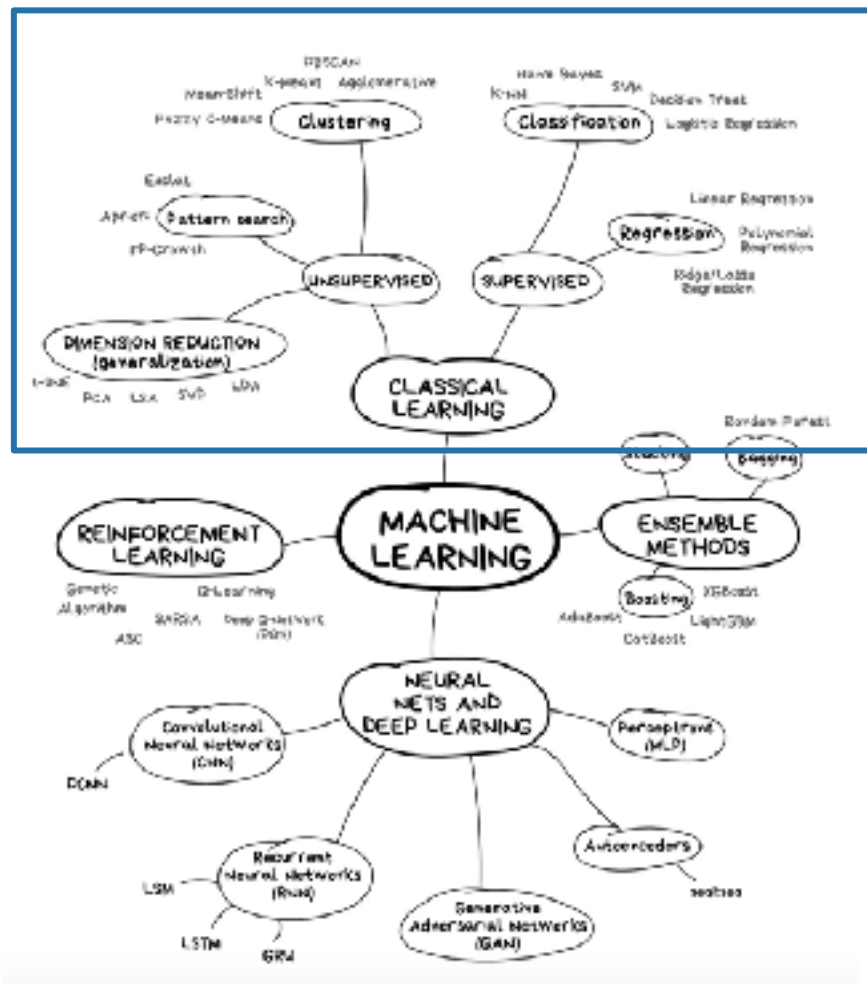
Clustering (categorical) & dimensionality reduction (continuous)

can automatically identify structure in data
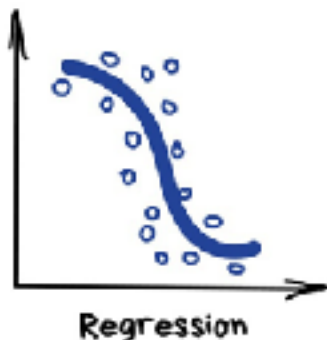
# CLASSICAL MACHINE LEARNING

Data is pre-categorized or numerical

Data is not labeled in any way

## SUPERVISED

## UNSUPERVISED

Predict a category

Predict a number

Divide by similarity

Identify sequences

### CLASSIFICATION
«Divide the socks by color»

### REGRESSION
«Divide the ties by length»

### CLUSTERING
«Split up similar clothing into stacks»

Find hidden dependencies

### ASSOCIATION
«Find what clothes I often wear together»

### DIMENSION REDUCTION (generalization)
«Make the best outfits from the given clothes»

# Regression

*"Draw a line through these dots. Yep, that's the machine learning"*

Today this is used for:

- Stock price forecasts
- Demand and sales volume analysis
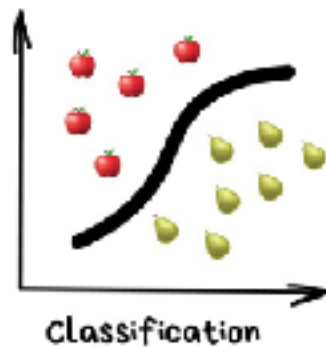- Medical diagnosis
- Any number-time correlations

Popular algorithms are <u>Linear</u> and <u>Polynomial</u> regressions.



Regression

# Classification

*"Splits objects based at one of the attributes known beforehand. Separate socks by based on color, documents based on language, music by genre"*

Today used for:

– Spam filtering
– Language detection
– A search of similar documents
– Sentiment analysis
– Recognition of handwritten characters and numbers
– Fraud detection

Popular algorithms: <u>Naive Bayes</u>, <u>Decision Tree</u>, <u>Logistic Regression</u>, <u>K-Nearest Neighbours</u>, <u>Support Vector Machine</u>



Classification

# **Regression**:
predicting <u>continuous</u> variables
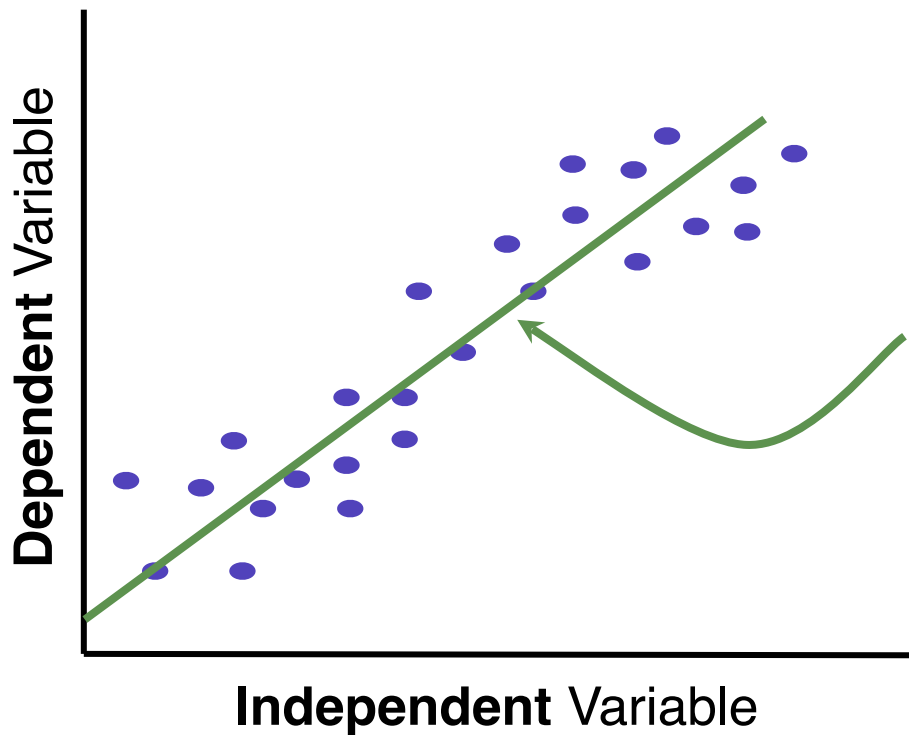(i.e. Age)

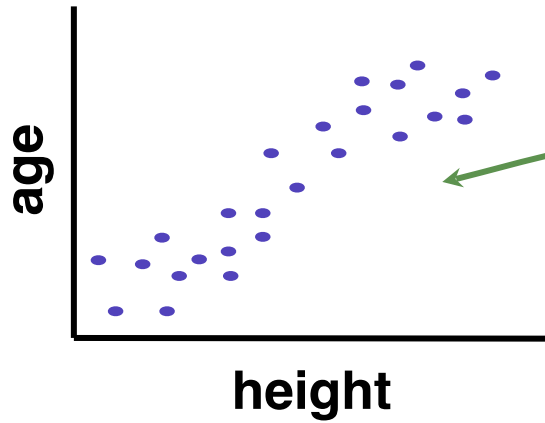<u>continuous</u> variable prediction

# **Classification**:
predicting <u>categorical</u> variables
(i.e. education level)

<u>categorical</u> variable prediction

Supervised Learning
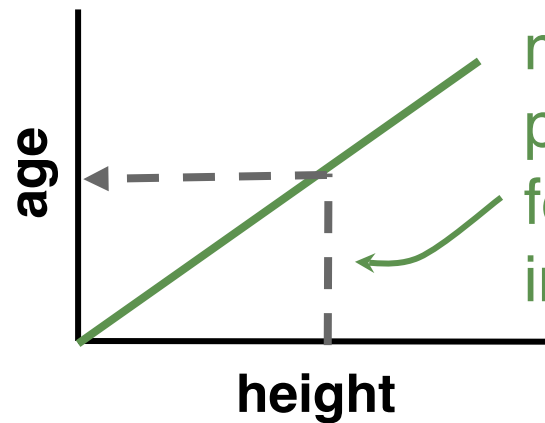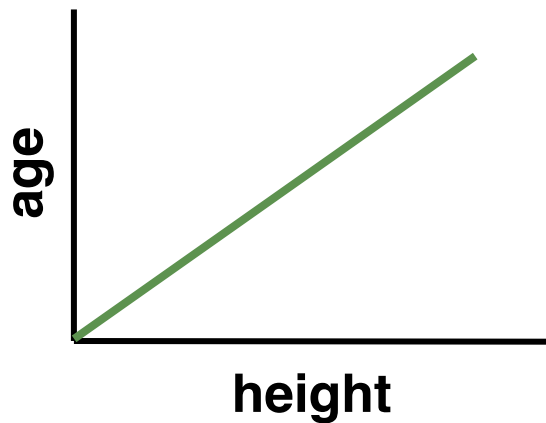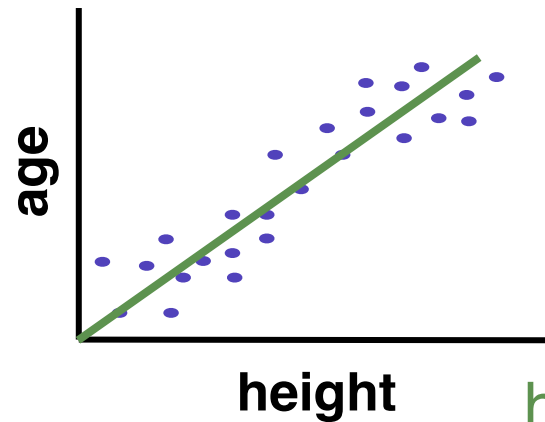
regression

continuous variable prediction

age

height

age

height

use linear regression to model the relationship

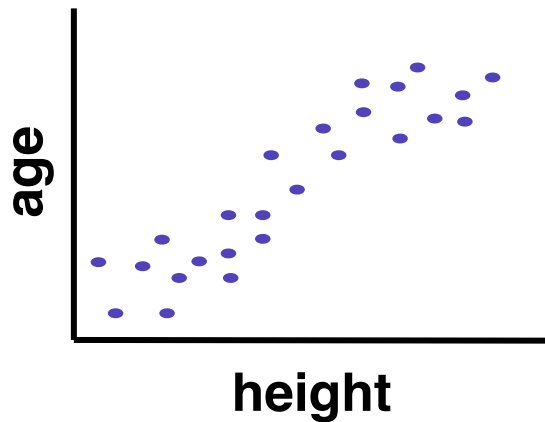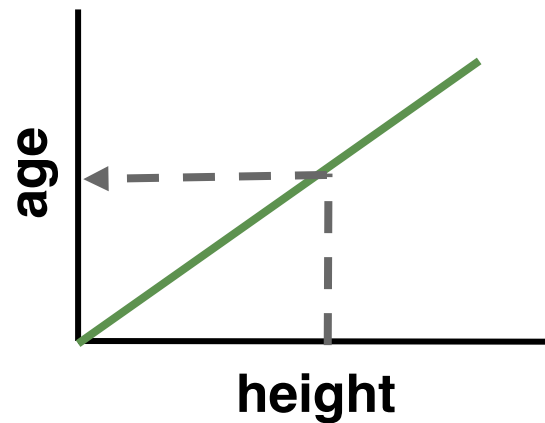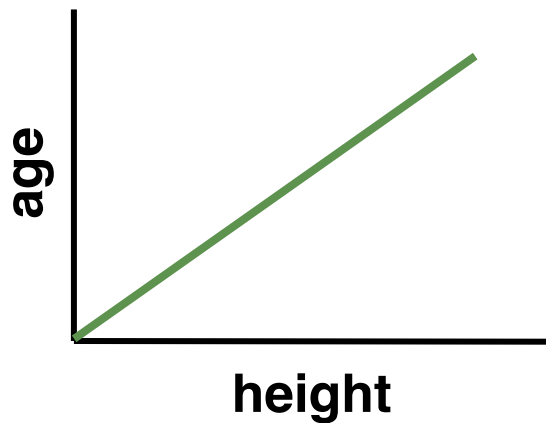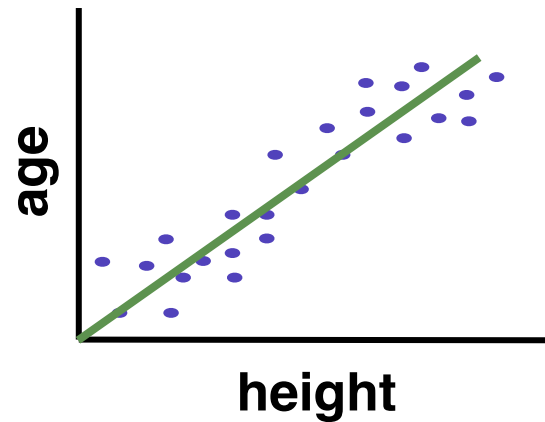For prediction, the individual values in the training data are *not* important. We only need the model.

regression

continuous variable prediction

PREDICT TRAFFIC JAMS

LINEAR
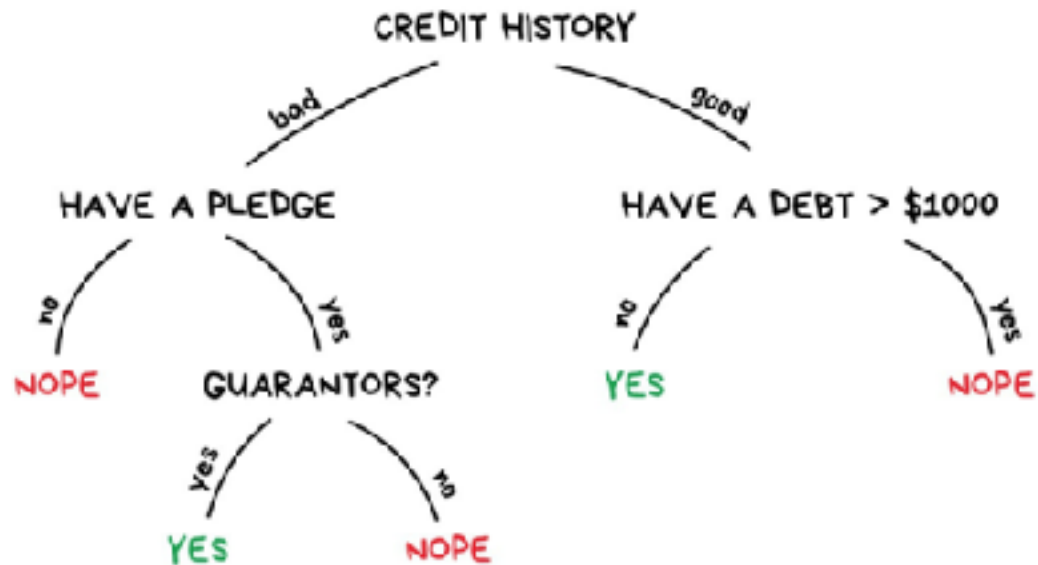
POLYNOMIAL

REGRESSION

**Regression**:
predicting continuous
variables
(i.e. Age)

**Classification**:
predicting categorical
variables
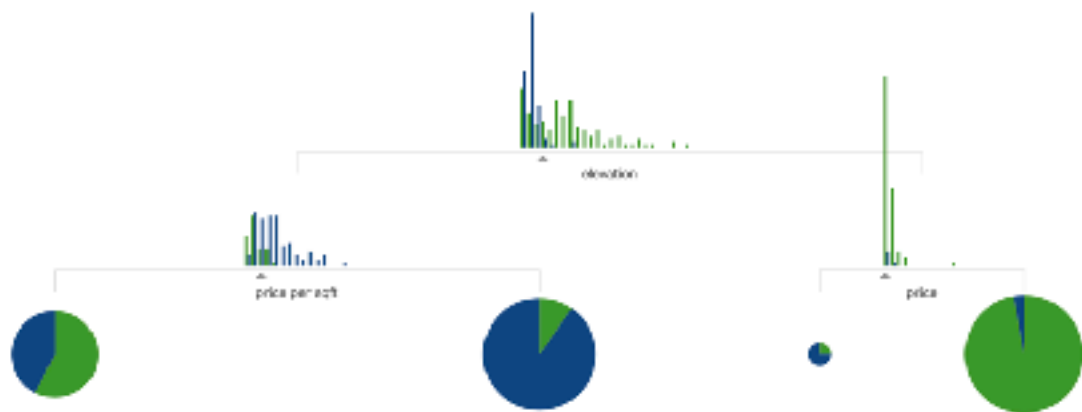(i.e. give a loan?)

GIVE A LOAN?

DECISION TREE

# Growing a tree

Additional forks will add new information that can increase a tree's **prediction accuracy.**

# K-nearest neighbors

1–Nearest Neighbor Classifier

15-Nearest Neighbor Classifier

Logistic regression

| Input Data | Nearest Neighbors | Linear SVM | RBF SVM | Gaussian Process | Decision Tree | Random Forest | Neural Net |
|---|---|---|---|---|---|---|---|
| | .97 | .88 | .97 | .97 | 95 | .95 | 90 |
| | .93 | .40 | .88 | .90 | 80 | .85 | 90 |
| | .93 | .93 | .95 | .93 | 95 | .95 | 95 |

https://scikit-learn.org/stable/auto_examples/classification/plot_classifier_comparison.html

# Unsupervised Learning

# To modes of machine learning



Supervised Learning

Unsupervised Learning

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

## Clustering

*"Divides objects based on unknown features. Machine chooses the best way"*

Nowadays used:

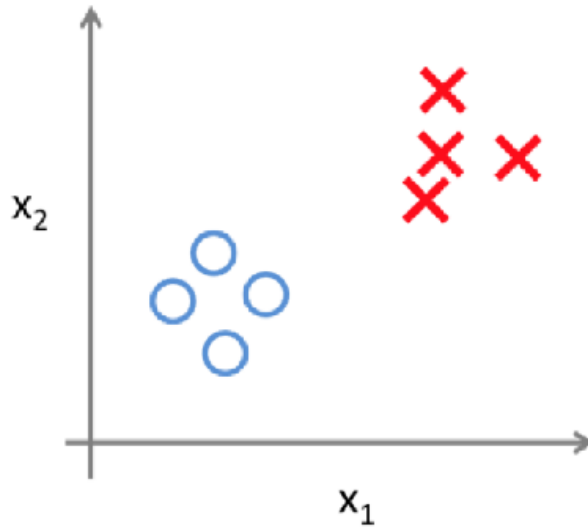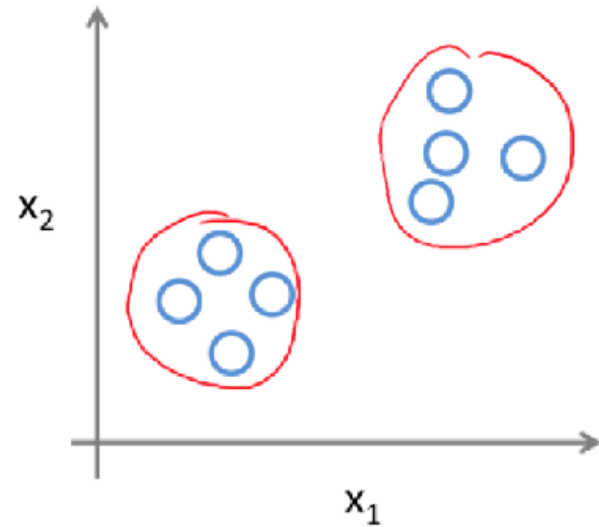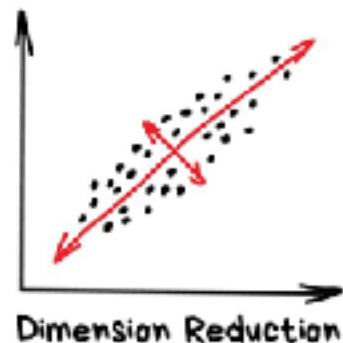- For market segmentation (types of customers, loyalty)
- To merge close points on a map
- For image compression
- To analyze and label new data
- To detect abnormal behavior

Popular algorithms: K-means_clustering, Mean-Shift, DBSCAN

## Dimensionality Reduction (Generalization)

*"Assembles specific features into more high-level ones"*

Nowadays is used for:

- Recommender systems (★)
- Beautiful visualizations
- Topic modeling and similar document search
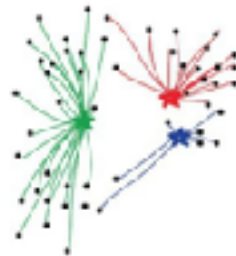- Fake image analysis
- Risk management

Popular algorithms: Principal Component Analysis (PCA), Singular Value Decomposition (SVD), Latent Dirichlet allocation (LDA), Latent Semantic Analysis (LSA, pLSA, GLSA), t-SNE (for visualization)
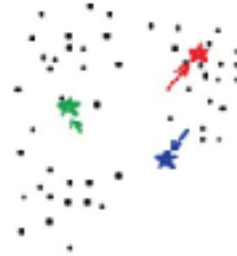
PUT KEBAB KIOSKS IN THE OPTIMAL WAY
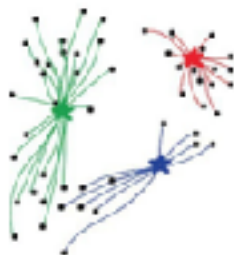(also illustrating the k-means method)

1. Put kebab kiosks in random places in city

2. Watch how buyers choose the nearest one

3. Move kiosks closer to the centers of their popularity

4. Watch and move again

5. Repeat a million times

6. Done! You're god of kebabs!

# Prediction Approach

You want to predict someone's emotion based on an image.

How would you approach this with machine learning?

A
Supervised,
Regression

B
Supervised,
Classification

C
Unsupervised
,
dimensionality
reduction

D
Unsupervised,
clustering

E
Unsupervised,
Neural
Network

"We have a thousand-layer network, dozens of video cards, but still no idea where to use it. Let's generate cat pics!"
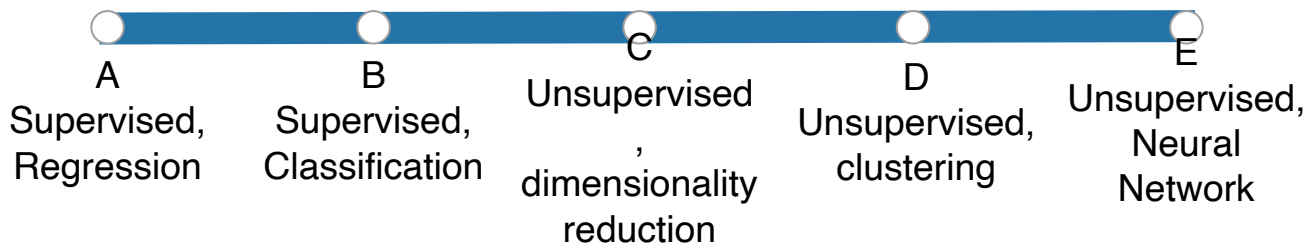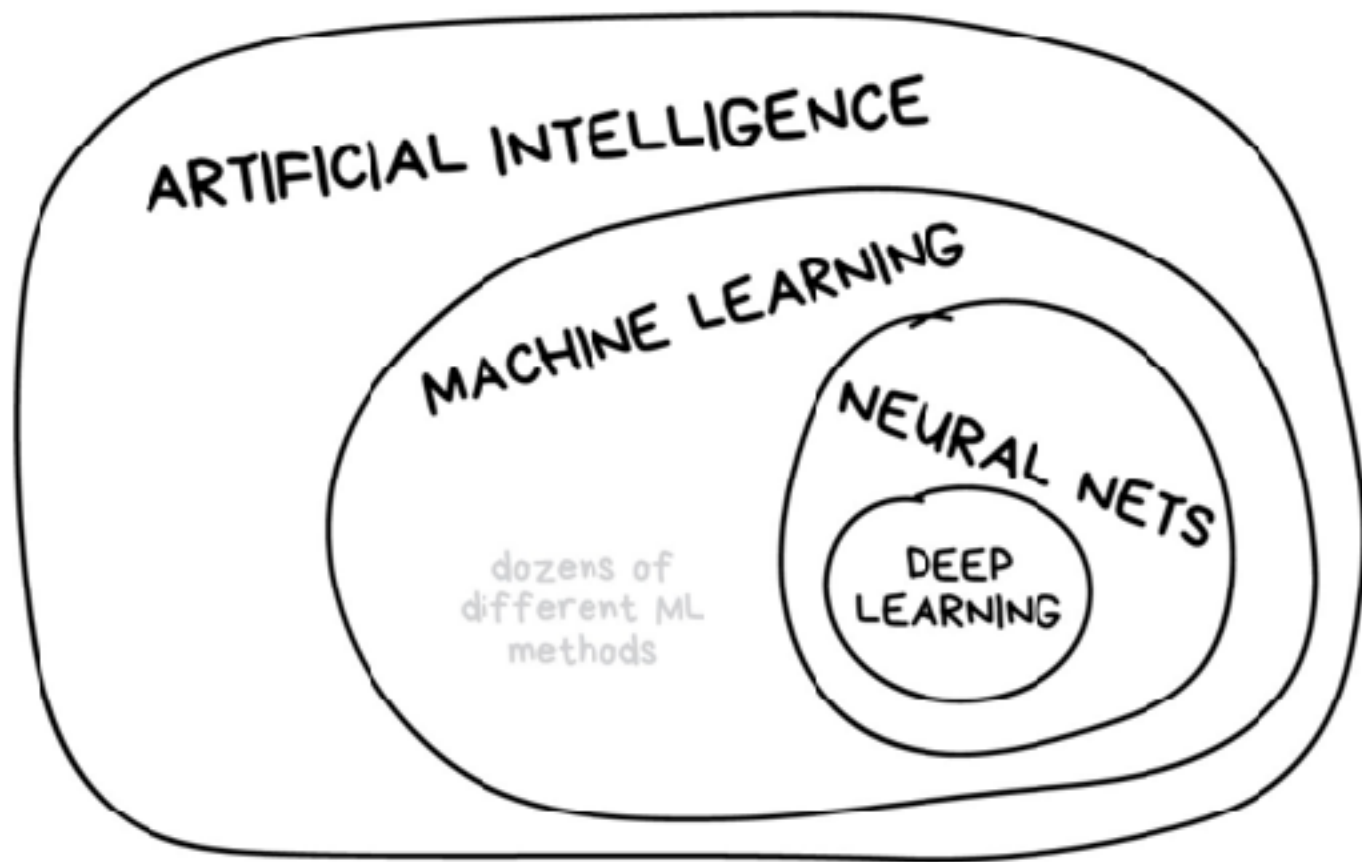
Used today for:

- Replacement of all algorithms above

- Object identification on photos and videos

- Speech recognition and synthesis

- Image processing, style transfer

- Machine translation



**Neural Networks**

Popular architectures: Perceptron, Convolutional Network (CNN), Recurrent Networks (RNN), Autoencoders

# WHAT IS A NEURON?



- Receives signal on synapse
- When trigger sends signal on axon

# MATHEMATICAL NEURON



$x_1$
$x_2$
$x_3$
$+1$

$h_{w,b}(x)$

- Mathematical abstraction, inspired by biological neuron
- Either on or off based on sum of input

This will likely not be the last time you see this (mostly unhelpful) neural net image

# Convolutional neural networks

These weights tell the neuron to respond more to one input and less to another. Weights are adjusted when training — that's how the network learns. Basically, that's all there is to it.

Inputs

Hidden layers

Outputs

0
1
2
3
4
5
6
7
8
9

4

1 pixel = 1 input

MULTILAYER PERCEPTRON (MLP)

# Manually labeling used to be the way...



Original image → Preliminary processing → Hand-crafted features → Neural Network → Result «cat»

# CNNs avoid manual features



☐▨▥▧ - convolution kernels

The neural network itself learns how to build
important features from the simple lines

«CAT»

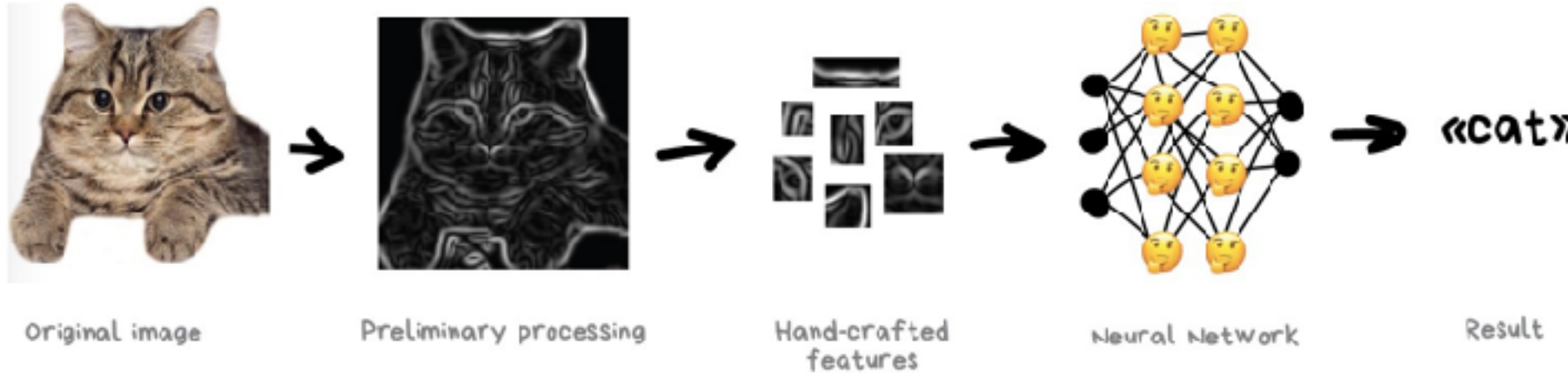Perceptron finds
important features
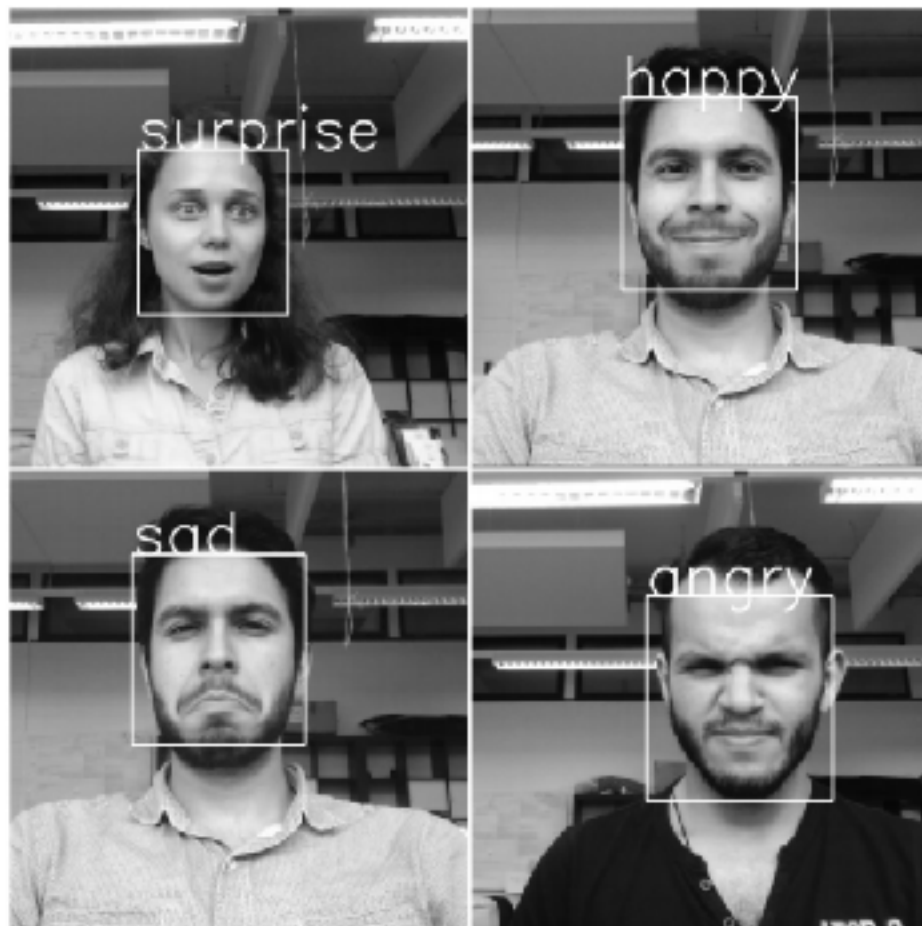specific to a cat

"CNNs are all the rage right now. They are used to search for objects on photos and in videos, face recognition, style transfer, generating and enhancing images, creating effects like slow-mo and improving image quality. Nowadays CNNs are used in all the cases that involve pictures and videos."

CONVOLUTIONAL NEURAL NETWORK (CNN)

Much of DL success comes from semi-supervised tricks to avoid large hand labelled datasets

## Masked LM

- **Solution**: Mask out $k\%$ of the input words, and then predict the masked words
    - We always use $k = 15\%$

```
                          store              gallon
                            ↑                  ↑
  the man went to the [MASK] to buy a [MASK] of milk
```

- Too little masking: Too expensive to train
- Too much masking: Not enough context

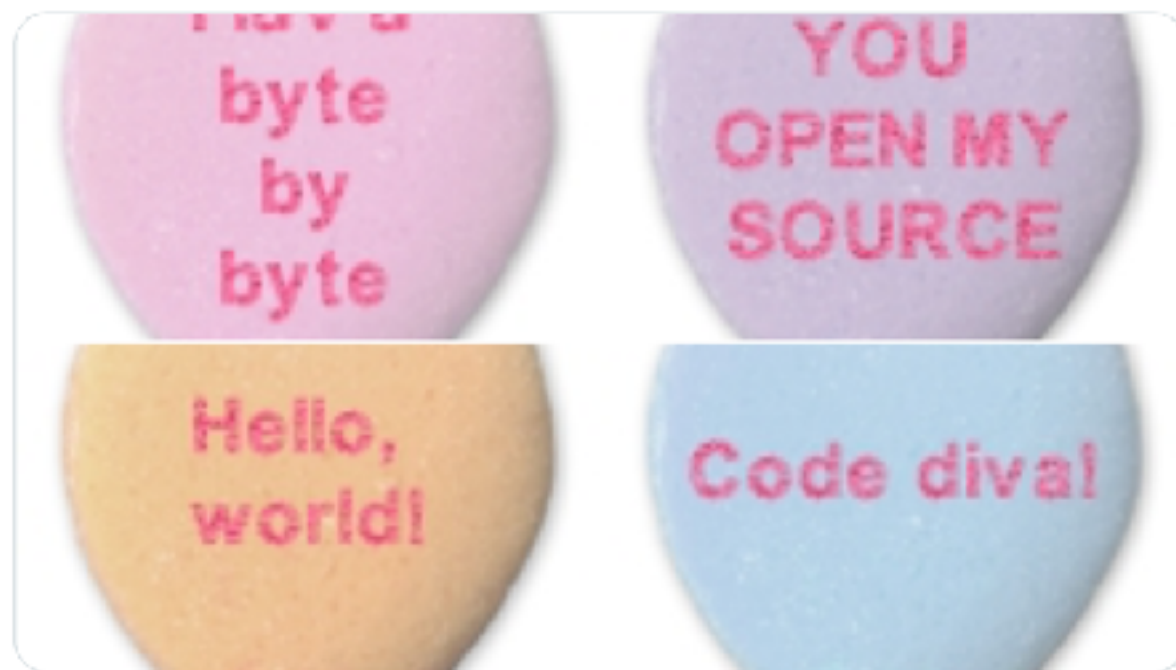https://neptune.ai/blog/unmasking-bert-transformer-model-performance

**all your base are belong to**
@jasongfleischer

I used @OpenAI GPT-3 to make some #programming themed candy hearts for you this Valentine's Day. Hope you feel the 100% computer generated 💕! Here's a selection of the one's I liked best (thread 1/3)

**all your base are belong to** @jasongfleischer · 1h

It's too long to fit on the candy heart generator but 🤣🤣💀 (bonus post, now the thread is over really)
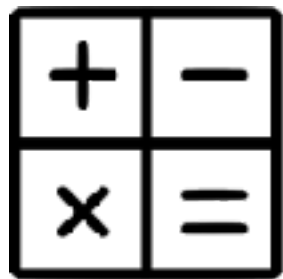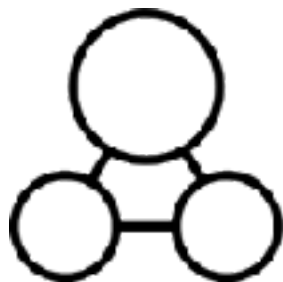
candy theme: normal
messages: BESTIE, CUTIE PIE, SOUL MATE, SWEET PEA, UR CUTE, YOU + ME, BE MINE, PICK ME, KISS ME, LOVE BIRDS, MARRY ME, OOO LA LA, TRUE LOVE, WINK WINK, XOXO
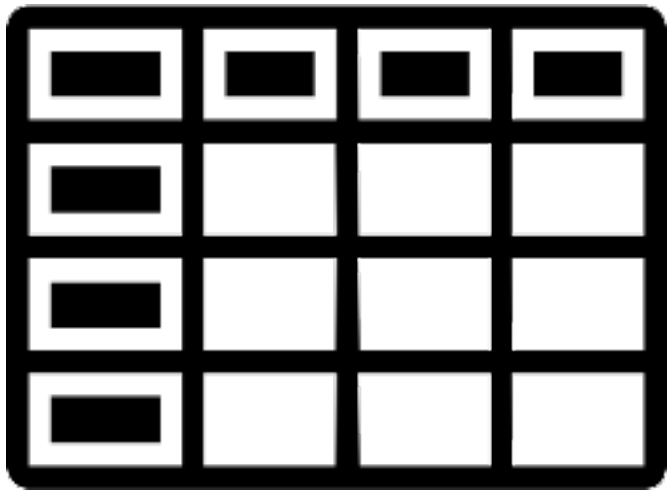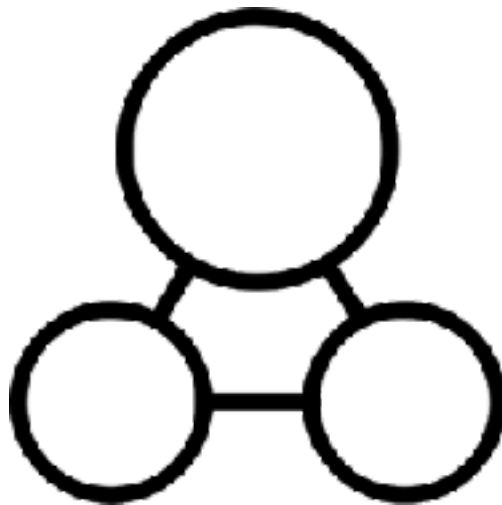candy theme: programming
messages:

(1) I <3 programmers, (2) Programmers do it better, (3) If you can't code it, you can't date it, (4) Code is poetry, (5) Programmers make the world go round, (6) If you can't code, you can't love, (7) Love is the language of the future, (8) I heart code, (9) Code is life, (10) If you don't code, you
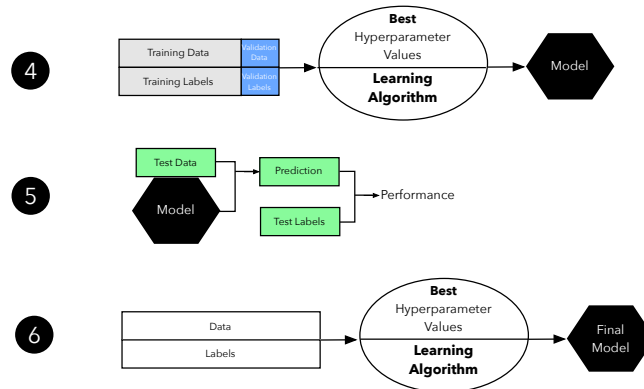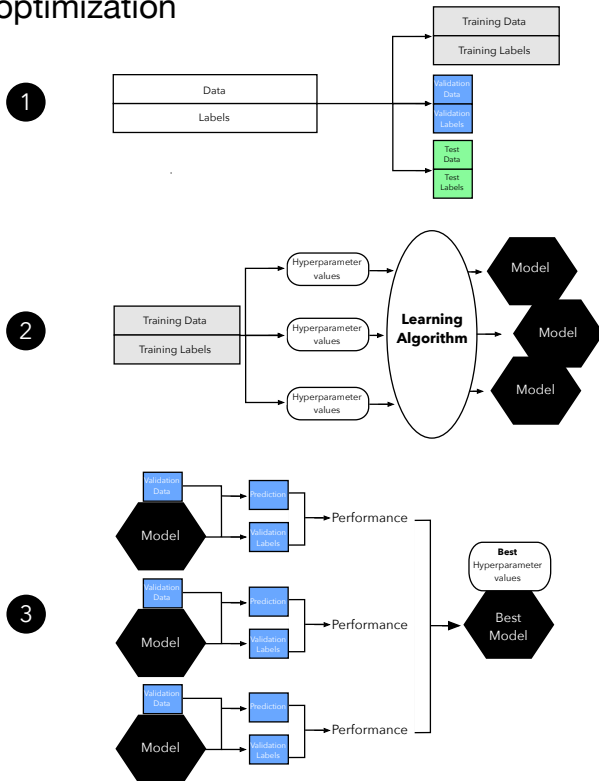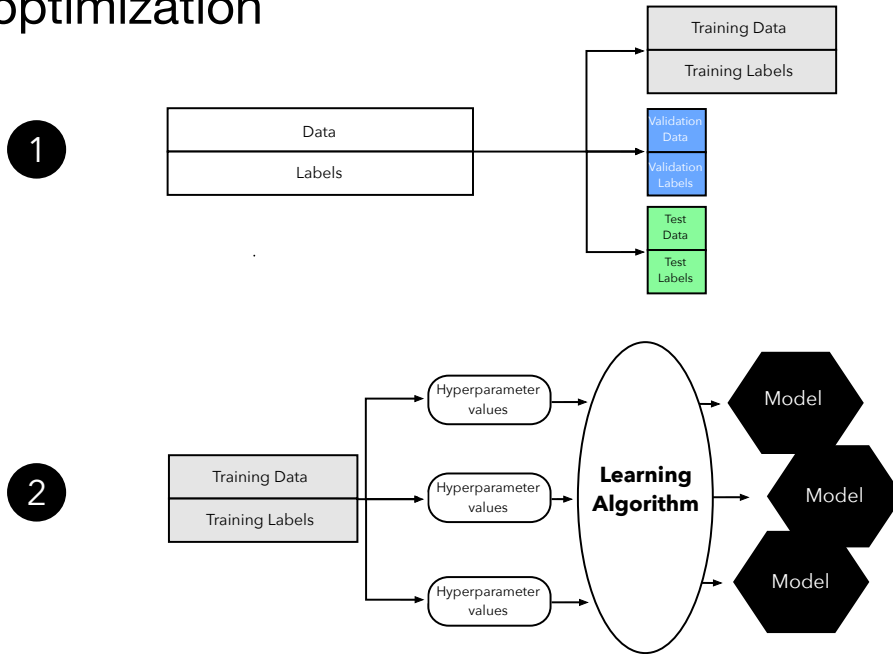
model selection

BIG
datasets

SIMPLE
models

# 3-Way Holdout

instead of "regular" holdout to avoid "data leakage" during hyperparameter optimization
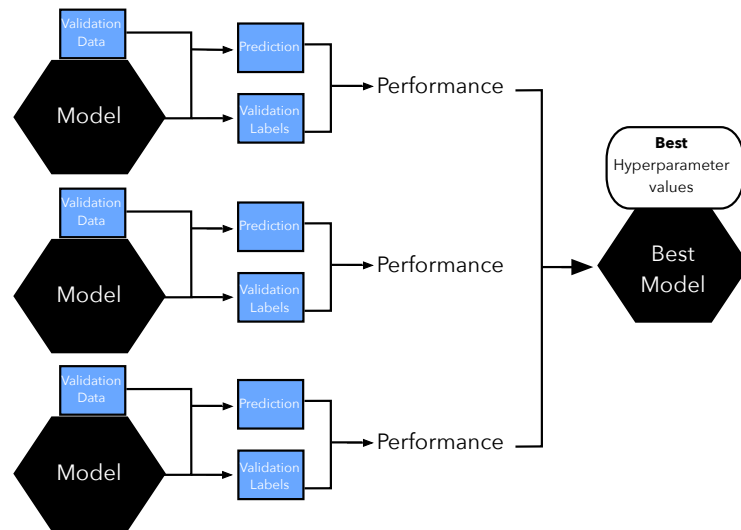
# 3-Way Holdout

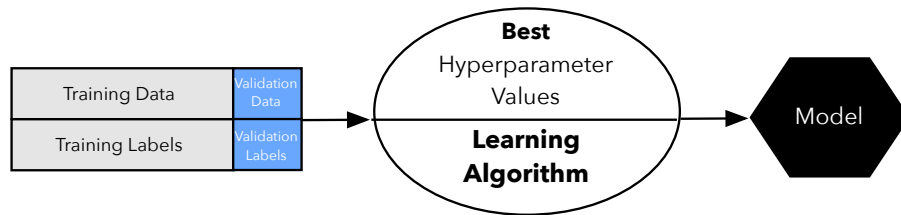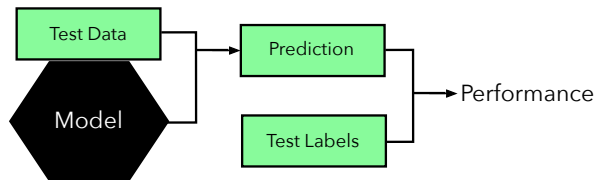instead of "regular" holdout to avoid "data leakage" during hyperparameter optimization

Sebastian Raschka                                    STAT 451: Intro to ML

# Model, feature, and hyperparameter selection

- Needs: Don't leak information from training/selection process into the test set!

- Trade-offs: Usually not enough data to have completely separate train, validation, test sets.  Which one do we prioritize?

    - Low training data -> bad fit

    - Low validation data -> bad selection of model/feature/hparam

    - Low test data -> poor estimate of generalization performance

SMALLER datasets        COMPLEX models

# k-Fold Cross-Validation

- Non-overlapping test data, overlapping training data among folds

- Small k -> biased pessimistic from small training data

- Variance increases with k; LOOCV is unbiased but hi variability

model assessment

# Root Mean Squared Error (RMSE)

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

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

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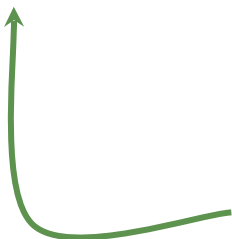
# Accuracy can mislead

- If classes are imbalanced, what would "chance performance" be?

- The classic example: detect cancer

  - 1/1000 actually have cancer

  - your prediction algorithm misdiagnoses 1% of healthy people

  - Do the math… your algorithm tells 10 people who are healthy they are sick for every 1 person who is actually sick

|  |  | Actual | |
|---|---|---|---|
|  |  | Positive | |
| Predicted | Positive | True Positive (TP) | False Positive |
| | Negative | False Negative | True Negative (TN) |

|  |  | 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

**Sensitivity**

$$\frac{TP}{TP + FN}$$



AUC = 0.94

Sensitivity

1 - Specificity

**Specificity**

$$\frac{TN}{TN + FP}$$

| Accuracy | What % were predicted correctly? |
|---|---|
| Sensitivity | Of those that *were* positives, what % were predicted to be positive? |
| Specificity | Of those that were *negatives*, what % were predicted to be negative? |

# Prediction Approach

You've been given a dataset with a number of features and have been asked to predict each individual's age.
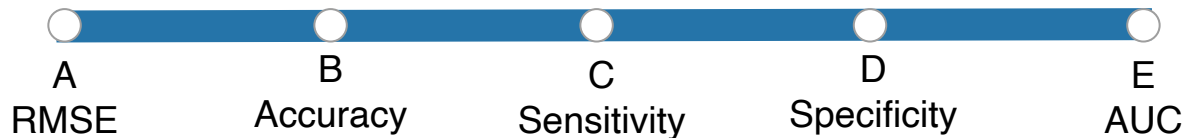
What prediction approach would you use?

A
regression
(supervised)

B
classification
(supervised)

C
clustering
(unsupervised)

D
dimensionality
reduction
(unsupervised)

# Prediction Approach

After predicting each person's age, how would you assess your model?

| A | B | C | D | E |
|---|---|---|---|---|
| RMSE | Accuracy | Sensitivity | Specificity | AUC |

# Prediction Approach

Which would be the error value you'd want from your model?



| A | B | C | D | E |
|---|---|---|---|---|
| 0.2 | 1.3 | 2.5 | 10.0 | 20.0 |