

# CMSI 5350/EECE 5998

# Machine Learning

Dr. Mandy Korpusik

The instructor gratefully acknowledges Alex Wong (Yale), Jessica Wu (HMC), and the many others who made their course materials freely available online.

# From last time... NumPy

```
import numpy as np
```

```
x = np.zeros((5, 4, 3))
print(x.shape)
```

```
x[:, :, 0] = 1 # Reminder: what did this do?
x[..., 1:3] = [2, 3] # And this?
x = np.reshape(x, (-1, 5)) # what's the shape?
```

```
def global_mean(x):
    return np.sum(x) / np.prod(x.shape)
```

```
global_mean(x) # what do we expect for the mean?
np.mean(x) # Let's check it
```

# Resources for NumPy

<https://numpy.org/doc/stable/user/quickstart>

The screenshot shows the NumPy User Guide website. The top navigation bar includes links for User Guide (which is underlined in blue), API reference, Building from source, Development, Release notes, Learn (with a dropdown arrow), and More (with a dropdown arrow). The main content area has a breadcrumb navigation path: Home > NumPy user guide > NumPy quickstart. The title of the page is "NumPy quickstart". Below the title is a section titled "Prerequisites" with the text: "You'll need to know a bit of Python. For a refresher, see the [Python tutorial](#). To work the examples, you'll need `matplotlib` installed in addition to NumPy." There is also a "Learner profile" section and a "Learning Objectives" section. On the left sidebar, there is a "Section Navigation" section with links to "Getting started", "What is NumPy?", "Installation", "NumPy quickstart" (which is highlighted with a blue border), "NumPy: the absolute basics for beginners", "Fundamentals and usage", "NumPy fundamentals", "NumPy for MATLAB users", "NumPy tutorials", "NumPy how-tos", "Advanced usage and interoperability", "Using NumPy C-API", and "F2PY user guide and reference manual".

# Scikit-learn library

```
import sklearn  
import sklearn.datasets as skdata  
  
housing = skdata.fetch_california_housing()  
x = housing.data  
y = housing.target  
feat_names = housing.feature_names
```

# Visualizing data

Let's plot our data using the matplotlib library!

```
from matplotlib import pyplot as plt
```

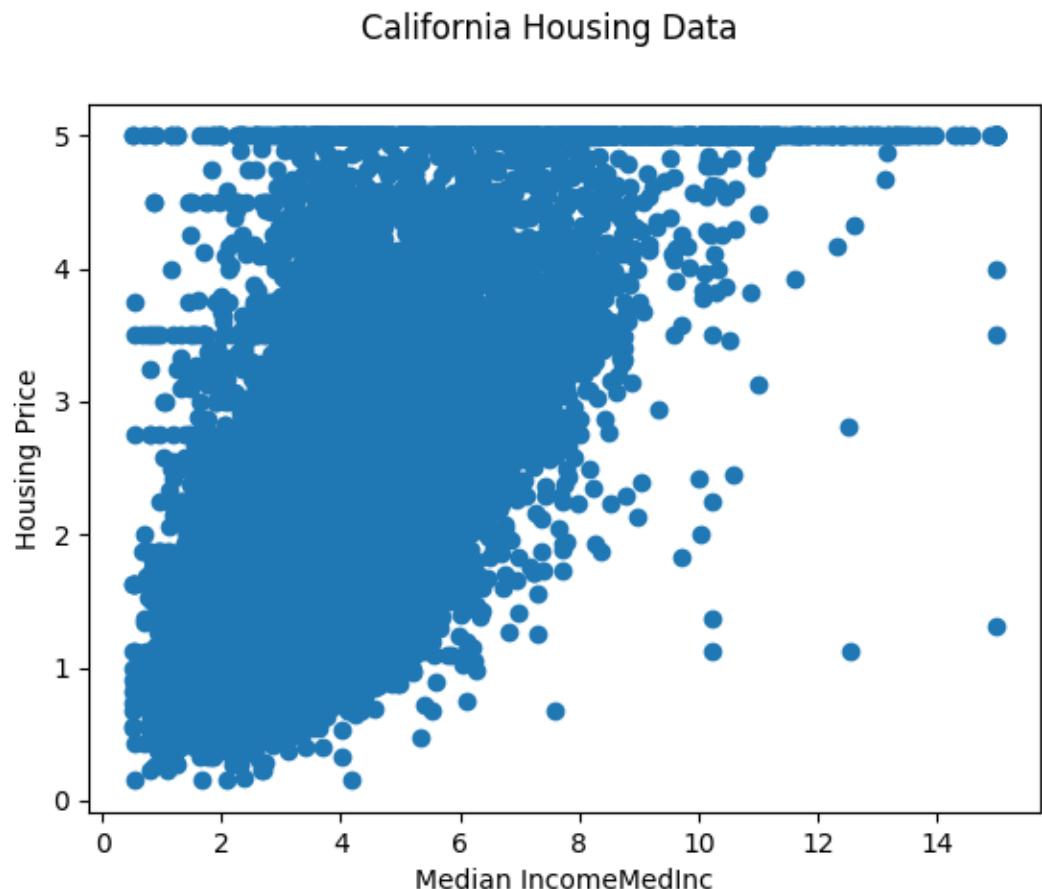
Let's start with how income varies with price.

```
income = x[:, 0] # MedInc is the 1st dimension  
  
fig = plt.figure()  
  
ax = fig.add_subplot(1, 1, 1) # (row, col, idx)  
ax.scatter(income, y)  
  
plt.show(block=True)
```

# Visualizing data

Let's add labels to the axes:

```
fig = plt.figure()  
  
ax = fig.add_subplot()  
ax.scatter(income, price)  
  
fig.suptitle('California Housing Data')  
ax.set_ylabel('Housing Price')  
ax.set_xlabel('Median Income')  
  
plt.show(block=True)
```



# Visualizing data

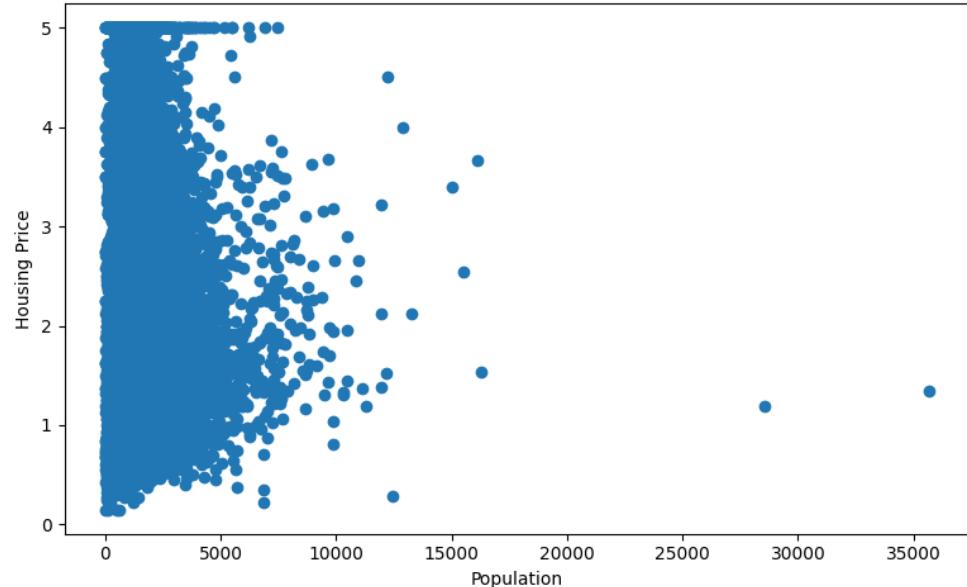
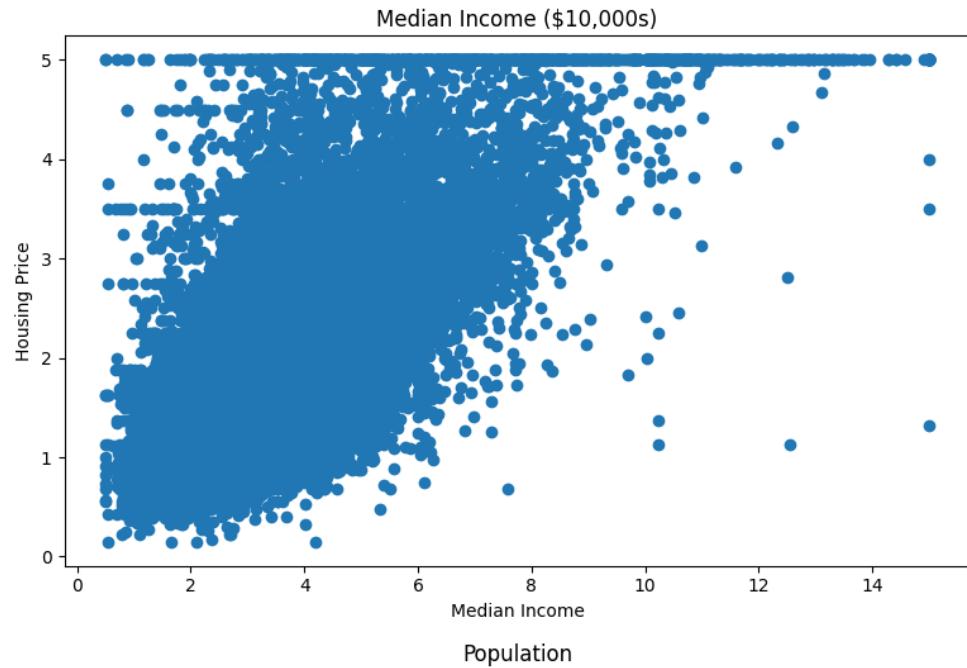
Let's see population vs

```
fig = plt.figure()
fig.suptitle('California')

ax1 = fig.add_subplot(2,
ax1.set_title('Median Inc')
ax1.set_ylabel('Housing P')
ax1.set_xlabel('Median In'
ax1.scatter(income, y)

population = x[:, 4]
ax2 = fig.add_subplot(2,
ax2.set_title('Population')
ax2.set_ylabel('Housing P')
ax2.set_xlabel(feat_names
ax2.scatter(population, y

plt.show(block=True)
```



# Visualizing data

Let's try overlaying the two sets of features:

```
fig = plt.figure()
fig.suptitle('California Housing Data')

ax = fig.add_subplot(1, 1, 1)
ax.set_ylabel('Housing Price')
ax.set_xlabel('Features')

obs = (income, population)
ys = (y, y)
clr = ('blue', 'red')
marker = ('o', '^')
lbl = (feat_names[0], feat_names[4])

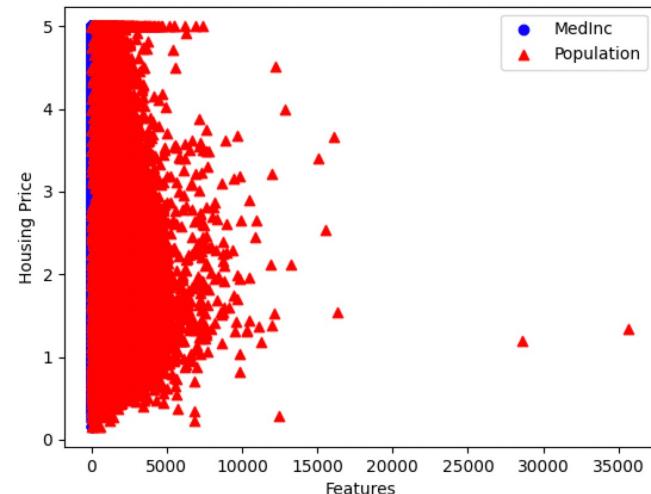
for o, yi, c, l, m in zip(obs, ys, clr, lbl, marker):
    ax.scatter(o, yi, c=c, label=l, marker=m)
ax.legend(loc='upper right')

plt.show(block=True)
```

# Visualizing data

The scales are too different... 😞

California Housing Data



```
print(np.min(income), np.max(income))
# MedInc: (0.4999, 15.0001)
```

```
print(np.min(population), np.max(population))
# Population: (3.0, 35682.0)
```

*Q: How can we get them on the same scale?*

# Tuesday's learning objectives

Students will be able to:

- ✓ Explore datasets in scikit-learn (sklearn) library
- ✓ Plot data with the matplotlib library
  - Implement min-max and standard normalization
  - Visualize data in 3D

# 1) Min-max normalization

Normalize using min and max values:

$$z = \frac{x - \min(x)}{\max(x) - \min(x)}$$

```
def min_max_norm(x):
    return (x - np.min(x)) / (np.max(x) - np.min(x))

# Let's apply it to income and population!
income_minmax = min_max_norm(income)
populn_minmax = min_max_norm(population)

print(np.min(income_minmax), np.max(income_minmax))
print(np.min(populn_minmax), np.max(populn_minmax))
```

# 1) Min-max normalization

Now let's overlay the two feature sets again:

```
fig = plt.figure()
fig.suptitle('California Housing Data')

ax = fig.add_subplot(1, 1, 1)
ax.set_ylabel('Housing Price')
ax.set_xlabel('Min-Max Norm Features')

obs = (income_minmax, popuLn_minmax)
for o, yi, c, l, m in zip(obs, ys, clr, lbl, marker):
    ax.scatter(o, yi, c=c, label=l, marker=m)
ax.legend(loc='upper right')

plt.show(block=True)
```

# 1) Min-max normalization

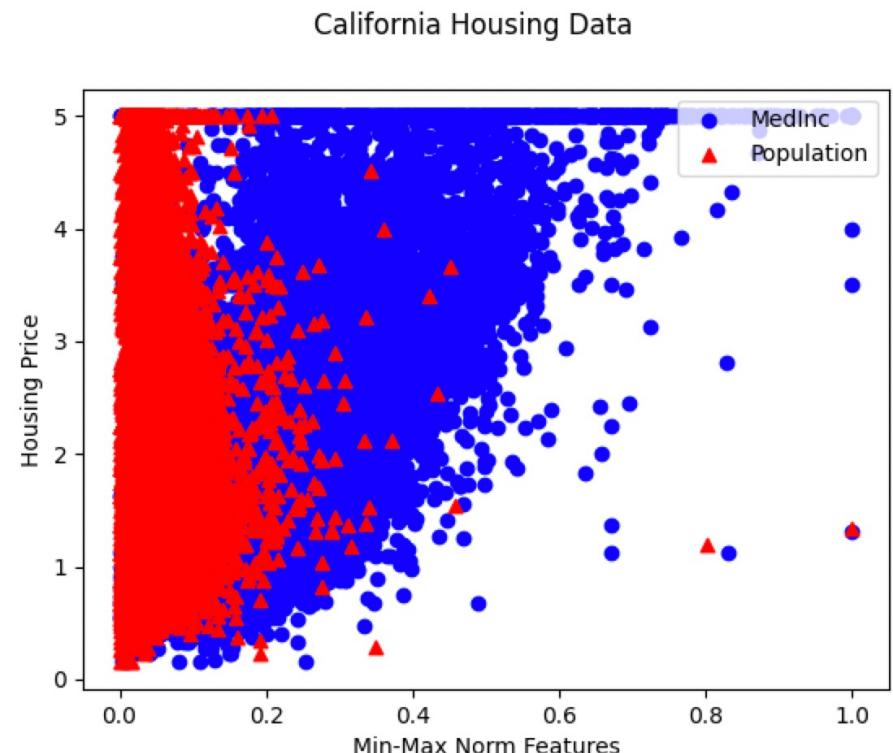
Now we can see trends in the data:

- Prices increase as median income increases
- Not a strong relationship between price and population

Q: Any issues with min-max?

Dividing population by max  
squishes observations btw. 0  
and 0.3, with a few outliers...

Q: Is there a way to normalize  
data so it's better distributed?



## 2) Standard normalization

Results in a mean of 0 and standard deviation of 1:

$$z = \frac{x - \mu}{\sigma}$$

```
def standard_norm(x):
    return (x - np.mean(x)) / (np.std(x))

# Let's apply it to income and population!
income_std = standard_norm(income)
populn_std = standard_norm(population)

print(np.min(income_std), np.max(income_std))
print(np.min(populn_std), np.max(populn_std))
```

## 2) Exercise: Standard normalization

Overlay the two features again. What do you notice?

```
fig = plt.figure()
fig.suptitle('California Housing Data')

ax = fig.add_subplot(1, 1, 1)
ax.set_ylabel('Housing Price')
ax.set_xlabel('Standard Norm Features')

obs = (income_std, popuLn_std)
for o, yi, c, l, m in zip(obs, ys, clr, lbl, marker):
    ax.scatter(o, yi, c=c, label=l, marker=m)
ax.legend(loc='upper right')

plt.show(block=True)
```

## 2) Standard normalization



- Min-max has a bounded range, but this will squish the data.
- Standard norm is not bounded, but retains info about outliers.

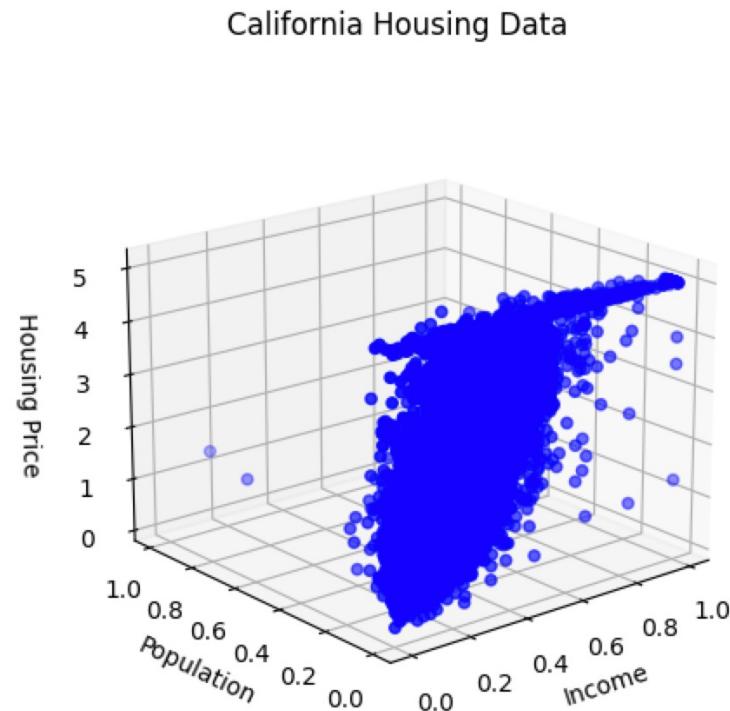
# Visualizing in 3D

Let's visualize income and population in 3D:

```
from mpl_toolkits.mplot3d import Axes3D  
  
fig = plt.figure()  
fig.suptitle('California Housing Data')  
  
ax = fig.add_subplot(1, 1, 1, projection='3d')  
ax.set_zlabel('Housing Price')  
ax.set_xlabel('Income')  
ax.set_ylabel('Population')  
ax.scatter(income_minmax, popuIn_minmax, y, c='blue',  
marker='o')  
  
plt.show(block=True)
```

# Visualizing in 3D

Can now see how the data is distributed together—how both features vary together with respect to housing price.



# Plotly

```
pip install plotly pandas
```

```
import plotly.express as px
import pandas as pd

data = pd.DataFrame(x, columns=feat_names)
print(data)

fig = px.scatter_geo(data, lat='Latitude',
                      lon='Longitude', scope='usa', color='MedInc')

# Try assigning different features to color!
fig.show()
```

# Exercise

Explore the UC Irvine Breast Cancer dataset:

```
import sklearn.datasets as skdata  
  
breast_cancer_data = skdata.load_breast_cancer()  
  
# Extract the data, feature names, and target  
# Pick a couple features to plot against the target!  
# Try normalizing the data
```

Target is binary (1 or 0) for benign or malignant.

# Today's learning objectives

Students will be able to:

- Define Machine Learning.
- List three types of Machine Learning algorithms.
- Give examples of supervised learning tasks.
- Name an unsupervised learning technique.
- Explain how reinforcement learning works.
- Describe the learning task with math.

# What is learning?

How do humans learn?

Observations → Learning → Action

How do machines learn?

Data (observations) → Learning → Decision (action)

# Reminder: Machine Learning is...

Machine learning, a branch of artificial intelligence, concerns the construction and study of systems that can learn from data.

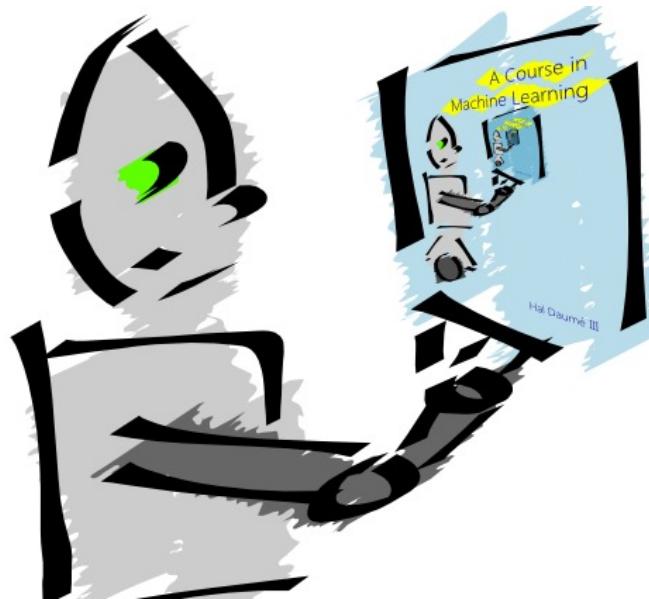


**WIKIPEDIA**  
The Free Encyclopedia

# Machine Learning is...

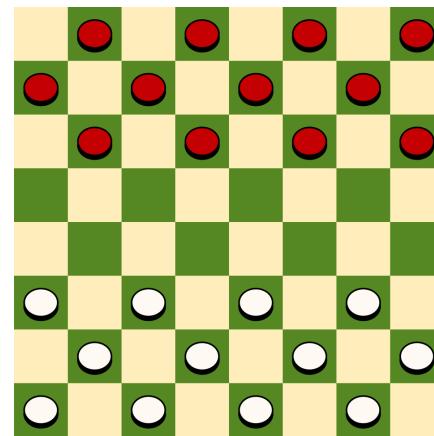
Machine learning is about **predicting the future** based on the past.

-- Hal Daume III



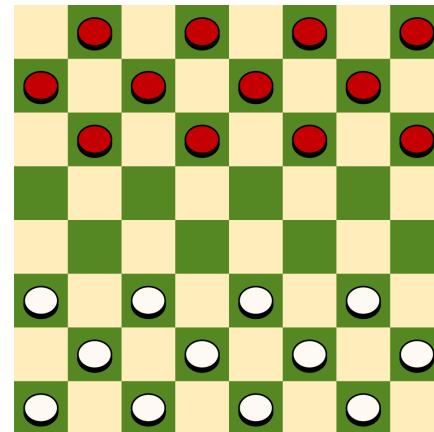
# Machine Learning is...

“the field of study that gives computers the ability to learn without being explicitly programmed.” -Arthur Samuel (1959)

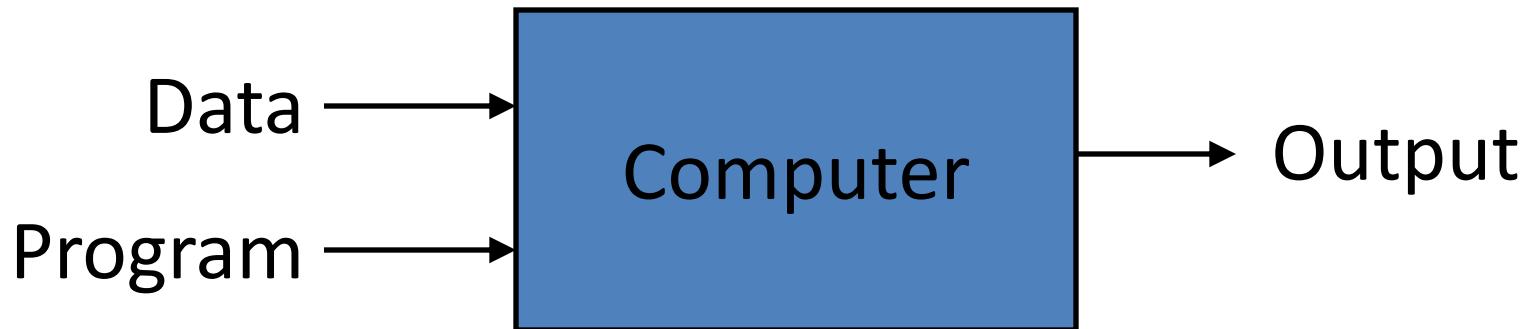


# Machine Learning is...

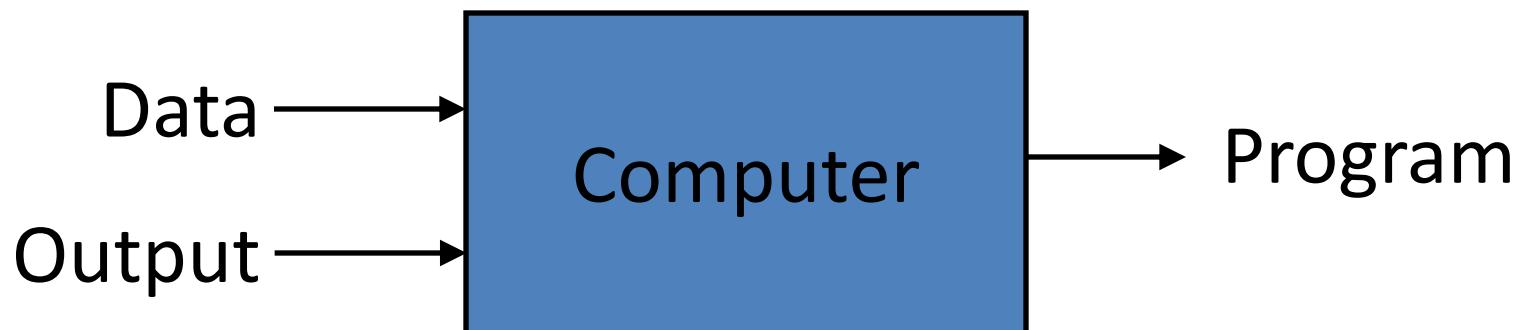
“a computer program is said to learn from experience E with respect to some class of tasks T and performance P, if its performance at tasks in T, as measured by P, improves with experience E.” -Tom Mitchell (1998)



## Traditional Programming



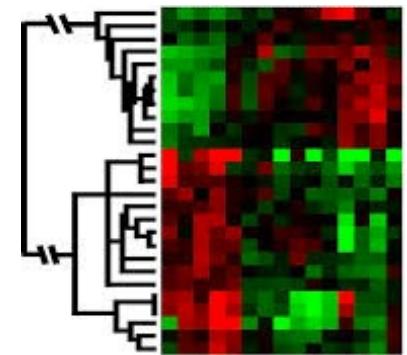
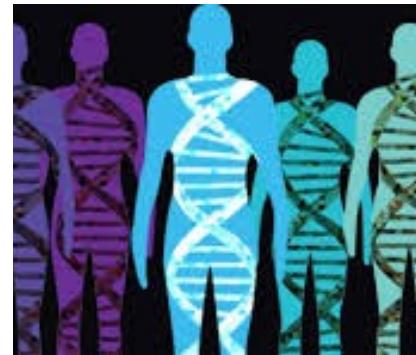
## Machine Learning



# When do we use Machine Learning?

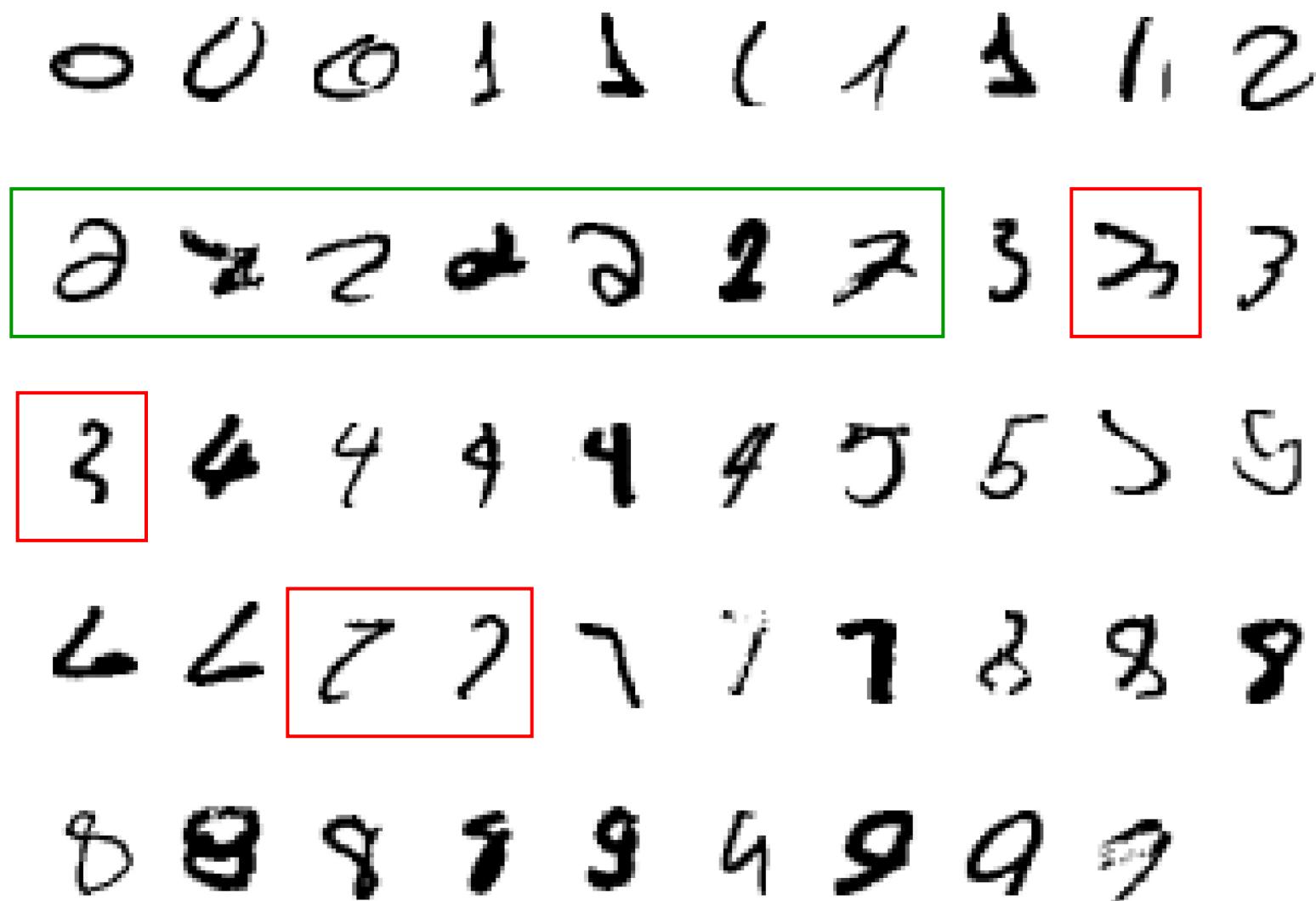
ML is used when:

- Human expertise does not exist
- Humans cannot explain their expertise
- Models must be customized
- Models are based on huge amounts of data



Learning is not always useful

A classic example of a task that requires machine learning:  
It is very hard to say what makes a 2



# Some examples of tasks that are best solved by using a learning algorithm

- Learning patterns
  - Predicting passenger survivability on the *Titanic*
  - Recognizing tweets as positive or negative
  - Clustering faces by identity
- State-of-the-art applications
  - Autonomous cars
  - Automatic speech recognition
  - Anomalies in credit card transactions
  - ChatGPT
- [Your favorite area here]

# ML/AI jobs on the rise

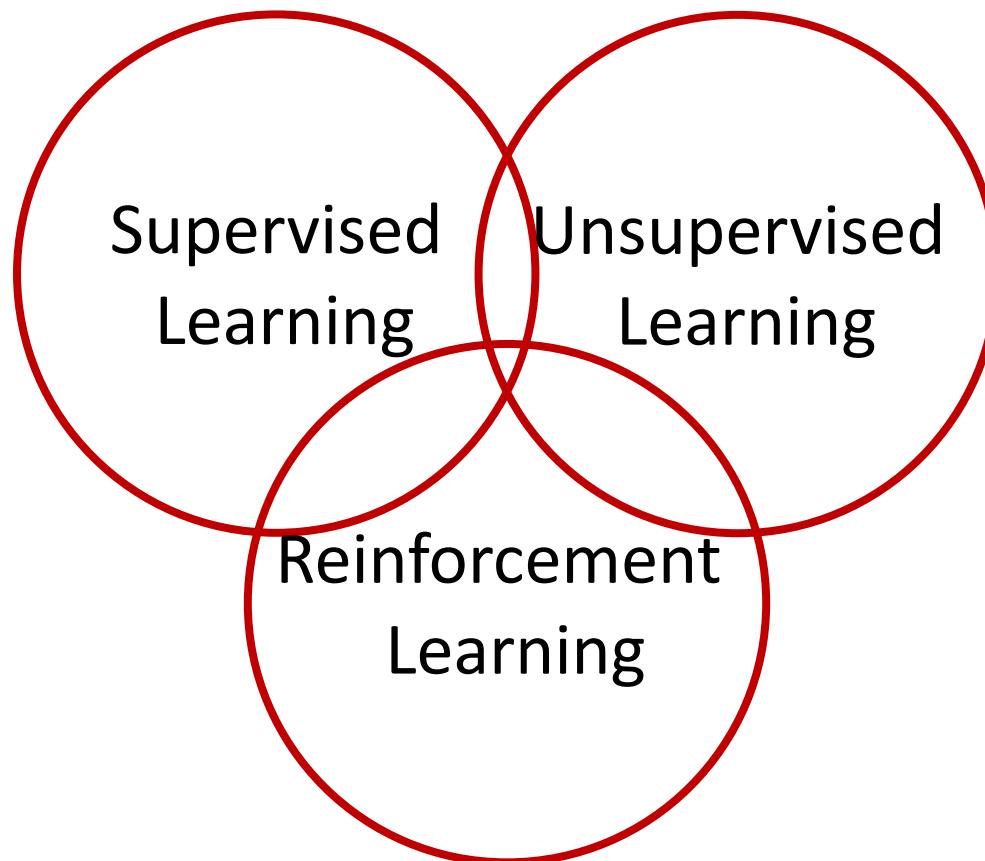
1. Chief Growth Officer
2. Government Program Analyst (**Data Analysis**)
3. Environment Health Safety Manager
4. Director of Revenue Operations
5. Sustainability Analyst (**Data Analysis**)
6. Advanced Practice Provider
7. VP of Diversity and Inclusion
8. **AI Consultant**
9. Recruiter
- 10. AI Engineer**

# Today's learning objectives

Students will be able to:

-  Define Machine Learning.
- List three types of Machine Learning algorithms.
- Give examples of supervised learning tasks.
- Name an unsupervised learning technique.
- Explain how reinforcement learning works.
- Describe the learning task with math.

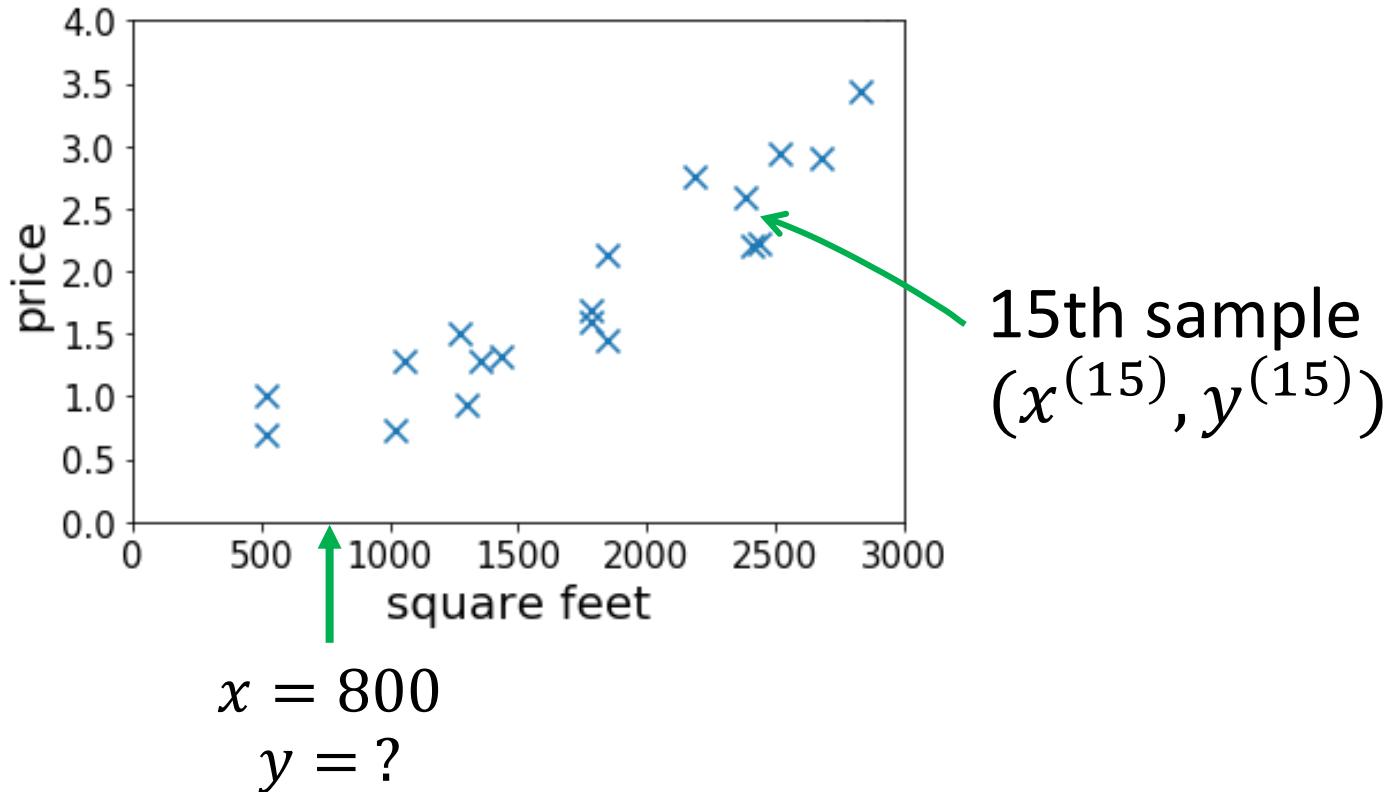
# Types of Machine Learning



# Supervised learning: Housing price

Given: a dataset of  $n$  samples  $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$

Task: if a residence has  $x$  sq ft, what is its price  $y$ ?



# More features

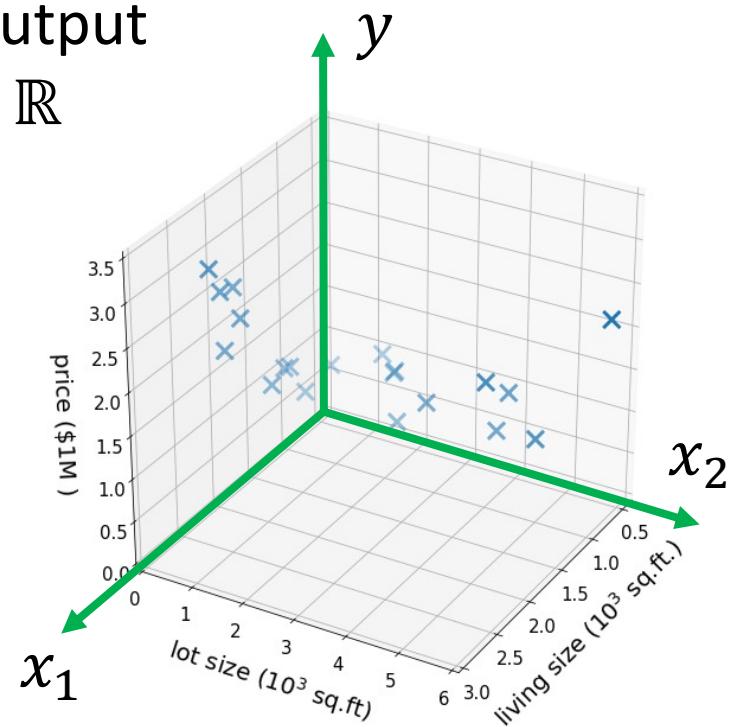
Suppose we also know the lot size

Task: find a function that maps

$(\text{size}, \text{lot size}) \rightarrow \text{price}$

features/input  
 $x \in \mathbb{R}^2$

label/output  
 $y \in \mathbb{R}$



# High-dimensional features

$x \in \mathbb{R}^d$  for large  $d$

Ex)

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ \vdots \\ \vdots \\ x_d \end{bmatrix} \begin{array}{l} \text{--- median income} \\ \text{--- house age} \\ \text{--- average rooms} \\ \vdots \end{array} \xrightarrow{\hspace{1cm}} y \text{ --- price}$$

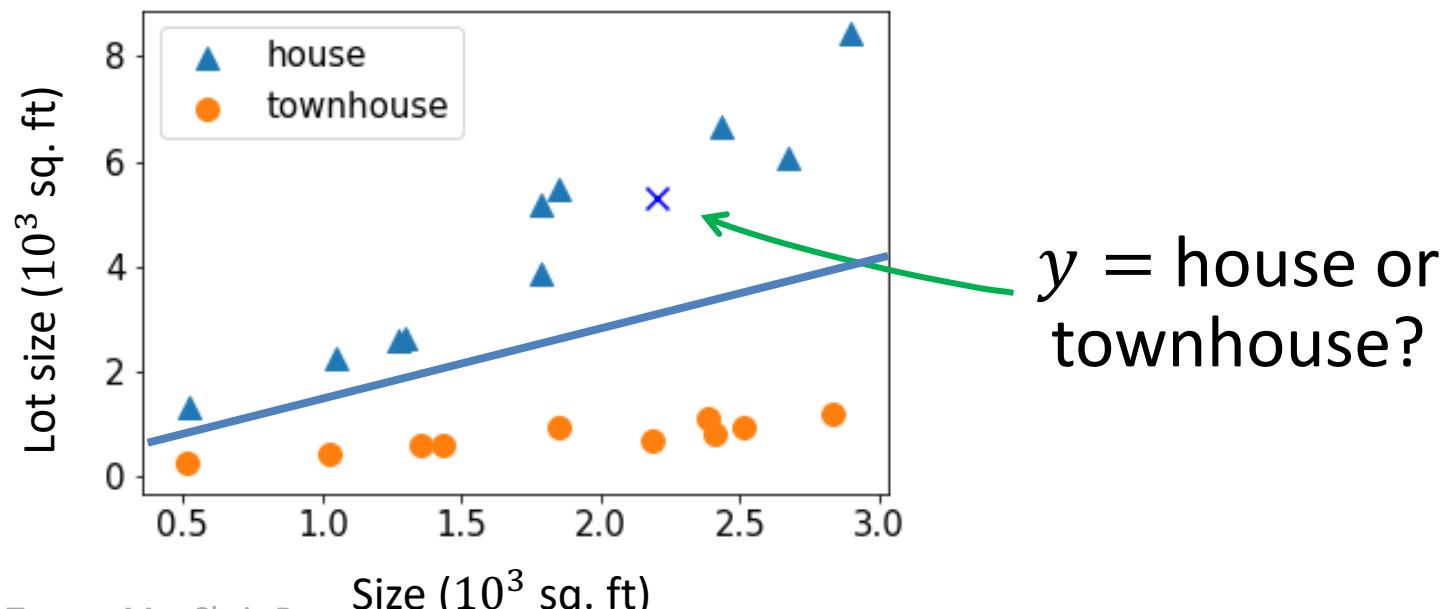
# Regression vs. classification

**Regression:** if  $y \in \mathbb{R}$  is a continuous variable

- e.g., price prediction

**Classification:** the label is a discrete variable

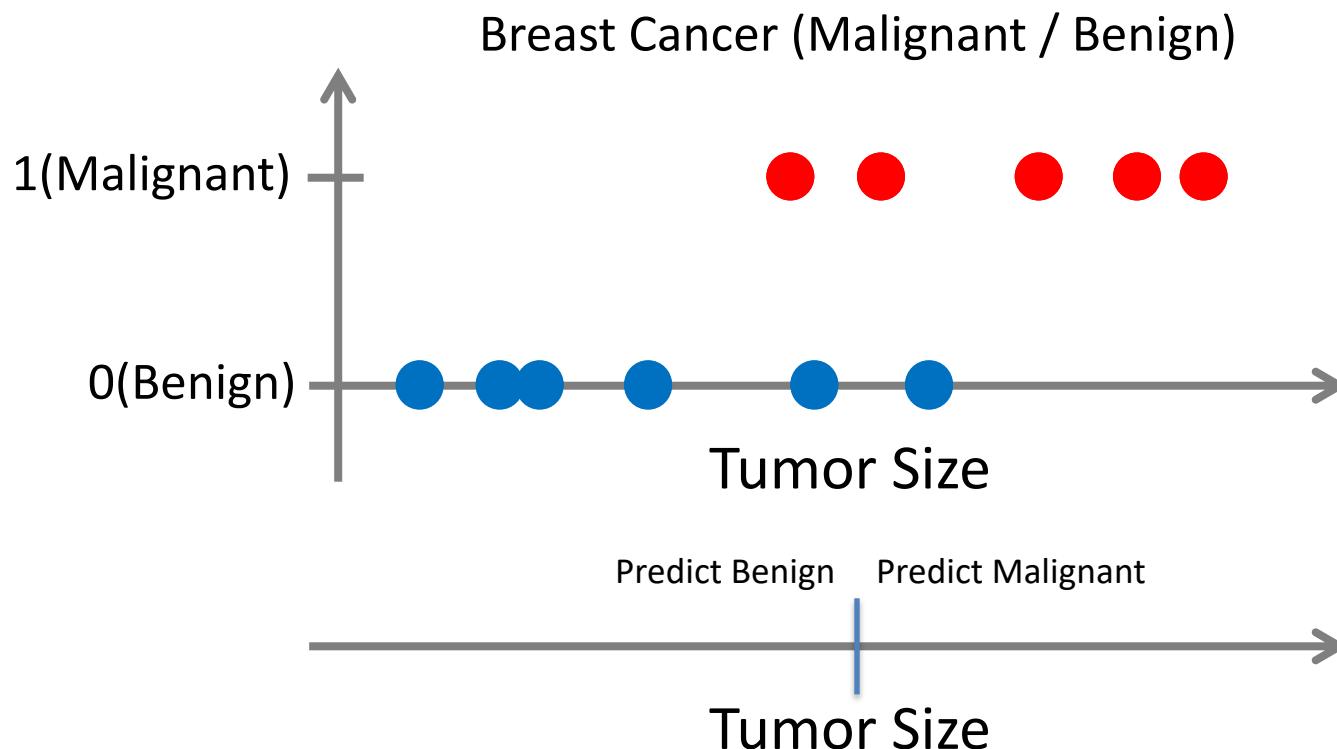
- e.g., the task of predicting the types of residence (size, lot size) → house or townhouse?



# Supervised Learning: Classification

Given: a dataset of  $n$  samples  $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$

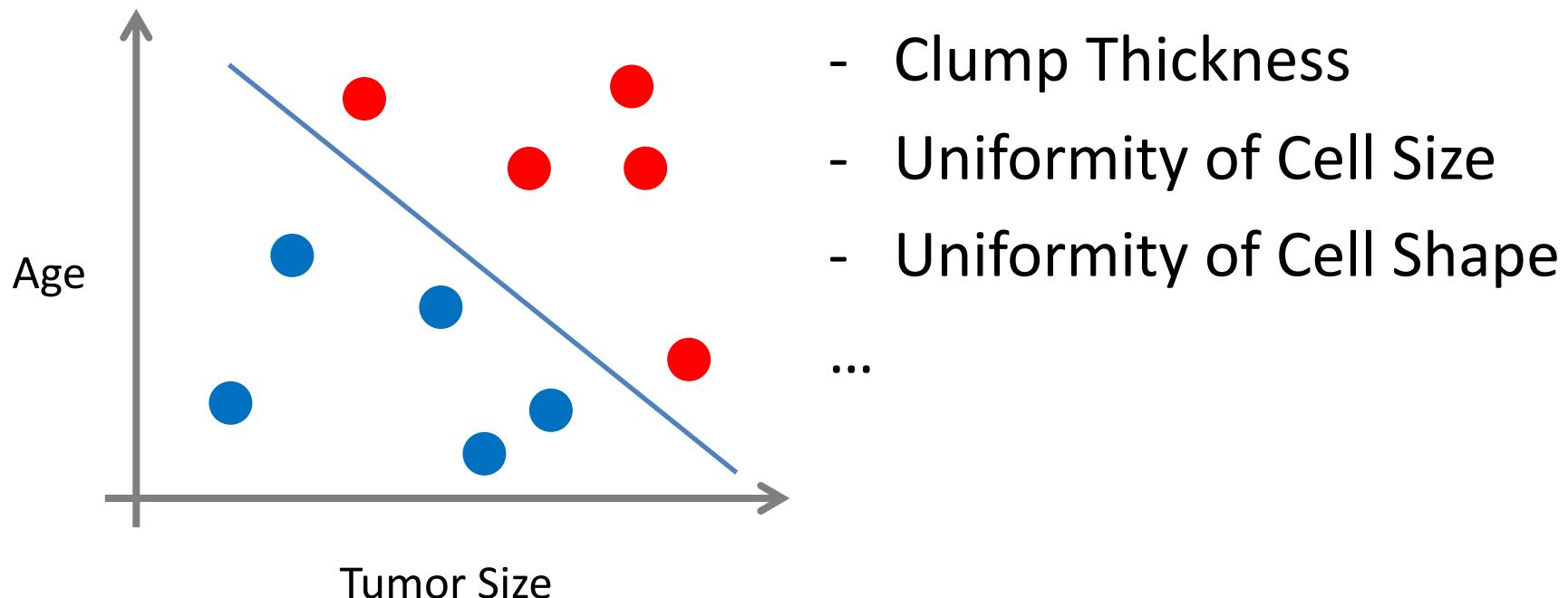
Task: given tumor size  $x$ , is it benign or malignant?



# Supervised Learning

$x$  can be multi-dimensional:

- each dimension corresponds to an attribute



- Clump Thickness
- Uniformity of Cell Size
- Uniformity of Cell Shape

...

# Supervised learning: Vision

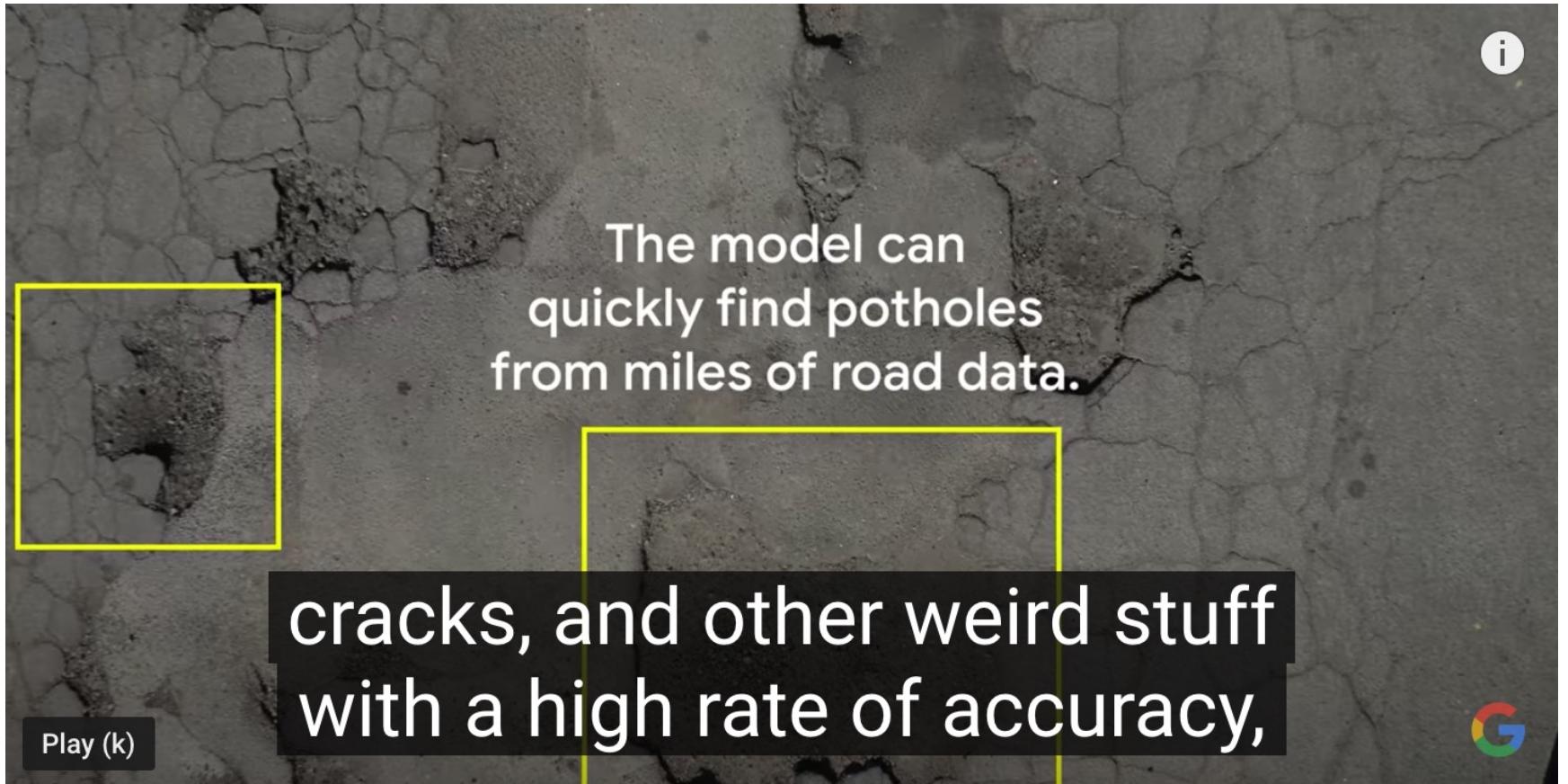
Image classification:

$x$  = raw pixels of the image

$y$  = the main object



# Example: Pothole detection



# Supervised learning: Natural Language Processing (NLP)

## Machine translation

Google Translate



The screenshot shows the Google Translate interface. The source text "Machine translation is a supervised learning problem" is in English (selected as the source language). The target text "机器翻译是一种有监督的学习问题" is in Chinese (Simplified) (selected as the target language). Below the text, there is a green arrow pointing from left to right, indicating the direction of translation. The input text has a character count of 52/5000. On the right side, there are icons for microphone, speaker, and sharing, along with a "Send feedback" link.

Machine translation is a supervised learning problem

机器翻译是一种有监督的学习问题

Jīqì fānyì shì yī zhǒng yǒu jiāndū de xuéxí wèntí

$x \rightarrow y$

**Note:** this course only covers the basics of NLP—take 5370 next semester for more detail!

# Today's learning objectives

Students will be able to:

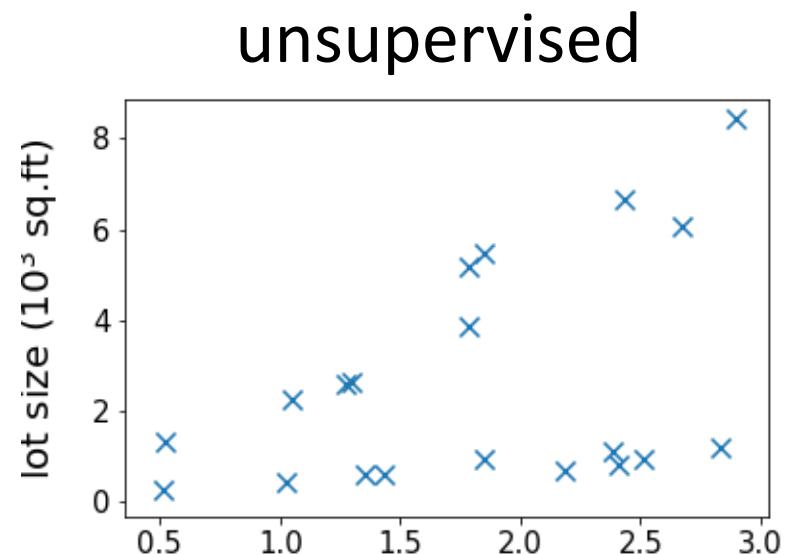
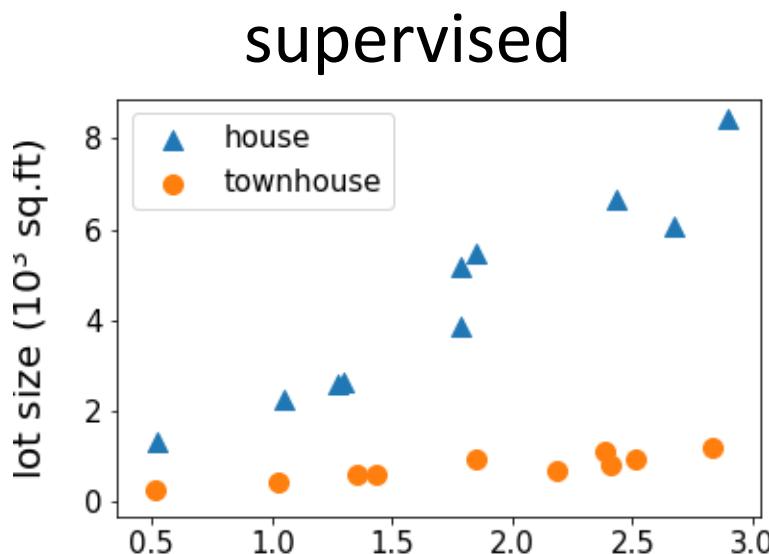
- ✓ Define Machine Learning.
- ✓ List three types of Machine Learning algorithms.
- ✓ Give examples of supervised learning tasks.
  - Name an unsupervised learning technique.
  - Explain how reinforcement learning works.
  - Describe the learning task with math.

# Unsupervised learning

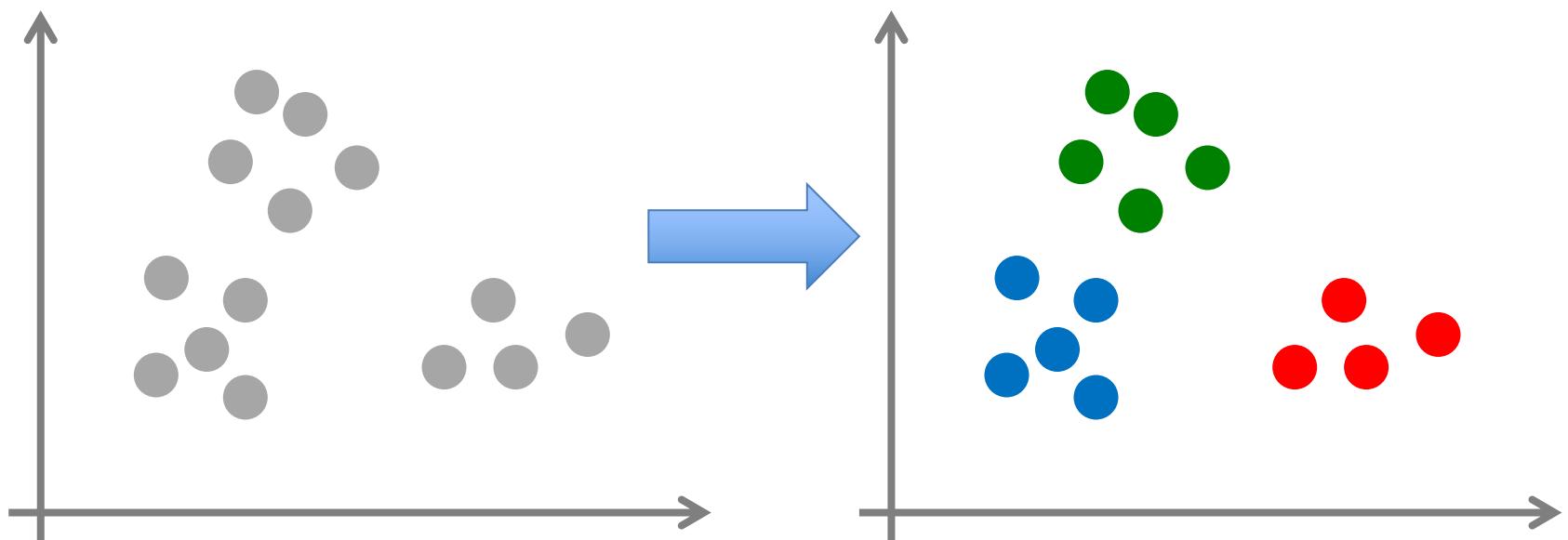
Dataset contains **no labels**:  $x^{(1)}, \dots, x^{(n)}$

**Goal:** find hidden structures in the data

- e.g., clustering

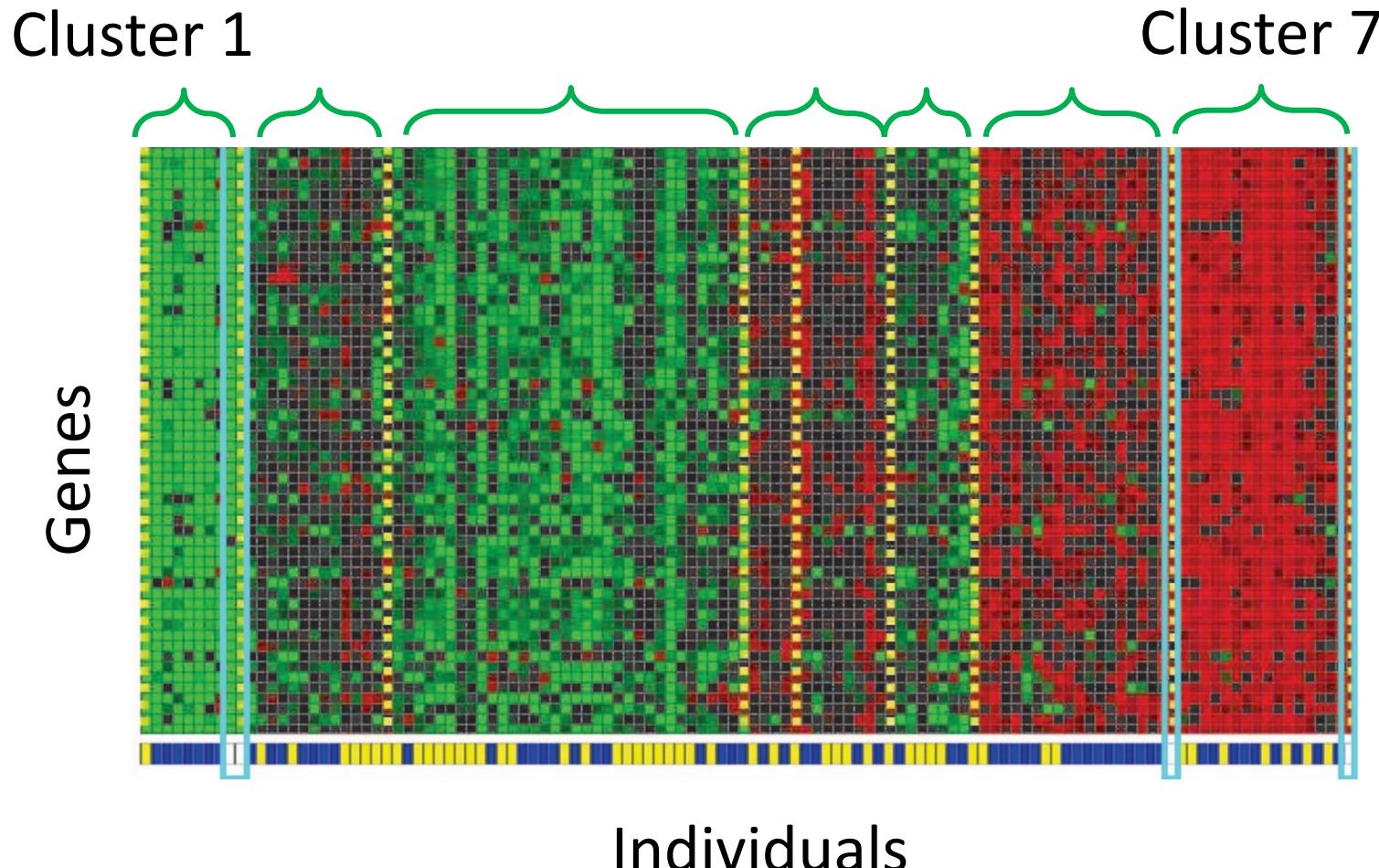


# Clustering



# Clustering genes

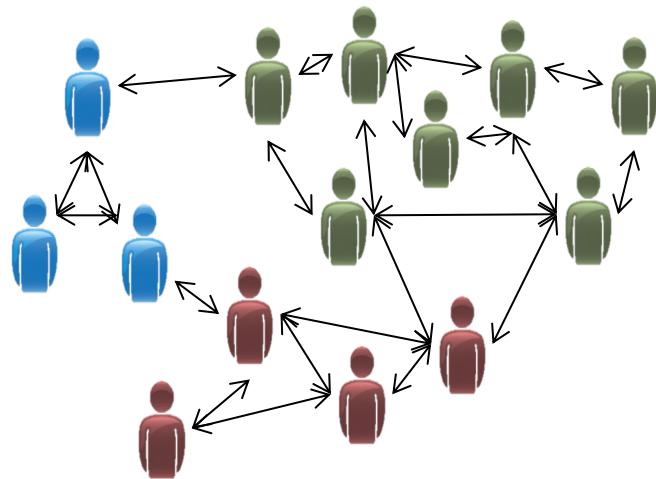
Genomics: group individuals by genetic similarity



# Clustering customer segments



Market segmentation



Social network analysis

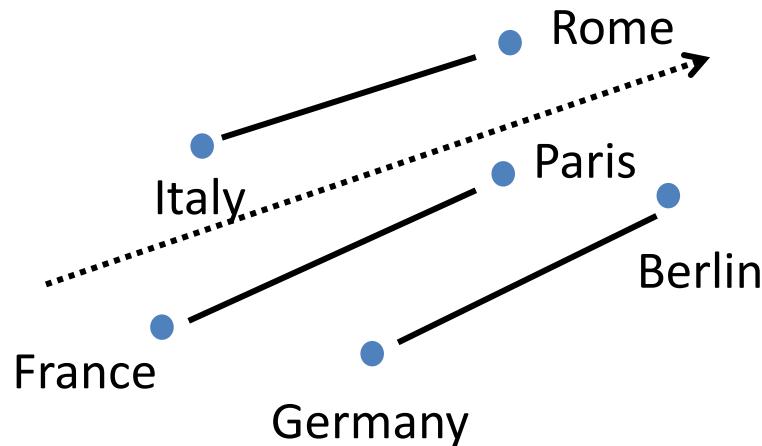
# Unsupervised NLP: Word embeddings

Represent words by vectors

- word       $\xrightarrow{\text{encode}}$       vector
- relation     $\xrightarrow{\text{encode}}$       direction



Unlabeled dataset



Word2vec [Mikolov et al'13]  
GloVe [Pennington et al'14]

# Today's learning objectives

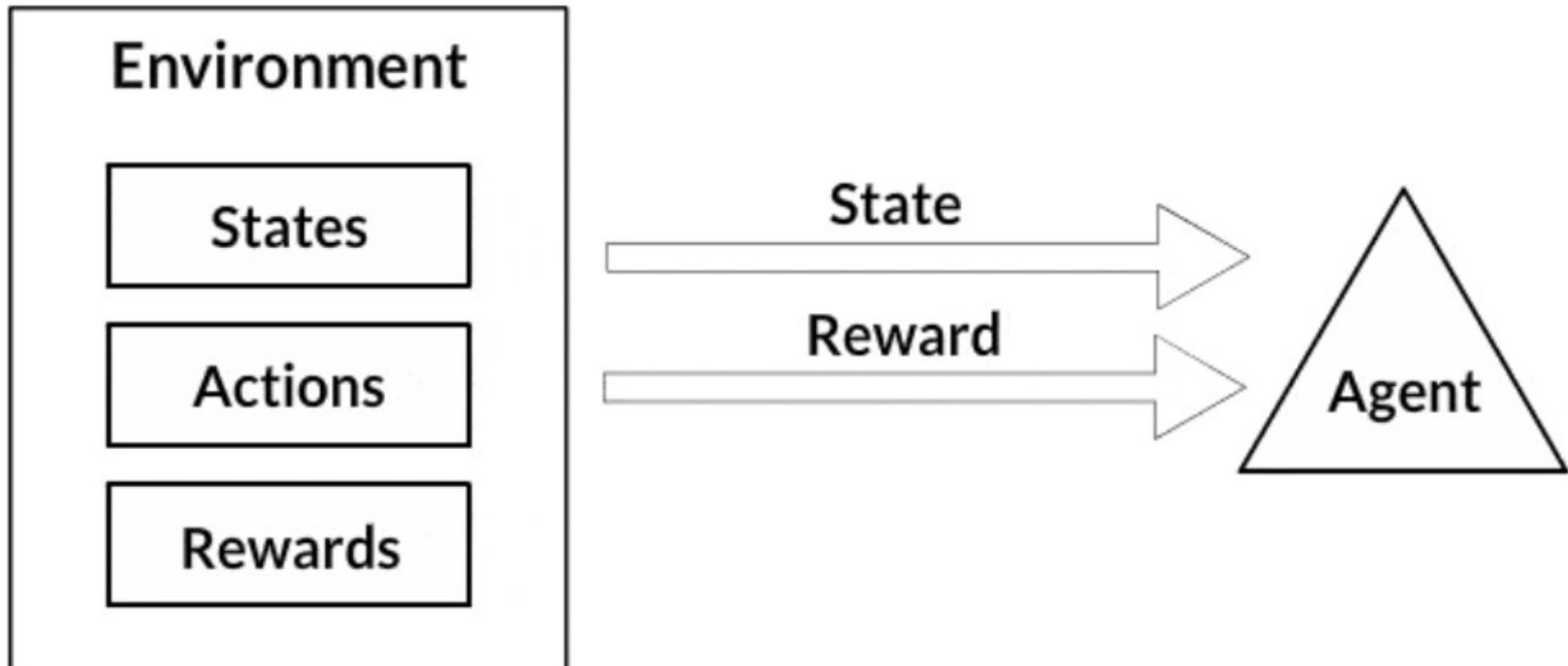
Students will be able to:

- ✓ Define Machine Learning.
- ✓ List three types of Machine Learning algorithms.
- ✓ Give examples of supervised learning tasks.
- ✓ Name an unsupervised learning technique.
- Explain how reinforcement learning works.
- Describe the learning task with math.

# Reinforcement learning

Collects data interactively

**Goal:** maximize the reward



# Reinforcement learning

Backgammon



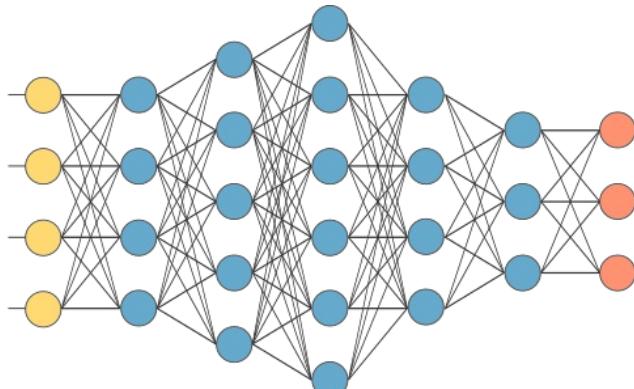
Given sequences of moves and whether or not the player won at the end, learn to make good moves

# Reinforcement learning



During training, the car explores with random actions.

# Brief Peek at Topics

- **Supervised learning**
  - Perceptron
  - Decision trees
  - K-nearest neighbors
  - Linear regression
  - Logistic regression
  - Support vector machines
  - Ensemble methods
- **Unsupervised learning**
  - Dimensionality reduction
  - Clustering
- **Reinforcement learning**
  - None: covered in-depth in AI (CMSI 3300)
- **Deep learning**
- **Bias and fairness in ML**

# Today's learning objectives

Students will be able to:

- ✓ Define Machine Learning.
- ✓ List three types of Machine Learning algorithms.
- ✓ Give examples of supervised learning tasks.
- ✓ Name an unsupervised learning technique.
- ✓ Explain how reinforcement learning works.
  - Describe the learning task with math.

# Framing a Learning Problem

Learning Goals

- List stages of ML pipeline
- Name some key issues in ML

# Defining the learning task

Improve on task T, with respect to  
performance metric P, based on experience E

T: Playing checkers

P: Percentage of games won against an arbitrary opponent

E: Playing practice games against itself

T: Recognizing hand-written words

P: Percentage of words correctly classified

E: Database of human-labeled images of handwritten words

T: Categorize email messages as spam or legitimate

P: Percentage of email messages correctly classified

E: Database of emails, some with human-given labels

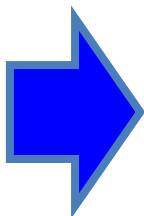
*Q: Which type of ML algorithm does each task use?*

# Representing examples

What is an example?

How is it represented?

examples



features

$\text{feat}_1, \text{feat}_2, \text{feat}_3, \text{feat}_4, \dots$

color, shape, leaf, weight, ...

red, round, leaf, 3oz, ...

green, round, no leaf, 4oz, ...

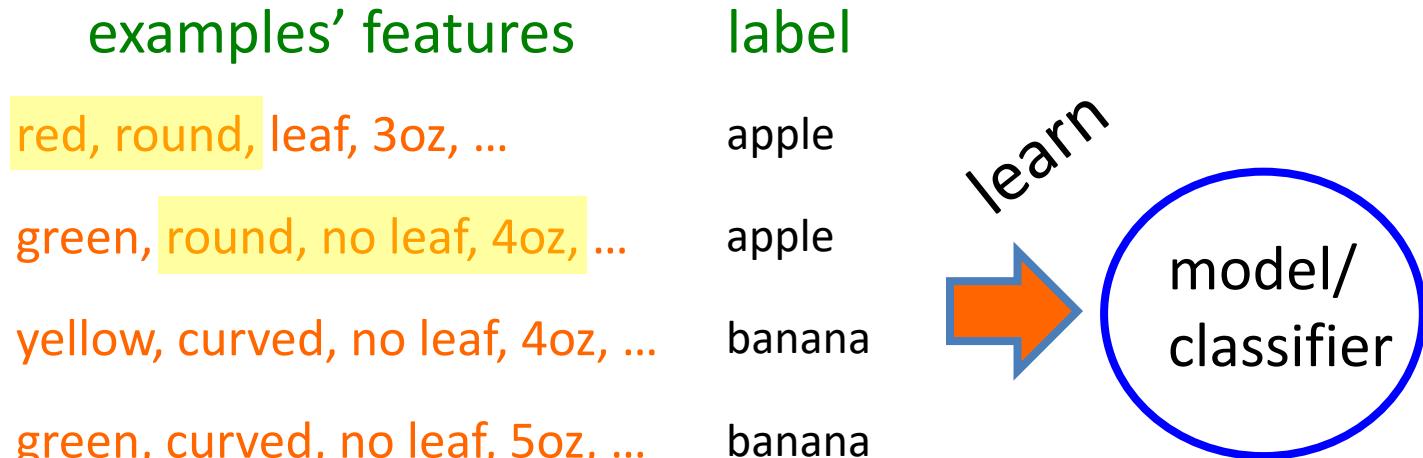
yellow, curved, no leaf, 4oz, ...

green, curved, no leaf, 5oz, ...

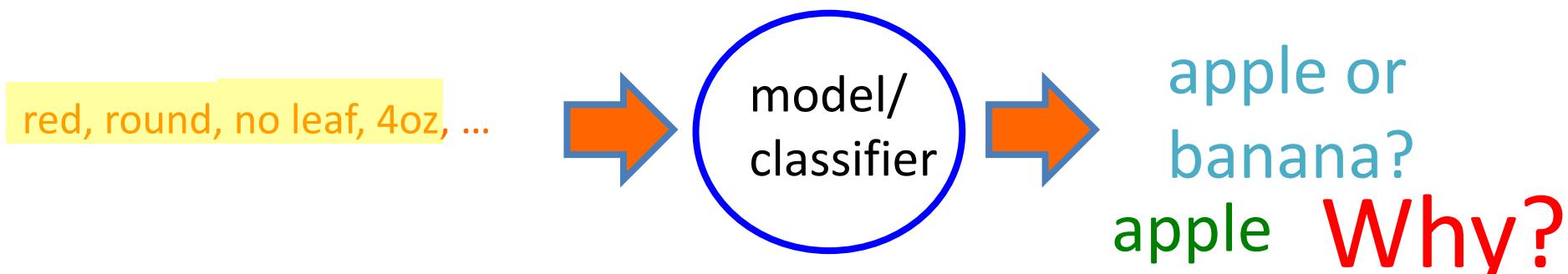
How our algorithms actually “view” the data

Features are the questions we can ask about the examples

# Learning system



During learning/training/induction, learn a model of what distinguishes apples and bananas *based on the features*



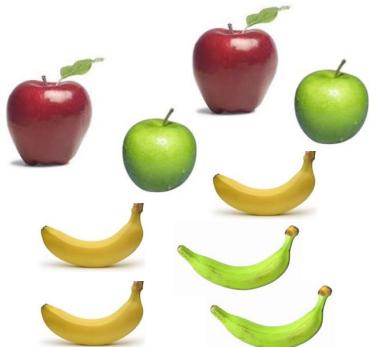
The model can then classify a new example *based on the features*

# Learning system

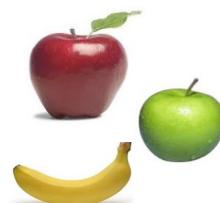
Learning is **generalizing** from training data.

What does this **assume** about train and test sets?

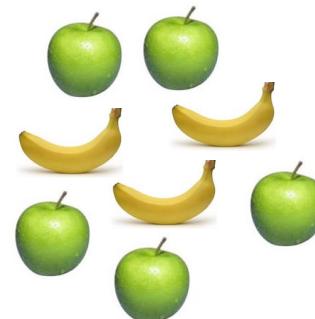
Training data



Test set



Training data



Test set



Not always the case, but  
we'll often assume it is!

# More technically...

We will use the ***probabilistic model*** of learning.

There is some (unknown) probability distribution over example/label pairs called the ***data generating distribution***.

**Both** the training data **and** the test set are generated based on this distribution.

- we call this i.i.d. (independent and *identically distributed*)

# ML as function approximation

## Problem Setting:

- Set of possible instances  $X$
- Set of possible labels  $Y$
- Unknown target function  $f: X \rightarrow Y$
- Set of function hypotheses  $H = \{h \mid h: X \rightarrow Y\}$

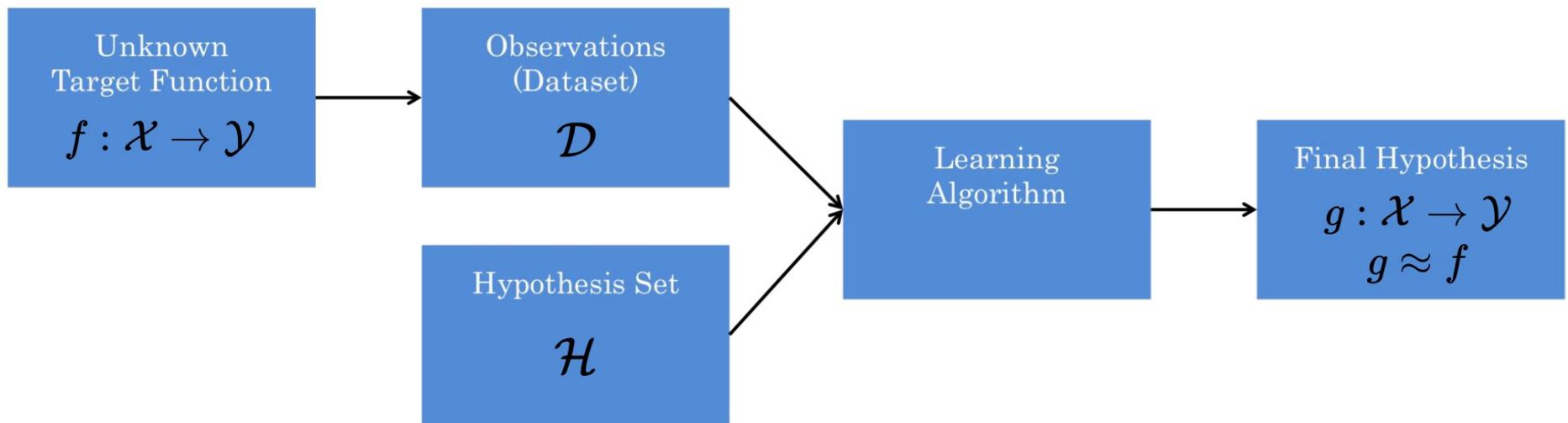
**Input:** Training examples of unknown target function  $f$

$$\left\{(x^{(i)}, y^{(i)})\right\}_{i=1}^n = \left\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\right\}$$

(superscript denotes example number)

**Output:** Hypothesis  $h \rightarrow H$  that best approximates  $f$

# The learning problem



# Stages of Machine Learning

**Given:** labeled training data  $X, Y = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$

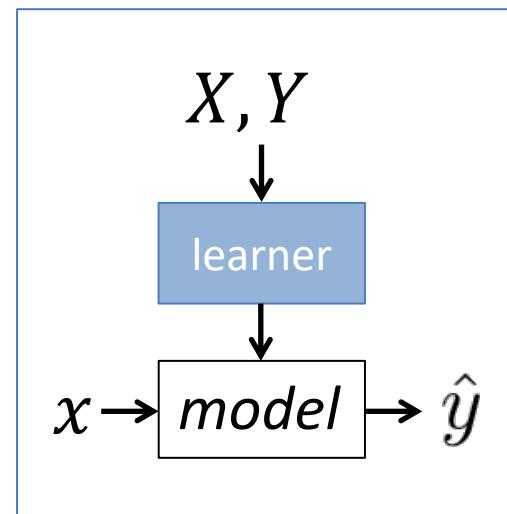
**Train model:**

$model \leftarrow learner.\text{train}(X, Y)$

**Apply model to new data:**

Given: new unlabeled instance  $x$

$\hat{y} \leftarrow model.\text{predict}(x)$



# Example: Regression

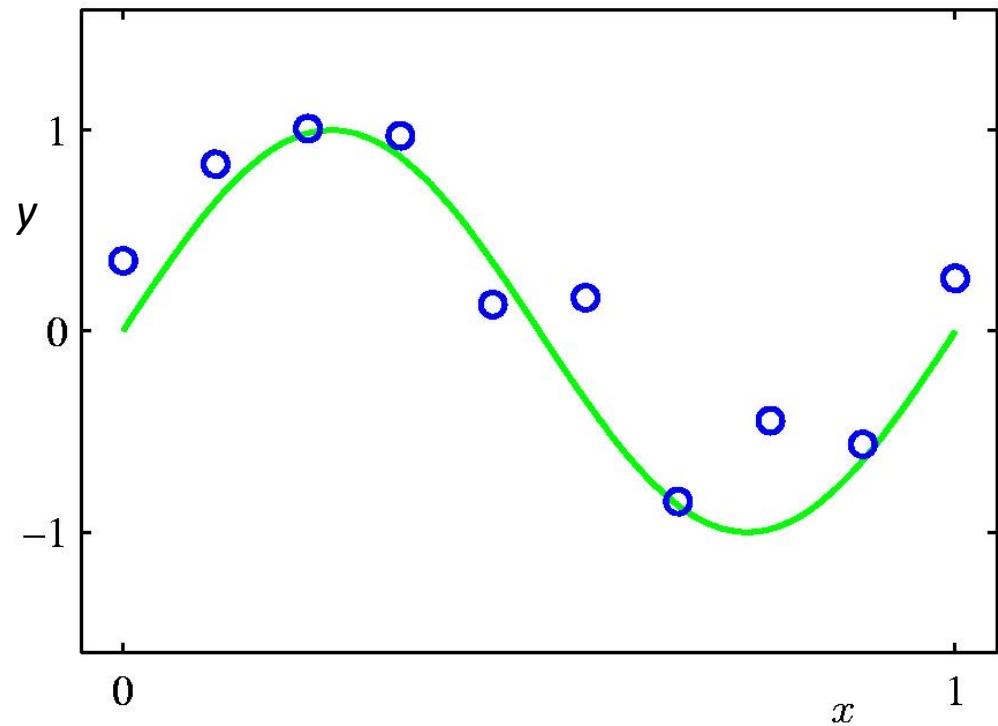
Consider regression:

- $f: X \rightarrow Y$
- $x, y \in \mathbb{R}$

**Question 1:** How should we pick the hypothesis space  $H$ ?

**Question 2:** How do we find the best  $h$  in this space?

Dataset: 10 points generated from sine function with noise

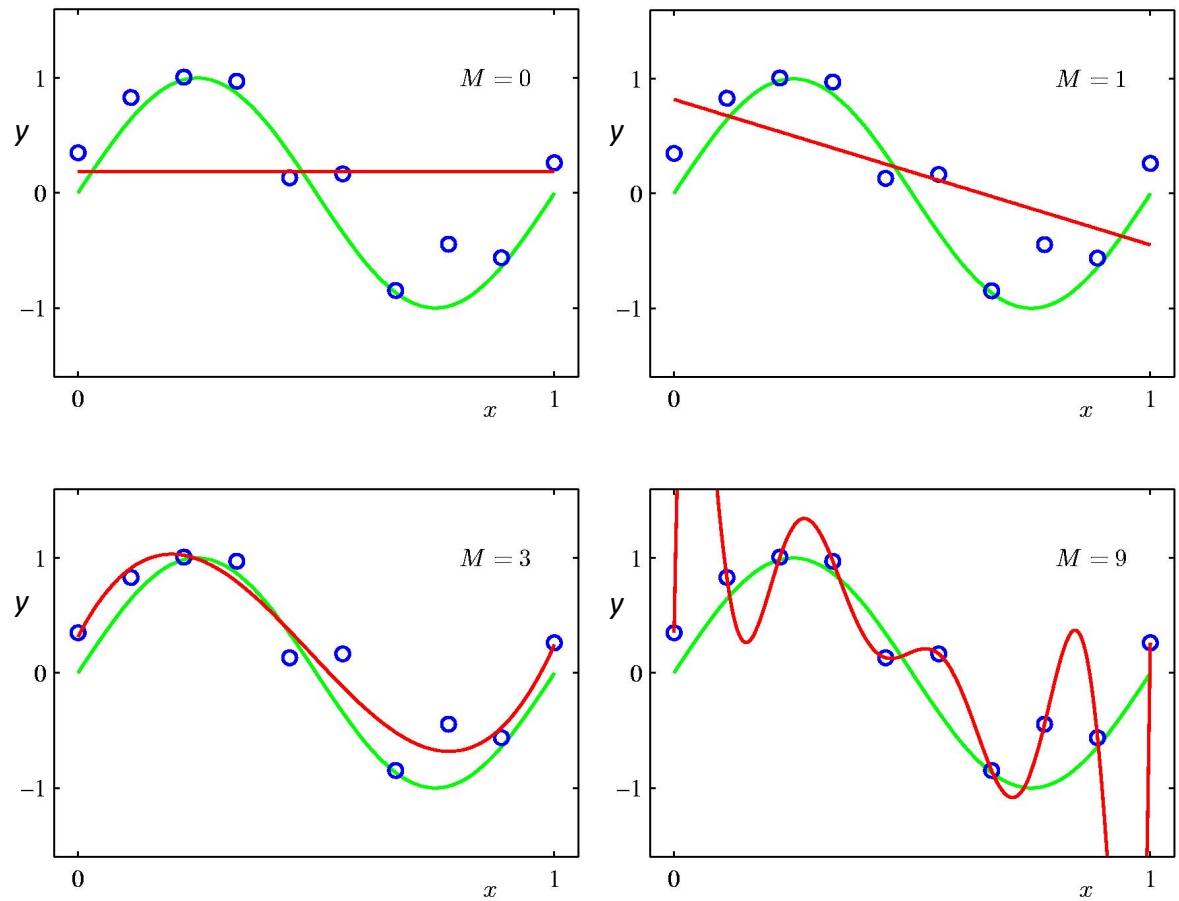


# Hypothesis space: Degree- $M$ polynomials

Infinitely many hypotheses

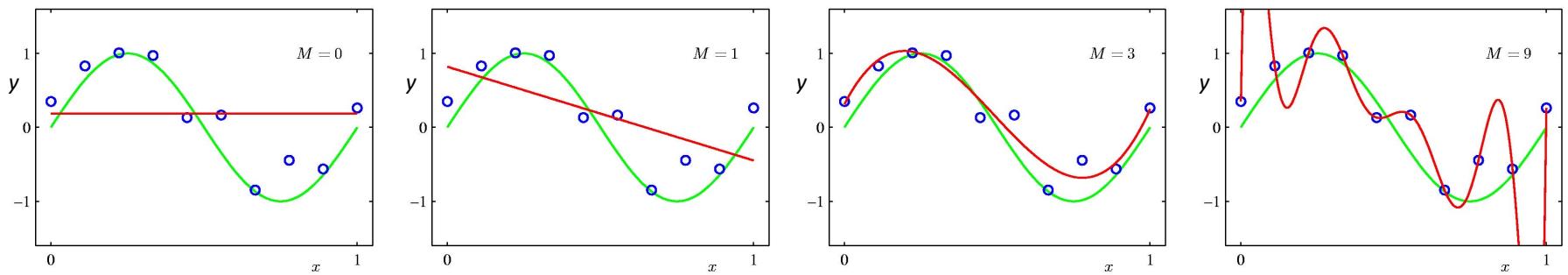
$M < 9$  inconsistent with dataset, but  $M \geq 9$  is consistent

Which one is **best**?



$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

# Hypothesis space: Degree-M polynomials



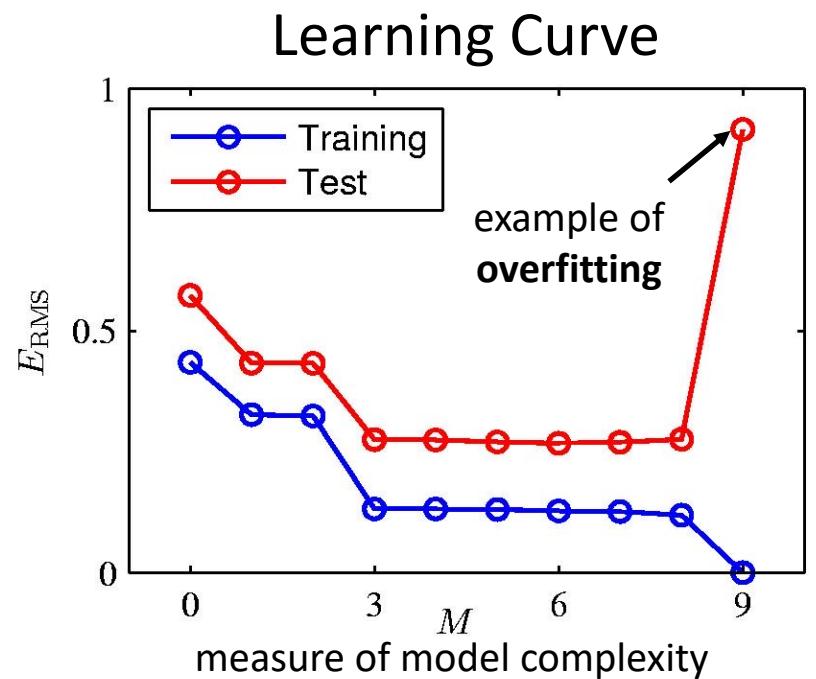
We measure error using a loss function  $L(y, \hat{y})$ .

For regression, common choice is squared loss:

$$L\left(y^{(i)}, h(x^{(i)})\right) = \left(y^{(i)} - h(x^{(i)})\right)^2$$

*Empirical loss* of function  $h$  applied to training data:

$$\frac{1}{n} \sum_{i=1}^n L\left(y^{(i)}, h(x^{(i)})\right) = \frac{1}{n} \sum_{i=1}^n \left(y^{(i)} - h(x^{(i)})\right)^2$$



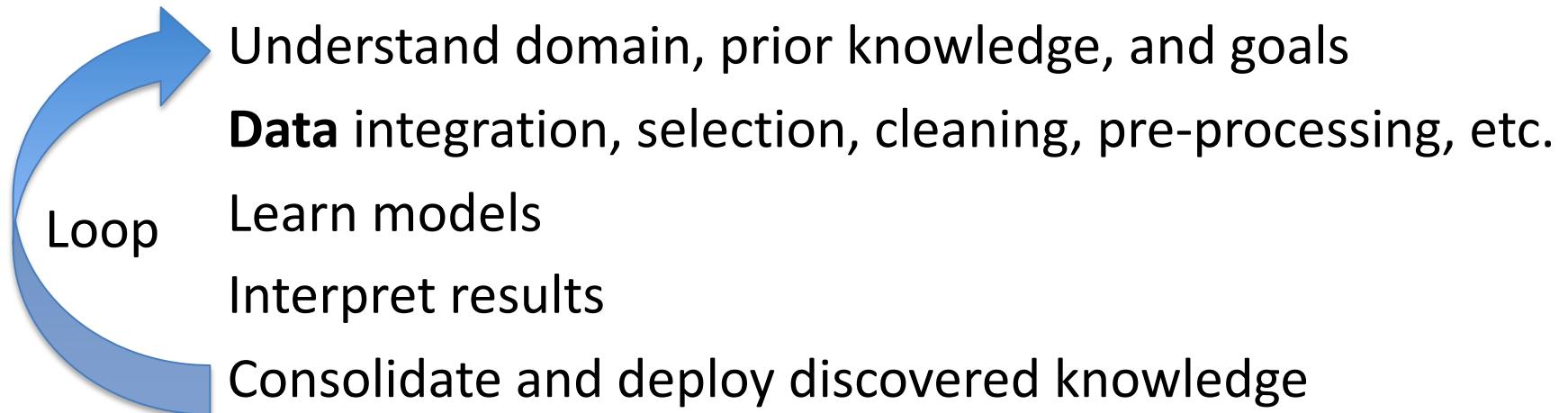
# Key issues in Machine Learning

**Representation:** How do we choose a hypothesis space?

**Optimization:** How do we find the *best* hypothesis?

**Evaluation:** How can we gauge the accuracy of a hypothesis on unseen testing data?

# Machine Learning in practice



# Today's learning objectives

Students will be able to:

- ✓ Define Machine Learning.
- ✓ List three types of Machine Learning algorithms.
- ✓ Give examples of supervised learning tasks.
- ✓ Name an unsupervised learning technique.
- ✓ Explain how reinforcement learning works.
- ✓ Describe the learning task with math.

# Exercise

Pick 2 applications and frame them as a learning problem with input (features)  $X$  and predicted output  $y$ .

- Web search
- Finance
- Healthcare
- E-commerce
- Natural language processing
- Robotics
- Social networks
- [Your favorite area]

# Quote of the day

“

Attitude is a little thing  
that makes a big difference.

WINSTON CHURCHILL

”

# Socrative: Optional feedback

<http://socrative.com>

Go to student login (it's anonymous!)



Student Login

Room Name

NYL5DRJJ

JOIN

