

#### 10-301/10-601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

# Machine Learning as Function Approximation

Matt Gormley & Henry Chai Lecture 2 Jan. 15, 2025

#### Reminders

- Background Test
  - Fri, Sep 1, in-class
- Homework 1: Background
  - Out: Mon, Jan 13
  - Due: Wed, Jan 22 at 11:59pm
  - Two parts:
    - 1. written part to Gradescope
    - 2. programming part to Gradescope
  - unique policies for this assignment:
    - 1. unlimited submissions for programming (i.e. keep submitting until you get 100%)
    - 2. we will grant (essentially) any and all extension requests

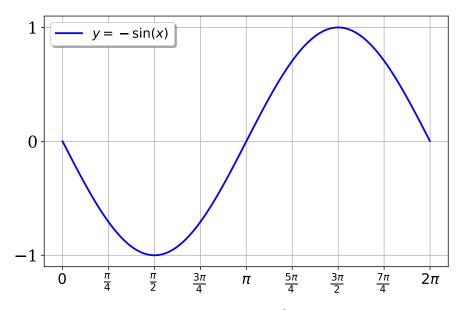
## Big Ideas

- 1. How to formalize a learning problem
- 2. How to learn an expert system (i.e. Decision Tree)
- 3. Importance of inductive bias for generalization
- 4. Overfitting

#### **FUNCTION APPROXIMATION**

## **Function Approximation**

**Quiz:** Implement a simple function which returns  $-\sin(x)$ .



#### A few constraints are imposed:

- 1. You can't call any other trigonometric functions
- You can call an existing implementation of sin(x) a few times (e.g. 100) to test your solution
- You only need to evaluate it for x in [0, 2\*pi]

#### **SUPERVISED MACHINE LEARNING**

## Medical Diagnosis

- Setting:
  - Doctor must decide whether or not patient is sick
  - Looks at attributes of a patient to make a medical diagnosis
  - (Prescribes treatment if diagnosis is positive)
- Key problem area for Machine Learning
- Potential to reshape health care

## Medical Diagnosis

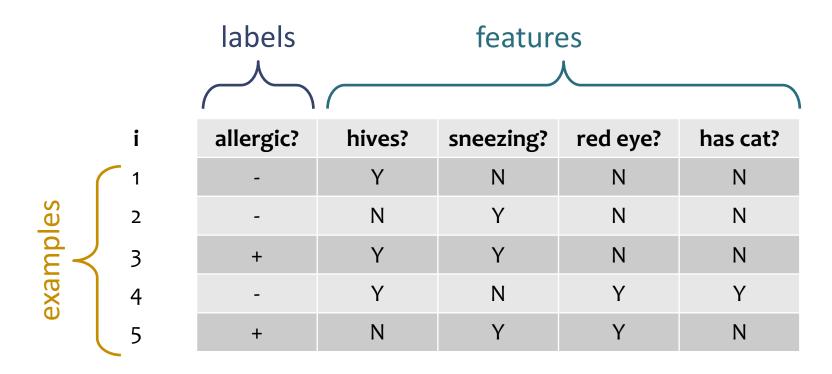
**Interview Transcript** 

**Date:** Jan. 15, 2023

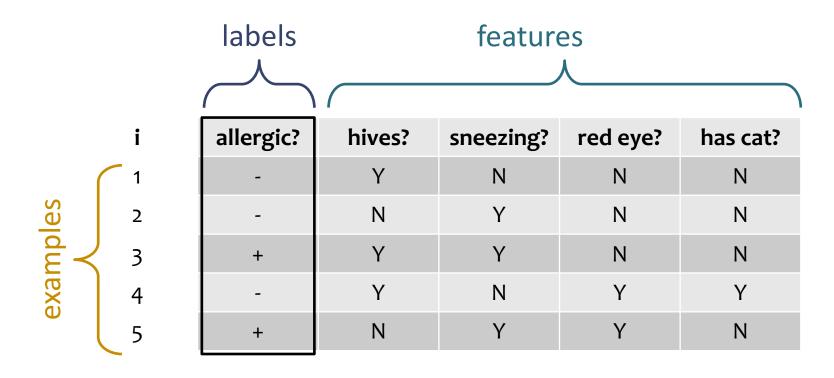
**Parties:** Matt Gormley and Doctor S.

**Topic:** Medical decision making

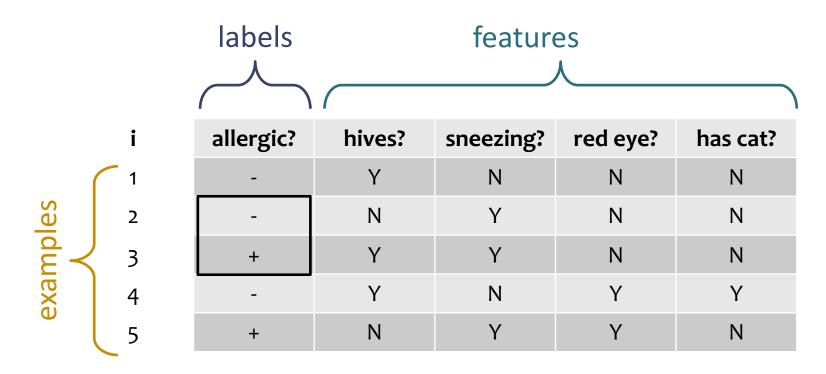
As a (supervised) binary classification task



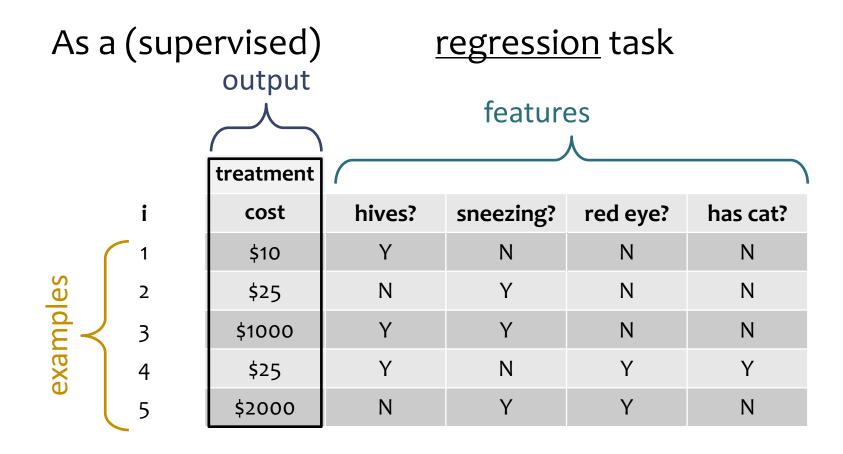
As a (<u>supervised</u>) binary classification task



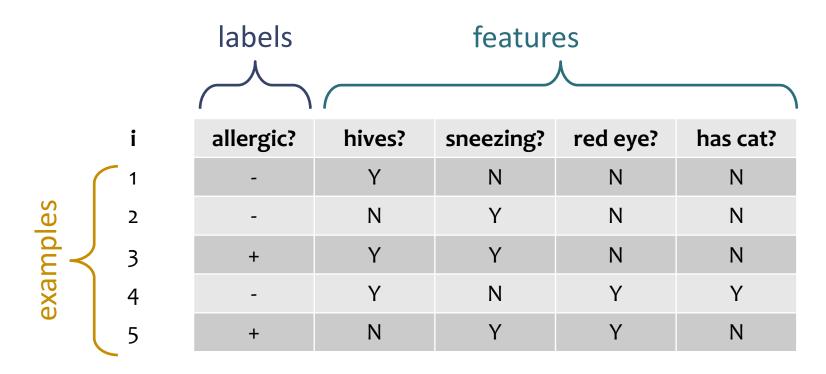
As a (supervised) binary classification task



As a (supervised) classification task labels features allergy sneezing? red eye? hives? has cat? Υ Ν Ν Ν none examples 2 Ν Ν Ν none Ν 3 Υ Ν dust Ν Υ 4 Υ Y none mold Ν Υ Υ Ν



As a (supervised) binary classification task



Doctor diagnoses the patient as sick or not  $y \in \{+, -\}$  based on attributes of the patient  $x_1, x_2, ..., x_M$ 

	у	$X_1$	X <sub>2</sub>	$X_3$	<b>X</b> <sub>4</sub>
i	allergic?	hives?	sneezing?	red eye?	has cat?
1	-	Υ	N	N	N

Doctor diagnoses the patient as sick or not  $y \in \{+, -\}$  based on attributes of the patient  $x_1, x_2, ..., x_M$ 

	у	$X_1$	$X_2$	$X_3$	<b>X</b> <sub>4</sub>
i	allergic?	hives?	sneezing?	red eye?	has cat?
1	-	Y	N	N	N
2	-	N	Υ	N	N
3	+	Υ	Υ	N	N
4	-	Υ	N	Υ	Υ
5	+	N	Y	Y	N

Doctor diagnoses the patient as sick or not  $y \in \{+, -\}$  based on attributes of the patient  $x_1, x_2, ..., x_M$ 

	у	$X_1$	$X_2$	$X_3$	$X_4$
i	allergic?	hives?	sneezing?	red eye?	has cat?
1	y <sup>(1)</sup> -	X <sub>1</sub> <sup>(1)</sup> Y	$x_2^{(1)} N$	x <sub>3</sub> <sup>(1)</sup> N	x <sub>4</sub> <sup>(1)</sup> N
2	y <sup>(2)</sup> -	$X_1^{(2)} N$	$X_2^{(2)} Y$	$X_3^{(2)} N$	$X_4^{(2)} N$
3	y <sup>(3)</sup> +	X <sub>1</sub> <sup>(3)</sup> Y	X <sub>2</sub> <sup>(3)</sup> Y	x <sub>3</sub> <sup>(3)</sup> N	x <sub>4</sub> <sup>(3)</sup> N
4	y <sup>(4)</sup> -	X <sub>1</sub> <sup>(4)</sup> Y	$X_2^{(4)} N$	x <sub>3</sub> <sup>(4)</sup> Y	x <sub>4</sub> <sup>(4)</sup> Y
5	y <sup>(5)</sup> +	X <sub>1</sub> <sup>(5)</sup> N	X <sub>2</sub> <sup>(5)</sup> Y	x <sub>3</sub> <sup>(5)</sup> Y	x <sub>4</sub> <sup>(5)</sup> N

Doctor diagnoses the patient as sick or not  $y \in \{+, -\}$  based on attributes of the patient  $x_1, x_2, ..., x_M$ 

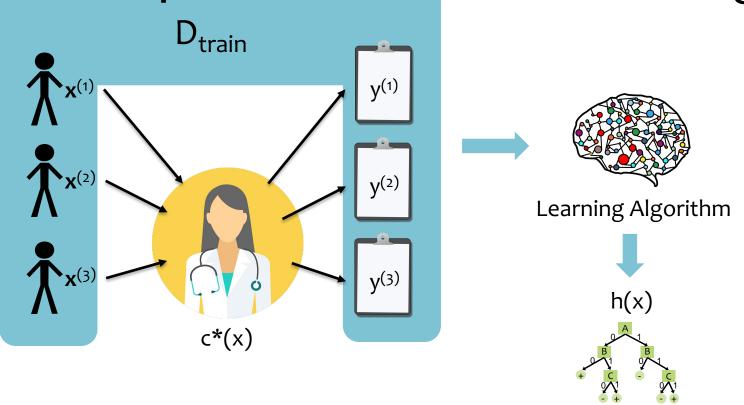
	У	X <sub>1</sub>	$X_2$	$X_3$	<b>X</b> <sub>4</sub>	
i	allergic?	hives?	sneezing?	red eye?	has cat?	
1	y <sup>(1)</sup> -	X <sub>1</sub> <sup>(1)</sup> Y	$X_2^{(1)} N$	x <sub>3</sub> <sup>(1)</sup> N	x <sub>4</sub> <sup>(1)</sup> N	X <sup>(1)</sup>
2	y <sup>(2)</sup> -	$X_1^{(2)} N$	$X_2^{(2)} Y$	$x_3^{(2)} N$	x <sub>4</sub> <sup>(2)</sup> N	$X^{(2)}$
3	y <sup>(3)</sup> +	Χ <sub>1</sub> <sup>(3)</sup> Υ	X <sub>2</sub> <sup>(3)</sup> Y	x <sub>3</sub> <sup>(3)</sup> N	x <sub>4</sub> <sup>(3)</sup> N	<b>X</b> (3)
4	y <sup>(4)</sup> -	X <sub>1</sub> <sup>(4)</sup> Y	x <sub>2</sub> <sup>(4)</sup> N	x <sub>3</sub> <sup>(4)</sup> Y	x <sub>4</sub> <sup>(4)</sup> Y	<b>X</b> <sup>(4)</sup>
5	y <sup>(5)</sup> +	X <sub>1</sub> <sup>(5)</sup> <b>N</b>	X <sub>2</sub> <sup>(5)</sup> Y	<b>x</b> <sub>3</sub> <sup>(5)</sup> <b>Y</b>	x <sub>4</sub> <sup>(5)</sup> N	<b>X</b> (5)

N = 5 training examples

M = 4 attributes

ML as Function Approximation

## Supervised Machine Learning



Doctor diagnoses the patient as sick or not  $y \in \{+, -\}$  based on attributes of the patient  $x_1, x_1, ..., x_M$ 

	y	X <sub>1</sub>	$X_2$	$X_3$	$X_4$	
i	allergic? <sub>C</sub>	hives?	sneezing?	red eye?	has cat?	
1	y <sup>(1)</sup> -	χ <sub>1</sub> <sup>(1)</sup> Υ	X <sub>2</sub> <sup>(1)</sup> N	x <sub>3</sub> <sup>(1)</sup> N	x <sub>4</sub> <sup>(1)</sup> N	X <sup>(1)</sup>
2	y <sup>(2)</sup>	x <sub>1</sub> <sup>(2)</sup> N	$X_2^{(2)} Y$	$x_3^{(2)} N$	x <sub>4</sub> <sup>(2)</sup> N	X <sup>(2)</sup>
3	y(3) #	χ <sub>1</sub> <sup>(3)</sup> Υ	X <sub>2</sub> <sup>(3)</sup> Y	x <sub>3</sub> <sup>(3)</sup> N	x <sub>4</sub> <sup>(3)</sup> N	<b>X</b> <sup>(3)</sup>
4	y(4)	X <sub>1</sub> <sup>(4)</sup> Y	x <sub>2</sub> <sup>(4)</sup> N	<b>X</b> <sub>3</sub> <sup>(4)</sup> <b>Y</b>	x <sub>4</sub> <sup>(4)</sup> Y	X <sup>(4)</sup>
5	y(5) 4	X <sub>1</sub> <sup>(5)</sup> <b>N</b>	X <sub>2</sub> <sup>(5)</sup> Y	x <sub>3</sub> <sup>(5)</sup> Y	x <sub>4</sub> <sup>(5)</sup> <b>N</b>	X <sup>(5)</sup>

N = 5 training examples M = 4 attributes

Example hypothesis function:
$$h(x) = \begin{cases} + & \text{if sneezing} = Y \\ - & \text{otherwise} \end{cases}$$

## Supervised Machine Learning

#### Problem Setting

- Set of possible inputs,  $\mathbf{x} \in \mathcal{X}$  (all possible patients)
- Set of possible outputs,  $y \in \mathcal{Y}$  (all possible diagnoses)
- Exists an unknown target function,  $c^* : \mathcal{X} \rightarrow \mathcal{Y}$  (the doctor's brain)
- Set,  $\mathcal{H}$ , of candidate hypothesis functions,  $h: \mathcal{X} \rightarrow \mathcal{Y}$  (all possible decision trees)
- Learner is given N training examples  $D = \{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), ..., (\mathbf{x}^{(N)}, \mathbf{y}^{(N)})\}$  where  $\mathbf{y}^{(i)} = \mathbf{c}^*(\mathbf{x}^{(i)})$  (history of patients and their diagnoses)
- Learner produces a hypothesis function,  $\hat{y} = h(x)$ , that best approximates unknown target function  $y = c^*(x)$  on the training data

## Supervised Machine Learning

#### Problem Setting

- Set of possible inputs,  $x \in \mathcal{X}$  (all possible patients)
- Set of possible outputs,  $y \in \mathcal{Y}$  (all possible diagnoses)
- Exists an unknown tar function,  $c^*: \mathcal{X} \rightarrow \mathcal{Y}$  (the doctor's brain)
- Set,  $\mathcal{H}$ , of candidate hypoth (all possible decision trees) consider:
- Learner is given N training
   D = {(x<sup>(1)</sup>, y<sup>(1)</sup>), (x<sup>(2)</sup>, y<sup>(2)</sup>), ...,
   where y<sup>(i)</sup> = c\*(x<sup>(i)</sup>)
   (history of patients and the
- Learner produces a hypoth approximates unknown tar

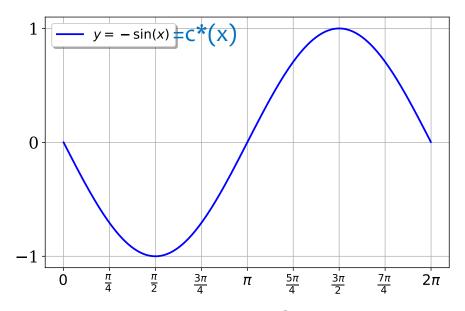
Two important settings we'll consider:

- Classification: the possible outputs are discrete
- 2. Regression: the possible outputs are real-valued

data

## **Function Approximation**

**Quiz:** Implement a simple function which returns  $-\sin(x)$ .



#### A few constraints are imposed:

- 1. You can't call any other trigonometric functions
- You can call an existing implementation of sin(x) a few times (e.g. 100) to test your solution
- You only need to evaluate it for x in [0, 2\*pi]

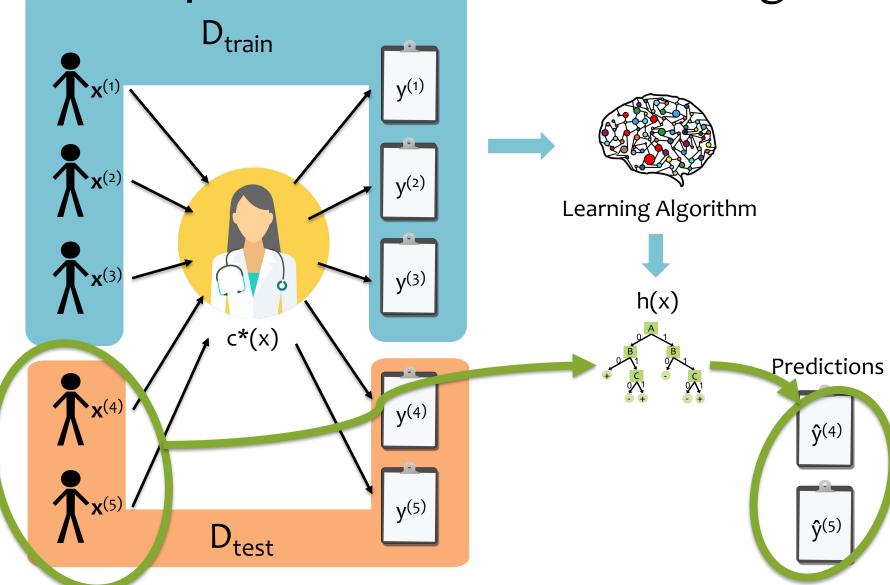
## Supervised Machine Learning

#### Problem Setting

- Set of possible inputs,  $x \in \mathcal{X}$  (all values in [0, 2\*pi])
- Set of possible outputs,  $y \in \mathcal{Y}$  (all values in [-1,1])
- Exists an unknown target function,  $c^*: \mathcal{X} \rightarrow \mathcal{Y}$  ( $c^*(x) = \sin(x)$ )
- Set,  $\mathcal{H}$ , of candidate hypothesis functions,  $h: \mathcal{X} \rightarrow \mathcal{Y}$  (all possible piecewise linear functions)
- Learner is given N training examples  $D = \{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), ..., (\mathbf{x}^{(N)}, \mathbf{y}^{(N)})\}$  where  $\mathbf{y}^{(i)} = \mathbf{c}^*(\mathbf{x}^{(i)})$  (true values of sin(x) for a few random x's)
- Learner produces a hypothesis function,  $\hat{y} = h(x)$ , that best approximates unknown target function  $y = c^*(x)$  on the training data

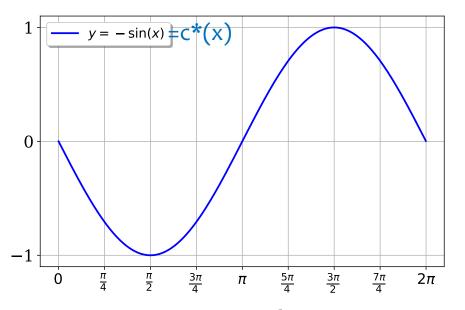
## **EVALUATION OF MACHINE LEARNING ALGORITHM**

## Supervised Machine Learning



## **Function Approximation**

**Quiz:** Implement a simple function which returns  $-\sin(x)$ .





How well does h(x) approximate c\*(x)?

#### A few constraints are imposed:

- 1. You can't call any other trigonometric functions
- 2. You can call an existing implementation of sin(x) a few times (e.g. 100) to test your solution
- You only need to evaluate it for x in [0, 2\*pi]

## Evaluation of ML Algorithms

- **Definition:** loss function,  $\ell: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ 
  - o Defines how "bad" predictions,  $\hat{y} = h(x)$ , are compared to the true labels,  $y = c^*(x)$
  - Common choices:
    - 1. Squared loss (for regression):  $\ell(y, \hat{y}) = (y \hat{y})^2$
    - 2. Binary or 0-1 loss (for classification):  $\ell(y, \hat{y}) = \mathbb{1}(y \neq \hat{y})$

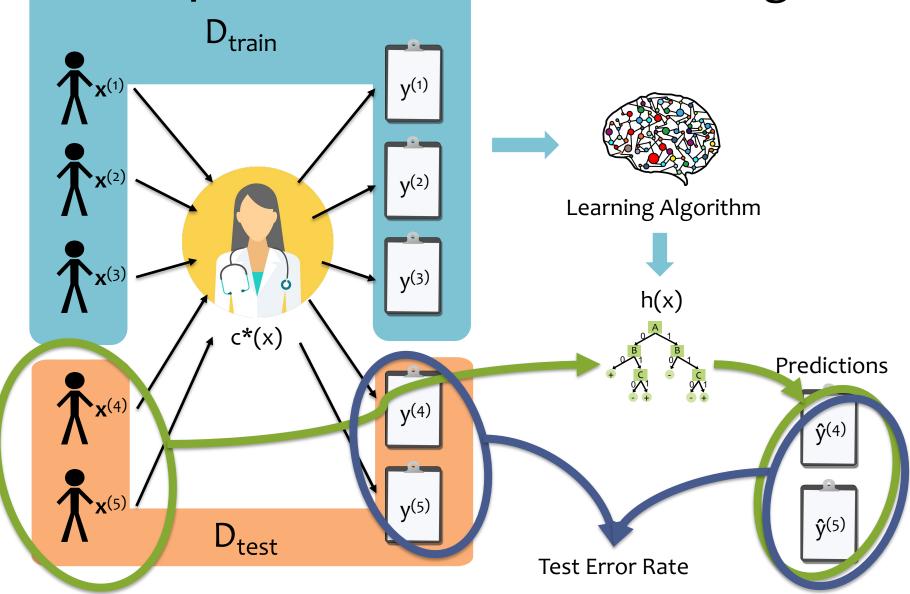
## Evaluation of ML Algorithms

- **Definition:** loss function,  $\ell: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ 
  - o Defines how "bad" predictions,  $\hat{y} = h(x)$ , are compared to the true labels,  $y = c^*(x)$
  - Common choices:
    - 1. Squared loss (for regression):  $\ell(y, \hat{y}) = (y \hat{y})^2$
    - 2. Binary or 0-1 loss (for classification):  $\ell(y, \hat{y}) = \mathbb{1}(y \neq \hat{y})$
- **Definition:** the error rate of a hypothesis h on a dataset  $\mathcal{D}$  is the average 0-1 loss:

$$error(h, \mathcal{D}) = \frac{1}{N} \sum_{n=1}^{N} \mathbb{1}(y^{(n)} \neq \hat{y}^{(n)})$$

- **Definition:** the mean squared error is the average squared loss (more on this later)
- Q: How do we evaluate a machine learning algorithm? A: Check its average loss on a separate test dataset,  $\mathcal{D}_{\text{test}}$ .

## Supervised Machine Learning



#### **Error Rate**

• Consider a hypothesis *h* its...

... error rate over all training data: error(h, D<sub>train</sub>)

... error rate over all test data: error(h,  $D_{test}$ )

... true error over all data: error<sub>true</sub>(h)



This is the quantity we care most about! But, in practice, error<sub>true</sub>(h) is **unknown**.

## Majority Vote Classifier Example

#### **Dataset:**

Output Y, Attributes A and B

Y	Α	В
-	1	0
-	1	0
+	1	0
+	1	0
+	1	1
+	1	1
+	1	1
+	1	1

#### **In-Class Exercise**

What is the **training error** (i.e. error rate on the training data) of the **majority vote classifier** on this dataset?

Choose one of: {0/8, 1/8, 2/8, ..., 8/8}

## LEARNING ALGORITHMS FOR SUPERVISED CLASSIFICATION

Algorithm 1 majority vote: predict the most common label in the training dataset

	у	$X_1$	$X_2$	$X_3$	$X_4$
predictions	allergic?	hives?	sneezing?	red eye?	has cat?
-	-	Y	N	N	N
-	-	N	Υ	N	N
-	+	Υ	Υ	N	N
-	-	Υ	N	Υ	Υ
-	+	N	Y	Y	N

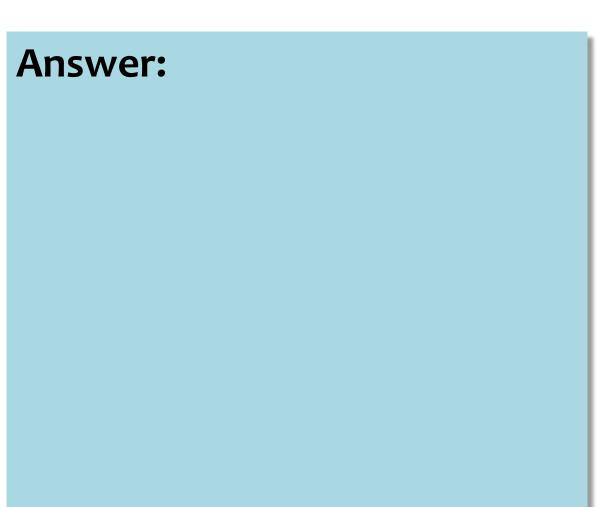
<u>Algorithm 2</u> memorizer: if a set of features exists in the training dataset, predict its corresponding label; otherwise, predict a random label

	у	$X_1$	$X_2$	$x_3$	$X_4$
predictions	allergic?	hives?	sneezing?	red eye?	has cat?
-	-	Y	N	N	N
-	-	N	Υ	N	N
+	+	Υ	Υ	N	N
-	-	Υ	N	Υ	Υ
+	+	N	Y	Υ	N

The memorizer always gets zero training error!

#### **Question:**

If we have 100 features, how many patients does the memorizer need to see to ensure zero test error?



## Algorithm 1: Majority Vote

Pseudocode

## Algorithm 2: Memorizer

Pseudocode

Algorithm 3 decision stump: based on a single feature,  $x_d$ , predict the most common label in the training dataset among all data points that have the same value for  $x_d$ 

	у	$X_1$	$X_2$	$X_3$	$X_4$
predictions	allergic?	hives?	sneezing?	red eye?	has cat?
-	-	Y	N	N	N
+	-	N	Υ	N	N
+	+	Υ	Υ	N	N
-	-	Υ	N	Υ	Υ
+	+	N	Υ	Υ	N

Nonzero training error, but perhaps still better than the memorizer

Example decision stump:  

$$h(x) = \begin{cases} + \text{ if sneezing} = Y \\ - \text{ otherwise} \end{cases}$$

## Algorithm 3: Decision Stump

Pseudocode

Algorithm 3 decision stump: based on a single feature,  $x_d$ , predict the most common label in the training dataset among all data points that have the same value for  $x_d$ 

#### **Questions:**

1. How do we pick which feature to split on?

2. Why stop at one feature?

## Algorithm 4: Decision Tree (preview)

Example

#### Tree to Predict C-Section Risk

Learned from medical records of 1000 women (Sims et al., 2000)

Negative examples are C-sections

```
[833+,167-] .83+ .17-
Fetal_Presentation = 1: [822+,116-] .88+ .12-
| Previous_Csection = 0: [767+,81-] .90+ .10-
  Primiparous = 0: [399+,13-] .97+ .03-
  Primiparous = 1: [368+,68-] .84+ .16-
| \ | \ | Fetal_Distress = 0: [334+,47-] .88+ .12-
| \ | \ | Birth_Weight >= 3349: [133+,36.4-] .78+
  | Fetal_Distress = 1: [34+,21-] .62+ .38-
| Previous_Csection = 1: [55+,35-] .61+ .39-
Fetal_Presentation = 2: [3+,29-] .11+ .89-
Fetal_Presentation = 3: [8+,22-] .27+ .73-
```