

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

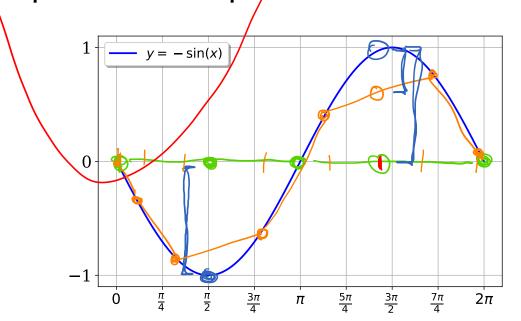
Big Ideas

- 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 - $\sin(x)$.



(1) chect and call sin(x)

anyway

(2) Faylor series approximation

(3) geometric solutions

- 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
 - 3. 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

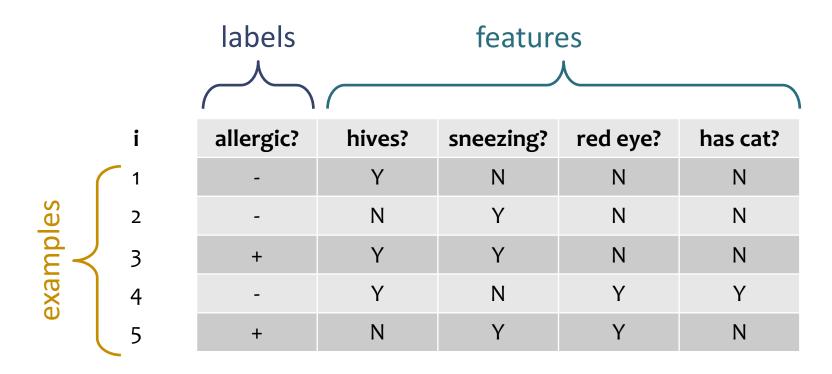
Interview Transcript

Date: Jan.7, 2025

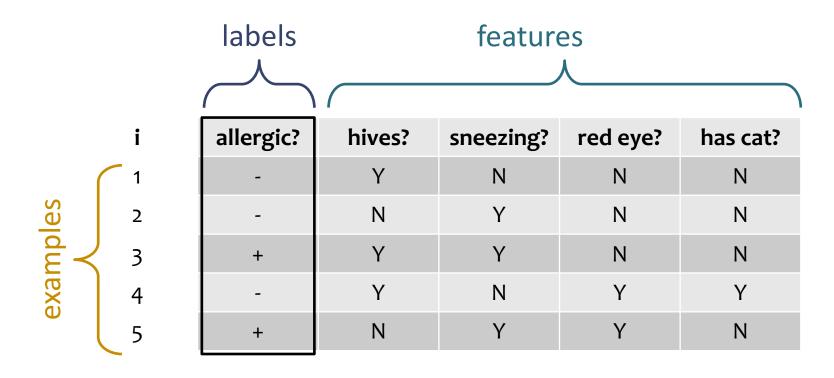
Parties: Matt Gormley and Doctor S.

Topic: Medical decision making

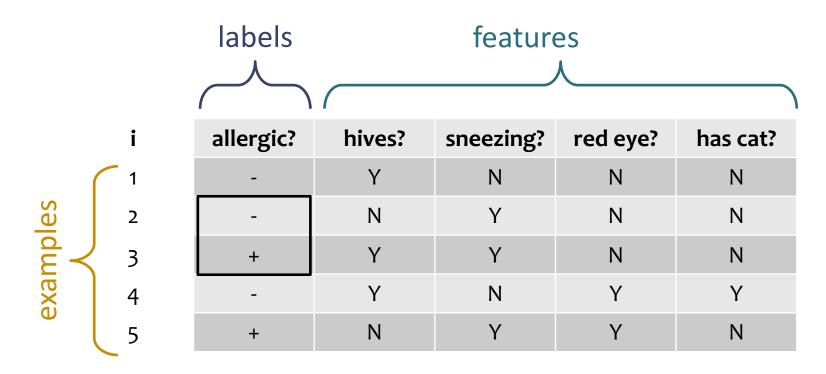
As a (supervised) binary classification task



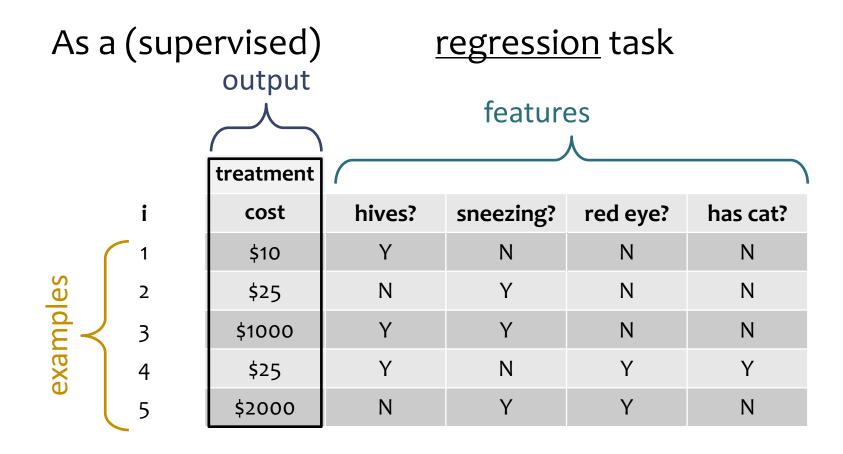
As a (<u>supervised</u>) binary classification task



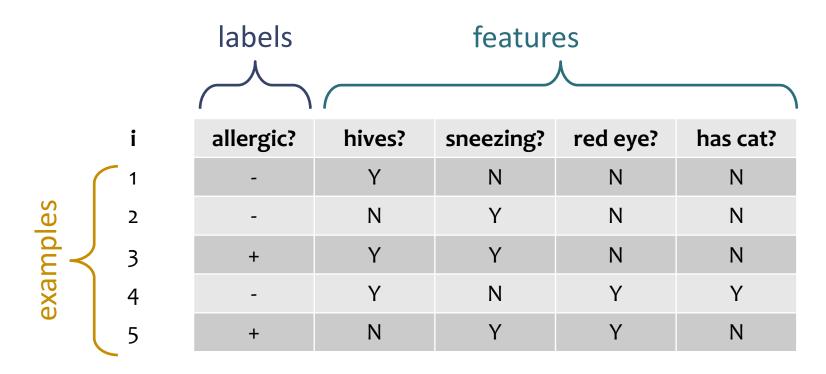
As a (supervised) binary classification task



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_2	X ₃	X ₄
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	X ₄
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 ⁽¹⁾ -	X ₁ ⁽¹⁾ Y	$x_2^{(1)} N$	x ₃ ⁽¹⁾ N	x ₄ ⁽¹⁾ N
2	y ⁽²⁾ -	$X_1^{(2)} N$	$X_2^{(2)} Y$	$X_3^{(2)} N$	$X_4^{(2)} N$
3	y ⁽³⁾ +	X ₁ ⁽³⁾ Y	X ₂ ⁽³⁾ Y	x ₃ ⁽³⁾ N	x ₄ ⁽³⁾ N
4	y ⁽⁴⁾ -	X ₁ ⁽⁴⁾ Y	$X_2^{(4)} N$	x ₃ ⁽⁴⁾ Y	x ₄ ⁽⁴⁾ Y
5	y ⁽⁵⁾ +	X ₁ ⁽⁵⁾ N	X ₂ ⁽⁵⁾ Y	x ₃ ⁽⁵⁾ Y	x ₄ ⁽⁵⁾ 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 ⁽¹⁾ -	X ₁ ⁽¹⁾ Y	$X_2^{(1)} N$	x ₃ ⁽¹⁾ N	x ₄ ⁽¹⁾ N	X ⁽¹⁾
2	y ⁽²⁾ -	$X_1^{(2)} N$	$X_2^{(2)} Y$	$x_3^{(2)} N$	x ₄ ⁽²⁾ N	$X^{(2)}$
3	y ⁽³⁾ +	Χ ₁ ⁽³⁾ Υ	X ₂ ⁽³⁾ Y	x ₃ ⁽³⁾ N	x ₄ ⁽³⁾ N	X (3)
4	y ⁽⁴⁾ -	X ₁ ⁽⁴⁾ Y	x ₂ ⁽⁴⁾ N	x ₃ ⁽⁴⁾ Y	x ₄ ⁽⁴⁾ Y	X ⁽⁴⁾
5	y ⁽⁵⁾ +	X ₁ ⁽⁵⁾ N	X ₂ ⁽⁵⁾ Y	x ₃ ⁽⁵⁾ Y	x ₄ ⁽⁵⁾ N	X (5)

N = 5 training examples

M = 4 attributes

ML as Function Approximation

troblem Setting:

- Set of possible inputs & (all possible feature vectors)
- Set at possible outputs of (all possible labels)
- Unknown teaget function c*:2-2
- Set of candidate hypotheses

Learner is given:

- Taining examples $\mathfrak{D} = \{ (\hat{x}^{(1)}, y^{(1)}), (\hat{x}^{(2)}, y^{(2)}), \dots, (\hat{x}^{(N)}, y^{(N)}) \}$ of unknown target Sunction y(i) = c*(x(i)), Yie \$1,..., Ns - N=# of examples M=# of features = 1×(i)

Learner produces:

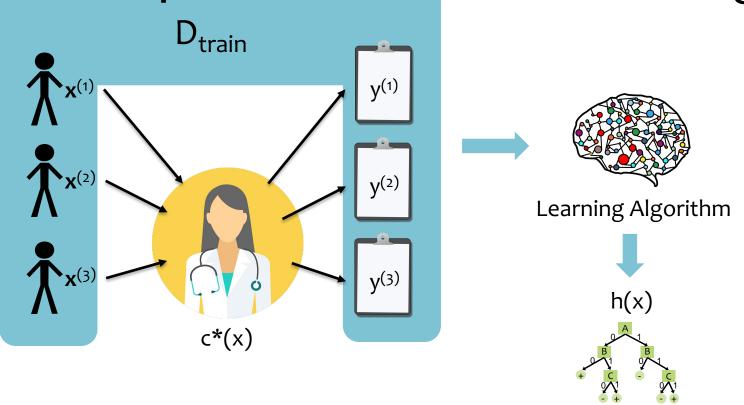
- Hypothesis he H that "best approximates" c* according to D

To Evaluate:

- Loss Function 1: yxy -> R measures how "bad" smaller for better pred. predictions $\hat{y} = h(\hat{x})$ are compared to true values $y^* = c^*(\hat{x})$
- Another dataset $\mathcal{D}_{test} = \{(\bar{x}^{(i)}, y^{(i)}), \dots, (\bar{x}^{(N')}, y^{(N')})\}$
- Evaluate the average loss of h(x) on Dtest

Aside: Function Types f(x,,x2,x3) = (x,x2)2 +x3 f: R³→R

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

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⁽¹⁾, y⁽¹⁾), (x⁽²⁾, y⁽²⁾), ...,
 where y⁽ⁱ⁾ = c*(x⁽ⁱ⁾)
 (history of patients and the
- Learner produces a hypoth approximates unknown tar

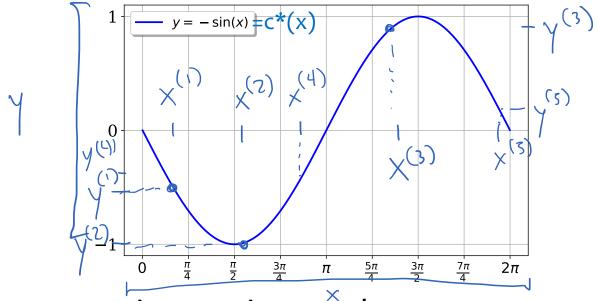
Two important settings we'll consider:

- 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
 - 3. You only need to evaluate it for x in [0, 2*pi]

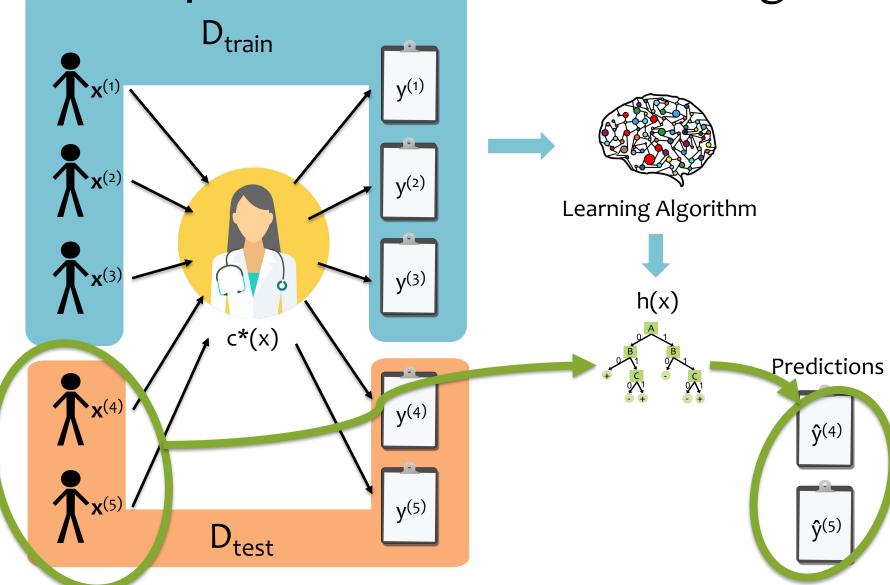
Supervised Machine Learning

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 $\sin(\mathbf{x})$ 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

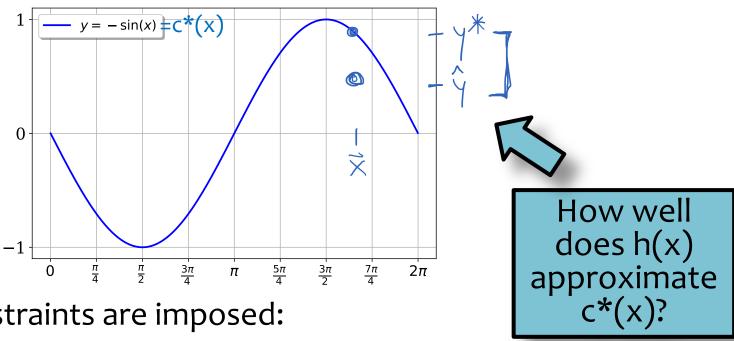
EVALUATION OF MACHINE LEARNING ALGORITHM

Supervised Machine Learning



Function Approximation

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



- A few constraints are imposed:
 - 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]

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:

2. Binary or 0-1 loss (for classification): $\ell(y,\hat{y}) = \mathbb{1}(y \neq \hat{y}) = \{ (y,\hat{y}) = \{ (y,\hat{y$

Evaluation of ML Algorithms

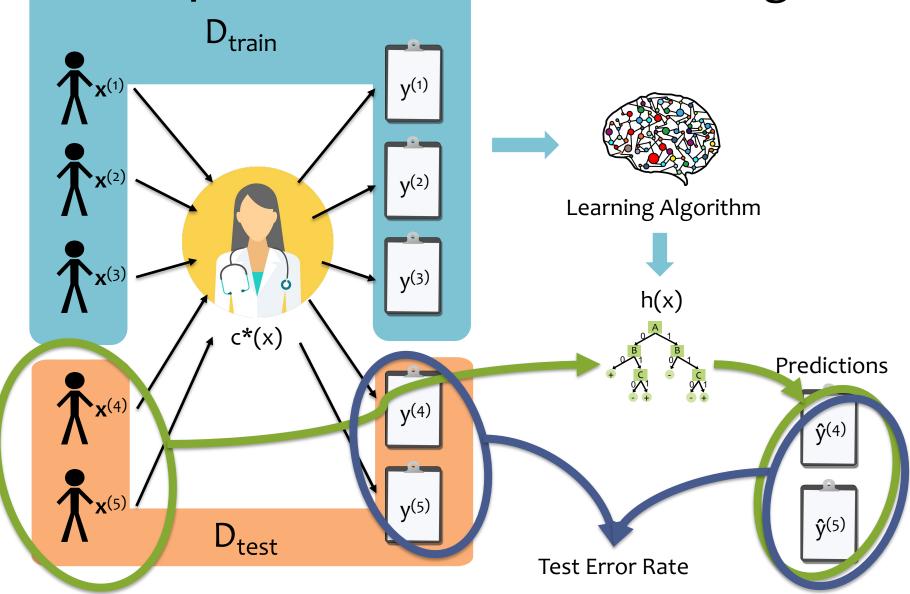
- **Definition:** loss function, $\ell: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$
 - o Defines how "bad" predictions, $\hat{y} = h(x)$, are compared to the true labels, $y = c^*(x)$
 - Common choices:
 - 1. Squared loss (for regression): $\ell(y, \hat{y}) = (y \hat{y})^2$
 - 2. Binary or 0-1 loss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

... true error over all data:

This is the quantity we care most about! But, in practice, error_{true}(h) is **unknown**.

Dataset:

Majority Vote Classifier Example

egad, we're already broken our supposed

In-Class Exercise

Dataset:

^	\wedge			
<u>Q</u>	$\overline{\gamma}$	Y	Α	В
	+	-	1	0
1	+	-	1	0
\bigcirc	+	+	1	0
\bigcirc	+	+	1	O
0	+	+	1	1
0	+	+	1	1
D	+	+	1	1
\bigcirc	+	+	1	1

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

Algorithms for 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	Υ	N

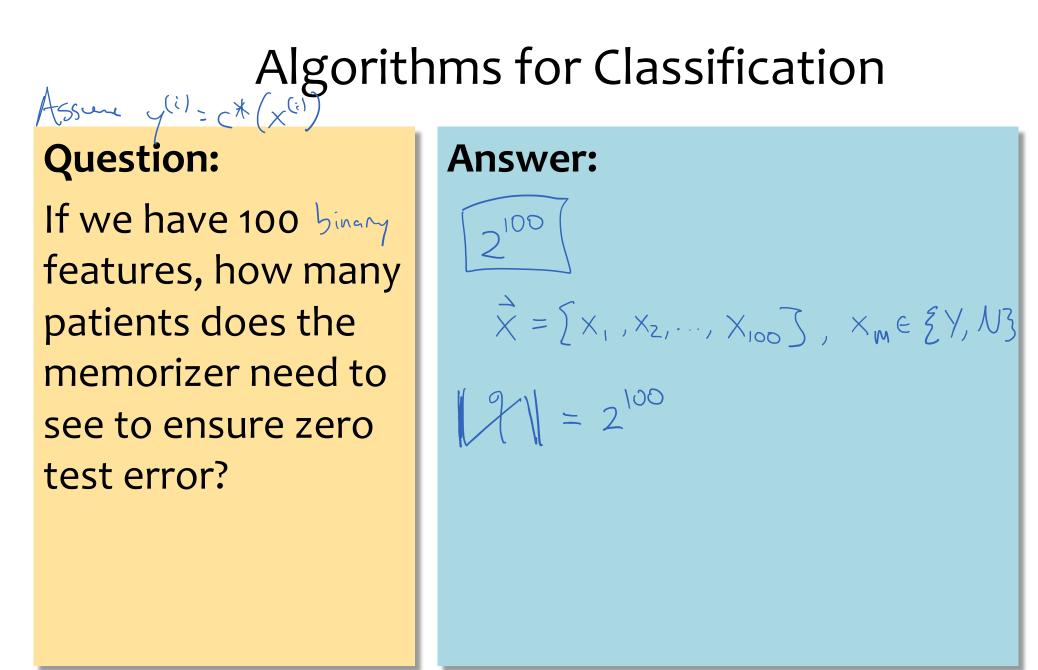
Algorithms for Classification

<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!

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



Algorithm 1: Majority Vote

Pseudocode

def train (D):

store
$$V = majority - vote(D)$$

= the dess yelf that appears most often in D

def $h(\vec{x})$:

return V

def predict (D test):

for $(\vec{x}^{(i)}, y^{(i)}) \in D_{Lest}$:

 $\hat{y}^{(i)} = h(\vec{x}^{(i)})$

Algorithm 2: Memorizer

Pseudocode

def train (D):
Store D
def
$$h(\dot{x})$$
:
if $J\dot{x}^{(i)} \in D$ s.t. $\dot{x}^{(i)} = \dot{x}$
return $y^{(i)}$
else:
return random $y \in \mathcal{Y}$

Algorithms for Classification

<u>Algorithm 3</u> **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?
V	-	<u>-</u>	Υ	N	N	N
X	+	-	N	Υ	N	N
\checkmark	+	+	Y	Υ	N	N
$\sqrt{}$	-	<u>-</u>	Υ	N	Y	Υ
\checkmark	+	+	N	Y	Y	N

Nonzero training error, but perhaps still better than the memorizer

Example decision stump:
$$h(\mathbf{x}) = \mathbf{1} + \mathbf{1} \text{ f sneezing} = \mathbf{Y}$$

$$\mathbf{1} \text{ otherwise } \mathbf{5} \mathbf{werry} = \mathbf{N}$$

Algorithm 3: Decision Stump

Assume Xm ∈ {0,13

Pseudocode

Jef fram (D):

(1) pick an attribute, M

(2) divide dataset on
$$\times M$$

$$D^{(0)} = \underbrace{X}(X^{(i)}, Y^{(i)}) \in D : X^{(i)} = 03$$

$$D^{(i)} = \underbrace{X}(X^{(i)}, Y^{(i)}) \in D : X^{(i)} = 13$$

(3) two votes
$$V^{(0)} = Majority - Vote(D^{(0)})$$

$$V^{(i)} = Majority - Vote(D^{(0)})$$

Jef
$$h(\vec{x})$$
:
if $x_m = 0$: return $V^{(0)}$
if $x_m = 1$: return $V^{(i)}$

Algorithms for Classification

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

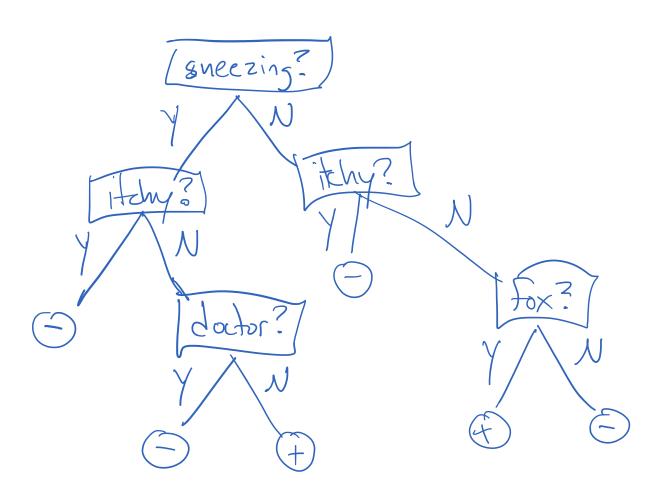
Questions:

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-
```