



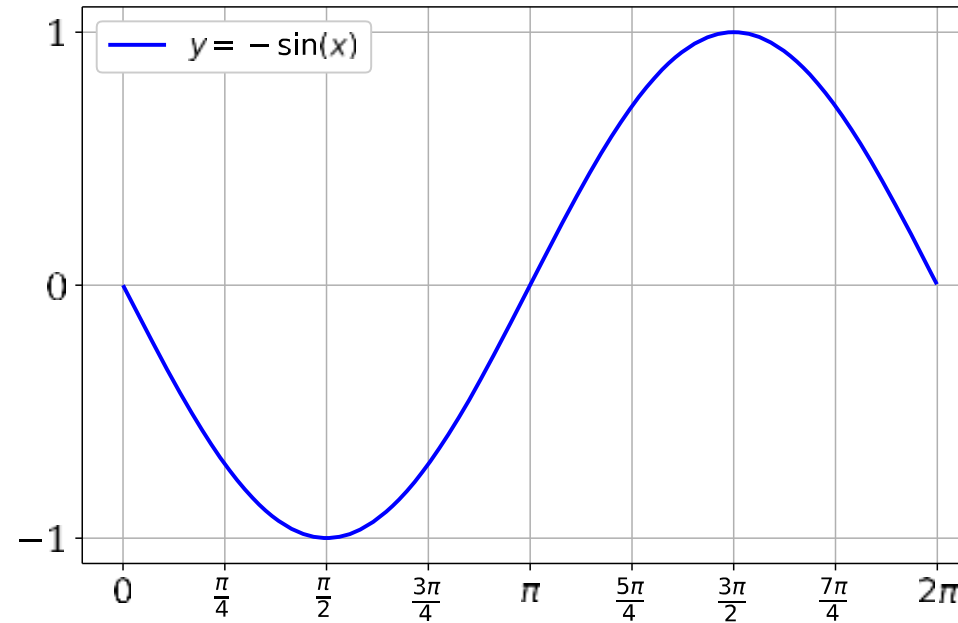
# CS182: Introduction to Machine Learning – ML as Function Approximation

Yujiao Shi  
SIST, ShanghaiTech  
Spring, 2025

## Warm-up Activity



- Challenge: implement a function that computes  $-\sin(x)$  for  $x \in [0, 2\pi]$



- You may not call any trigonometric functions
- You may call an existing implementation of  $\sin(x)$  a few times (e.g., 100) to check your work

# Our first Machine Learning Task

- Learning to diagnose heart disease  
as a **(supervised) binary classification task**

features			labels
Family History	Resting Blood Pressure	Cholesterol	Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes

# Our first Machine Learning Task

- Learning to diagnose heart disease  
as a **(supervised) binary classification** task

features			labels
Family History	Resting Blood Pressure	Cholesterol	Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes

# Our first Machine Learning Task

- Learning to diagnose heart disease

as a **(supervised)** classification task

	features			labels
	Family History	Resting Blood Pressure	Cholesterol	Risk
data points	Yes	Low	Normal	Low Risk
	No	Medium	Normal	Low Risk
	No	Low	Abnormal	Medium Risk
	Yes	Medium	Normal	High Risk
	Yes	High	Abnormal	High Risk

# Our first Machine Learning Task

- Learning to diagnose heart disease

as a (supervised)

regression task

features			targets
Family History	Resting Blood Pressure	Cholesterol	Medical Costs
Yes	Low	Normal	\$0
No	Medium	Normal	\$20
No	Low	Abnormal	\$30
Yes	Medium	Normal	\$100
Yes	High	Abnormal	\$5000

## Notation

- Feature space,  $\mathcal{X}$
- Label space,  $\mathcal{Y}$
- (Unknown) Target function,  $c^*: \mathcal{X} \rightarrow \mathcal{Y}$
- Training dataset:

$$\mathcal{D} = \{(\mathbf{x}^{(1)}, c^*(\mathbf{x}^{(1)}) = y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}) \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$$

- Example:  $(\mathbf{x}^{(n)}, y^{(n)}) = (x_1^{(n)}, x_2^{(n)}, \dots, x_D^{(n)}, y^{(n)})$
- Hypothesis space:  $\mathcal{H}$
- Goal: find a classifier,  $h \in \mathcal{H}$ , that best approximates  $c^*$

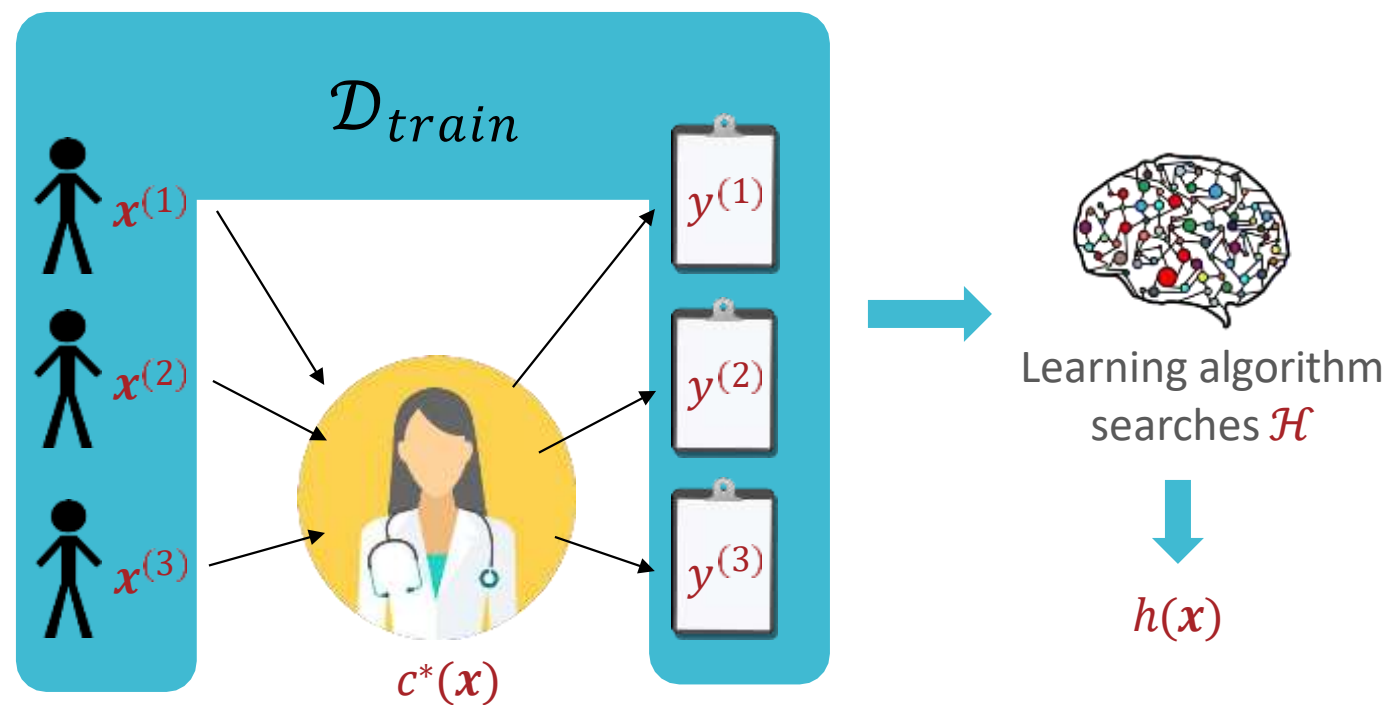
## Notation: Example

- $N = 5$  and  $D = 3$
- $\mathbf{x}^{(2)} = (x_1^{(2)} = \text{"No"}, x_2^{(2)} = \text{"Medium"}, x_3^{(2)} = \text{"Normal"})$

$x_1$ Family History	$x_2$ Resting Blood Pressure	$x_3$ Cholesterol	$y$ Heart Disease?	$y$ Predictions
Yes	Low	Normal	No	Yes
$\mathbf{x}^{(2)}$ No	Medium	Normal	No	Yes
No	Low	Abnormal	Yes	Yes
Yes	Medium	Normal	Yes	Yes
Yes	High	Abnormal	Yes	Yes



# Our first Machine Learning Task



# Evaluation



- Loss function,  $P: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ 
  - 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):
$$\ell(y, \hat{y}) = \begin{cases} 1 & \text{if } y \neq \hat{y} \\ 0 & \text{otherwise} \end{cases}$$

# Evaluation



- Loss function,  $P: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ 
  - 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):

$$\ell(y, \hat{y}) = \mathbb{1}(y \neq \hat{y})$$

- Error rate:

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

## Different Kinds of Error

- Training error rate =  $err(h, \mathcal{D}_{train})$
- Test error rate =  $err(h, \mathcal{D}_{test})$
- True error rate =  $err(h)$ 
  - = the error rate of  $h$  on all possible examples
- In machine learning, this is the quantity that we care about but, in most cases, it is unknowable.

# Our second Machine Learning Task

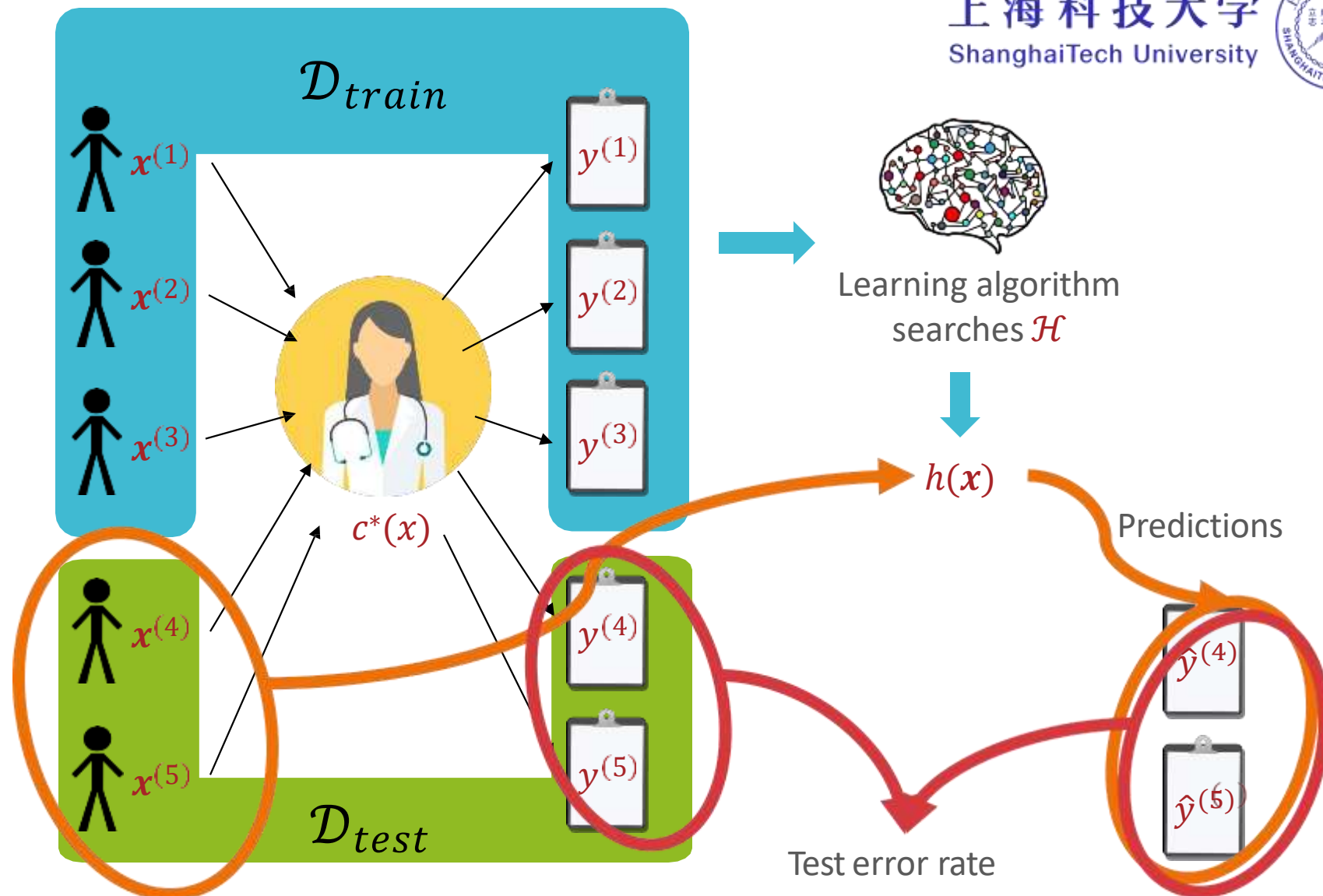


Figure courtesy of Matt Gormley

# Recall: Our first Machine Learning Classifier

- A **classifier** is a function that takes feature values as input and outputs a label
- Majority vote classifier: always predict the most common label in the **training** dataset

features				labels
Family History	Resting Blood Pressure	Cholesterol	Heart Disease?	Predictions
Yes	Low	Normal	No	Yes
No	Medium	Normal	No	Yes
No	Low	Abnormal	Yes	Yes
Yes	Medium	Normal	Yes	Yes
Yes	High	Abnormal	Yes	Yes

data points

# Our first Machine Learning Classifier: Pseudocode

- Majority vote classifier:

```
def train( $\mathcal{D}_{train}$ ):  
    store  $v = \text{mode}(y^{(1)}, y^{(2)}, \dots, y^{(N)})$   
  
def h( $x'$ ):  
    return  $v$   
  
def predict( $\mathcal{D}_{test}$ ):  
    for  $(x^n, y^n) \in \mathcal{D}_{test}$ :  
         $\hat{y}^{(n)} = h(x^{(n)})$ 
```

Test your  
understanding

$x_1$	$x_2$	$y$
1	0	-
1	0	-
1	0	+
1	0	+
1	1	+
1	1	+
1	1	+
1	1	+

- What is the **training error** of the **majority vote classifier** on this dataset?



# Our first Machine Learning Classifier

- A **classifier** is a function that takes feature values as input and outputs a label
- Majority vote classifier: always predict the most common label in the **training** dataset

features			labels	
Family History	Resting Blood Pressure	Cholesterol	Heart Disease?	Predictions
Yes	Low	Normal	No	Yes
No	Medium	Normal	No	Yes
No	Low	Abnormal	Yes	Yes
Yes	Medium	Normal	Yes	Yes
Yes	High	Abnormal	Yes	Yes

- This classifier completely ignores the features...

## Our second Machine Learning Classifier



- A **classifier** is a function that takes feature values as input and outputs a label
- Memorizer: if a set of features exists in the **training** dataset, predict its corresponding label; otherwise, predict the majority vote

Family History	Resting Blood Pressure	Cholesterol	Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes

## Our second Machine Learning Classifier



- A **classifier** is a function that takes feature values as input and outputs a label
- Memorizer: if a set of features exists in the **training** dataset, predict its corresponding label; otherwise, predict the majority vote

Family History	Resting Blood Pressure	Cholesterol	Heart Disease?	Predictions
Yes	Low	Normal	No	No
No	Medium	Normal	No	No
No	Low	Abnormal	Yes	Yes
Yes	Medium	Normal	Yes	Yes
Yes	High	Abnormal	Yes	Yes

- The training error rate is 0!

## Our second Machine Learning Classifier

- A **classifier** is a function that takes feature values as input and outputs a label
- Memorizer: if a set of features exists in the **training** dataset, predict its corresponding label; otherwise, predict the majority vote

Family History	Resting Blood Pressure	Cholesterol	Heart Disease?	Predictions
Yes	Low	Normal	No	No
No	Medium	Normal	No	No
No	Low	Abnormal	Yes	Yes
Yes	Medium	Normal	Yes	Yes
Yes	High	Abnormal	Yes	Yes

- The training error rate is 0...

# Is the memorizer “learning”?



- A **classifier** is a function that takes feature values as input and outputs a label
- Memorizer: if a set of features exists in the **training** dataset, predict its corresponding label; otherwise, predict the majority vote

Family History	Resting Blood Pressure	Cholesterol	Heart Disease?	Predictions
Yes	Low	Normal	No	No
No	Medium	Normal	No	No
No	Low	Abnormal	Yes	Yes
Yes	Medium	Normal	Yes	Yes
Yes	High	Abnormal	Yes	Yes

- The training error rate is 0...

## Our second Machine Learning Classifier



- A **classifier** is a function that takes feature values as input and outputs a label
- Memorizer: if a set of features exists in the **training** dataset, predict its corresponding label; otherwise, predict the majority vote
- The memorizer (typically) does not **generalize** well, i.e., it does not perform well on unseen data points
- In some sense, good generalization, i.e., the ability to make accurate predictions given a small training dataset, is the whole point of machine learning!

## Our second Machine Learning Classifier: Pseudocode



- Memorizer:

```
def train( $\mathcal{D}$ ):  
    store  $\mathcal{D}$   
  
def h( $\mathbf{x}'$ ):  
    if  $\exists \mathbf{x}^{(n)} \in \mathcal{D}$  s.t.  $\mathbf{x}' = \mathbf{x}^{(n)}$ :  
        return  $y^{(n)}$   
  
    else  
        return mode( $y^{(1)}, y^{(2)}, \dots, y^{(N)}$ )
```

## Our third Machine Learning Classifier

- Alright, let's actually (try to) extract a pattern from the data

$x_1$ Family History	$x_2$ Resting Blood Pressure	$x_3$ Cholesterol	$y$ Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes

- 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  $x_d$



## Our third Machine Learning Classifier

- Alright, let's actually (try to) extract a pattern from the data

$x_1$ Family History	$x_2$ Resting Blood Pressure	$x_3$ Cholesterol	$y$ Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes

- Decision stump on  $x_1$ :

$$h(\mathbf{x}') = h(x'_1, \dots, x'_D) = \begin{cases} ??? & \text{if } x'_1 = \text{"Yes"} \\ ??? & \text{otherwise} \end{cases}$$

# Our third Machine Learning Classifier



- Alright, let's actually (try to) extract a pattern from the data

$x_1$ Family History	$x_2$ Resting Blood Pressure	$x_3$ Cholesterol	$y$ Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes

- Decision stump on  $x_1$ :

$$h(\mathbf{x}') = h(x'_1, \dots, x'_D) = \begin{cases} \text{"Yes"} & \text{if } x'_1 = \text{"Yes"} \\ ??? & \text{otherwise} \end{cases}$$

# Our third Machine Learning Classifier



- Alright, let's actually (try to) extract a pattern from the data

$x_1$ Family History	$x_2$ Resting Blood Pressure	$x_3$ Cholesterol	$y$ Heart Disease?	$y$ Predictions
Yes	Low	Normal	No	Yes
No	Medium	Normal	No	No
No	Low	Abnormal	Yes	No
Yes	Medium	Normal	Yes	Yes
Yes	High	Abnormal	Yes	Yes

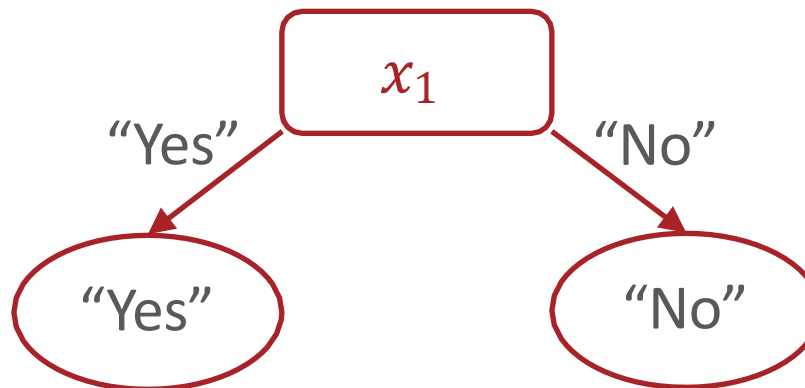
- Decision stump on  $x_1$ :

$$h(\mathbf{x}') = h(x'_1, \dots, x'_D) = \begin{cases} \text{"Yes"} & \text{if } x'_1 = \text{"Yes"} \\ \text{"No"} & \text{otherwise} \end{cases}$$

# Our third Machine Learning Classifier

- Alright, let's actually (try to) extract a pattern from the data

$x_1$ Family History	$x_2$ Resting Blood Pressure	$x_3$ Cholesterol	$y$ Heart Disease?	$y$ Predictions
Yes	Low	Normal	No	Yes
No	Medium	Normal	No	No
No	Low	Abnormal	Yes	No
Yes	Medium	Normal	Yes	Yes
Yes	High	Abnormal	Yes	Yes



## Decision Stumps: Pseudocode

```
def train( $\mathcal{D}$ ):
```

1. pick a feature,  $x_d$
2. split  $\mathcal{D}$  according to  $x_d$

for  $v$  in  $V(x_d)$ , all possible values of  $x_d$ :

$$\mathcal{D}_v = \{(\mathbf{x}^{(n)}, y^{(n)}) \in \mathcal{D} \mid x_d^{(n)} = v\}$$

3. Compute the majority vote for each split

for  $v$  in  $V(x_d)$ :

$$\hat{y}_v = \text{mode}(\text{labels in } \mathcal{D}_v)$$

```
def h( $\mathbf{x}'$ ):
```

for  $v$  in  $V(x_d)$ :

if  $x'_d = v$ : return  $\hat{y}_v$

# Decision Stumps: Questions



1. How can we pick which feature to split on?
2. Why stop at just one feature?
  - a) If we split on more than one feature, how do we decide the order to split on?