

Fall 2022 | Lecture 17

Decision Trees

Ariel Procaccia | Harvard University

SUPERVISED LEARNING

- We are given a **training set** of n examples $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$ where each pair was generated by an unknown function $y = f(x)$
- The goal is to find a **hypothesis** $h \in \mathcal{H}$ that approximates f , where \mathcal{H} is called the **hypothesis space**
- h is chosen to be a **best-fit function** for which each $h(x^{(i)})$ is “close” to $y^{(i)}$
- h **generalizes** well if it gives accurate predictions on a fresh **test set**

CLASSIFICATION

- **Classification** is the task of learning f whose range is a discrete, finite set
- Such a function is called a **classifier**
- When the cardinality of the range is 2 then the task is known as **binary classification**, otherwise it's called **multi-class classification**

EXAMPLE: IMAGE CLASSIFICATION



$(x^{(1)}, \text{dog})$



$(x^{(2)}, \text{cat})$



$(x^{(3)}, \text{dog})$

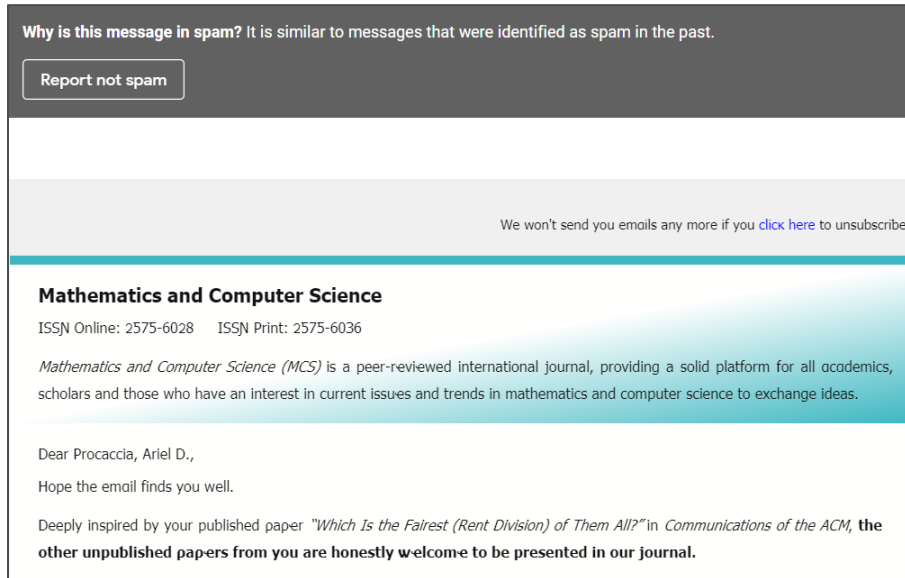


$(x^{(4)}, \text{dog})$

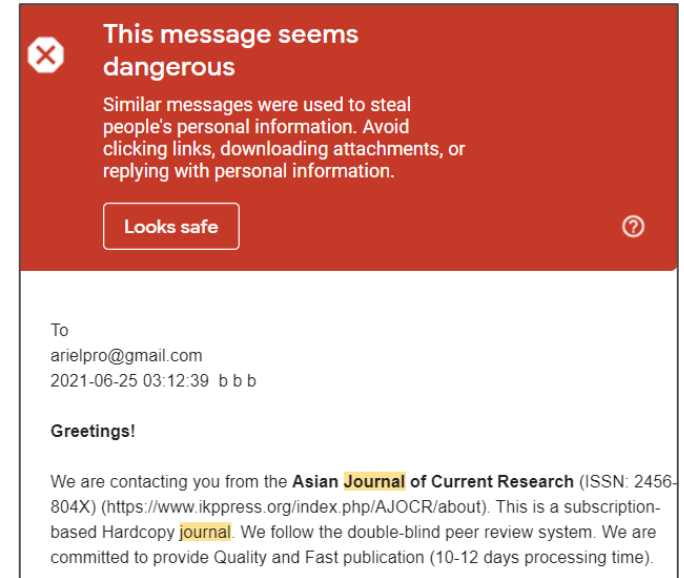


$f(x) = ?$

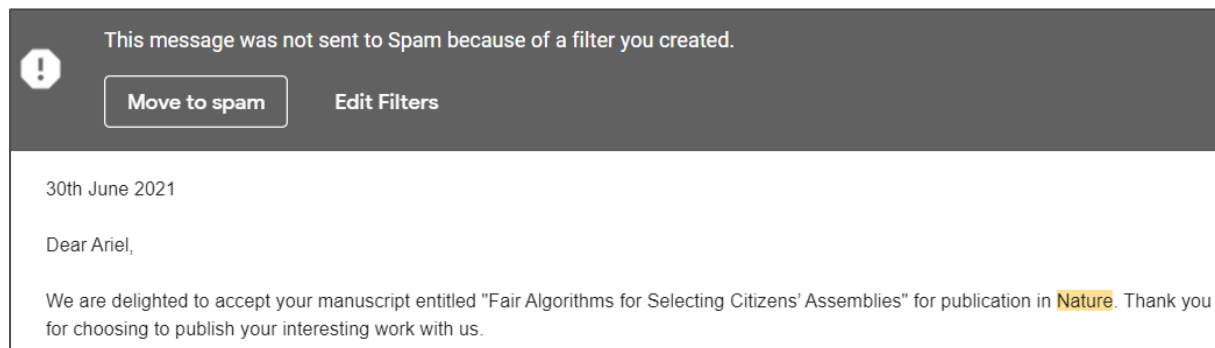
EXAMPLE: SPAM FILTER



$(x^{(1)}, spam)$



$(x^{(2)}, spam)$

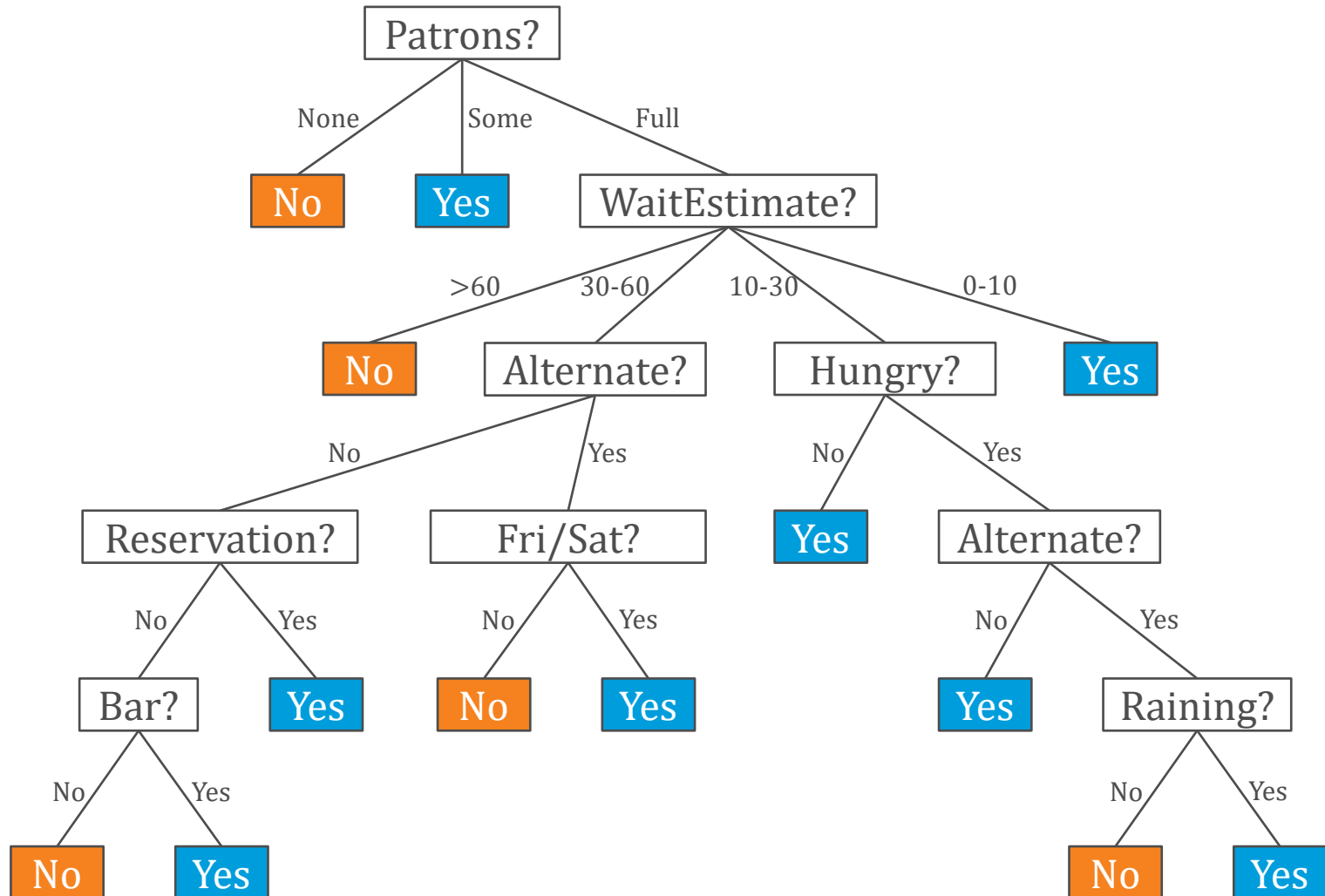


$f(x) = ?$

EXAMPLE: RESTAURANT WAITING

Example	Input Features										Output
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>Wait</i>
$\mathbf{x}^{(1)}$	<i>Y</i>	<i>N</i>	<i>N</i>	<i>Y</i>	<i>Some</i>	<i>\$\$\$</i>	<i>N</i>	<i>Y</i>	<i>French</i>	0-10	$y^{(1)} = Y$
$\mathbf{x}^{(2)}$	<i>Y</i>	<i>N</i>	<i>N</i>	<i>Y</i>	<i>Full</i>	<i>\$</i>	<i>N</i>	<i>N</i>	<i>Thai</i>	30-60	$y^{(2)} = N$
$\mathbf{x}^{(3)}$	<i>N</i>	<i>Y</i>	<i>N</i>	<i>N</i>	<i>Some</i>	<i>\$</i>	<i>N</i>	<i>N</i>	<i>Burger</i>	0-10	$y^{(3)} = Y$
$\mathbf{x}^{(4)}$	<i>Y</i>	<i>N</i>	<i>Y</i>	<i>Y</i>	<i>Full</i>	<i>\$</i>	<i>Y</i>	<i>N</i>	<i>Thai</i>	10-30	$y^{(4)} = Y$
$\mathbf{x}^{(5)}$	<i>Y</i>	<i>N</i>	<i>Y</i>	<i>N</i>	<i>Full</i>	<i>\$\$\$</i>	<i>N</i>	<i>Y</i>	<i>French</i>	>60	$y^{(5)} = N$
$\mathbf{x}^{(6)}$	<i>N</i>	<i>Y</i>	<i>N</i>	<i>Y</i>	<i>Some</i>	<i>\$\$</i>	<i>Y</i>	<i>Y</i>	<i>Italian</i>	0-10	$y^{(6)} = Y$
$\mathbf{x}^{(7)}$	<i>N</i>	<i>Y</i>	<i>N</i>	<i>N</i>	<i>None</i>	<i>\$</i>	<i>Y</i>	<i>N</i>	<i>Burger</i>	0-10	$y^{(7)} = N$
$\mathbf{x}^{(8)}$	<i>N</i>	<i>N</i>	<i>N</i>	<i>Y</i>	<i>Some</i>	<i>\$\$</i>	<i>Y</i>	<i>Y</i>	<i>Thai</i>	0-10	$y^{(8)} = Y$
$\mathbf{x}^{(9)}$	<i>N</i>	<i>Y</i>	<i>Y</i>	<i>N</i>	<i>Full</i>	<i>\$</i>	<i>Y</i>	<i>N</i>	<i>Burger</i>	>60	$y^{(9)} = N$
$\mathbf{x}^{(10)}$	<i>Y</i>	<i>Y</i>	<i>Y</i>	<i>Y</i>	<i>Full</i>	<i>\$\$\$</i>	<i>N</i>	<i>Y</i>	<i>Italian</i>	0-30	$y^{(10)} = N$
$\mathbf{x}^{(11)}$	<i>N</i>	<i>N</i>	<i>N</i>	<i>N</i>	<i>None</i>	<i>\$</i>	<i>N</i>	<i>N</i>	<i>Thai</i>	0-10	$y^{(11)} = N$
$\mathbf{x}^{(12)}$	<i>Y</i>	<i>Y</i>	<i>Y</i>	<i>Y</i>	<i>Full</i>	<i>\$</i>	<i>N</i>	<i>N</i>	<i>Burger</i>	30-60	$y^{(12)} = Y$

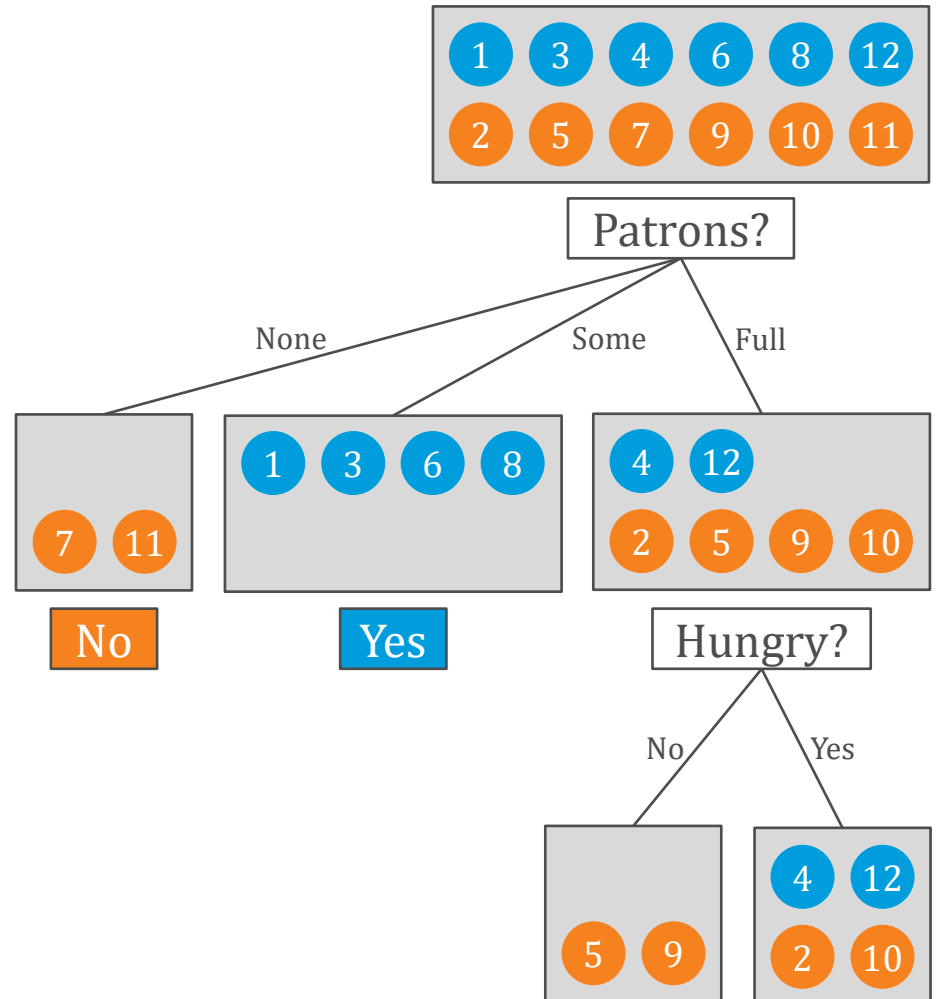
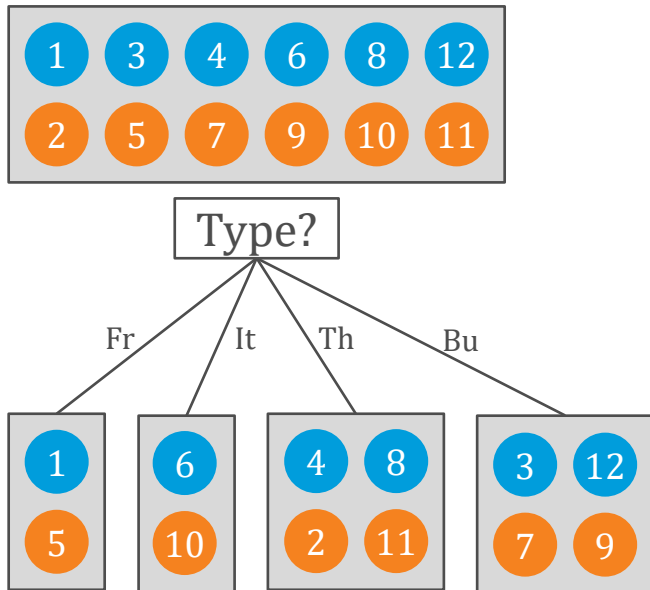
DECISION TREES



DECISION TREES

- A **decision tree** reaches an output (in the leaves) through a sequence of tests on the input attributes (in internal nodes)
- Decision trees can represent any classifier, but some may require a large tree
- **Poll 1:** Which of the following Boolean functions can be represented via a tree of size linear in the number of features?
 - Unanimity
 - Parity
 - Majority
 - None of the above

SPLITTING ON FEATURES



LEARNING DECISION TREES

```
function LEARN-DT(examples, features, parent_examples)
  if examples =  $\emptyset$  then return PLURALITY-VALUE(parent_examples)
  else if all examples have the same label then return
    that label
  else if features =  $\emptyset$  then return PLURALITY-VALUE(examples)
  else
     $A \leftarrow \operatorname{argmax}_{a \in \text{features}} \text{IMPORTANCE}(a, \text{examples})$ 
    tree  $\leftarrow$  new decision tree with root test A
    for each value v of A do
      new_examples  $\leftarrow \{e \in \text{examples} : e.A = v\}$ 
      subtree  $\leftarrow$  LEARN-DT(new_examples, features  $\setminus \{A\}$ , examples)
      add branch to tree with label A = v and subtree
  return tree
```



Claude Shannon

1916-2001

Mathematician and electrical engineer,
father of information theory.



INFORMATION GAIN

- To instantiate the IMPORTANCE function we will use the notion of **entropy**, which is measured in bits
- The entropy of random variable V that takes each value v with probability $P(v)$ is

$$H(V) = \sum_v P(v) \log \frac{1}{P(v)} = - \sum_v P(v) \log P(v)$$

- The entropy of a fair coin flip is 1 bit:

$$H(Fair) = -(0.5 \log 0.5 + 0.5 \log 0.5) = 1$$

- The entropy of a biased coin with 99% heads is:

$$H(Biased) = -(0.99 \log 0.99 + 0.01 \log 0.01) \approx 0.08$$

- Denote the entropy of a Bernoulli random variable that is true with probability q by

$$B(q) = -(q \log q + (1 - q) \log(1 - q))$$

INFORMATION GAIN

- If a training set contains p positive examples and n negative examples, the entropy of the output variable is

$$H(\text{Output}) = B\left(\frac{p}{p+n}\right)$$

- Feature A with d values divides the training set into d subsets, each with p_k positive examples and n_k negative examples
- The entropy after testing A is

$$\text{Remainder}(A) = \sum_{k=1}^d \frac{p_k + n_k}{p + n} B\left(\frac{p_k}{p_k + n_k}\right)$$

- The **information gain** from testing A is

$$\text{Gain}(A) = H(\text{Output}) - \text{Remainder}(A)$$

- In LEARN-DT, we can measure IMPORTANCE based on information gain

INFORMATION GAIN: EXAMPLE



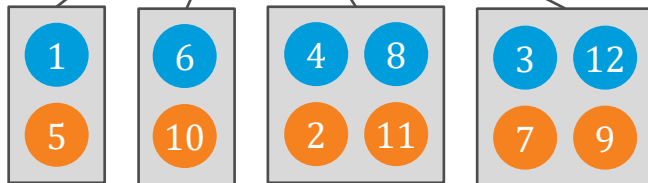
Type?

Fr

It

Th

Bu



Patrons?

None

Some

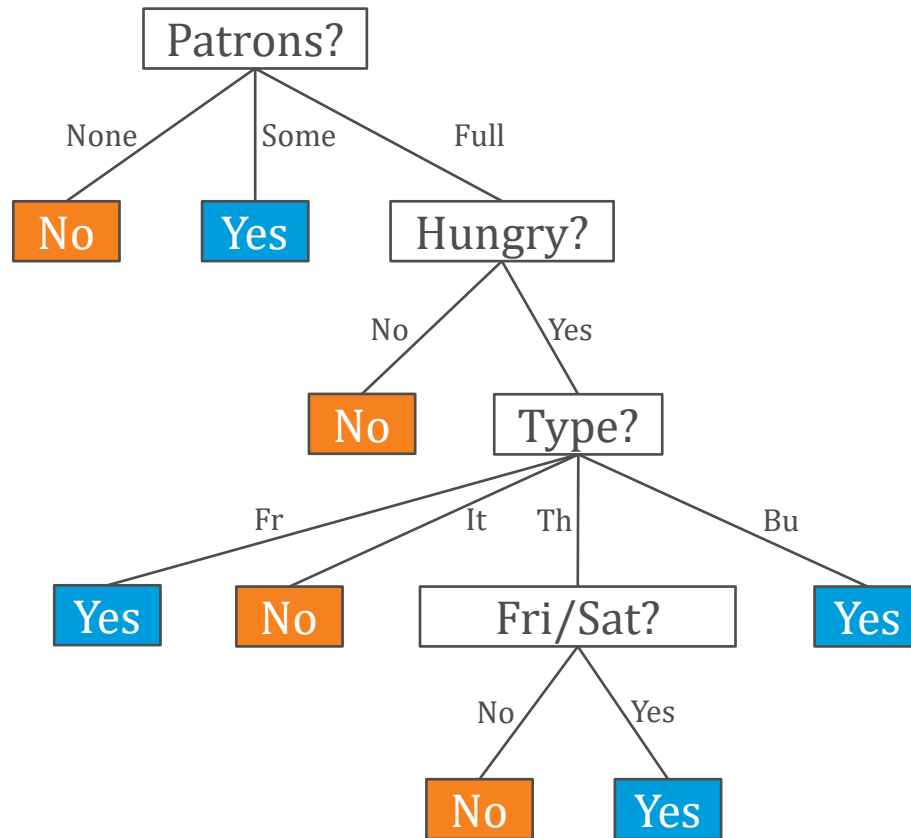
Full



$$Gain(Type) = 1 - \left[\frac{2}{12} B\left(\frac{1}{2}\right) + \frac{2}{12} B\left(\frac{1}{2}\right) + \frac{4}{12} B\left(\frac{2}{4}\right) + \frac{4}{12} B\left(\frac{2}{4}\right) \right] = 0$$

$$Gain(Patrons) = 1 - \left[\frac{2}{12} B\left(\frac{0}{2}\right) + \frac{4}{12} B\left(\frac{4}{4}\right) + \frac{6}{12} B\left(\frac{2}{6}\right) \right] \approx 0.541$$

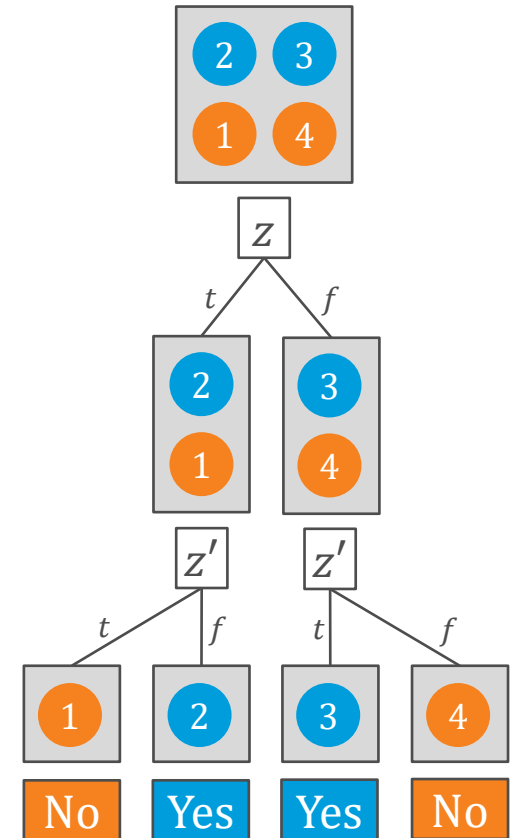
LEARNING DECISION TREES: EXAMPLE



The output of LEARN-DT is simpler than the original tree!

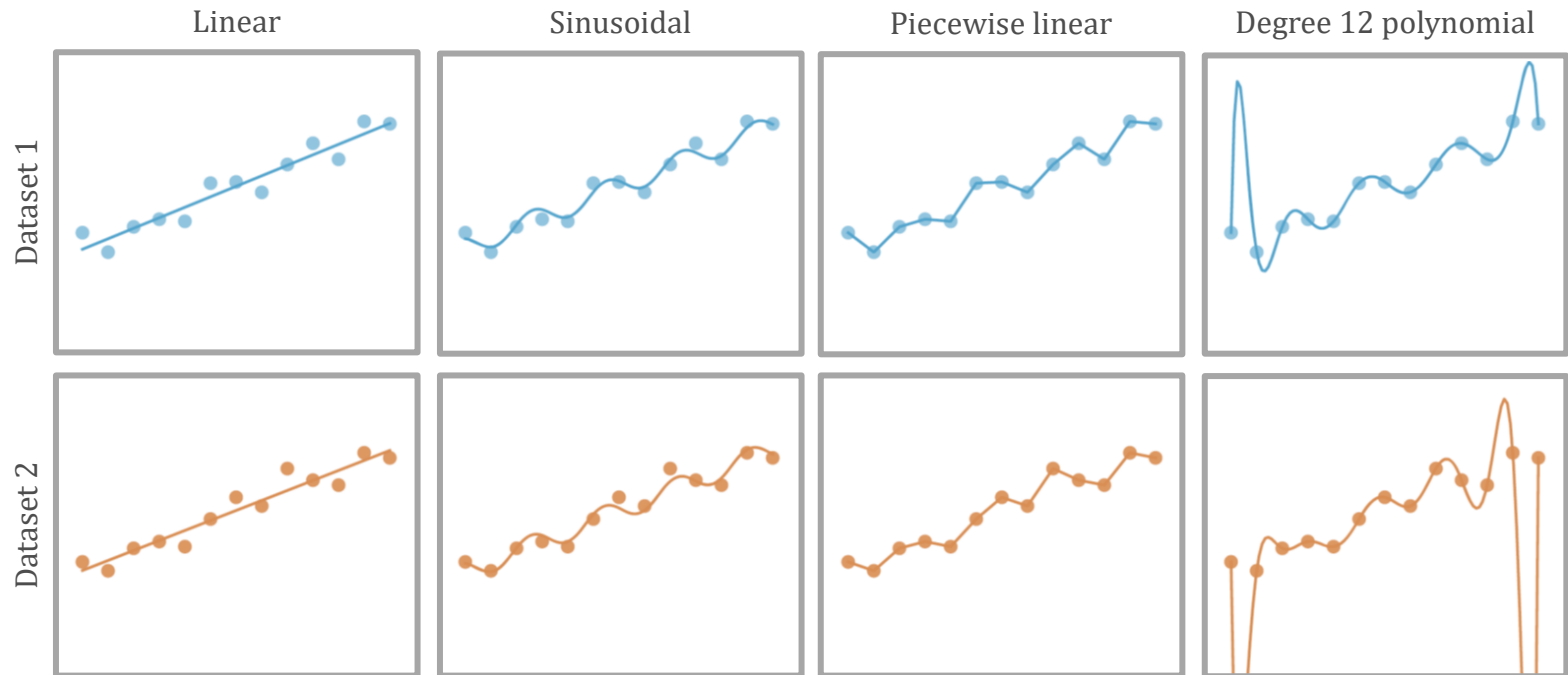
EARLY STOPPING

Example	Input features		Output
	z	z'	$z \oplus z'$
$x^{(1)}$	t	t	f
$x^{(2)}$	t	f	t
$x^{(3)}$	f	t	t
$x^{(4)}$	f	f	f



Should we stop LEARN-DT when the information gain is low? We may miss situations where combos of features are informative!

MODEL SELECTION



Degree 12 polynomials exhibit **overfitting**

MODEL SELECTION

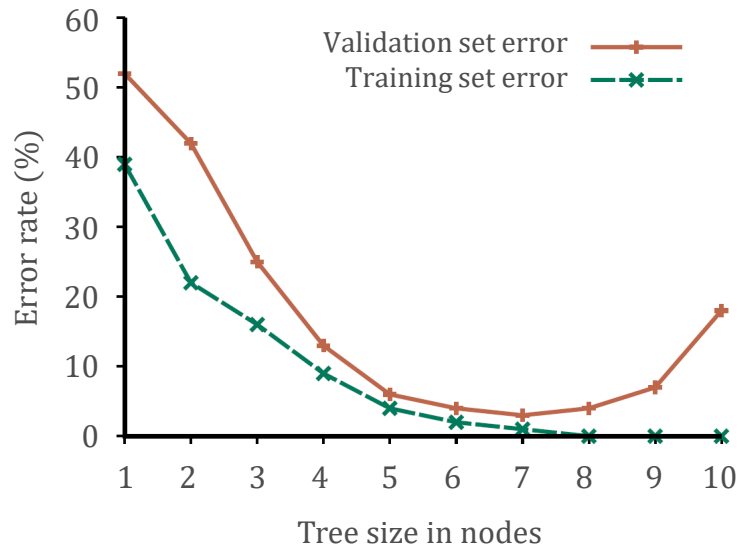
- Let us think of the quality of the “fit” of hypothesis h as **error rate**, i.e., the probability that $h(\mathbf{x}) \neq f(\mathbf{x})$
- We divide the data into three sets:
 1. **Training set** to train candidate models
 2. **Validation set** to choose among different models or hypothesis classes
 3. **Test set** to perform an unbiased evaluation of the best model

MODEL SELECTION

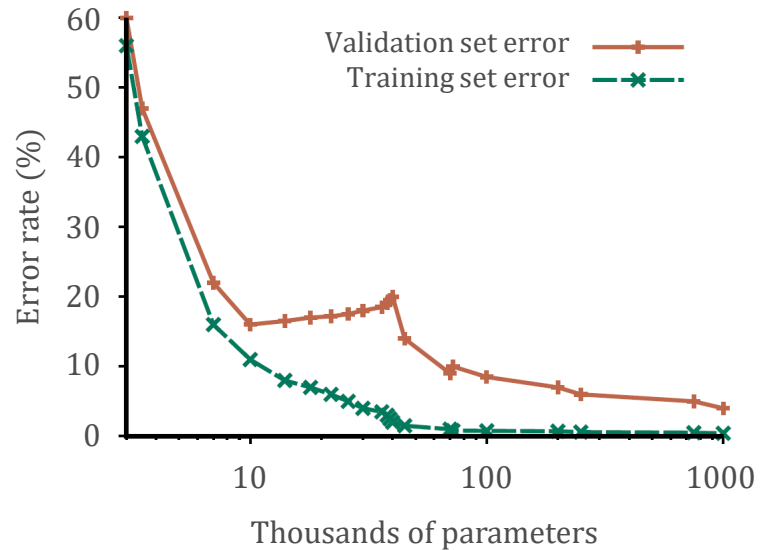
- We refer to the “complexity” of the hypothesis class (e.g., number of nodes in a decision tree) as the model size
- **Poll 2:** As the model size grows (check all possible options):
 - Training error decreases, validation error decreases
 - Training error decreases, validation error increases
 - Training error increases, validation error decreases
 - Training error increases, validation error increases

MODEL SELECTION

Decision Trees



Convolutional neural networks



The graph on the right is an example of one of the great mysteries of deep learning!