# Machine Learning (CSE 446):
## Decision Trees

Noah Smith
© 2017

University of Washington
nasmith@cs.washington.edu

September 29, 2017

# Features

Let $\phi$ be a function that maps from inputs $(x)$ to values.

# Features

Let $\phi$ be a function that maps from inputs $(x)$ to values.

- If $\phi$ maps to $\{0, 1\}$, we call it a "binary feature (function)."

# Features

Let $\phi$ be a function that maps from inputs ($x$) to values.

- If $\phi$ maps to $\{0, 1\}$, we call it a "binary feature (function)."
- If $\phi$ maps to $\mathbb{R}$, we call it a "real-valued feature (function)."

# Features

Let $\phi$ be a function that maps from inputs $(x)$ to values.

- If $\phi$ maps to $\{0, 1\}$, we call it a "binary feature (function)."
- If $\phi$ maps to $\mathbb{R}$, we call it a "real-valued feature (function)."
- Feature functions can map to categorical values, ordinal values, integers, and more.

# Features

Data derived from https://archive.ics.uci.edu/ml/datasets/Auto+MPG

| mpg | cylinders | displacement | horsepower | weight | acceleration | year | origin |
|------|-----------|--------------|------------|--------|--------------|------|--------|
| 18.0 | 8 | 307.0 | 130.0 | 3504. | 12.0 | 70 | 1 |
| 15.0 | 8 | 350.0 | 165.0 | 3693. | 11.5 | 70 | 1 |
| 18.0 | 8 | 318.0 | 150.0 | 3436. | 11.0 | 70 | 1 |
| 16.0 | 8 | 304.0 | 150.0 | 3433. | 12.0 | 70 | 1 |
| 17.0 | 8 | 302.0 | 140.0 | 3449. | 10.5 | 70 | 1 |
| 15.0 | 8 | 429.0 | 198.0 | 4341. | 10.0 | 70 | 1 |
| 14.0 | 8 | 454.0 | 220.0 | 4354. | 9.0 | 70 | 1 |
| 14.0 | 8 | 440.0 | 215.0 | 4312. | 8.5 | 70 | 1 |
| 14.0 | 8 | 455.0 | 225.0 | 4425. | 10.0 | 70 | 1 |
| 15.0 | 8 | 390.0 | 190.0 | 3850. | 8.5 | 70 | 1 |
| 15.0 | 8 | 383.0 | 170.0 | 3563. | 10.0 | 70 | 1 |
| 14.0 | 8 | 340.0 | 160.0 | 3609. | 8.0 | 70 | 1 |
| 15.0 | 8 | 400.0 | 150.0 | 3761. | 9.5 | 70 | 1 |
| 14.0 | 8 | 455.0 | 225.0 | 3086. | 10.0 | 70 | 1 |
| 24.0 | 4 | 113.0 | 95.00 | 2372. | 15.0 | 70 | 3 |
| 22.0 | 6 | 198.0 | 95.00 | 2833. | 15.5 | 70 | 1 |
| 18.0 | 6 | 199.0 | 97.00 | 2774. | 15.5 | 70 | 1 |
| 21.0 | 6 | 200.0 | 85.00 | 2587. | 16.0 | 70 | 1 |
| 27.0 | 4 | 97.00 | 88.00 | 2130. | 14.5 | 70 | 3 |
| 26.0 | 4 | 97.00 | 46.00 | 1835. | 20.5 | 70 | 2 |
| 25.0 | 4 | 110.0 | 87.00 | 2672. | 17.5 | 70 | 2 |
| 24.0 | 4 | 107.0 | 90.00 | 2430. | 14.5 | 70 | 2 |

Goal: predict whether mpg is $< 23$ ("bad" $= 0$) or above ("good" $= 1$) given other attributes (other columns).

201 "good" and 197 "bad"; guessing the most frequent class (good) will get 50.5% accuracy.

# Contingency Table

| values of $y$ | values of feature $\phi$ | | | |
|---|---|---|---|---|
| | $v_1$ | $v_2$ | $\cdots$ | $v_K$ |
| 0 | | | | |
| 1 | | | | |

# Decision Stump Example

| $y$ | maker | | |
|---|---|---|---|
| | america | europe | asia |
| 0 | 174 | 14 | 9 |
| 1 | 75 | 56 | 70 |
| | ↓ | ↓ | ↓ |
| | 0 | 1 | 1 |

# Decision Stump Example

| $y$ | maker | | |
|---|---|---|---|
| | america | europe | asia |
| 0 | 174 | 14 | 9 |
| 1 | 75 | 56 | 70 |
| | ↓ | ↓ | ↓ |
| | 0 | 1 | 1 |

# Decision Stump Example



| $y$ | maker | | |
|---|---|---|---|
| | america | europe | asia |
| 0 | 174 | 14 | 9 |
| 1 | 75 | 56 | 70 |
| | ↓ | ↓ | ↓ |
| | 0 | 1 | 1 |

Errors: $75 + 14 + 9 = 98$   (about 25%)

# Decision Stump Example

# Decision Stump Example
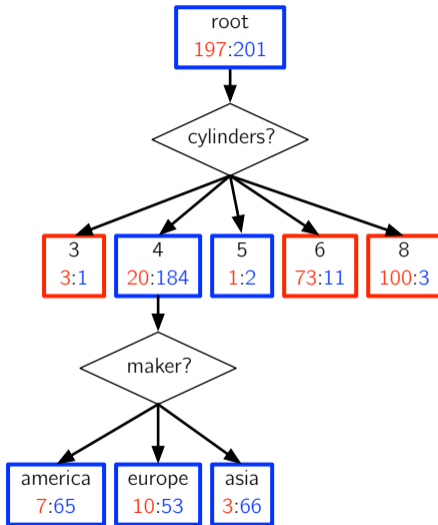


Errors: $1 + 20 + 1 + 11 + 3 = 36$ (about 9%)

# Key Idea: Recursion

A single feature **partitions** the data.

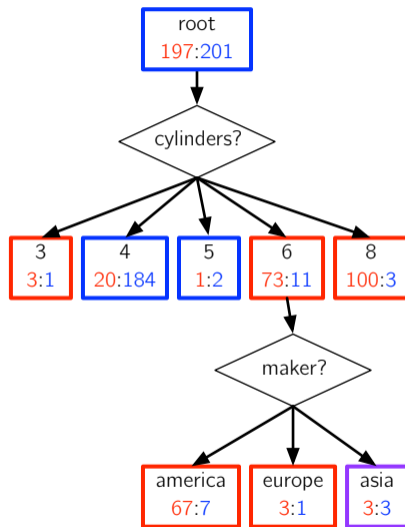For each partition, we could choose another feature and partition further.

Applying this recursively, we can construct a **decision tree**.
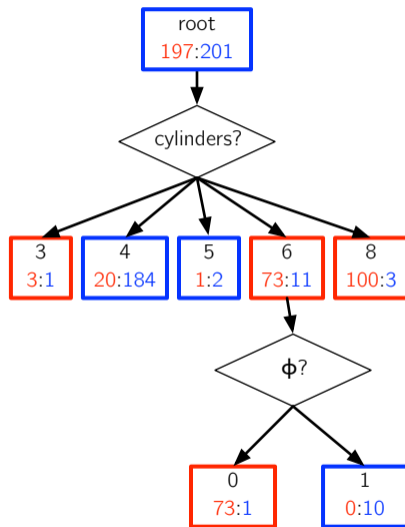
# Decision Tree Example



Error reduction compared to the cylinders stump?

# Decision Tree Example



Error reduction compared to the cylinders stump?
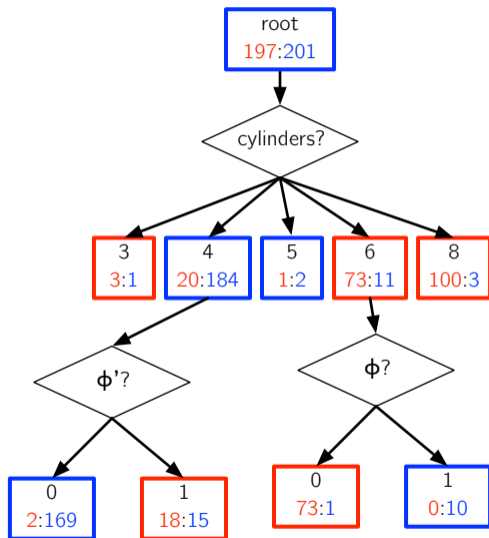
# Decision Tree Example



Error reduction compared to the cylinders stump?
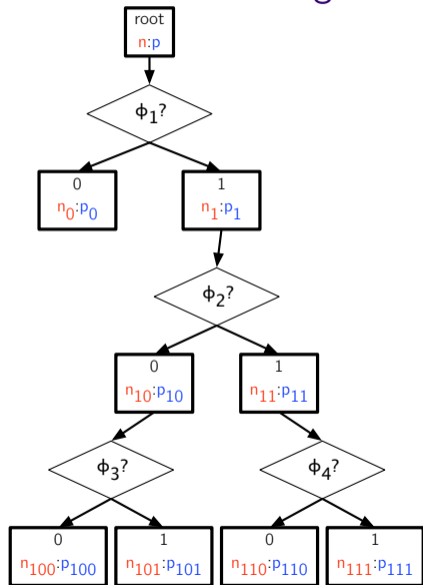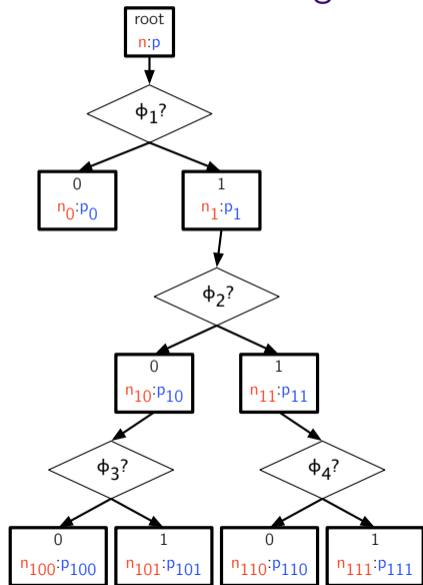
# Decision Tree Example



Error reduction compared to the cylinders stump?

# Decision Tree: Making a Prediction

# Decision Tree: Making a Prediction



**Data**: decision tree $t$, input example $x$
**Result**: predicted class
**if** *$t$ has the form* $\text{LEAF}(y)$ **then**
  return $y$;
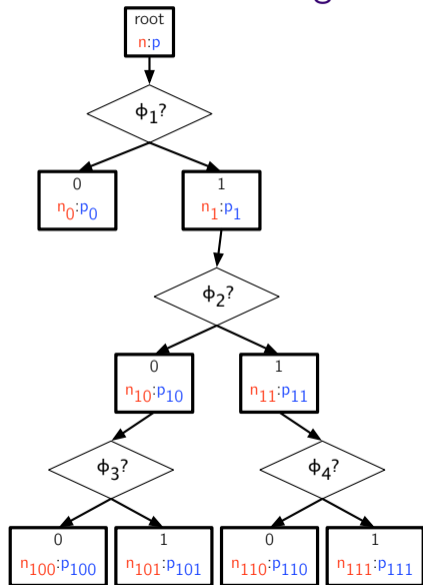**else**
  # $t.\phi$ is the feature associated with $t$;
  # $t.\text{child}(v)$ is the subtree for value $v$;
  return $\text{DTREETEST}(t.\text{child}(t.\phi(x)),\ x)$;
**end**

**Algorithm 1:** $\text{DTREETEST}$

# Decision Tree: Making a Prediction



Equivalent boolean formulas:

$$(\phi_1 = 0) \Rightarrow [\![ \mathsf{n}_0 < \mathsf{p}_0 ]\!]$$
$$(\phi_1 = 1) \wedge (\phi_2 = 0) \wedge (\phi_3 = 0) \Rightarrow [\![ \mathsf{n}_{100} < \mathsf{p}_{100} ]\!]$$
$$(\phi_1 = 1) \wedge (\phi_2 = 0) \wedge (\phi_3 = 1) \Rightarrow [\![ \mathsf{n}_{101} < \mathsf{p}_{101} ]\!]$$
$$(\phi_1 = 1) \wedge (\phi_2 = 1) \wedge (\phi_4 = 0) \Rightarrow [\![ \mathsf{n}_{110} < \mathsf{p}_{110} ]\!]$$
$$(\phi_1 = 1) \wedge (\phi_2 = 1) \wedge (\phi_4 = 1) \Rightarrow [\![ \mathsf{n}_{111} < \mathsf{p}_{111} ]\!]$$

# Tangent: How Many Formulas?

Assume we have $D$ binary features.

# Tangent: How Many Formulas?

Assume we have $D$ binary features.

Each feature could be set to 0, or set to 1, or excluded (wildcard/don't care).

# Tangent: How Many Formulas?

Assume we have $D$ binary features.

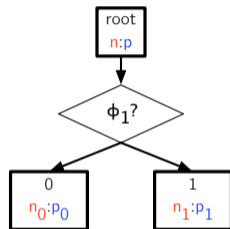Each feature could be set to 0, or set to 1, or excluded (wildcard/don't care).

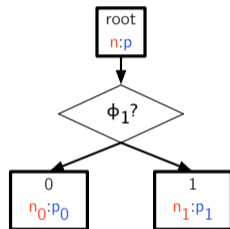$3^D$ formulas.

# Growing a Decision Tree

```
root
n:p
```

# Growing a Decision Tree
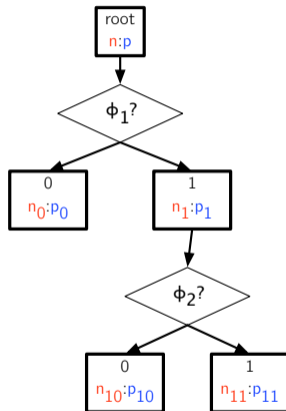


We chose feature $\phi_1$. Note that $n = n_0 + n_1$ and $p = p_0 + p_1$.

# Growing a Decision Tree



We chose not to split the left partition. Why not?

# Growing a Decision Tree



A decision tree diagram. The root node is labeled "root" with $n{:}p$. An arrow points down to a decision node $\phi_1?$. The left branch (0) leads to a node labeled "0" with $n_0{:}p_0$. The right branch (1) leads to a node labeled "1" with $n_1{:}p_1$, which points down to another decision node $\phi_2?$. Its left branch (0) leads to a node labeled "0" with $n_{10}{:}p_{10}$, and its right branch (1) leads to a node labeled "1" with $n_{11}{:}p_{11}$.

# Growing a Decision Tree
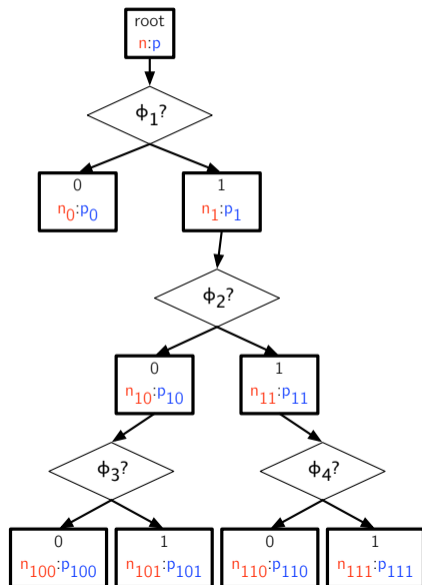
# Growing a Decision Tree

# Greedily Building a Decision Tree (Binary Features)

**Data**: data $D$, feature set $\Phi$

**Result**: decision tree

**if** *all examples in $D$ have the same label $y$, or $\Phi$ is empty and $y$ is the best guess* **then**

    | return $\text{LEAF}(y)$;

**else**

    | **for** *each feature $\phi$ in $\Phi$* **do**

        | | partition $D$ into $D_0$ and $D_1$ based on $\phi$-values;

        | | let mistakes$(\phi)$ = (non-majority answers in $D_0$) + (non-majority answers in $D_1$);
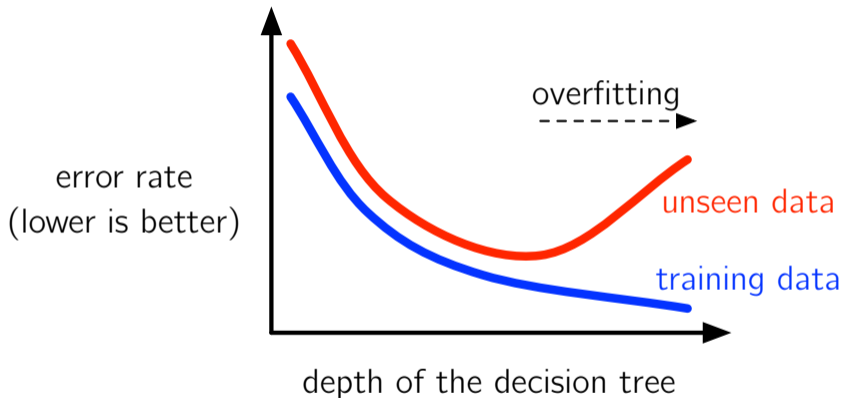
    | **end**

    | let $\phi^*$ be the feature with the smallest number of mistakes;

    | return $\text{NODE}(\phi^*, \{0 \to \text{DTREETRAIN}(D_0, \Phi \setminus \{\phi^*\}), 1 \to \text{DTREETRAIN}(D_1, \Phi \setminus \{\phi^*\})\})$;

**end**

**Algorithm 2:** $\text{DTREETRAIN}$

# Danger: Overfitting



error rate (lower is better)

overfitting

unseen data

training data

depth of the decision tree

# Detecting Overfitting

If you use all of your data to train, you won't be able to draw the red curve on the preceding slide!

# Detecting Overfitting

If you use all of your data to train, you won't be able to draw the red curve on the preceding slide!

Solution: hold some out. This data is called **development data**. More terms:

- Decision tree max depth is an example of a **hyperparameter**
- "I used my development data to **tune** the max-depth hyperparameter."

# Detecting Overfitting

If you use all of your data to train, you won't be able to draw the red curve on the preceding slide!

Solution: hold some out. This data is called **development data**. More terms:

- Decision tree max depth is an example of a **hyperparameter**
- "I used my development data to **tune** the max-depth hyperparameter."

Better yet, hold out two subsets, one for tuning and one for a true, honest-to-science **test**.

# Detecting Overfitting

If you use all of your data to train, you won't be able to draw the red curve on the preceding slide!

Solution: hold some out. This data is called **development data**. More terms:

- Decision tree max depth is an example of a **hyperparameter**
- "I used my development data to **tune** the max-depth hyperparameter."

Better yet, hold out two subsets, one for tuning and one for a true, honest-to-science **test**.

Splitting your data into training/development/test requires careful thinking. Starting point: randomly shuffle examples with an 80%/10%/10% split.