# CS 189/289

Today's lecture
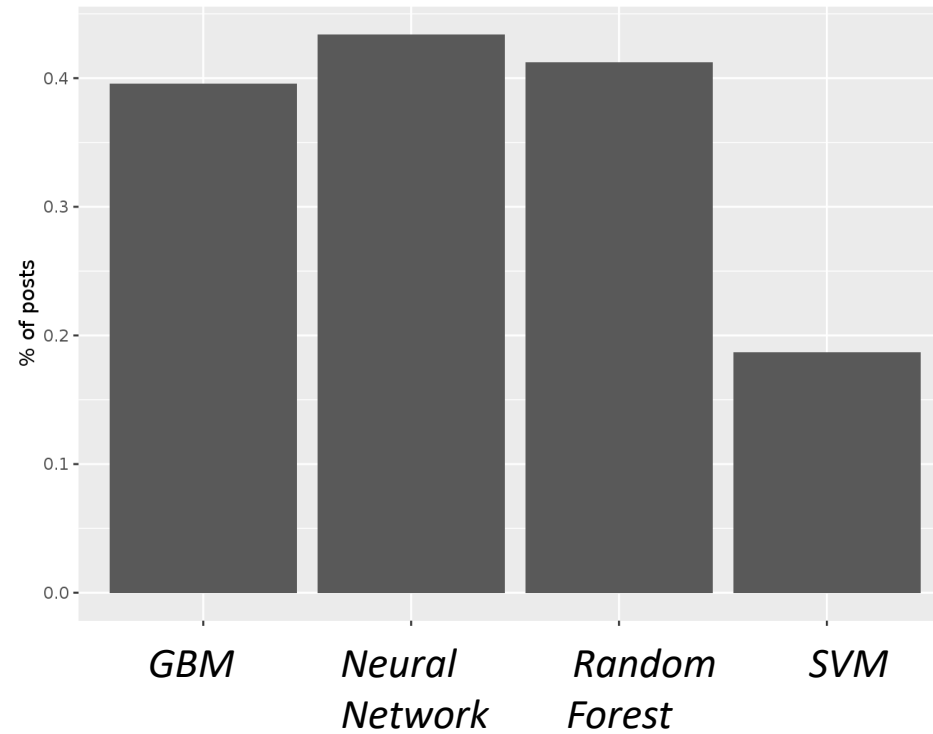
1. Decision Trees
2. Ensemble approaches

# Classes of supervised models so far

1. Linear regression (regression)
2. Logistic regression (classification)
3. Class conditional gaussians (Gaussian Discriminant Analysis)
4. Neural networks (regression & classification)
5. Today: the "game of 20 questions" approach to modeling

# Decision trees underlie state-of-the-art methods

- Competitive model classes for classification and regression rely on decision trees (DT) (*random forest* and *gradient-boosted trees*).

- Kaggle competitions 2016:

- DT-based models work well with default hyper-parameters.

- Are highly interpretable.

- NNs may generally perform better.



GBM= "Gradient Boosted Machine" (*ensemble* of Decision Trees)

# The Game of 20 Questions

| mpg | cylinders | displacement | horsepower | weight | acceleration | modelyear | maker |
|-----|-----------|--------------|------------|--------|--------------|-----------|-------|
| good | 4 | low | low | low | high | 75to78 | asia |
| bad | 6 | medium | medium | medium | medium | 70to74 | america |
| bad | 4 | medium | medium | medium | low | 75to78 | europe |
| bad | 8 | high | high | high | low | 70to74 | america |

$Y$        $X$

20Q.net Inc.

← → C 🔲 www.20q.net

▦ Apps 📁 NYC Stuff 📁 Research 📁 Paris Stuff Ⓜ Кино онлайн беспл... 📁 teaching 📁 B

**20Q** *the neural-net on the internet.*

**Play**

? **Play 20Q**
**About Us**
**Products**
**More ...**

games played online
87,115,788

I Can Read Your ...

**Q11. I am guessing that it is a pomegranate?**
**Right, Wrong, Close**

10. Can you order it at a restaurant? **No**.
9. Does it open? **Yes**.
8. Does it have lots of seeds? **Yes**.
7. Is it red? **Yes**.
6. Does it come from a plant? **Yes**.
5. Does it have bumpy skin? **No**.
4. Does it grow on a tree? **Yes**.
3. Is it made in many different styles? **No**.
2. Does it have leaves? **No**.
1. It is classified as **Vegetable**.

- How might you use something like this as a classifier?
- How would a sample get classified?
- How would you train this model?

# Decision Tree running example: classify fuel efficiency

- 40 data points

- Goal: predict MPG

- Need to find: $f : X \rightarrow Y$

- Discrete data (for now)

| mpg | cylinders | displacement | horsepower | weight | acceleration | modelyear | maker |
|-----|-----------|--------------|------------|--------|--------------|-----------|-------|
|     |           |              |            |        |              |           |       |
| good | 4 | low | low | low | high | 75to78 | asia |
| bad | 6 | medium | medium | medium | medium | 70to74 | america |
| bad | 4 | medium | medium | medium | low | 75to78 | europe |
| bad | 8 | high | high | high | low | 70to74 | america |
| bad | 6 | medium | medium | medium | medium | 70to74 | america |
| bad | 4 | low | medium | low | medium | 70to74 | asia |
| bad | 4 | low | medium | low | low | 70to74 | asia |
| bad | 8 | high | high | high | low | 75to78 | america |
| : | : | : | : | : | : | : | : |
| : | : | : | : | : | : | : | : |
| : | : | : | : | : | : | : | : |
| bad | 8 | high | high | high | low | 70to74 | america |
| good | 8 | high | medium | high | high | 79to83 | america |
| bad | 8 | high | high | high | low | 75to78 | america |
| good | 4 | low | low | low | low | 79to83 | america |
| bad | 6 | medium | medium | medium | high | 75to78 | america |
| good | 4 | medium | low | low | low | 79to83 | america |
| good | 4 | low | low | medium | high | 79to83 | america |
| bad | 8 | high | high | high | low | 70to74 | america |
| good | 4 | low | medium | low | medium | 75to78 | europe |
| bad | 5 | medium | medium | medium | medium | 75to78 | europe |

$Y$

$X$

From the UCI repository (thanks to Ross Quinlan)

# Hypotheses: decision trees $f : X \to Y$

| mpg | cylinders | displacement | horsepower | weight | acceleration | modelyear | maker |
|-----|-----------|--------------|------------|--------|--------------|-----------|-------|
| good | 4 | low | low | low | high | 75to78 | asia |
| bad | 6 | medium | medium | medium | medium | 70to74 | america |
| bad | 4 | medium | medium | medium | low | 75to78 | europe |
| bad | 8 | high | high | high | low | 70to74 | america |

$Y$=fuel efficiency (mpg)
$\in \{good, bad\}$

$Y$     $X$

- Each internal node tests an attribute $x_i$

- Each branch "picks off" an attribute value $x_i = v$

- Each leaf assigns a class $y$

- To classify input $x$: traverse the tree from root to leaf, output the labeled $y$
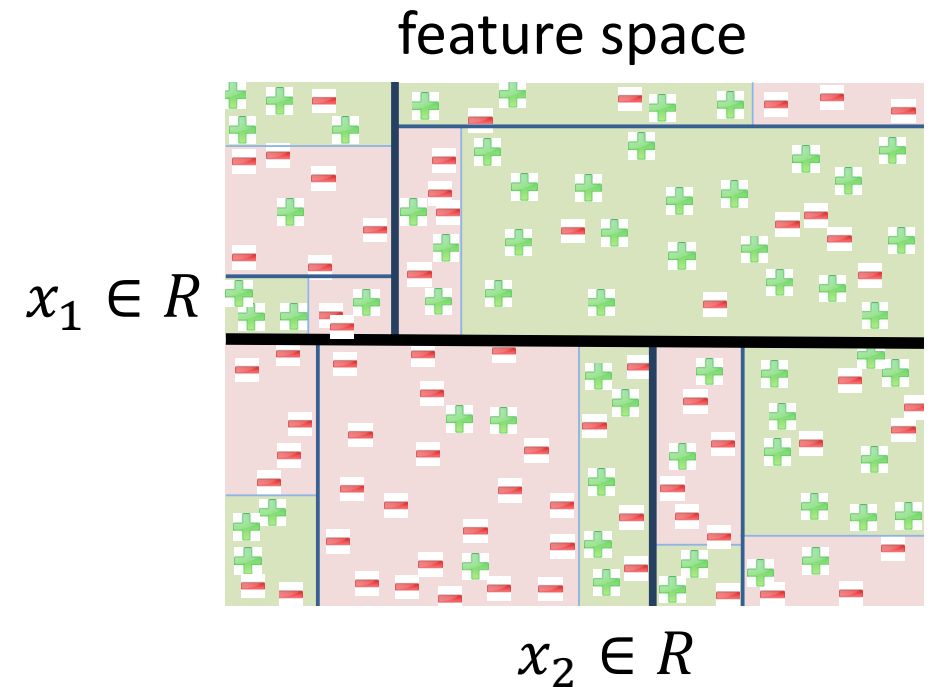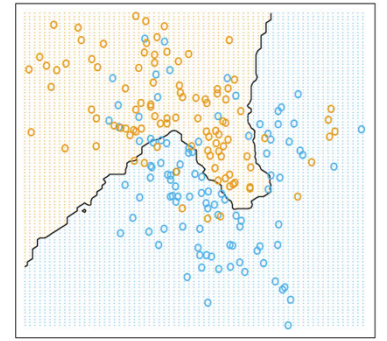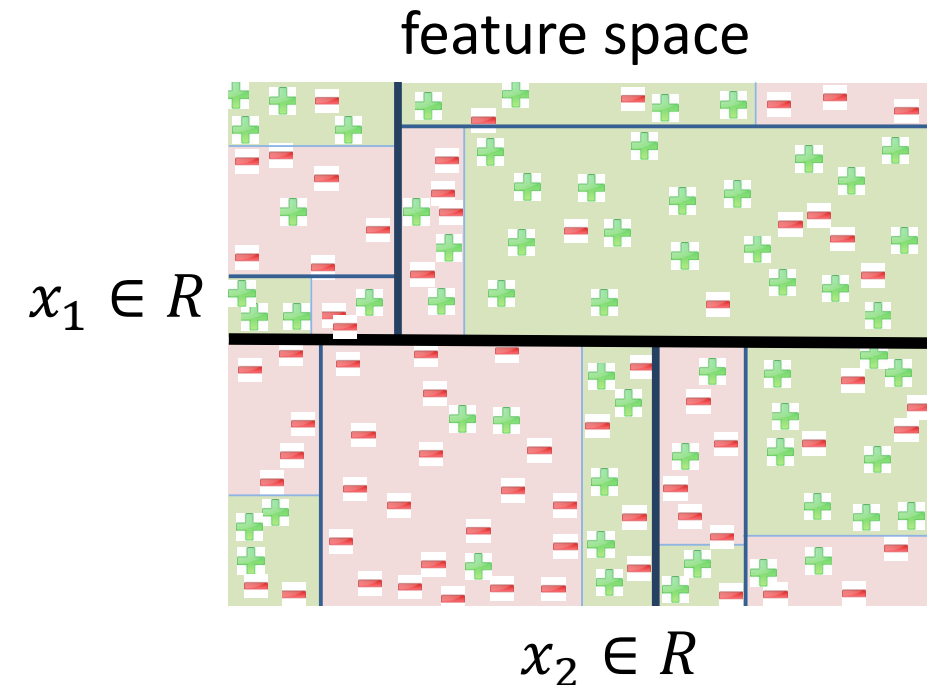
# DT classification boundaries

What would this kind of diagram look like for a decision tree?



(recall 15-NN classifier)

- The partitions would be …axis-aligned
  i.e., no diagonals boundaries

feature space



$x_1 \in R$

$x_2 \in R$

figure credit: Yisong Yue

# DT classification boundaries

What would this kind of diagram look like for a decision tree?



(recall 15-NN classifier)

- The partitions would be …axis-aligned
    i.e., no diagonals boundaries
- It would be a "piecewise static" function class
    i.e., each partition has a static prediction.



feature space
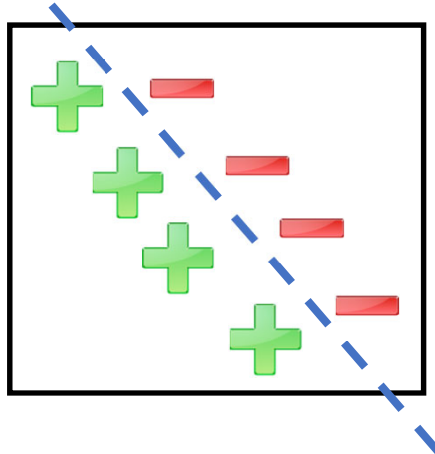
$x_1 \in R$

$x_2 \in R$

fitgure credit: Yisong Yue

# Decision Trees vs Linear Models
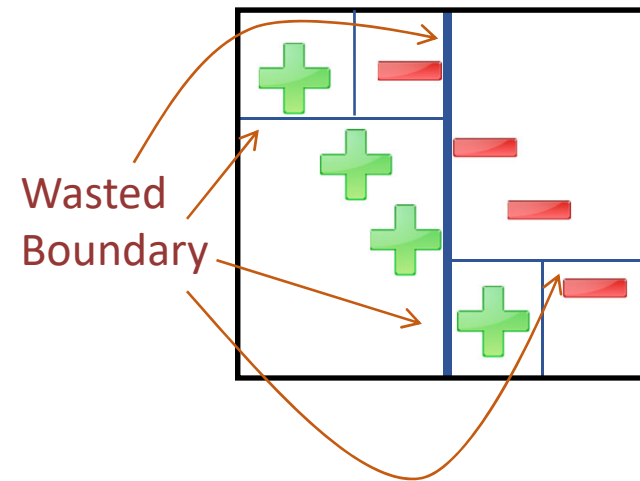
Decision Trees are AXIS-ALIGNED!

- Cannot easily model diagonal boundaries
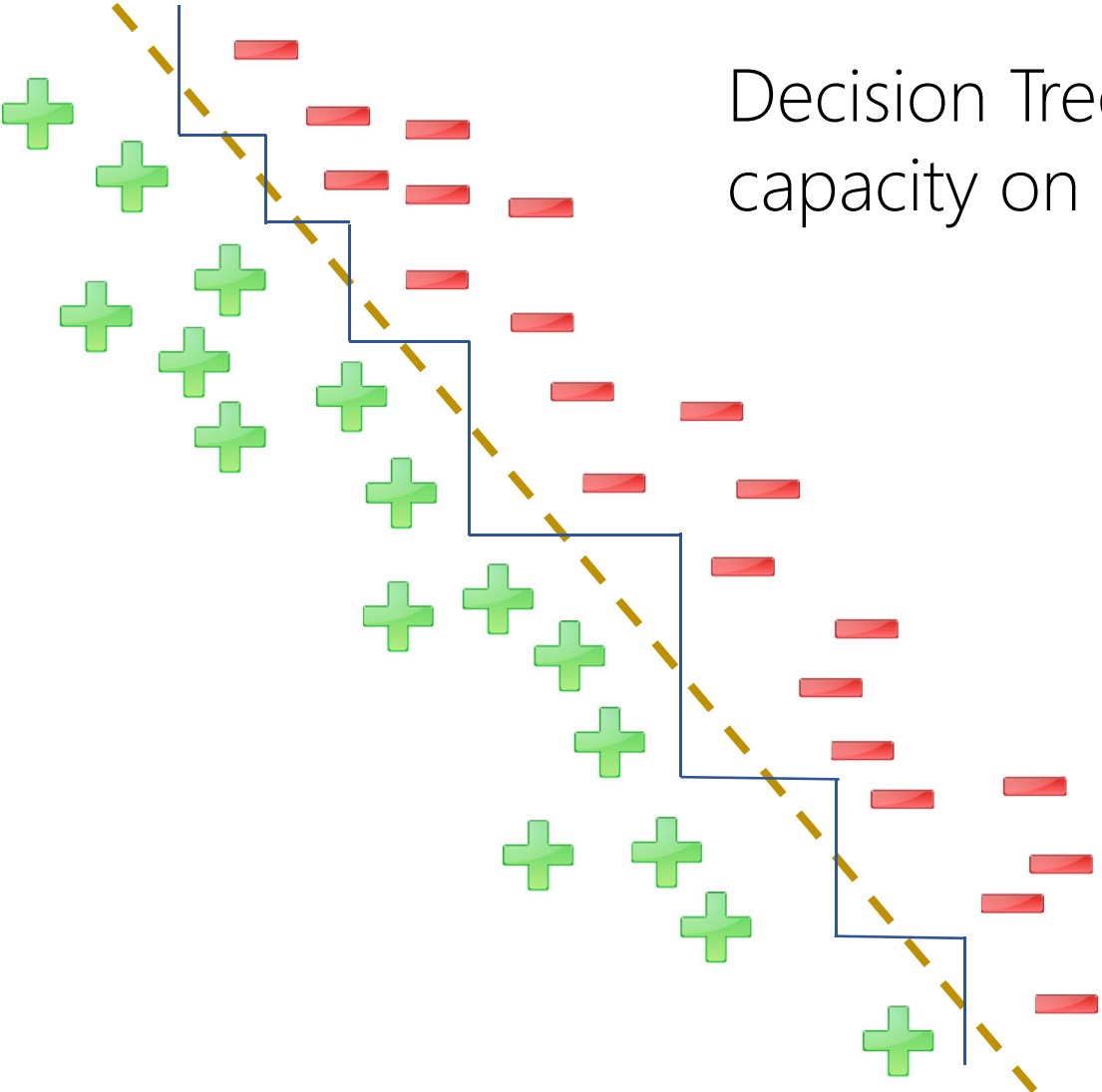
Example:



Logistic regression can have arbitrary linear boundary

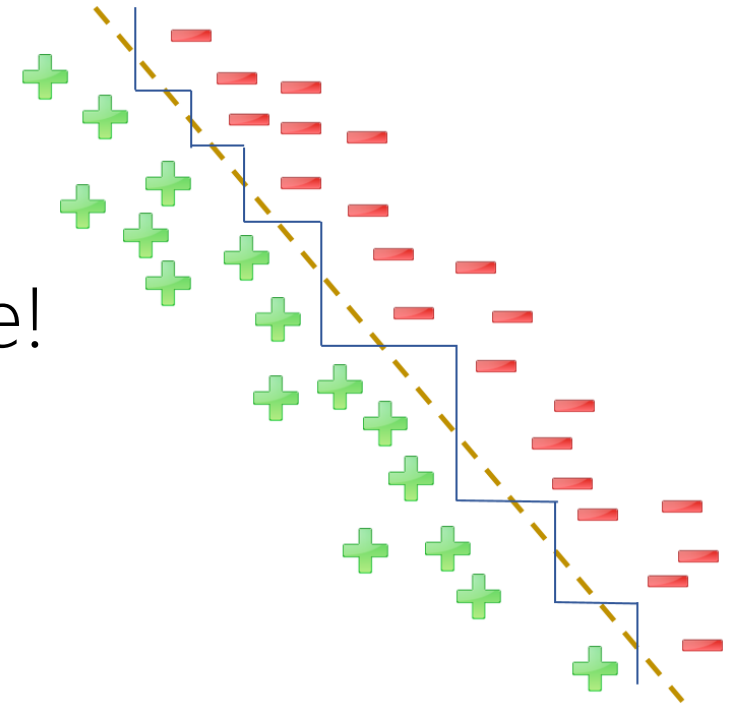Decision Trees Require Complex Axis-Aligned Partitioning

Wasted Boundary

# More Extreme Example

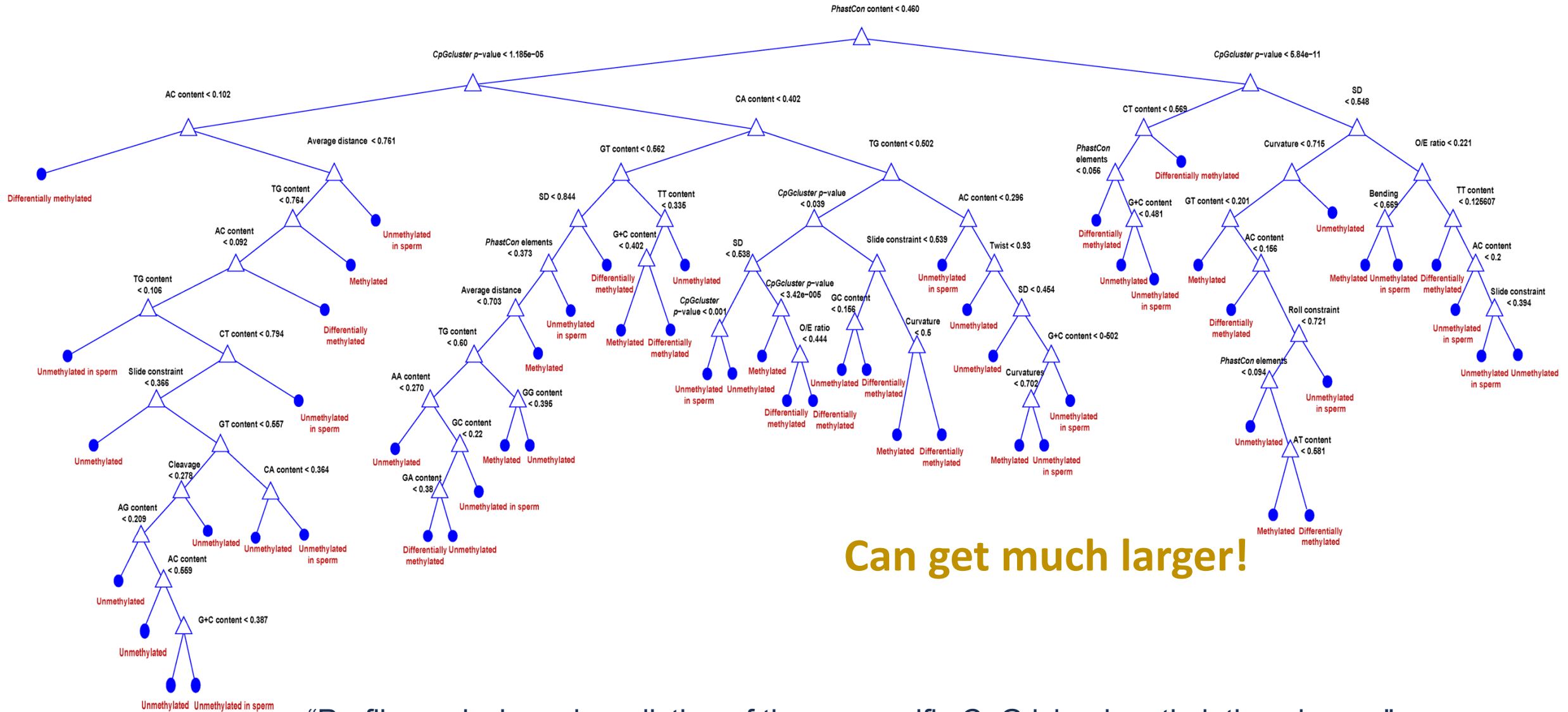Decision Tree can waste a lot of model capacity on useless boundaries.

# Decision Trees vs Linear Models

- Still, Decision Trees are often more accurate!
- Their non-linearity powerful.
- Catch: requires sufficient training data.
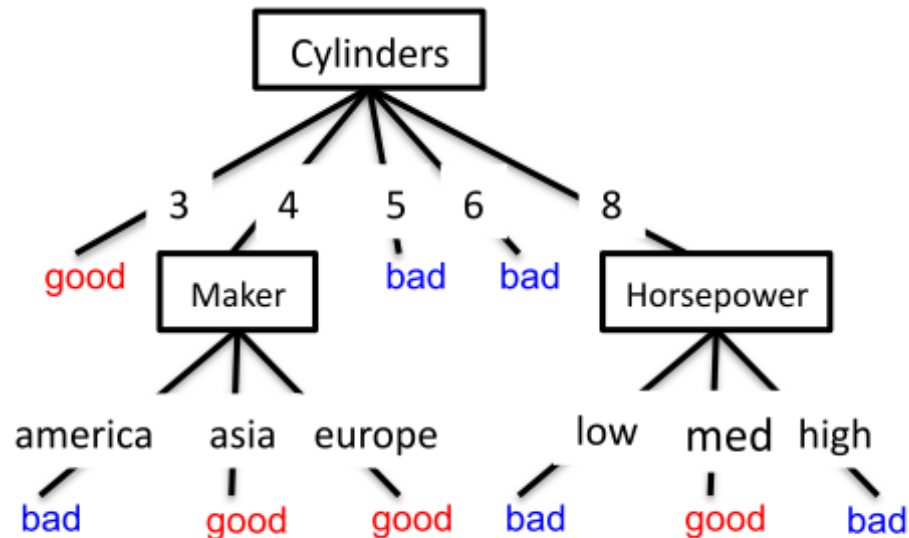
# Real Decision Trees



**Can get much larger!**

"Profile analysis and prediction of tissue-specific CpG island methylation classes"
Image Source: http://www.biomedcentral.com/1471-2105/10/116

# Model space

- How many possible models?

- What functions can be represented?

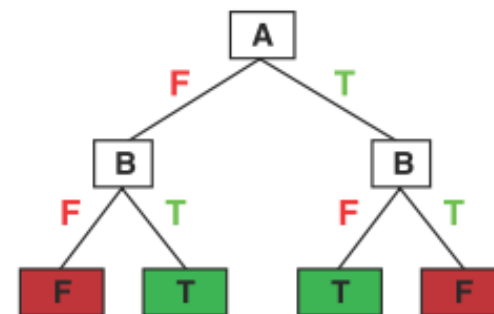| mpg | cylinders | displacement | horsepower | weight | acceleration | modelyear | maker |
|-----|-----------|--------------|------------|--------|--------------|-----------|-------|
| | | | | | | | |
| good | 4 | low | low | low | high | 75to78 | asia |
| bad | 6 | medium | medium | medium | medium | 70to74 | america |
| bad | 4 | medium | medium | medium | low | 75to78 | europe |
| bad | 8 | high | high | high | low | 70to74 | america |
| bad | 6 | medium | medium | medium | medium | 70to74 | america |
| bad | 4 | low | medium | low | medium | 70to74 | asia |
| bad | 4 | low | medium | low | low | 70to74 | asia |
| bad | 8 | high | high | high | low | 75to78 | america |
| : | : | : | : | : | : | : | : |
| : | : | : | : | : | : | : | : |
| : | : | : | : | : | : | : | : |
| bad | 8 | high | high | high | low | 70to74 | america |
| good | 8 | high | medium | high | high | 79to83 | america |
| bad | 8 | high | high | high | low | 75to78 | america |
| good | 4 | low | low | low | low | 79to83 | america |
| bad | 6 | medium | medium | medium | high | 75to78 | america |
| good | 4 | medium | low | low | low | 79to83 | america |
| good | 4 | low | low | medium | high | 79to83 | america |
| bad | 8 | high | high | high | low | 70to74 | america |
| good | 4 | low | medium | low | medium | 75to78 | europe |
| bad | 5 | medium | medium | medium | medium | 75to78 | europe |

$Y$  $X$



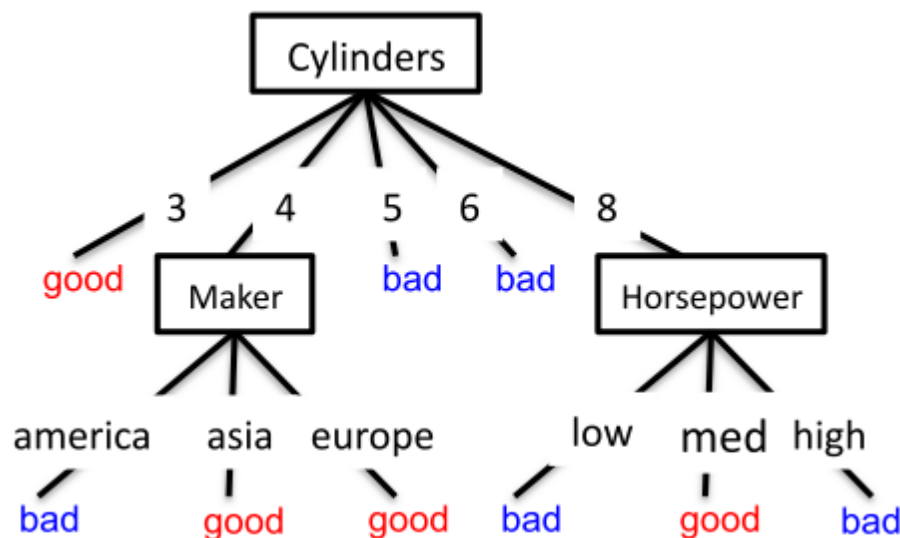Can you think of a function that a decision tree could not represent?

# What functions can be represented?

- Decision trees can represent any function of the input attributes!

- For Boolean functions, path to leaf gives truth table row

- But, could require exponentially many nodes...

| A | B | A xor B |
|---|---|---------|
| F | F | F |
| F | T | T |
| T | F | T |
| T | T | F |



(Figure from Stuart Russell)



cyl=3 ∨ (cyl=4 ∧ (maker=asia ∨ maker=europe)) ∨ …

# Hypothesis space

- How many possible lmodels?

- What functions can be represented?

- How many will be consistent with a given dataset?

- How will we choose the best one?



Lets first consider how to build a tree from training data, then revisit this question.

# What is the Simplest Tree?

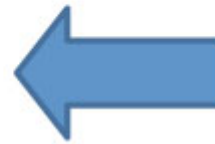| mpg | cylinders | displacement | horsepower | weight | acceleration | modelyear | maker |
|-----|-----------|--------------|------------|--------|--------------|-----------|-------|
| good | 4 | low | low | low | high | 75to78 | asia |
| bad | 6 | medium | medium | medium | medium | 70to74 | america |
| bad | 4 | medium | medium | medium | low | 75to78 | europe |
| bad | 8 | high | | | | 70to74 | america |
| bad | 6 | medium | | | | to74 | america |
| bad | 4 | low | | | | to74 | asia |
| bad | 4 | low | | | | to74 | asia |
| bad | 8 | high | | | | to78 | america |
| : | : | : | | | | : | : |
| : | : | : | | | | : | : |
| : | : | : | | | | : | : |
| bad | 8 | high | | | | to74 | america |
| good | 8 | high | | | | to83 | america |
| bad | 8 | high | | | | to78 | america |
| good | 4 | low | low | low | low | 79to83 | america |
| bad | 6 | medium | medium | medium | high | 75to78 | america |
| good | 4 | medium | low | low | low | 79to83 | america |
| good | 4 | low | low | medium | high | 79to83 | america |
| bad | 8 | high | high | high | low | 70to74 | america |
| good | 4 | low | medium | low | medium | 75to78 | europe |
| bad | 5 | medium | medium | medium | medium | 75to78 | europe |

**root**

**22**   **18**

**pchance = 0.001**

$Y$                          $X$

root node only
predict majority class:
    mpg=bad

## Is this a good tree?

[22+, 18-]  ⬅  Means:
correct on 22 examples
incorrect on 18 examples

# Next simplest tree: A Decision Stump
## (one feature splitting node)

| mpg | cylinders | displacement | horsepower | weight | acceleration | modelyear | maker |
|-----|-----------|--------------|------------|--------|--------------|-----------|-------|
| good | 4 | low | low | low | high | 75to78 | asia |
| bad | 6 | medium | medium | medium | medium | 70to74 | america |
| bad | 4 | medium | medium | medium | low | 75to78 | europe |
| bad | 8 | high | high | high | low | 70to74 | america |

$Y$  $X$

mpg values:  bad  good

root

22  18

pchance = 0.001

| cylinders = 3 | cylinders = 4 | cylinders = 5 | cylinders = 6 | cylinders = 8 |
|---------------|---------------|---------------|---------------|---------------|
| 0  0 | 4  17 | 1  0 | 8  0 | 9  1 |
| Predict bad | Predict good | Predict bad | Predict bad | Predict bad |

# Increase complexity by recursive partitioning

# Increase complexity by recursive partitioning

# Now have second level of tree



mpg values:  bad  good

root
22  18
pchance = 0.001

cylinders = 3
0  0
Predict bad

cylinders = 4
4  17
pchance = 0.135

cylinders = 5
1  0
Predict bad

cylinders = 6
8  0
Predict bad

cylinders = 8
9  1
pchance = 0.085

maker = america
0  10
Predict good

maker = asia
2  5
Predict good

maker = europe
2  2
Predict bad

horsepower = low
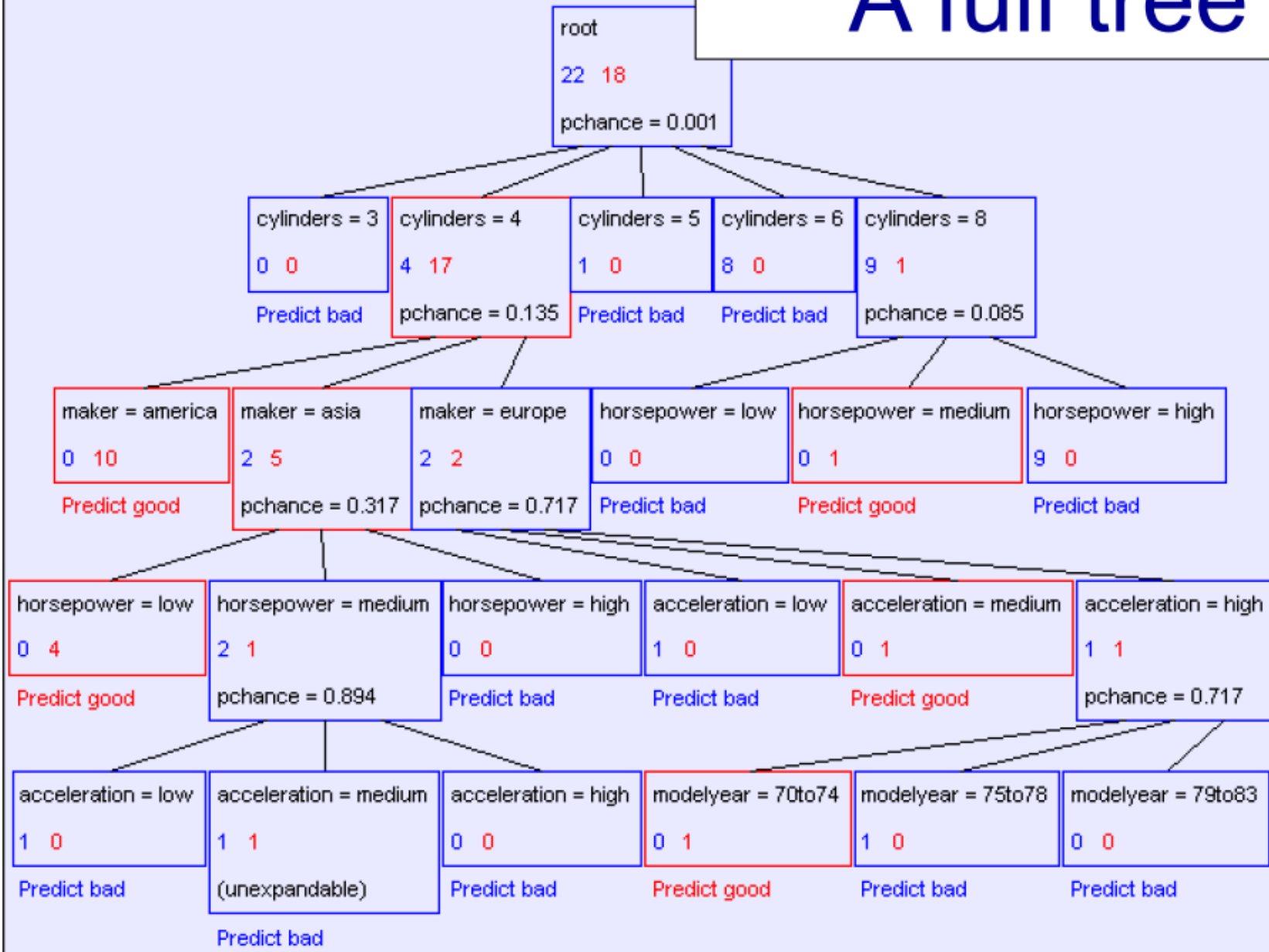0  0
Predict bad

horsepower = medium
0  1
Predict good

horsepower = high
9  0
Predict bad

Recursively build a tree from the seven
records in which there are four cylinders
and the maker was based in Asia

mpg values: bad good

# A full tree

root
22 18
pchance = 0.001

cylinders = 3
0 0
Predict bad

cylinders = 4
4 17
pchance = 0.135

cylinders = 5
1 0
Predict bad

cylinders = 6
8 0
Predict bad

cylinders = 8
9 1
pchance = 0.085

maker = america
0 10
Predict good

maker = asia
2 5
pchance = 0.317

maker = europe
2 2
pchance = 0.717

horsepower = low
0 0
Predict bad

horsepower = medium
0 1
Predict good

horsepower = high
9 0
Predict bad

horsepower = low
0 4
Predict good

horsepower = medium
2 1
pchance = 0.894

horsepower = high
0 0
Predict bad

acceleration = low
1 0
Predict bad

acceleration = medium
0 1
Predict good

acceleration = high
1 1
pchance = 0.717

acceleration = low
1 0
Predict bad

acceleration = medium
1 1
(unexpandable)
Predict bad

acceleration = high
0 0
Predict bad

modelyear = 70to74
0 1
Predict good

modelyear = 75to78
1 0
Predict bad

modelyear = 79to83
0 0
Predict bad
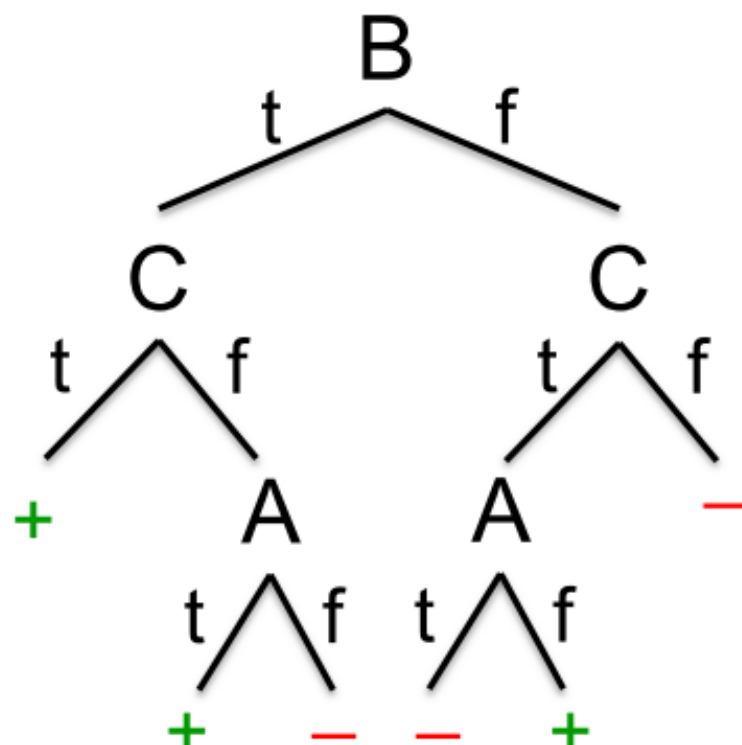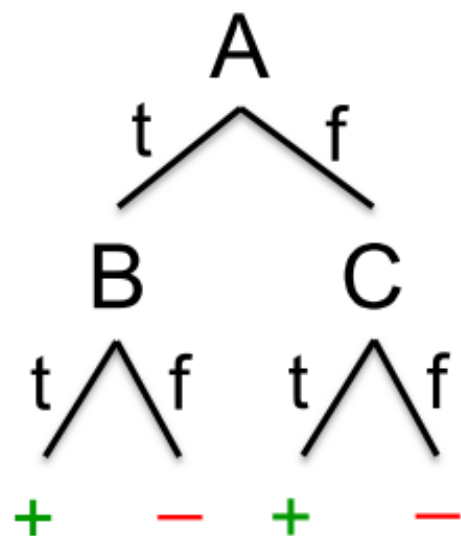
- Each leaf has only one example, or is "unexpandable".
- i.e., the value for any unused features is constant.

- Many trees can represent the same concept
- But, not all trees will have the same size!
  - e.g., $\phi = (A \wedge B) \vee (\neg A \wedge C)$



Which tree do we prefer?

# "Occam's razor"

- "Nunquam ponenda est pluralitis sin necesitate"



*English Franciscan friar William of Ockham (c. 1287–1347)*

- "Entities should not be multiplied beyond necessity"

- "when you have two competing theories that make exactly the same predictions, the simpler one is the better"

# Why is Occam's razor a reasonable heuristic for decision tree learning?

- there are fewer short models (i.e. small trees) than long ones

- a short model is unlikely to fit the training data well by chance

- a long model is more likely to fit the training data well coincidentally

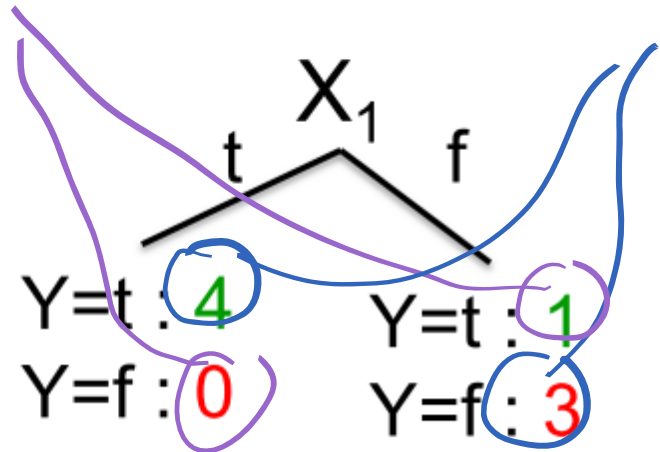# Learning simplest decision tree is hard

- Learning the simplest (smallest) decision tree is an NP-complete problem [Hyafil & Rivest '76]

- Resort to a greedy heuristic:
  - Start from empty decision tree
  - Split on **next best attribute (feature)**
  - Recurse
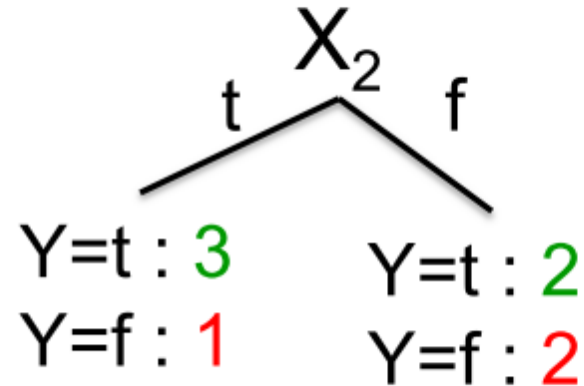
*but what does this mean?*

# How to pick a good feature to split on?

Would we prefer to split on $X_1$ or $X_2$?

incorrect

correct

$X_1$

t        f

Y=t : 4        Y=t : 1
Y=f : 0        Y=f : 3

$X_2$

t        f

Y=t : 3        Y=t : 2
Y=f : 1        Y=f : 2

| $X_1$ | $X_2$ | Y |
|-------|-------|---|
| T | T | T |
| T | F | T |
| T | T | T |
| T | F | T |
| F | T | T |
| F | F | F |
| F | T | F |
| F | F | F |

- Idea: use counts at leaves to define $p(Y)$.
- Goal: find distributions where we're not surprised to see the label.
- i.e. those with low entropy.

# Entropy: a measure of expected surprise

Think about a flipping a coin once, and how surprised you would be at observing a head.

$p(head) = 0$

$p(head) = 0.5$

$p(head) = 0.01$

$p(head) = 1$

# Entropy: a measure of expected surprise

- The "surprise" of observing that a discrete random variable (RV) $Y$ takes on value $k$ is:

$$log \frac{1}{P(Y = k)} = -\log(P(Y = k))$$

- As $P(Y = k) \to 0$, the surprise of observing $k$ approaches $\infty$.
- As $P(Y = k) \to 1$, the surprise of observing $k$ approaches $0$.
- The entropy of the distribution of $Y$ is the *expected surprise*:

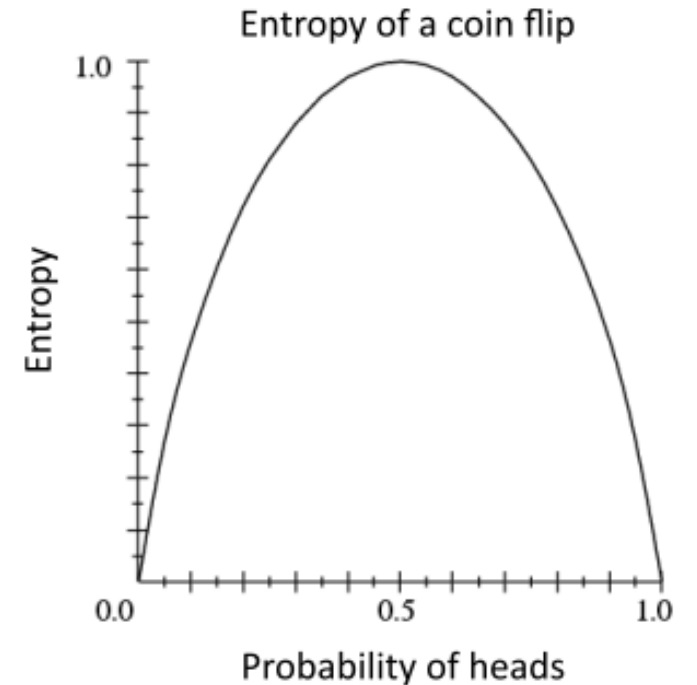$$H(Y) \equiv E_Y[-\log P(Y = k)] = \sum_k P(Y = k) \log P(Y = k)$$

# Entropy example: flipping a coin

$$H(Y) = -\sum_{i=1}^{\kappa} P(Y = y_i) \log_2 P(Y = y_i)$$

P(Y=t) = 5/6

P(Y=f) = 1/6

H(Y) = - 5/6 log$_2$ 5/6 - 1/6 log$_2$ 1/6

= 0.65

Entropy of a coin flip

Entropy

0.0          0.5          1.0
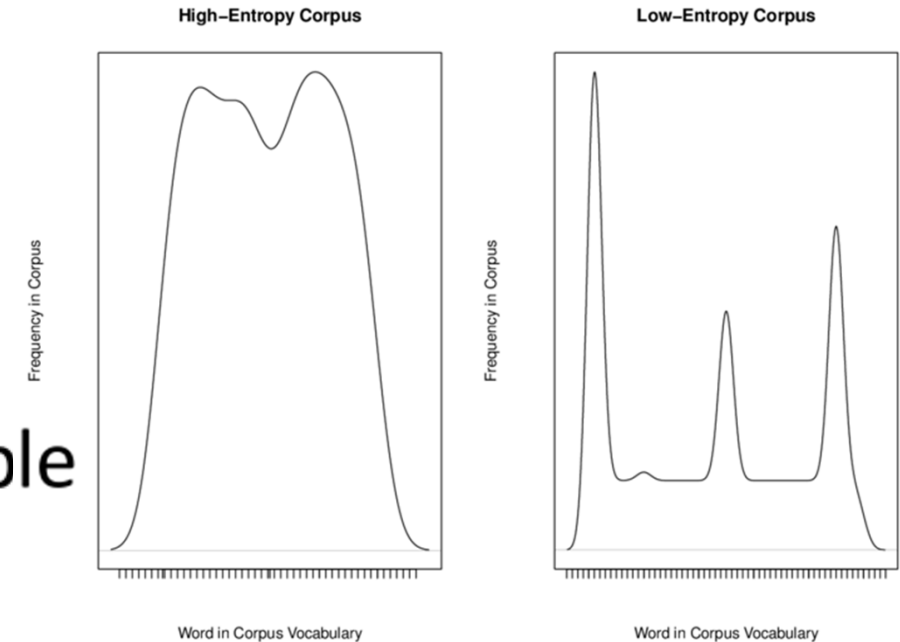
Probability of heads

# Entropy of a random variable $Y$:

"**High Entropy**"

– Y is from a uniform like distribution

– Flat histogram

– Values sampled from it are less predictable

"**Low Entropy**"

– Y is from a varied (peaks and valleys) distribution

– Histogram has          lows and highs
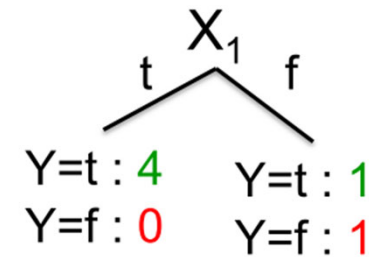
– Values sampled from it are more predictable



High–Entropy Corpus

Low–Entropy Corpus

Frequency in Corpus

Frequency in Corpus

Word in Corpus Vocabulary

Word in Corpus Vocabulary

*https://www.researchgate.net/figure/Hypothetical-distributions-of-term-frequency-in-high-and-low-entropy-corpora_fig1_305417514*

(Slide from Vibhav Gogate)

# Entropy for choosing feature to split next

| $X_1$ | $X_2$ | Y |
|---|---|---|
| T | T | T |
| T | F | T |
| T | T | T |
| T | F | T |
| F | T | T |
| F | F | F |

- Cannot control entropy of RV $Y$ (label).

- Can control the entropy of $p(Y|X)$, i.e. after splitting on feature(s).

- Goal: recursively reduce the *conditional entropy* at each node split until expected surprise at leaf nodes=0.

$X_1$

t     f

Y=t : 4     Y=t : 1
Y=f : 0     Y=f : 1

# Conditional entropy (for a tree split)

*weighted average of entropy on each side of the split of feature into $x_j < v$ and $x_j \geq v$*

$$H(Y|X_{j,v}) := P(X_{j,v} = 1)H(Y|X_{j,v} = 1) + P(X_{j,v} = 0)H(Y|X_{j,v} = 0)$$

| $X_1$ | $X_2$ | Y |
|-------|-------|---|
| T | T | T |
| T | F | T |
| T | T | T |
| T | F | T |
| F | T | T |
| F | F | F |

# Conditional entropy (for a tree split)

*weighted average of entropy on each side of the split of feature into $x_j < v$ and $x_j \geq v$*

$$H(Y|X_{j,v}) := P(X_{j,v} = 1)H(Y|X_{j,v} = 1) + P(X_{j,v} = 0)H(Y|X_{j,v} = 0)$$

**Example:**

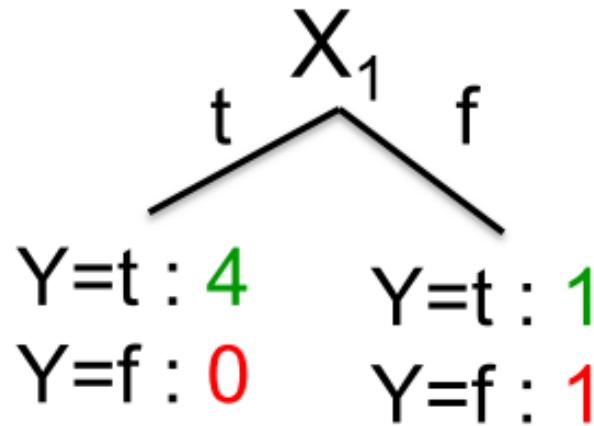$$H(Y) = -\sum_{i=1}^{\kappa} P(Y = y_i) \log_2 P(Y = y_i)$$

$X_1$

t    f

Y=t : 4     Y=t : 1
Y=f : 0     Y=f : 1

$P(X_1=t) = 4/6$

$P(X_1=f) = 2/6$

$H(Y|X_1) = - 4/6 \ (1 \log_2 1 + 0 \log_2 0)$

$\qquad\qquad\quad - 2/6 \ (1/2 \log_2 1/2 + 1/2 \log_2 1/2)$

$\qquad = 2/6$

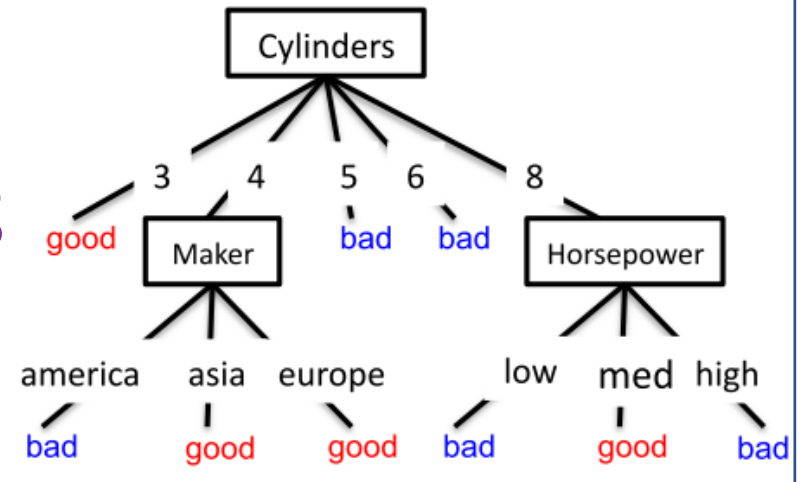| $X_1$ | $X_2$ | Y |
|-------|-------|---|
| T | T | T |
| T | F | T |
| T | T | T |
| T | F | T |
| F | T | T |
| F | F | F |

# Conditional entropy (for a tree split)

*weighted average of entropy on each side of the split of feature into $x_j < v$ and $x_j \geq v$*

$$H(Y|X_{j,v}) := P(X_{j,v} = 1)H(Y|X_{j,v} = 1) + P(X_{j,v} = 0)H(Y|X_{j,v} = 0)$$

Goal: at each recursion, find the feature (and split) which minimizes the conditional entropy.



$H(Y|X_1) = -4/6 (1 \log_2 1 + 0 \log_2 0)$

$\qquad - 2/6 (1/2 \log_2 1/2 + 1/2 \log_2 1/2)$

$\qquad = 2/6$

# Equivalently: maximize the *information gain*

(also called the *mutual information*)

$$\text{maximize } I(X_{j,v}; Y) := H(Y) - H(Y|X_{j,v})$$

constant
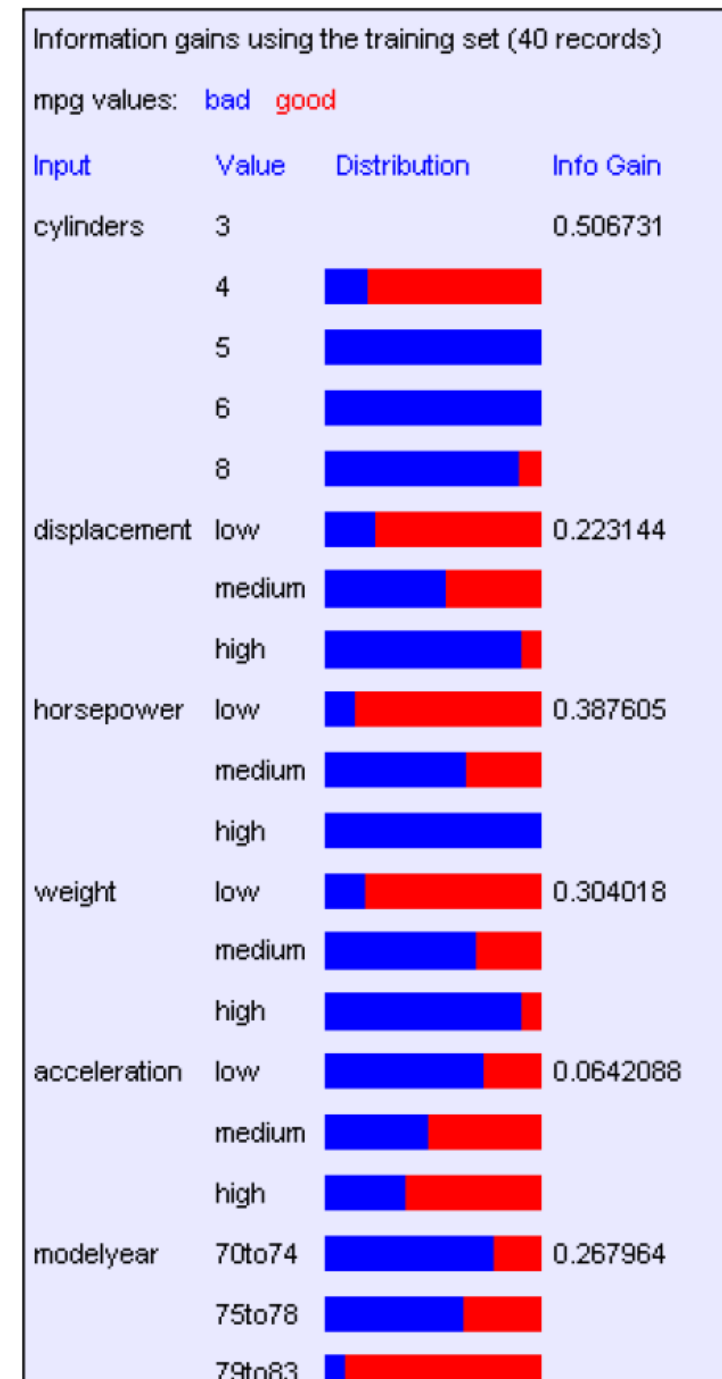
# "Learning" Decision Trees (aka building from data)

- Start from empty decision tree

- Split on **next best attribute (feature)**

  - Use, for example, information gain to select attribute:

$$\arg\max_i IG(X_i) = \arg\max_i H(Y) - H(Y \mid X_i)$$
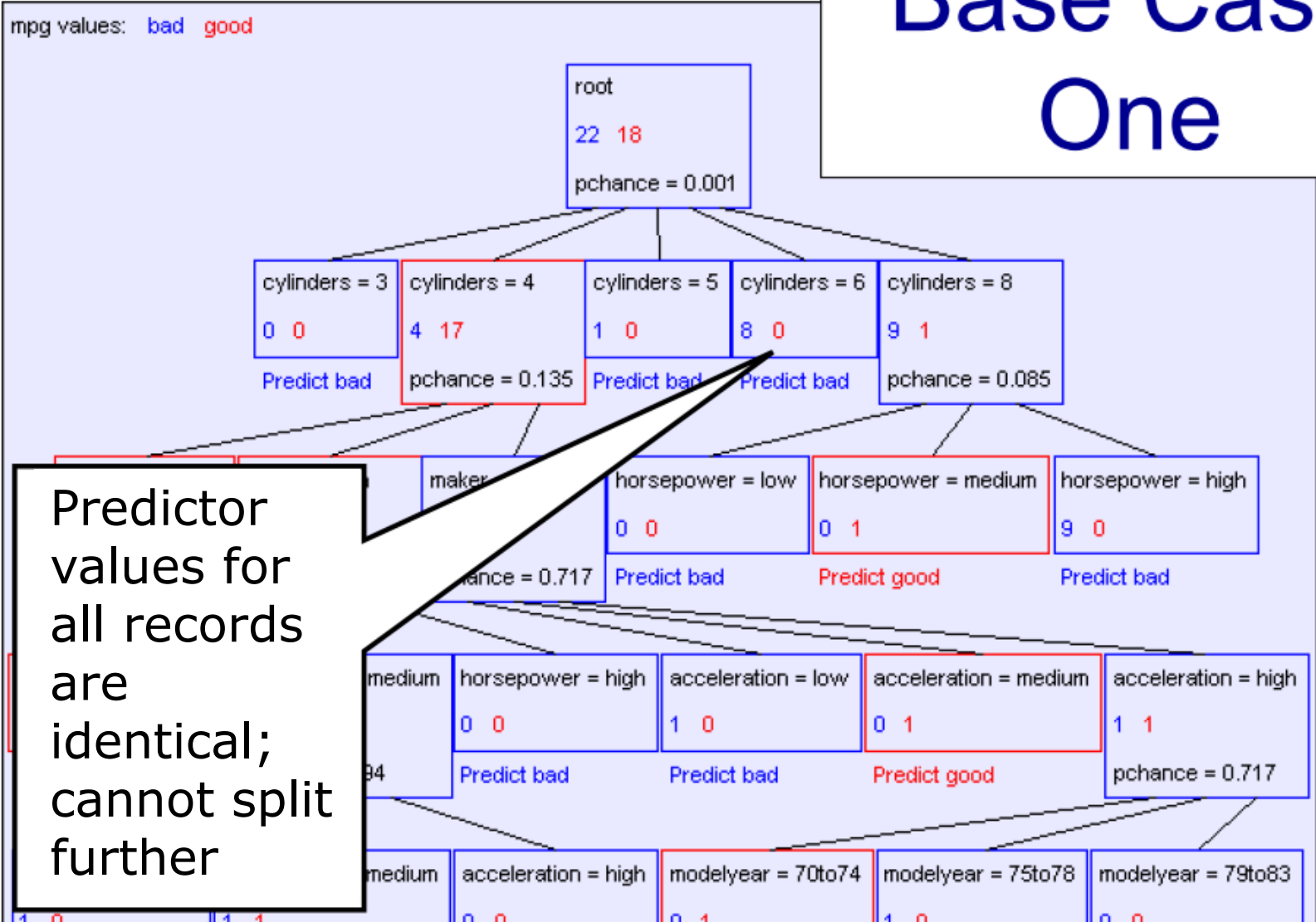
- Recurse

# Back to the "mpg" example

- Compute information gain for every possible split at every node.
- Categorical features: split into each category, or could do one vs. rest.
- Real-valued features: pick a threshold to split on; try different thresholds.



Information gains using the training set (40 records)

mpg values: bad good

| Input | Value | Distribution | Info Gain |
|---|---|---|---|
| cylinders | 3 | | 0.506731 |
| | 4 | | |
| | 5 | | |
| | 6 | | |
| | 8 | | |
| displacement | low | | 0.223144 |
| | medium | | |
| | high | | |
| horsepower | low | | 0.387605 |
| | medium | | |
| | high | | |
| weight | low | | 0.304018 |
| | medium | | |
| | high | | |
| acceleration | low | | 0.0642088 |
| | medium | | |
| | high | | |
| modelyear | 70to74 | | 0.267964 |
| | 75to78 | | |
| | 79to83 | | |

# When to stop recursing on a node?



Base Case One

mpg values: bad good

root
22 18
pchance = 0.001

| cylinders = 3 | cylinders = 4 | cylinders = 5 | cylinders = 6 | cylinders = 8 |
|---|---|---|---|---|
| 0 0 | 4 17 | 1 0 | 8 0 | 9 1 |
| Predict bad | pchance = 0.135 | Predict bad | Predict bad | pchance = 0.085 |

maker

horsepower = low
0 0
Predict bad

horsepower = medium
0 1
Predict good

horsepower = high
9 0
Predict bad

...ance = 0.717

medium

| horsepower = high | acceleration = low | acceleration = medium | acceleration = high |
|---|---|---|---|
| 0 0 | 1 0 | 0 1 | 1 1 |
| Predict bad | Predict bad | Predict good | pchance = 0.717 |

medium

acceleration = high

modelyear = 70to74

modelyear = 75to78

modelyear = 79to83

Predictor values for all records are identical; cannot split further

# When to stop recursing on a node?



Information gains using the training set (2 records)

mpg values: bad good

| Input | Value | Distribution | Info Gain |
|-------|-------|--------------|-----------|
| cylinders | 3 | | 0 |
| | 4 | | |
| | 5 | | |
| | 6 | | |
| | 8 | | |
| displacement | low | | 0 |
| | medium | | |
| | high | | |
| horsepower | low | | 0 |
| | medium | | |
| | high | | |
| weight | low | | 0 |
| | medium | | |
| | high | | |
| acceleration | low | | 0 |
| | medium | | |
| | high | | |
| modelyear | 70to74 | | 0 |
| | 75to78 | | |
| | 79to83 | | |
| maker | america | | 0 |
| | asia | | |
| | europe | | |

**Base Case Two**

Don't split a node if data points are identical on remaining attributes

mpg values: bad good

root
22 18
pchance = 0.001

cylinders = 3 — 0 0 — Predict bad
cylinders = 4 — 4 17 — pchance = 0.135
cylinders = 5 — 1 0 — Predict bad
cylinders = 6 — 8 0 — Predict bad
cylinders = 8 — 9 1 — pchance = 0.085

maker = america — 0 10 — Predict good
maker = asia — 2 5 — pchance = 0.317
maker = europe — 2 2 — pchance = 0.717
horsepower = low — 0 0 — Predict bad
horsepower = ... — 0 1 — Predict good

horsepower = low — 0 4 — Predict good
horsepower = medium — 2 1 — pchance = 0.894
horsepower = high — 0 0 — Predict bad
acceleration = low — 1 0 — Pr...

acceleration = low — 1 0 — Predict bad
acceleration = medium — 1 1 — (unexpandable) — Predict bad
...ation = high — 0 0 — Predict bad
modelyear = 70to74 — 0 1 — Predict good
modelyear = 75to78 — 1 0 — Predict bad
modelyear = 79to83 — 0 0 — Predict bad

# When to stop recursing on a node?

Base Case One: If all records in current data subset have the same output then don't recurse

Base Case Two: If all records have exactly the same set of input attributes then don't recurse

Proposed Base Case 3:
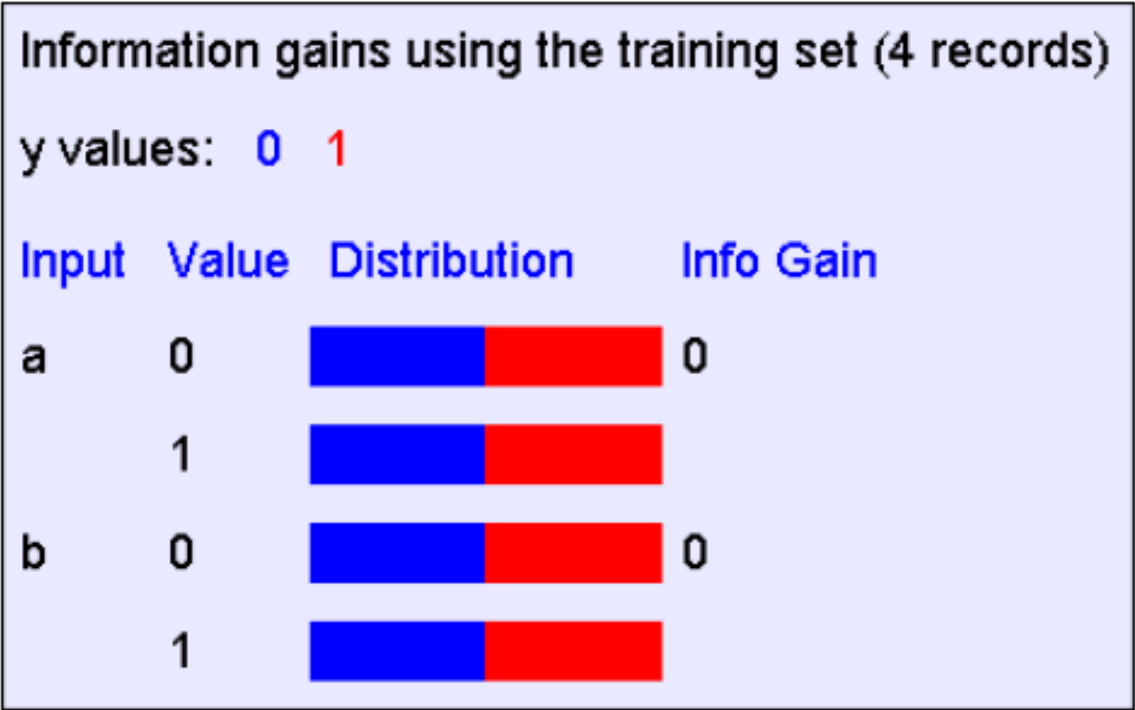If all attributes have zero information gain then don't recurse

- *Is this a good idea?*

# The problem with Base Case 3 for stopping

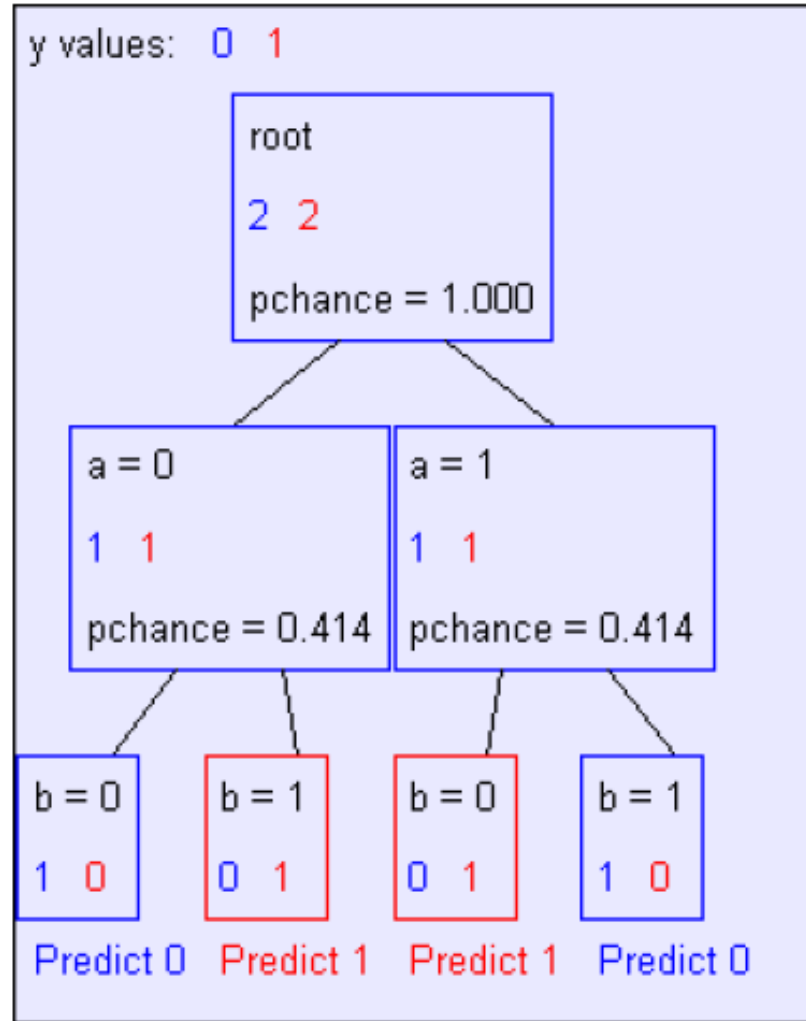We are using a greedy heuristic.

y = a XOR b

| a | b | y |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

Information gains using the training set (4 records)

y values:  0  1

| Input | Value | Distribution | Info Gain |
|-------|-------|--------------|-----------|
| a | 0 | | 0 |
|   | 1 | | |
| b | 0 | | 0 |
|   | 1 | | |

# If we omit Base Case 3 for stopping:

The resulting decision tree:



$y = a\ XOR\ b$

| a | b | y |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

Instead, perform *pruning* after building the tree using:
1. Statistical test
2. Hold out performance.

# Summary: building Decision Trees

*X*   *Y*

BuildTree(*DataSet,Output*)

- If all output values are the same in *DataSet*, return a leaf node that says "predict this unique output" **BASE CASE 1**

- If all input values are the same, return a leaf node that says **BASE CASE** "predict the majority output"

- Else find attribute *X* with highest Info Gain

- Suppose *X* has $n_X$ distinct values (i.e. X has arity $n$

  - Create a non-leaf node with $n_X$ children.
  - The *i*'th child should be built by calling

    BuildTree($DS_i$,*Output*)

    Where $DS_i$ contains the records in DataSet where X =

- What fundamental flaw does this have?

- It will keep going until the data are *perfectly* split/labelled.

# This algorithm will always overfit the data

Naive decision trees have no learning bias

– Training set error is always zero!

  • (If there is no label noise)

– Must introduce some bias towards simpler trees

# How to regularize tree-building

- Limit on depth of the tree.
- Min # data points at a node.
- Backward-greedy pruning (greedily remove node & descendants that most improves validation performance).
- Early stopping according to any number of criteria, such as a statistical test (but then subject to the problem of Base Case 3).
- Take an *ensemble* of small trees.

# Decision Tree for regression

- Prediction can be mean of those in the leaf "bucket".
- Or linear regression within a leaf bucket.
- Now choose nodes to split on by minimizing the sum of variances after a split:

$$\sum_{i:\ x_i \in R_1(j,s)} (y_i - \hat{y}_{R_1})^2 + \sum_{i:\ x_i \in R_2(j,s)} (y_i - \hat{y}_{R_2})^2$$

# CS 189/289

Today's lecture

1. Decision Trees
2. Ensemble approaches

# Bolstering Decision Trees

1. Decision trees can easily overfit if we don't regularize considerably.

2. Slightly different samples can lead to very different trees (high variance models).

3. If we average several randomized trees, we tend to do better: *Random Forests*.

4. Instantiation of broader, commonly used idea of *ensembling* models.

# Creating *ensemble* of DTs

1. Train $M$ models using $M$ samples of size $n'$ with replacement, from $n$ total data points.

(called *bootstrap replicates* of the data)

These trees will be weakly correlated.



Original Data

Bootstrapping

Aggregating

Bagging

Classifier

Classifier

Classifier

...

Ensemble classifier

# Creating *ensemble* of DTs

1. Train $M$ models using $M$ samples of size $n'$ with replacement, from $n$ total data points.

2. When building these $M$ DTs, allow use of only a random subset of the features (each time you recurse to create a new node), e.g. 2/3 of features.

3. 1 + 2 + averaging $M$ predictions is called a *Random Forest*.

4. 1 + averaging is called *Bagging* ("bootstrap aggregation")

   (These DTs are not uncorrelated, but randomization helps to make their errors uncorrelated, which provides the advantage in model averaging approaches.

# Empirical Comparison of Decision Boundaries

# Other forms of ensemble modelling

- So far we have considered *averaging* several models together to reduce the variance, $y = \frac{1}{m}\sum_{m=1}^{M} G_m(x)$.

- Why not consider a weighted average,
  $$y = \sum_{m=1}^{M} \alpha_m G_m(x)?$$

- Given a set of trained models of any kind (e.g. neural network, decision tree, ...), $\{G_m(x)\}$, if we then optimize a loss for the weights $\{\alpha_m\}$, this is called stacking.

- What about jointly optimizing $\{G_m(x)\}$ and $\{\alpha_m\}$?

- Difficult optimization problem, but...

# Boosting

- Too hard to jointly optimize $\{G_m(x)\}$ and $\{\alpha_m\}$ in
$$y = \sum_{m=1}^{M} \alpha_m G_m(x).$$

- Instead, lets use a greedy approximation wherein we sequentially train the next $G_m(x)$ conditioned on all previous learned base models $\{G_{m-1}, G_{m-2}, \dots, G_1\}$, and their corresponding weights, $\{\alpha_{m-1}, \alpha_{m-2}, \dots, \alpha_1\}$.

- The intuition is that at each iteration, we will re-weight the training points to focus on those that are not correctly classified.

# Boosting at an intuitive level (decision stump)



$$G(x) = \sum_{m=1}^{M} \alpha_m G_m(x)$$

# Boosting at an intuitive level (decision stump)
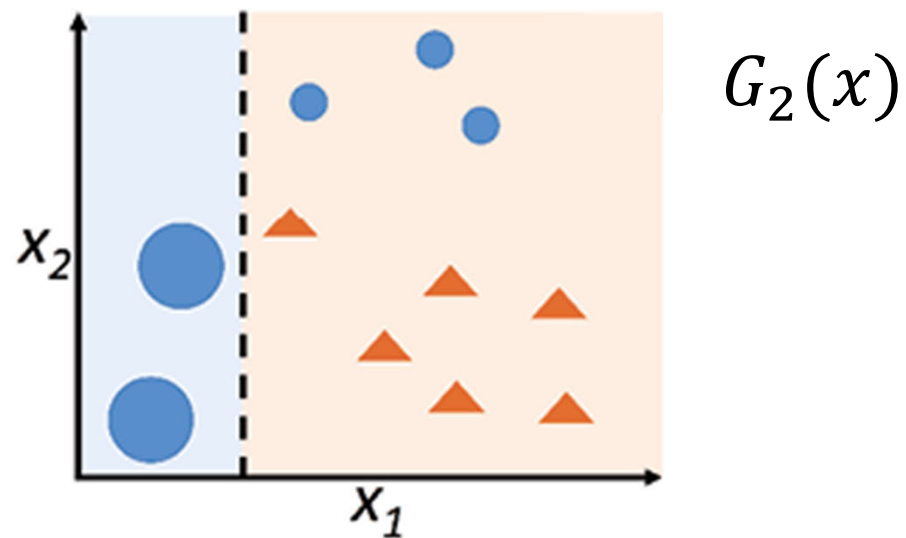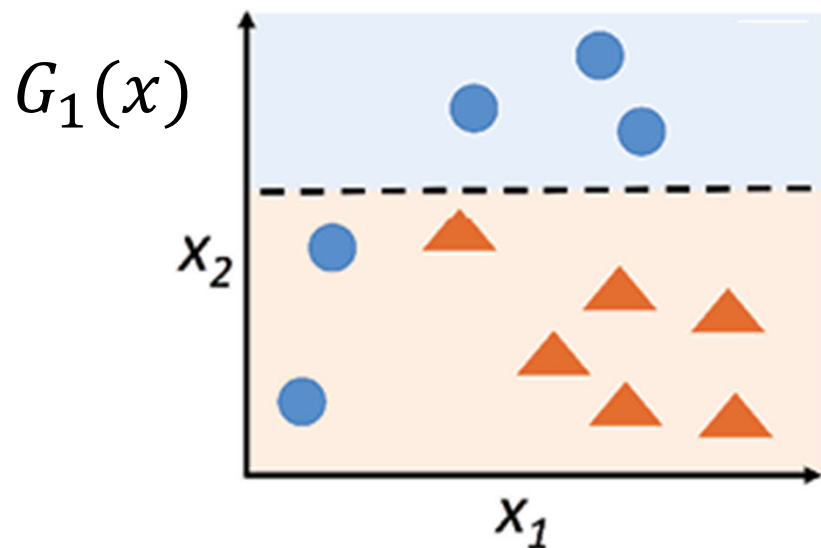


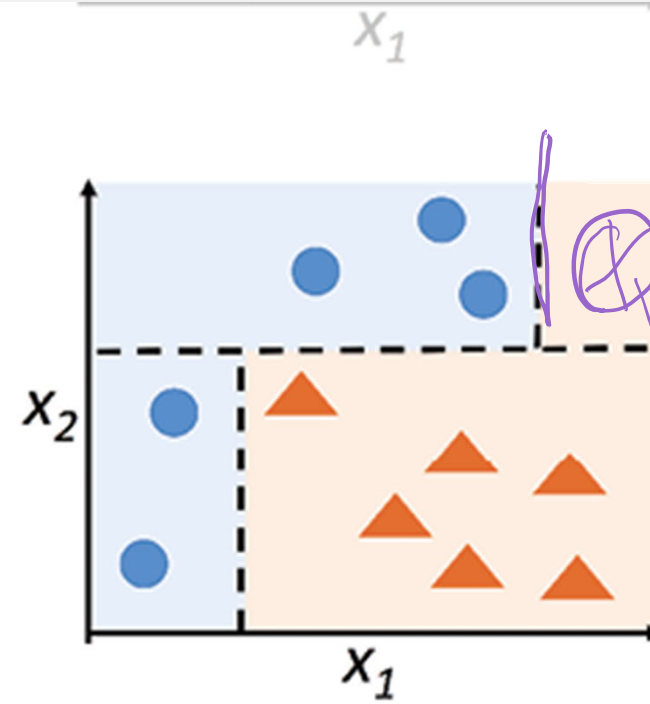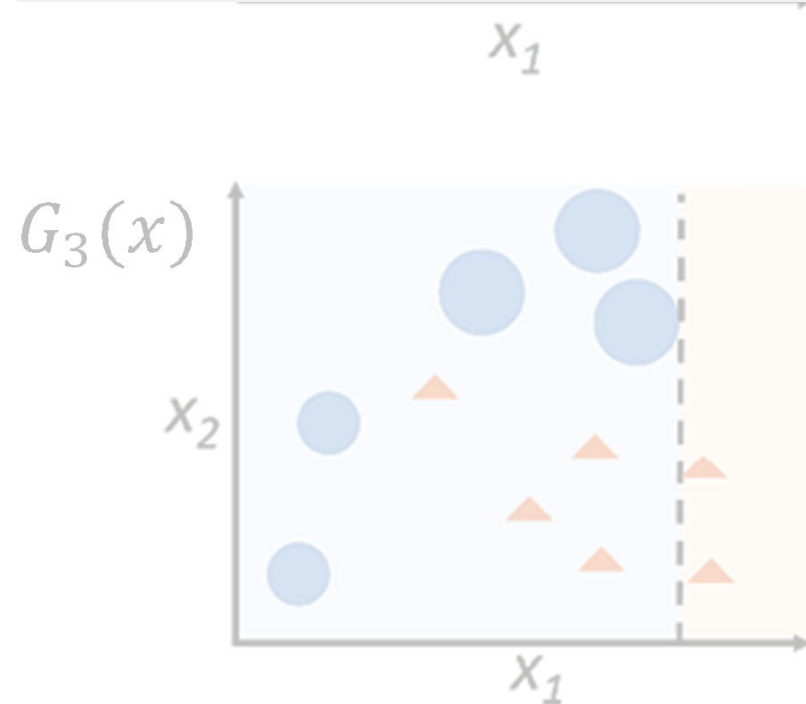$G_1(x)$

$G_2(x)$

$$G(x) = \sum_{m=1}^{M} \alpha_m G_m(x)$$

# Boosting at an intuitive level (decision stump)



$$G(x) = \sum_{m=1}^{M} \alpha_m G_m(x)$$

# Boosting at an intuitive level (decision stump)



$G_1(x)$

$G_2(x)$

$G_3(x)$

$$G(x) = \sum_{m=1}^{M} \alpha_m G_m(x)$$

# Boosting at an intuitive level (decision stump)

- Final boosted model looks a lot like a decision tree.
- Can you spot anything that makes you think a decision tree could not have come up with this?
- (Hint: try to reconstruct what a DT algorithm would do)



$$G(x) = \sum_{m=1}^{M} \alpha_m G_m(x)$$

DT would not split here because no data (already 100% pure)