

CS 189/289

Today's lecture

1. Decision Trees
2. Ensemble approaches

Assigned reading: none

Announcement: today's office hour ends at 3:55pm

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

State of the art (SOTA) Machine Learning models in 2022

| Competition Name | Numer of competing teams | Data Type | Machine learning problem class | Winning (SOTA) Machine Learning model |
|--|--------------------------|-----------------------------|--------------------------------|---|
| Feedback Prize - English Language Learning | 2654 | Text | Regression | Neural networks: - microsoft-deberta-v3-base - deberta-v3-large - deberta-v2-xlarge - roberta-large - distilbert-base-uncased |
| Open Problems - Multimodal Single-Cell Integration | 1220 | Single-cell multiomics data | Multimodal + Correlation | Neural networks - custom architecture |
| RSNA 2022 Cervical Spine Fracture Detection | 883 | Images | Classification | Neural networks: - resnet 18d unet - effv2 - convnext tiny - convnext nano - convnext pico - convnext tiny - nfnet lo |



Kaggle

[Article](#) [Talk](#)

From Wikipedia, the free encyclopedia

Kaggle is a [data science competition platform](#) and online community for [data scientists](#) and [machine learning](#) practitioners under [Google LLC](#). Kaggle enables users to find and publish datasets, explore and build models in a web-based data science environment, work with other data scientists and machine learning engineers, and enter competitions to solve data science challenges.^[1]

Decision trees underlie state-of-the-art methods

| | | | | |
|---------------------------------------|------|---------|--------------------------------|--|
| Google Universal Image Embedding | 1022 | Images | Creating image representations | Neural Networks: - ViT-H - ViT-L |
| Mayo Clinic - STRIP AI | 888 | Images | Classification | Neural Networks: - swin_large_patch4_window12_384 |
| HuBMAP + HPA - Hacking the Human Body | 1175 | Images | Segmentation | Neural networks: -1x SegFormer mit-b3 - 2x SegFormer mit-b4 - 1x SegFormer mit-b5 - SegFormer mit-b5 |
| American Express - Default Prediction | 4874 | Tabular | Classification | Gradient Boosted Decision Trees: - LGBM |

Decision trees underlie state-of-the-art methods

| | | | | |
|--|------|---------|----------------|--|
| Feedback Prize - Predicting Effective Arguments | 1557 | Text | Classification | Neural networks: - deberta-(v3)-large |
| Google AI4Code - Understand Code in Python Notebooks | 1135 | Text | Regression | Neural networks: - deberta-(v3)-large |
| UW-Madison GI Tract Image Segmentation | 1548 | Images | Segmentation | Neural networks: 4 unets, with efficientnet b4, b5, b6, b7 backbones |
| Foursquare - Location Matching | 1079 | Tabular | Classification | Gradient Boosted Decision Trees + Neural Networks: - LGBM - Bert |

Decision trees underlie state-of-the-art methods

| | | | | |
|--|------|----------------|---|--|
| U.S. Patent Phrase to Phrase Matching | 1889 | Text | Regression | ahotrod/electra_large _discriminator_squad 2_512 - Yanhao/simcse-bert- for-patent - funnel- transformer/large |
| March Machine Learning Mania 2022 - Men's | 930 | Tabular | Classification | Gradient Boosted Decision Trees: - XGB |
| March Machine Learning Mania 2022 - Women's | 651 | Tabular | Classification | Gradient Boosted Decision Trees: - LGBM |
| BirdCLEF 2022 | 807 | Sound | Classification | Neural Networks: - tf_efficientnet_b3_ns - eca_nfnet_lo |
| H&M Personalized Fashion Recommendations | 2952 | Events' stream | Recommendation system (classification) | Gradient Boosted Decision Trees: - LGBM |

Tabular data are the same as categorical data

| shoes-sparse-500.csv | | | | |
|-----------------------|----------|---------------|--------------|------------------------------|
| Air Max 90 | Nike | Running | White | The Air Max 90 is a classic |
| Chuck Taylor All Star | Converse | Casual | Black | The Chuck Taylor All Star is |
| Superstar | Adidas | Athletic | White | The Superstar is a versatil |
| Old Skool | Vans | Skateboarding | Black | The Old Skool is a skateboa |
| Air Jordan 1 | Jordan | Basketball | Red | The Air Jordan 1 is a legen |
| Stan Smith | Adidas | Tennis | Green | The Stan Smith is a classic |
| Ultra Boost | Adidas | Running | Black | The Ultra Boost is a popula |
| Classic Slip-On | Vans | Casual | Checkerboard | The Classic Slip-On is a la |
| Air Force 1 | Nike | Athletic | White | The Air Force 1 is a legend |
| Chuck 70 | Converse | Casual | White | The Chuck 70 is a modern ta |

Decision trees underlie state-of-the-art methods

- DT-based models work well with default hyper-parameters.
- Work natively with categorical features (*i.e.*, tabular data).
- Are highly interpretable (unless « ensemble » them).

(none of these things are true for neural networks)

The Game of 20 Questions

The screenshot shows a web browser window for 20Q.net Inc. The URL is www.20q.net. The page displays a game session titled "Q11. I am guessing that it is a pomegranate?". The user has asked 11 questions, and the responses are as follows:

- 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.**

The sidebar on the left includes links for "Play 20Q", "About Us", "Products", "More...", and statistics: "games played online 87,115,788". The bottom right corner features a circular graphic with the text "I Can Read Your Mind".

| 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

- 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 |

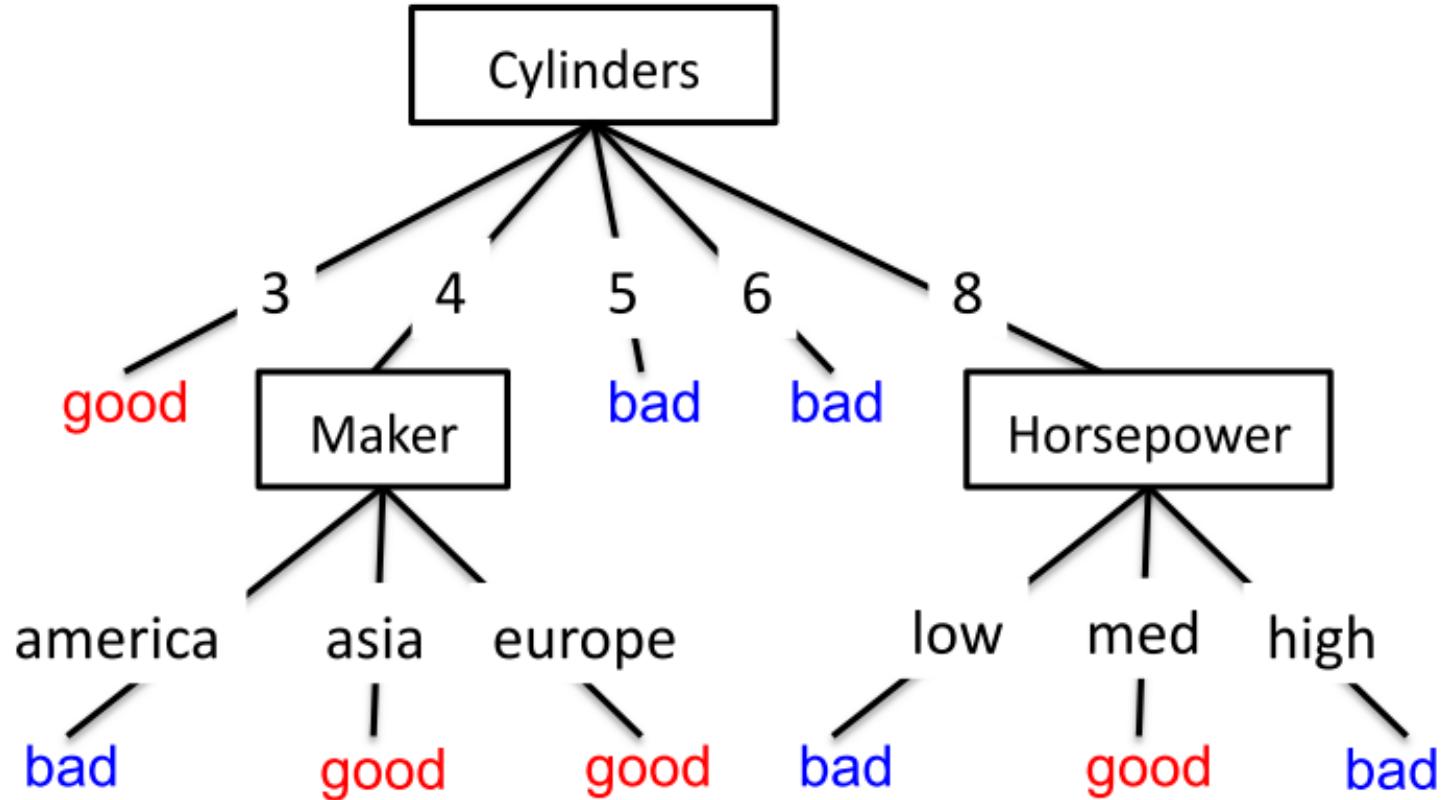
From the UCI repository (thanks to Ross Quinlan)

Hypotheses: decision trees $f : X \rightarrow Y$

- 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

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

| 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 |



Decision trees underlie state-of-the-art methods

- DT-based models work well with default hyper-parameters.
- Work natively with categorical features (*i.e.*, tabular data).
- Are highly interpretable (unless « ensemble » them).

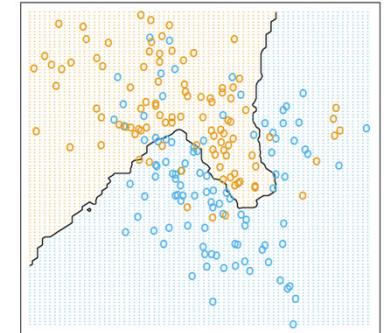
(none of these things are true for neural networks)

- Lack of smoothness in decision boundary
- Difficulty capturing additive signal
- Unstable---small change in data can make big changes to tree.

DT classification boundaries

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

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



(recall 15-NN classifier)

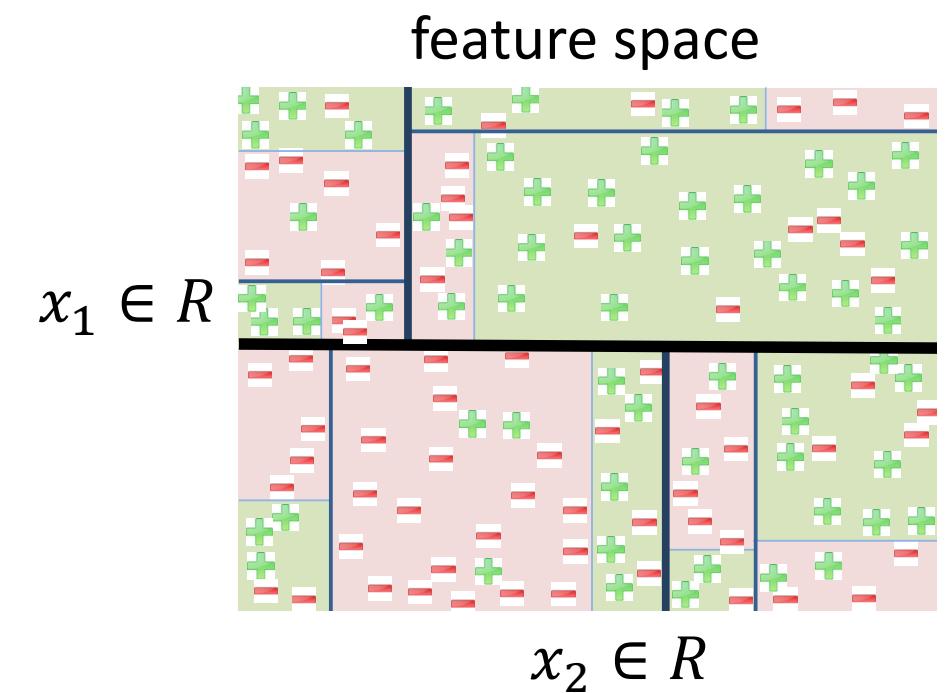
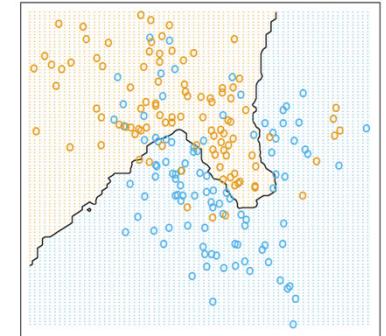


figure credit: Yisong Yue

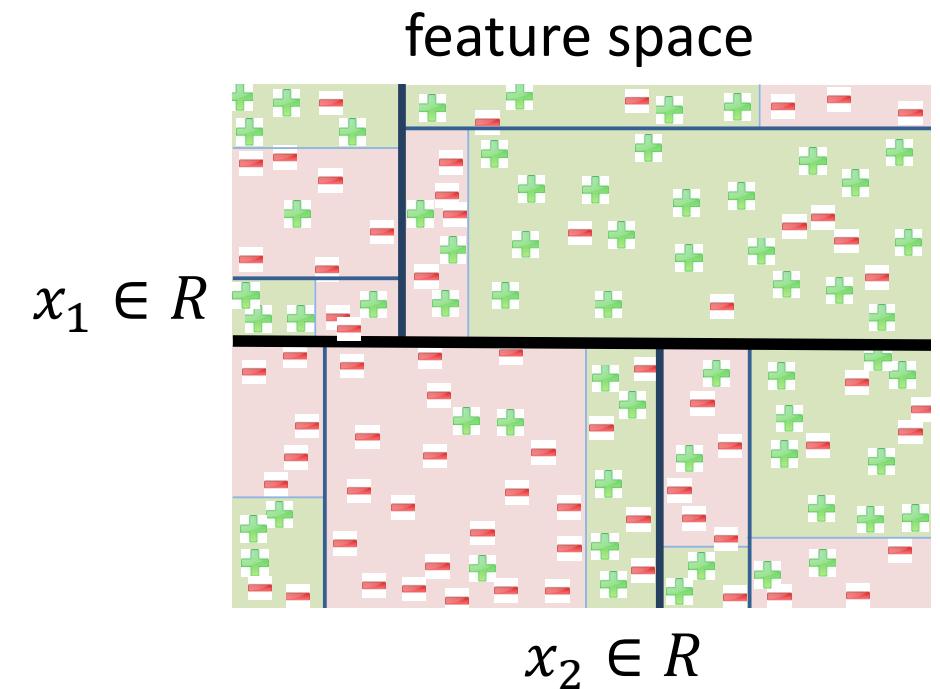
DT classification boundaries

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

- 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.



(recall 15-NN classifier)



fitgure credit: Yisong Yue

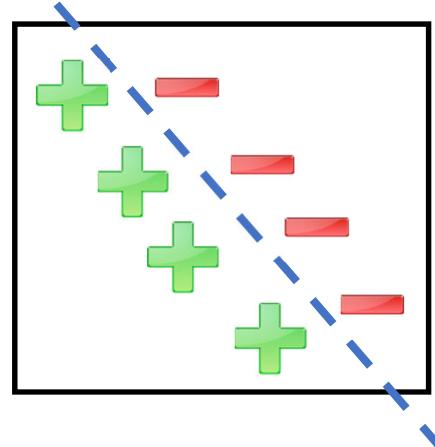
Decision Trees vs Linear Models

Decision Trees are AXIS-ALIGNED!

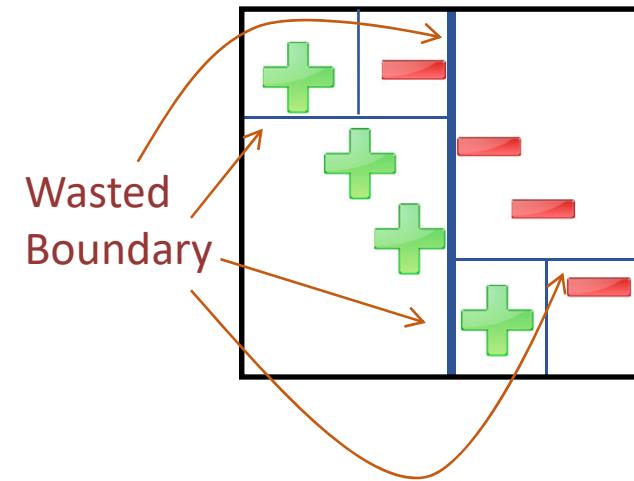
- Cannot easily model diagonal boundaries

Example:

Logistic regression
can have arbitrary
linear boundary



Decision Trees Require
Complex Axis-Aligned
Partitioning



slide credit: Yisong Yue

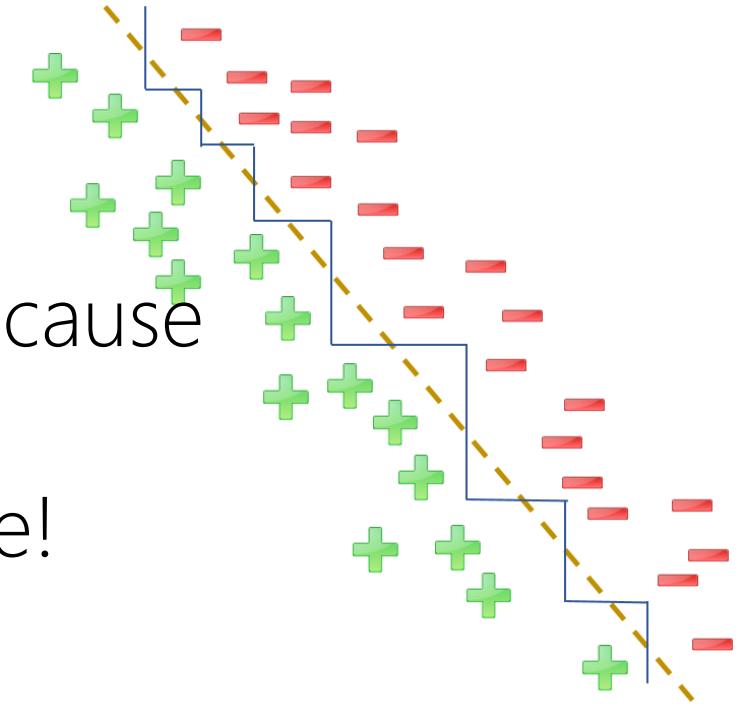
More Extreme Example



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

Decision Trees vs Linear Models

- Not optimal use of limited training data because so much potential wasted boundary.
- Still, Decision Trees are often more accurate!
- Their non-linearity is powerful.



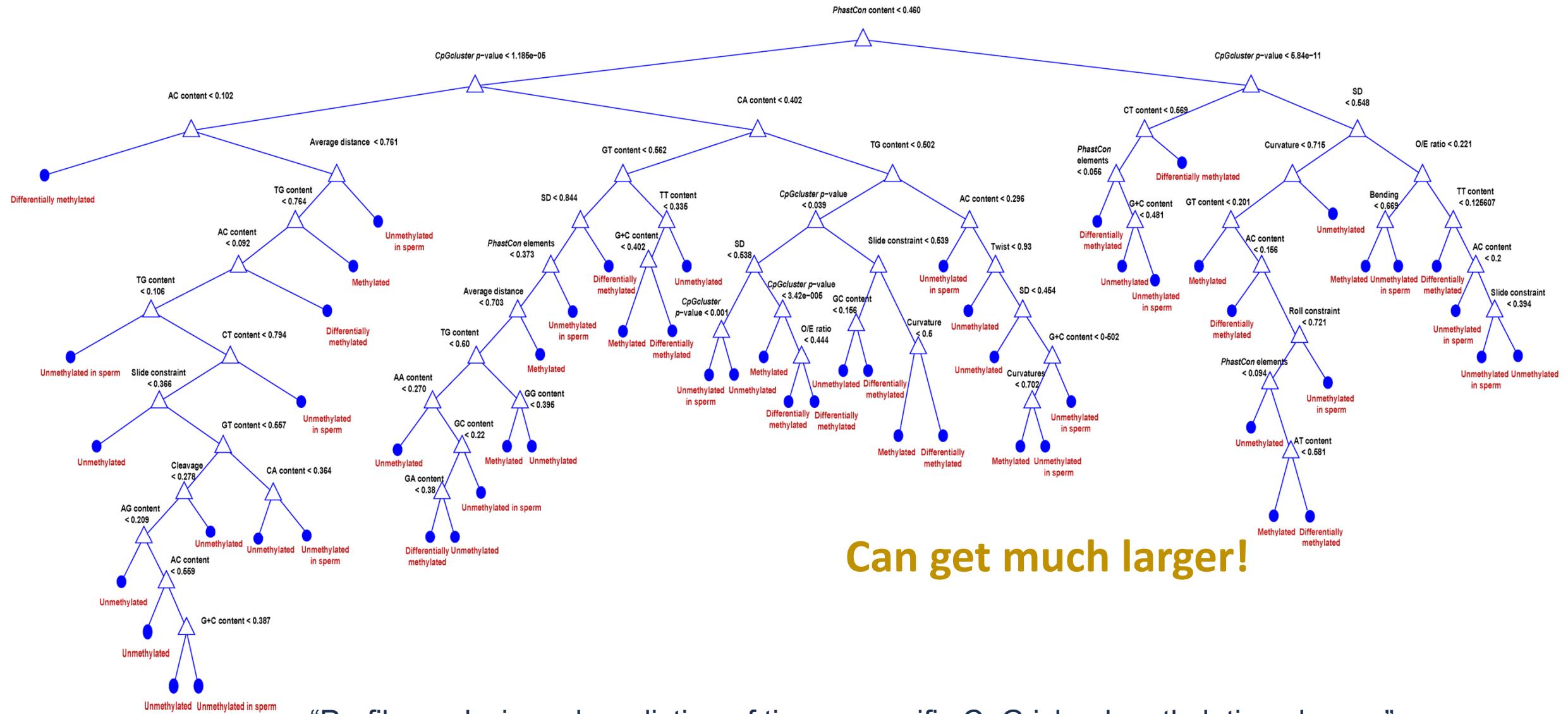
Decision trees underlie state-of-the-art methods

- DT-based models work well with default hyper-parameters.
- Work natively with categorical features (*i.e.*, tabular data).
- Are highly interpretable (unless « ensemble » them).

(none of these things are true for neural networks)

- Lack of smoothness in decision boundary
- Difficulty capturing additive signal
- Unstable---small change in data can make big changes to tree.

Real Decision Trees



Can get much larger!

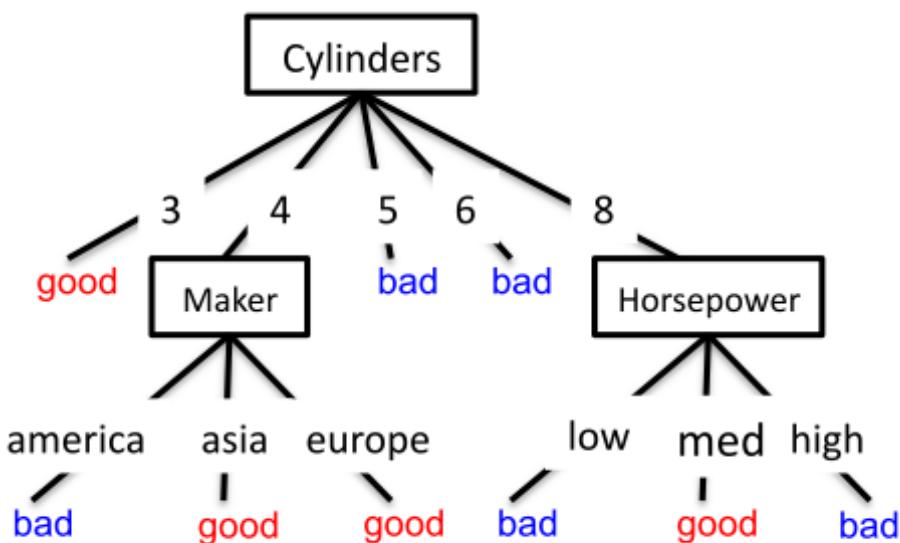
“Profile analysis and prediction of tissue-specific CpG island methylation classes”

Image Source: <http://www.biomedcentral.com/1471-2105/10/116>

slide credit: Yisong Yue

Model space

- How many possible models?
- What functions can be represented?



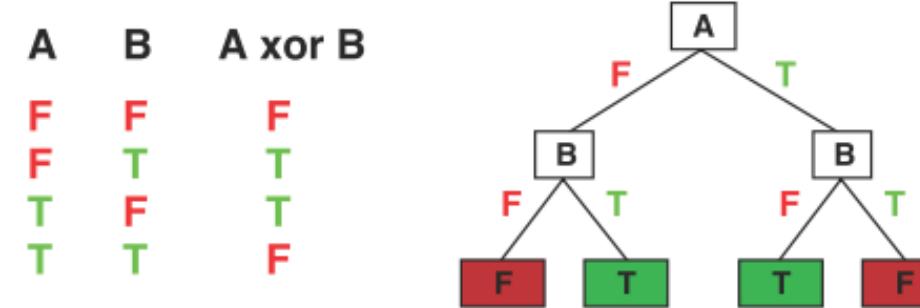
| mpg | cylinders | displacement | horsepower | weight | acceleration | modelyear | maker |
|------|-----------|--------------|------------|--------|--------------|-----------|---------|
| good | 4 | low | low | low | high | 75to78 | asia |
| bad | 6 | medium | medium | medium | medium | 70to74 | america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

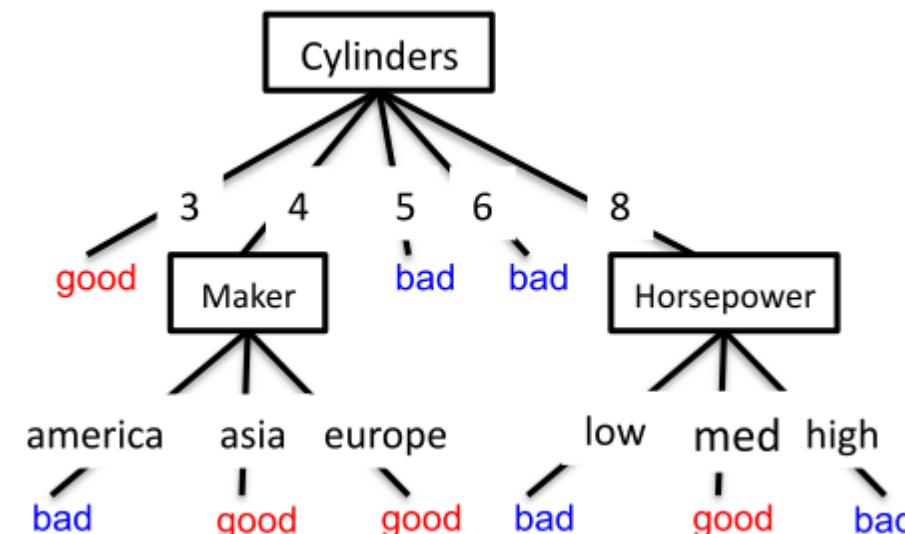
Can you think of a function that a decision tree could not represent?

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...



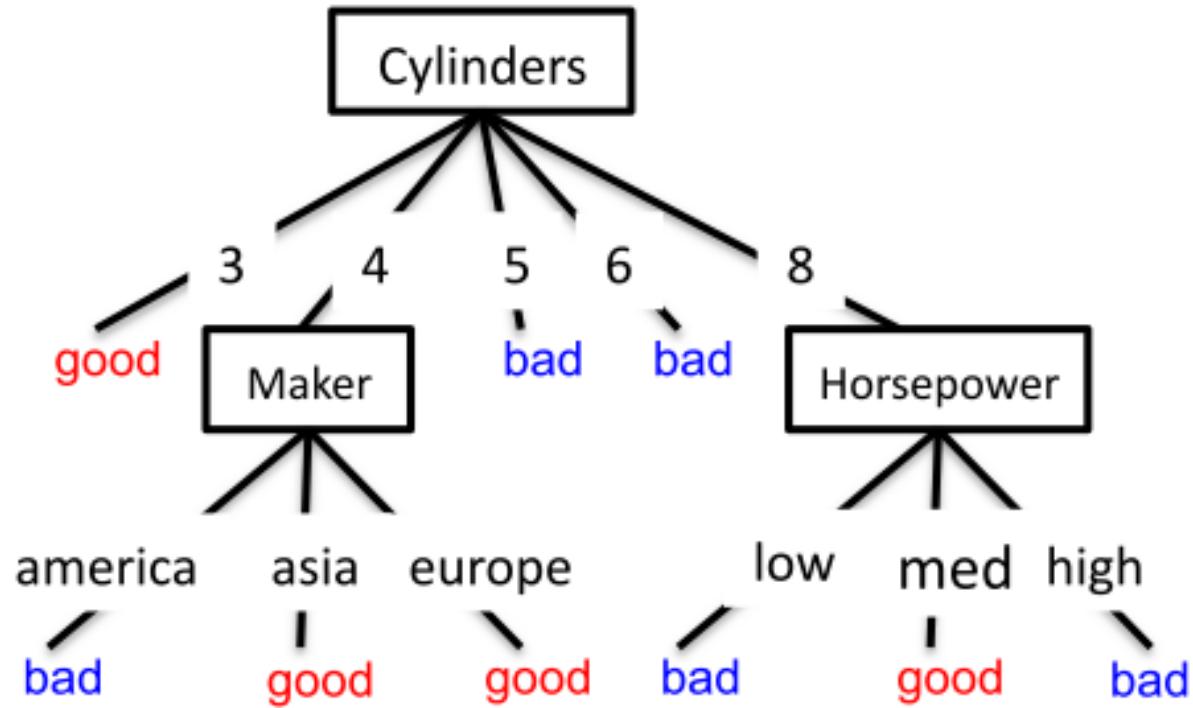
(Figure from Stuart Russell)



cyl=3 \vee (cyl=4 \wedge (maker=asia \vee maker=europe)) \vee ...

Hypothesis space

- How many possible models?
- What functions can be represented?
- How many will be consistent with a given dataset?
- How will we choose the best one?



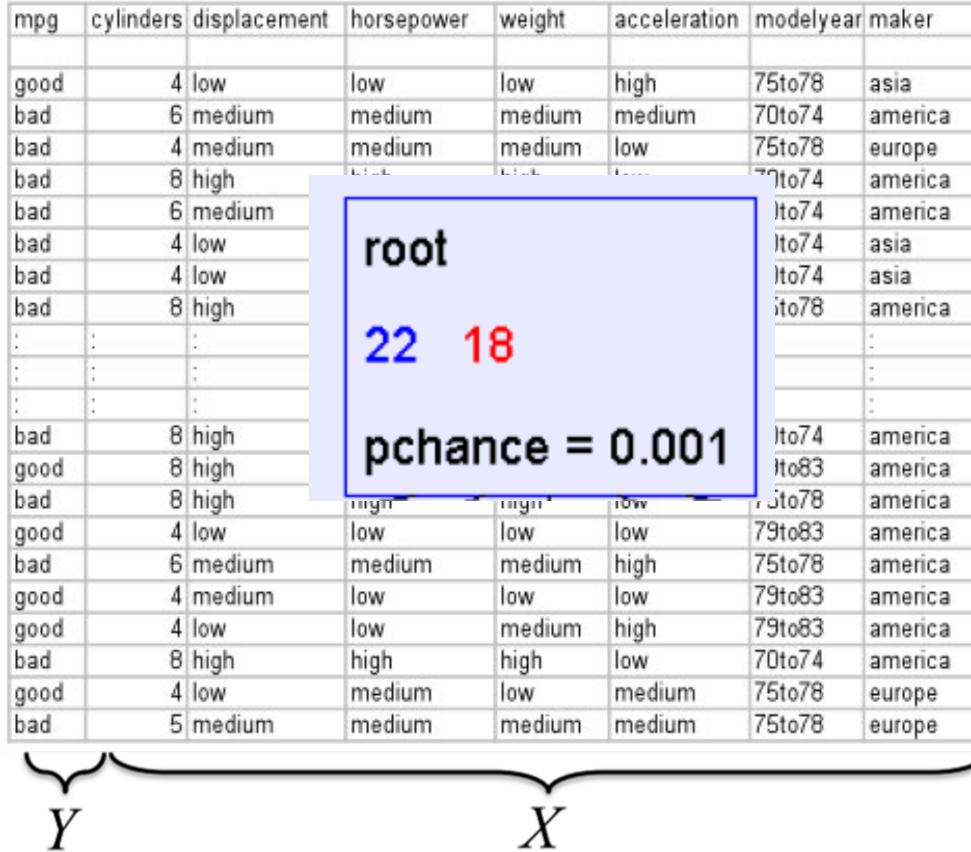
Lets first consider how to build a tree from training data, then revisit this question.

What is the Simplest Tree?

root node only
predict majority class:
 $\text{mpg}=\text{bad}$

Is this a good tree?

[22+, 18-]

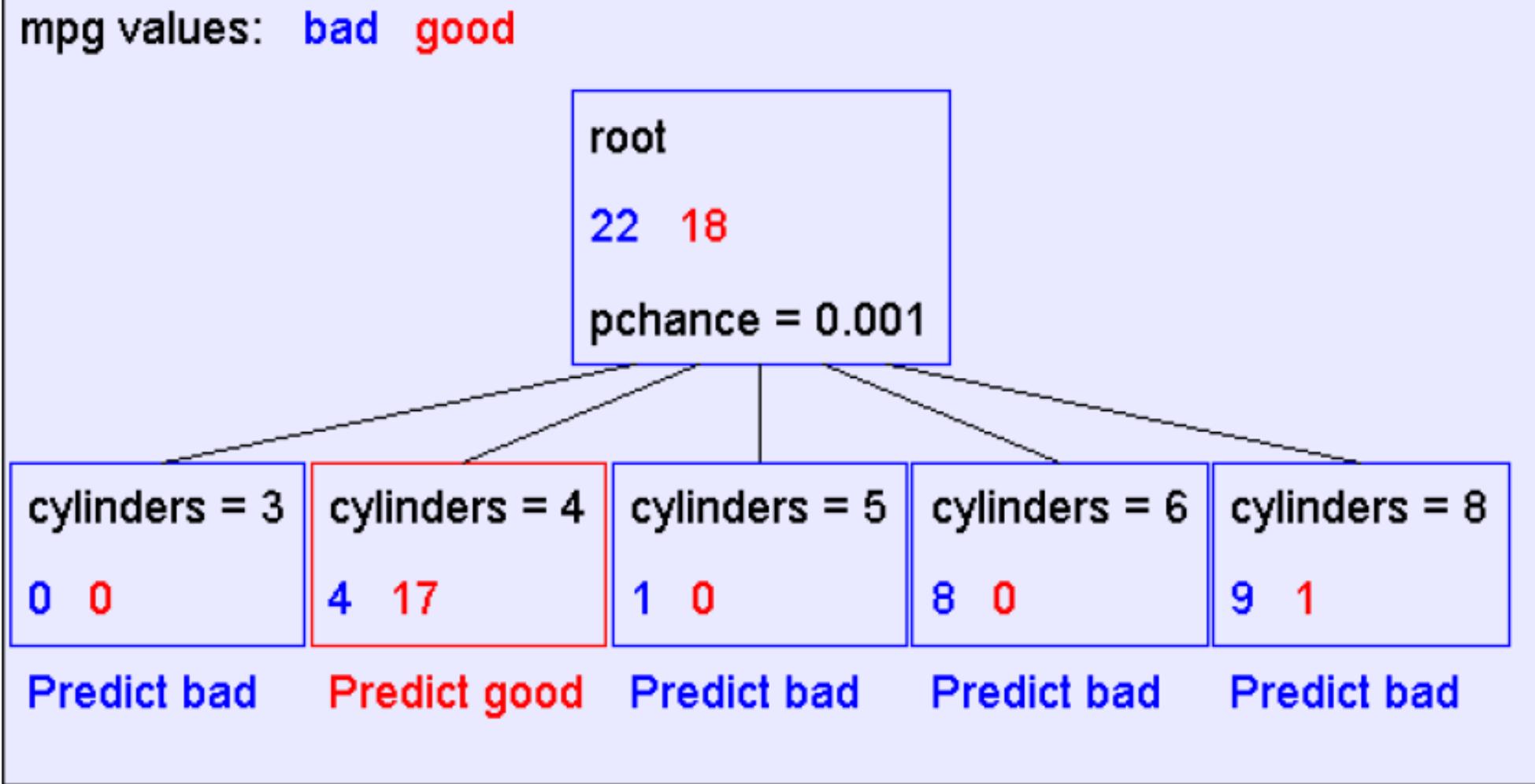


Means:
correct on 22 examples
incorrect on 18 examples

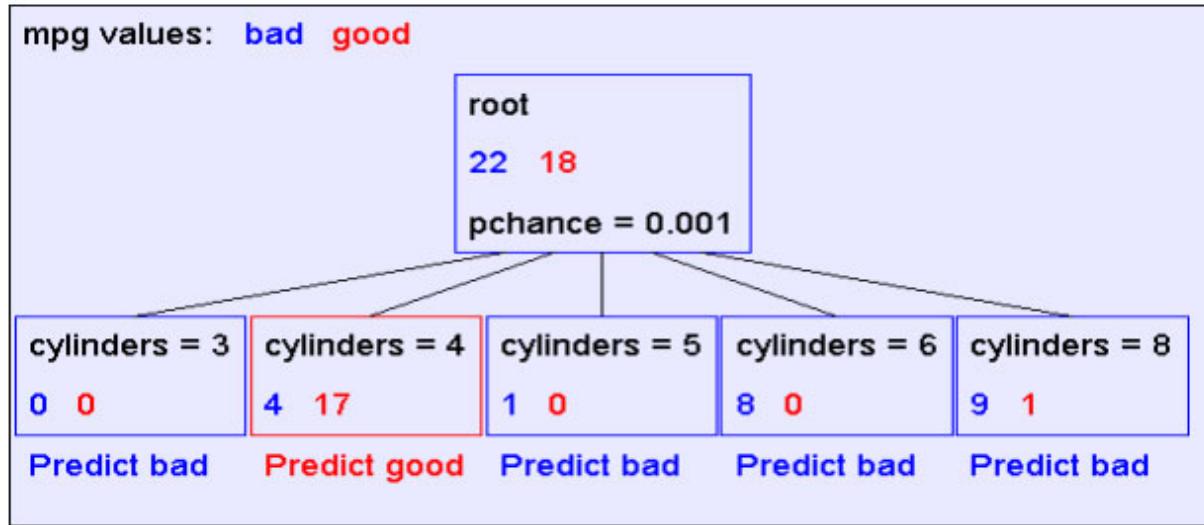
Next simplest tree: A Decision Stump (one feature splitting node)

| 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 |

\brace{Y} \brace{X}



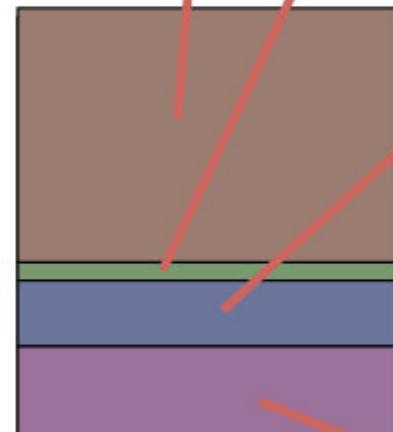
Increase complexity by recursive partitioning



Take the Original Dataset..



And partition it according to the value of the attribute we split on



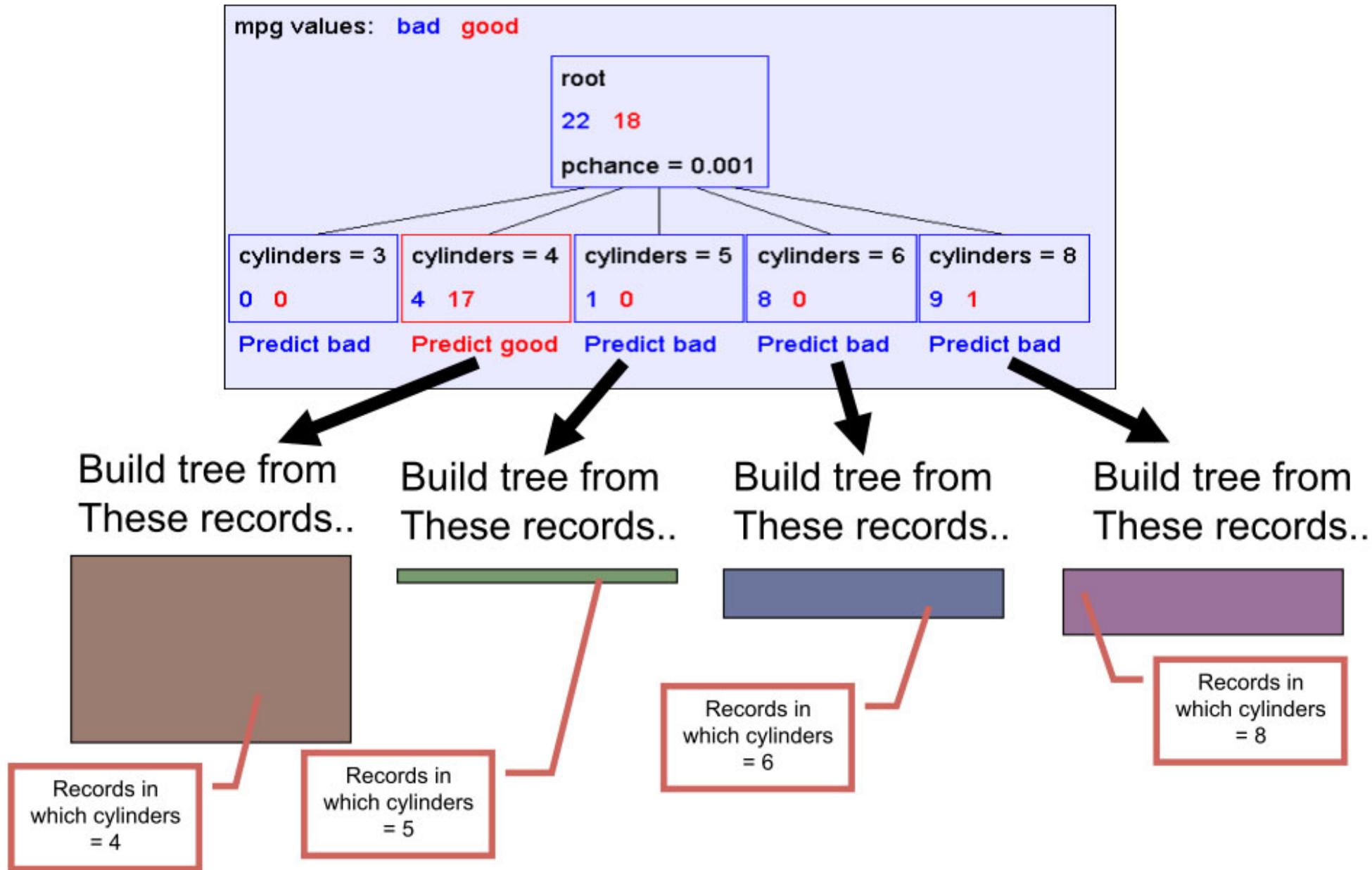
Records in which cylinders = 4

Records in which cylinders = 5

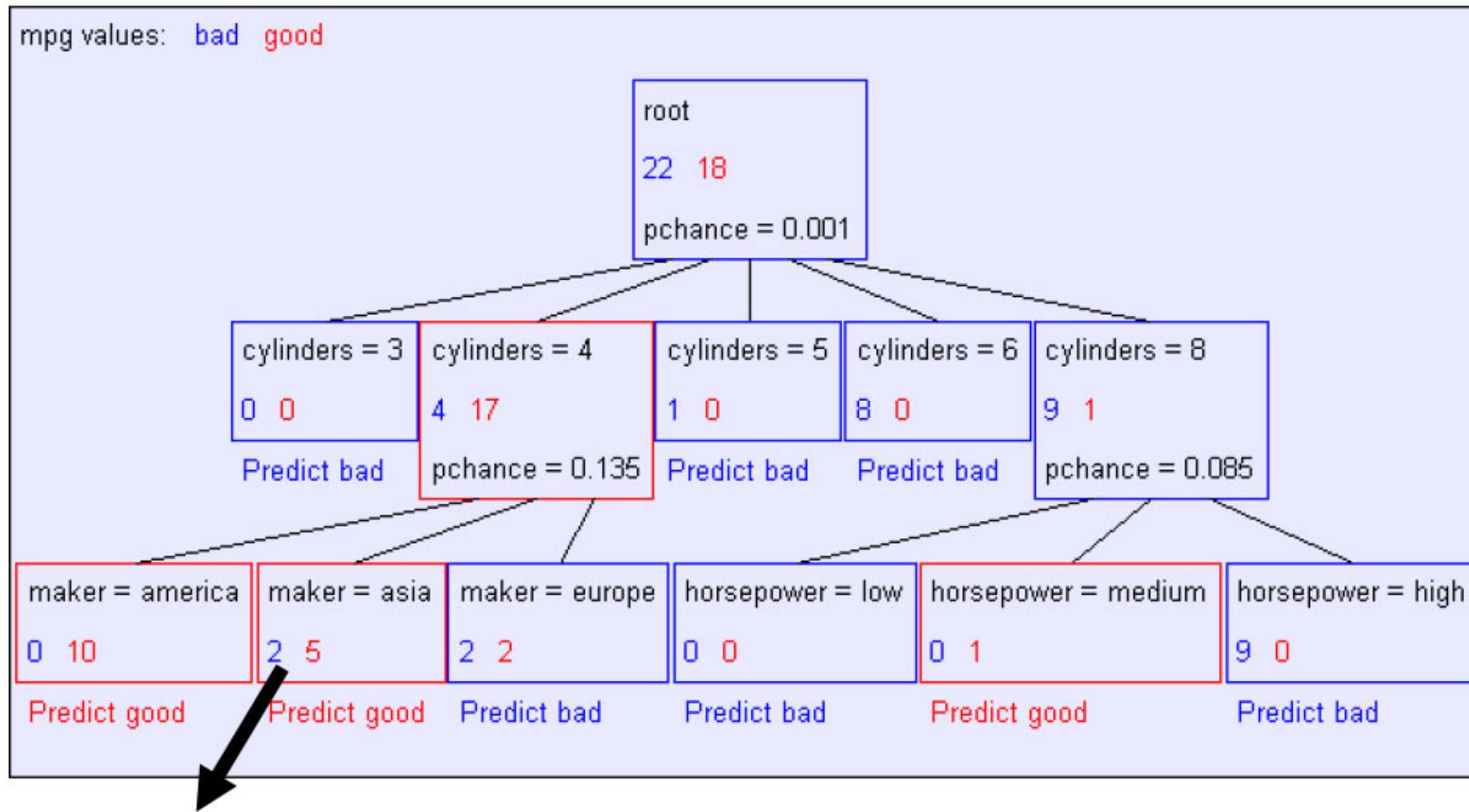
Records in which cylinders = 6

Records in which cylinders = 8

Increase complexity by recursive partitioning

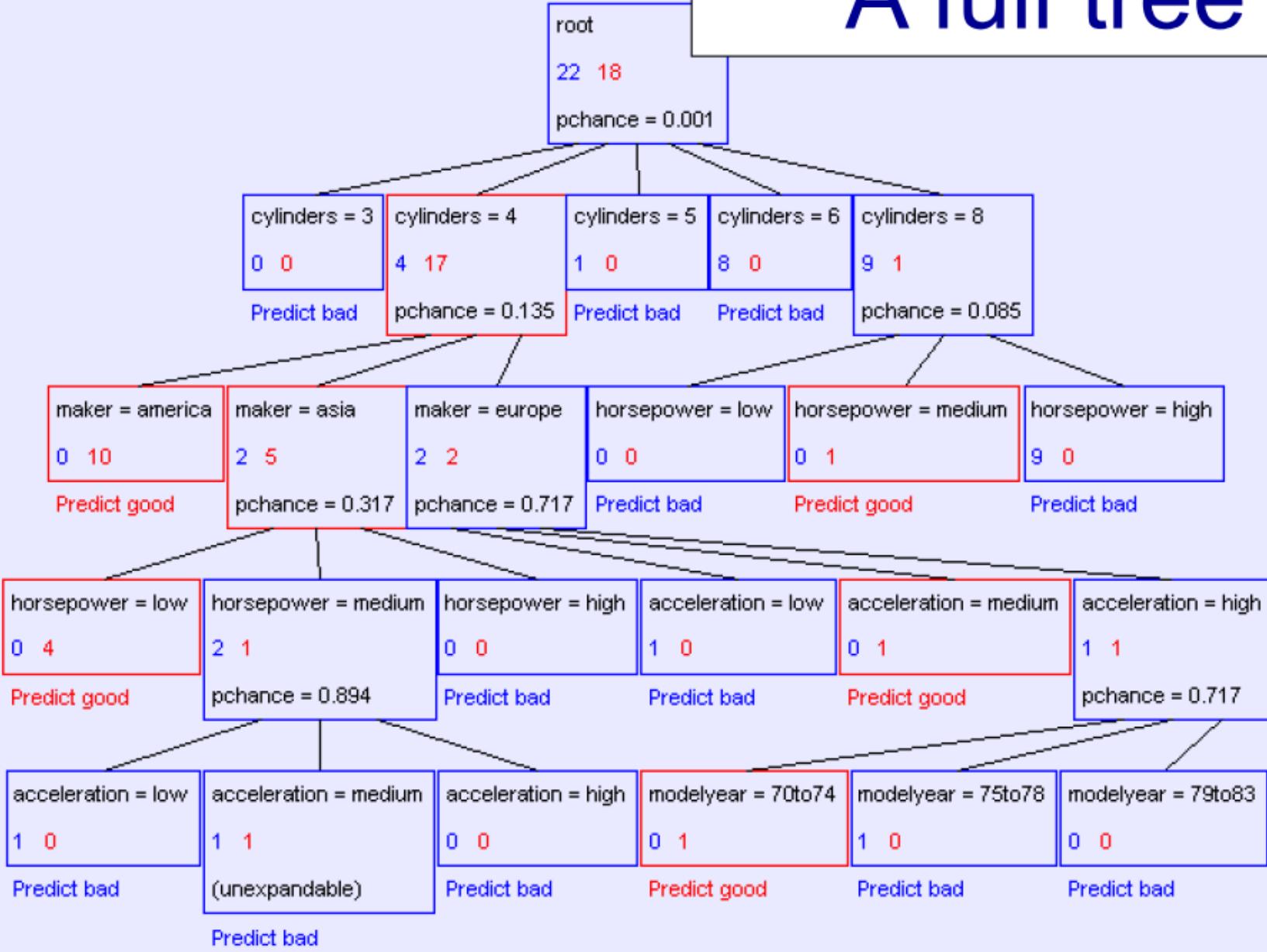


Now have second level of tree



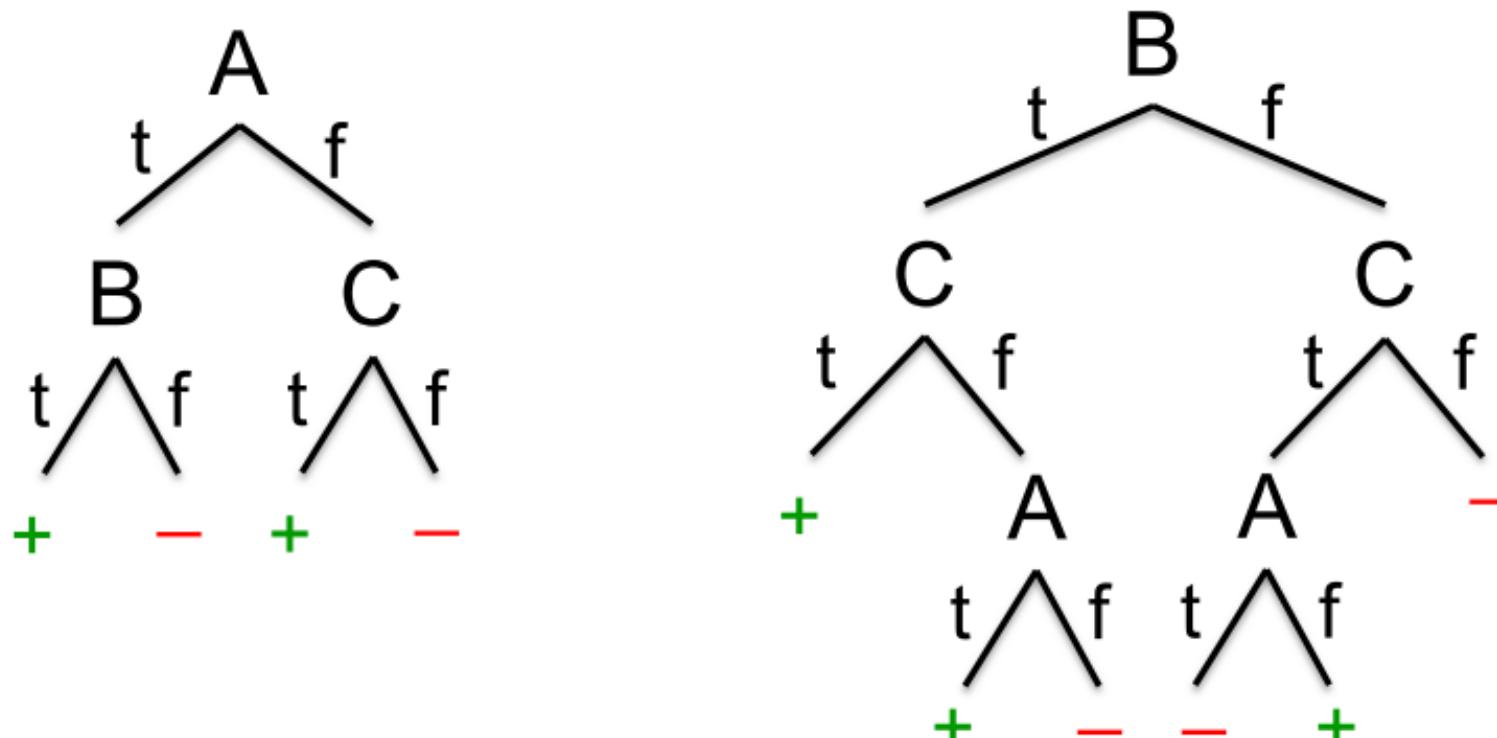
mpg values: bad good

A full tree



- 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 necessitate”



- “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”

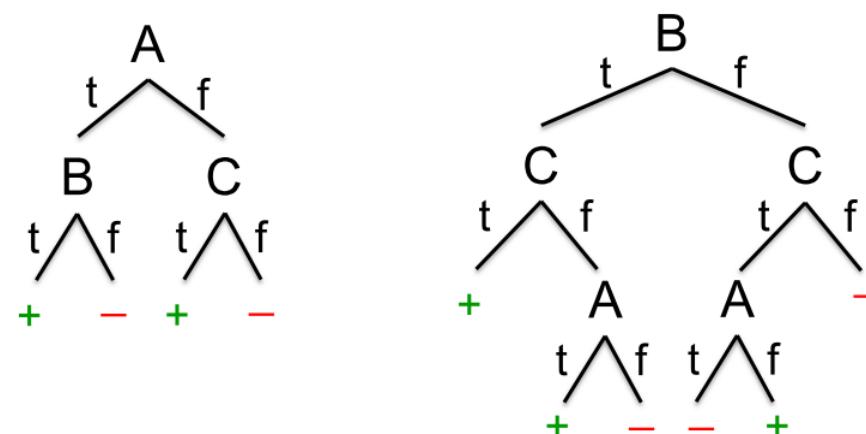


English Franciscan friar William of Ockham (c. 1287–1347)

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

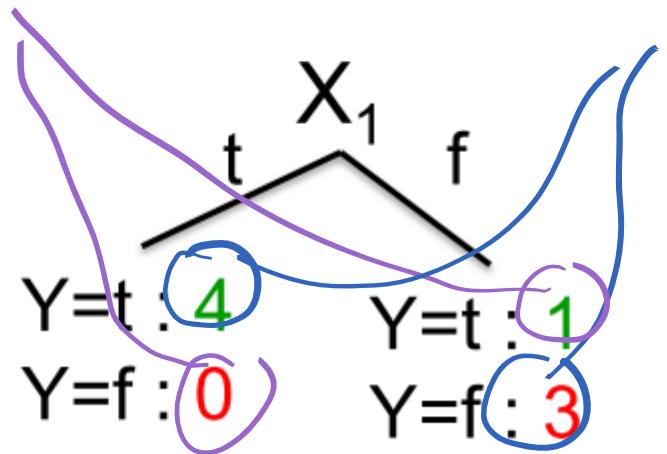
- Learning the simplest (smallest) decision tree is an NP-complete problem [Hyafil & Rivest '76]
- Resort to a greedy heuristic:
 - Start from empty decision tree
 - Split on **next best attribute (feature)**
 - Recurse

but what does
this mean?

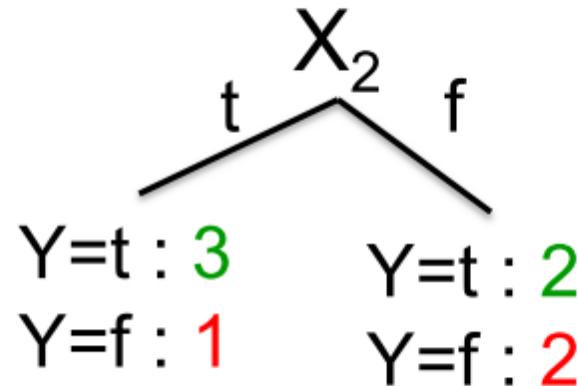
How to pick a good feature to split on?

Would we prefer to split on X_1 or X_2 ?

incorrect



correct



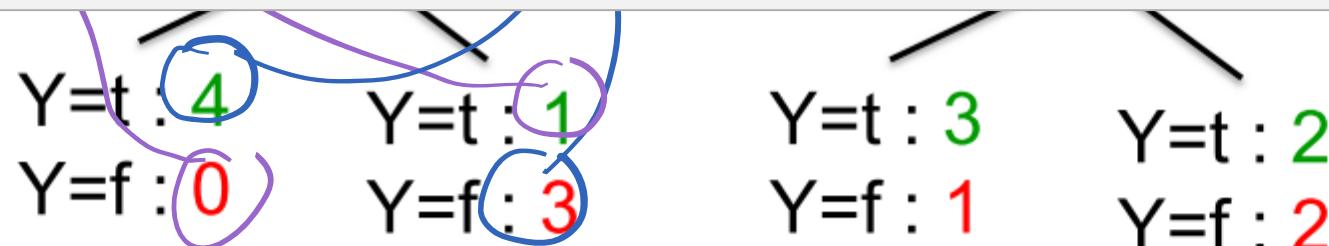
| X_1 | X_2 | Y |
|-------|-------|---|
| T | T | T |
| T | F | T |
| T | T | T |
| T | F | T |
| F | T | T |
| F | F | F |
| F | T | F |
| F | F | F |

- Goal: want leaf nodes that are a “pure split” (all belong to one class).
- Want to keep increasing the “purity” with each split.

How to pick a good feature to split on?

Would we prefer to split on X_1 or X_2 ?

- The more pure a node/feature, the less surprised we are to find a misclassified example there.
- So want to minimize surprise==min entropy!



- Goal: want leaf nodes that are a “pure split” (all belong to one class).
- Want to keep increasing the “purity” with each split.

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

Entropy: a measure of expected surprise

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

$$p(\text{head}) = 0.5$$



$$p(\text{head}) = 0$$

$$p(\text{head}) = 1$$



$$p(\text{head}) = 0.01$$



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

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

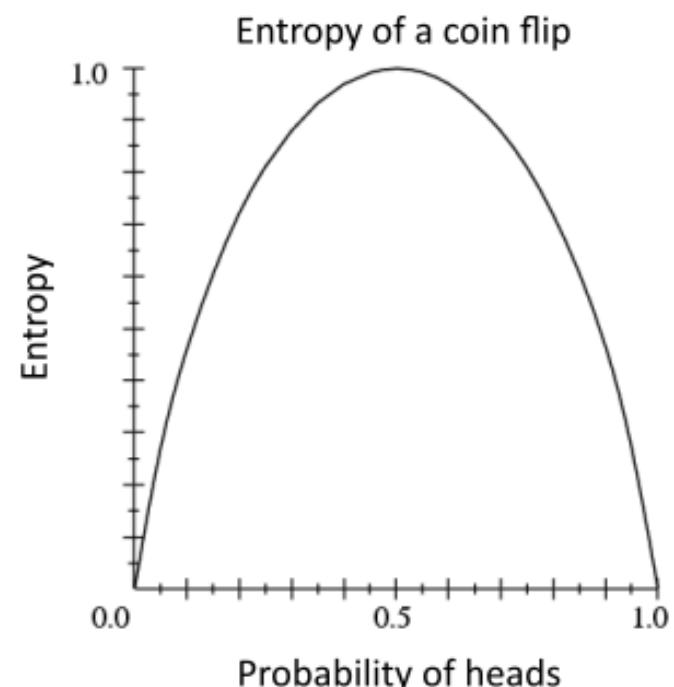
Entropy example: flipping a coin

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

$$P(Y=t) = 5/6$$

$$P(Y=f) = 1/6$$

$$\begin{aligned} H(Y) &= - 5/6 \log_2 5/6 - 1/6 \log_2 1/6 \\ &= 0.65 \end{aligned}$$



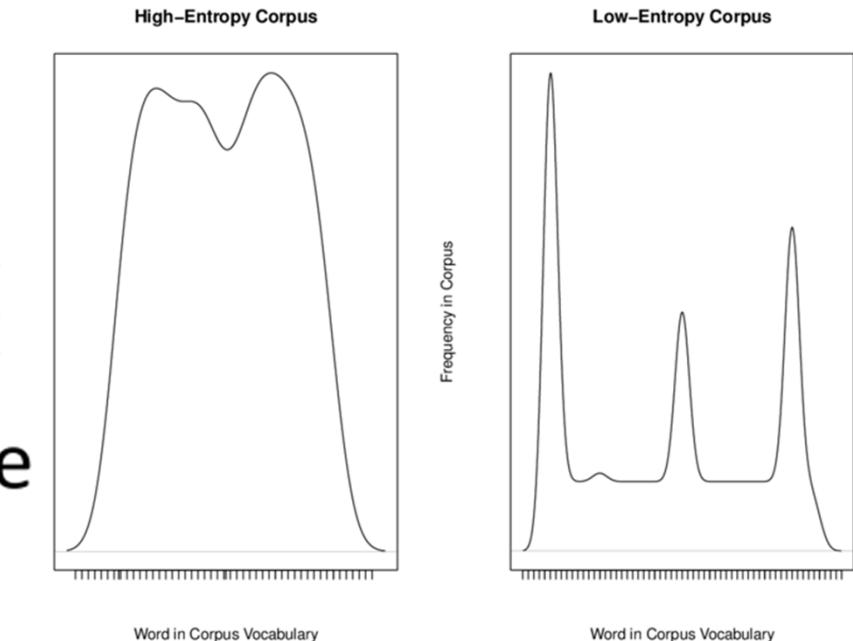
Entropy of a random variable Y :

“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

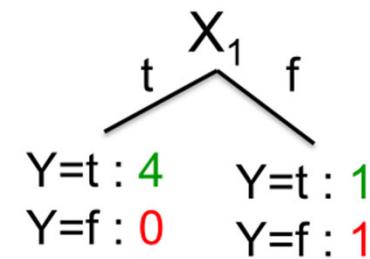


https://www.researchgate.net/figure/Hypothetical-distributions-of-term-frequency-in-high-and-low-entropy-corpora_fig1_305417514

Entropy for choosing feature to split next

- 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 | 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 $x_j \geq v$

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

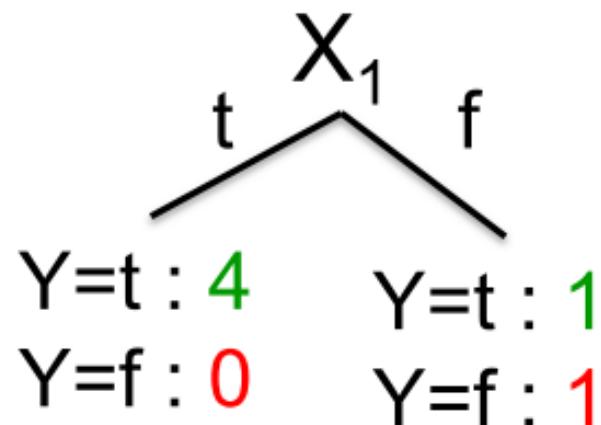
| X_1 | X_2 | Y |
|-------|-------|---|
| T | T | T |
| T | F | T |
| T | T | T |
| T | F | T |
| F | T | T |
| F | F | F |

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:



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

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

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

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

$$= 2/6$$

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

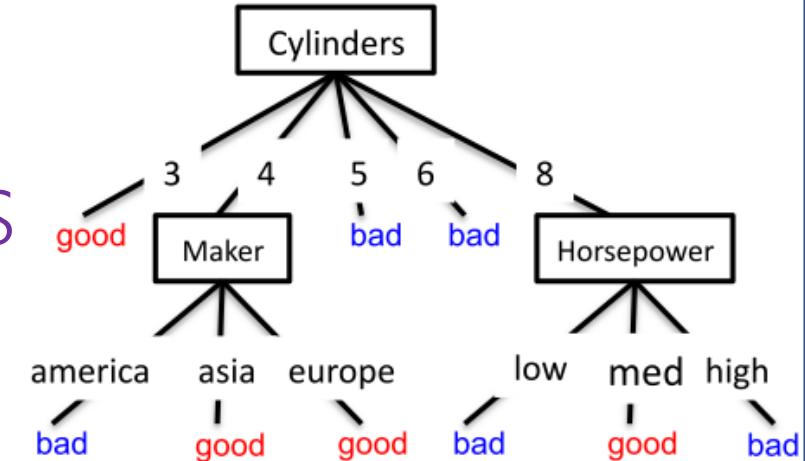
| X ₁ | X ₂ | Y |
|----------------|----------------|---|
| T | T | T |
| T | F | T |
| T | T | T |
| T | F | T |
| F | T | T |
| F | F | F |

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.



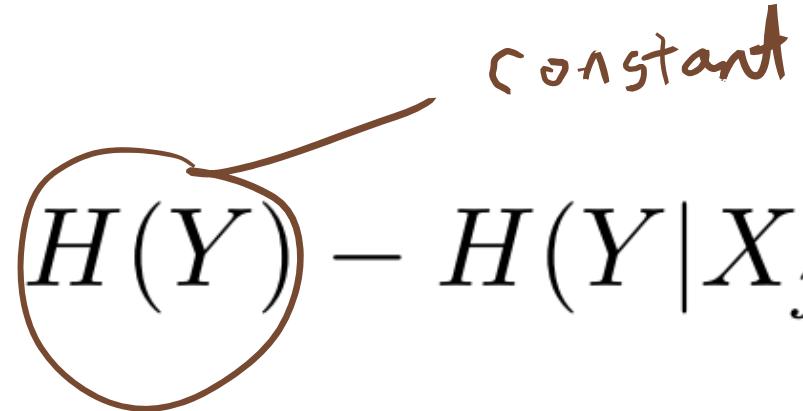
$$\begin{aligned} H(Y|X_1) &= -\frac{4}{6} (1 \log_2 1 + 0 \log_2 0) \\ &\quad - \frac{2}{6} (1/2 \log_2 1/2 + 1/2 \log_2 1/2) \\ &= 2/6 \end{aligned}$$

| | | |
|---|---|---|
| T | T | T |
| T | F | T |
| F | T | T |
| F | F | F |

Equivalently: maximize the *information gain*

(also called the *mutual information*)

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



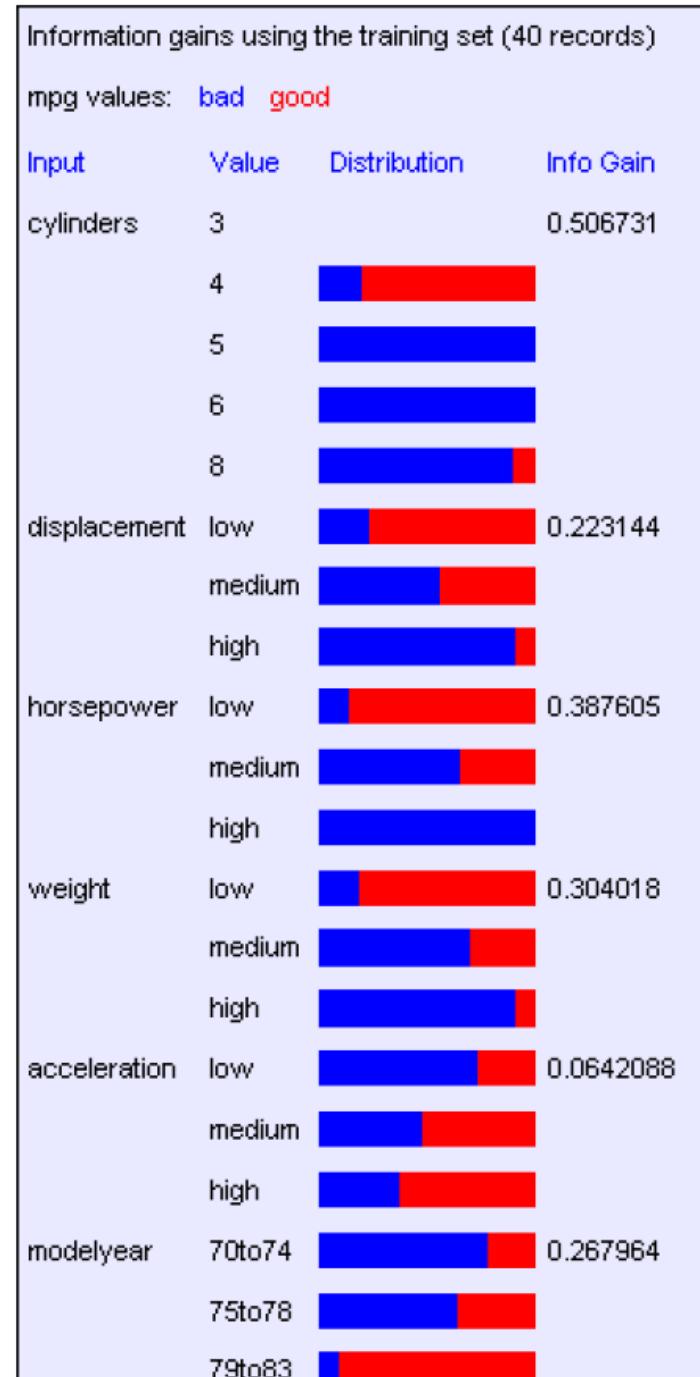
constant

“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:
$$\arg \max_i IG(X_i) = \arg \max_i H(Y) - H(Y | X_i)$$
- Recurse

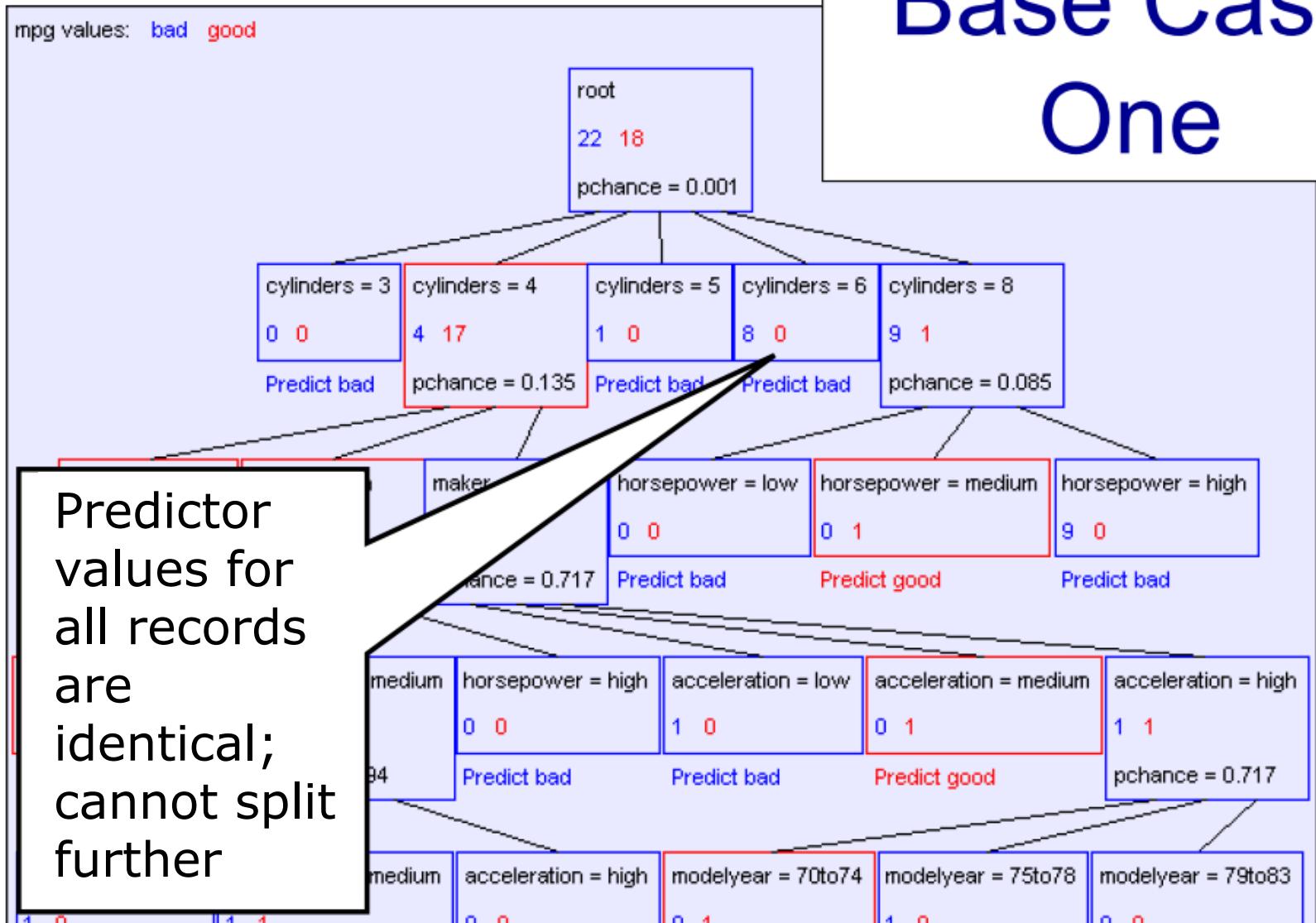
Back to the “mpg” example

- 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.



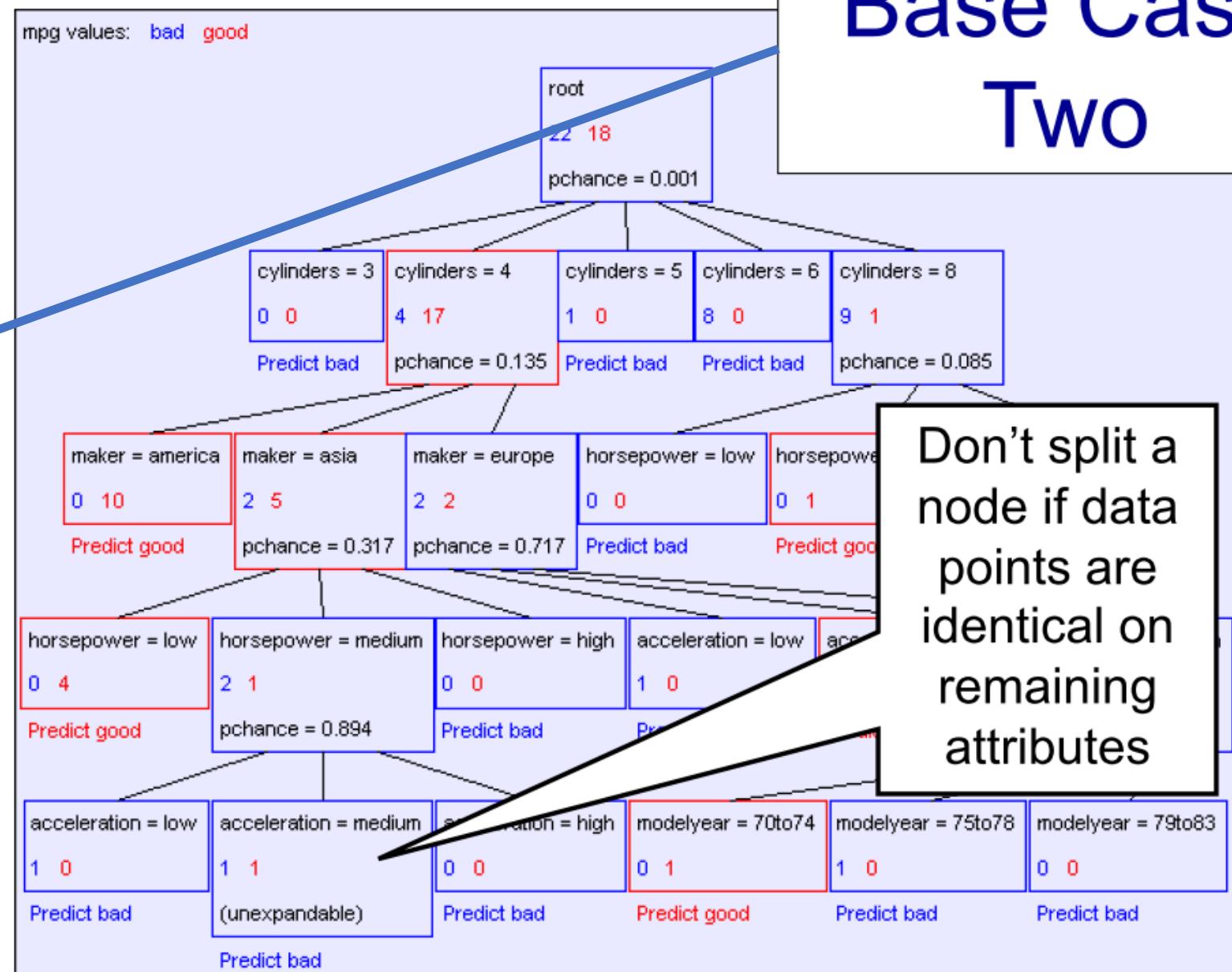
When to stop recursing on a node?

Base Case One



When to stop recursing on a node?

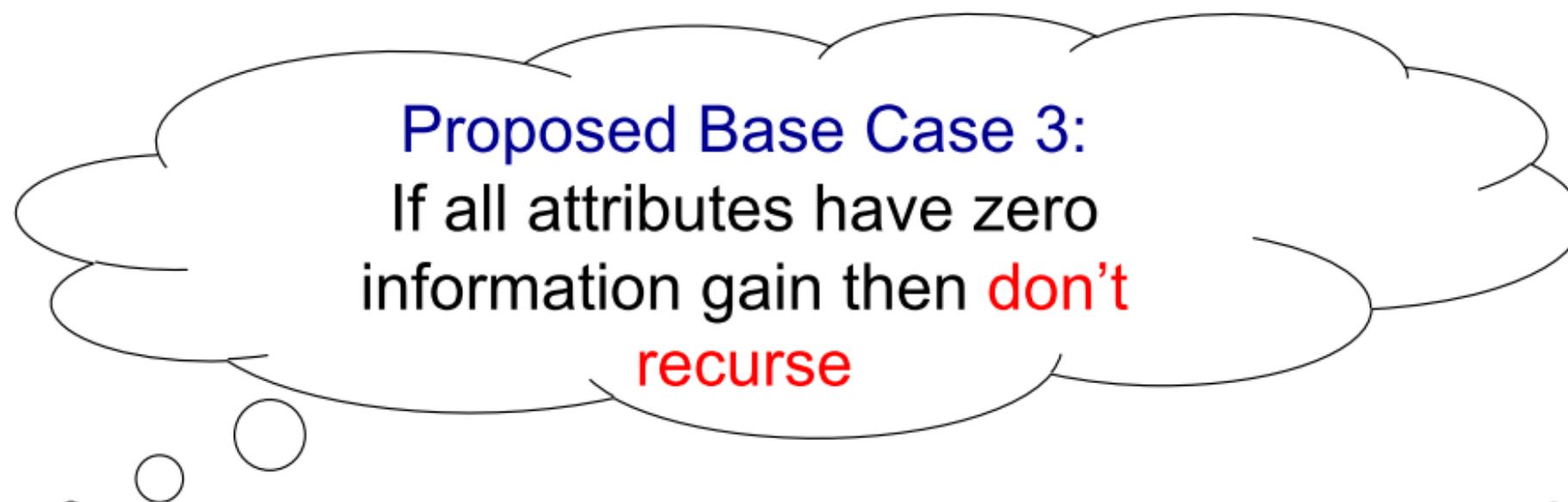
| Information gains using the training set (2 records) | | | |
|--|---------|--------------|-----------|
| mpg values: | | bad | good |
| Input | Value | Distribution | Info Gain |
| cylinders | 3 | 0 | |
| | 4 | 0.5 0.5 | |
| | 5 | | |
| | 6 | | |
| | 8 | | |
| displacement | low | 0.5 0.5 | 0 |
| | medium | | |
| | high | | |
| horsepower | low | 0 | |
| | medium | 0.5 0.5 | |
| | high | | |
| weight | low | 0.5 0.5 | 0 |
| | medium | | |
| | high | | |
| acceleration | low | 0 | |
| | medium | 0.5 0.5 | |
| | high | | |
| modelyear | 70to74 | 0.5 0.5 | 0 |
| | 75to78 | | |
| | 79to83 | | |
| maker | america | 0 | |
| | asia | 0.5 0.5 | |
| | europe | | |



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**



- *Is this a good idea?*

The problem with Base Case 3 for stopping

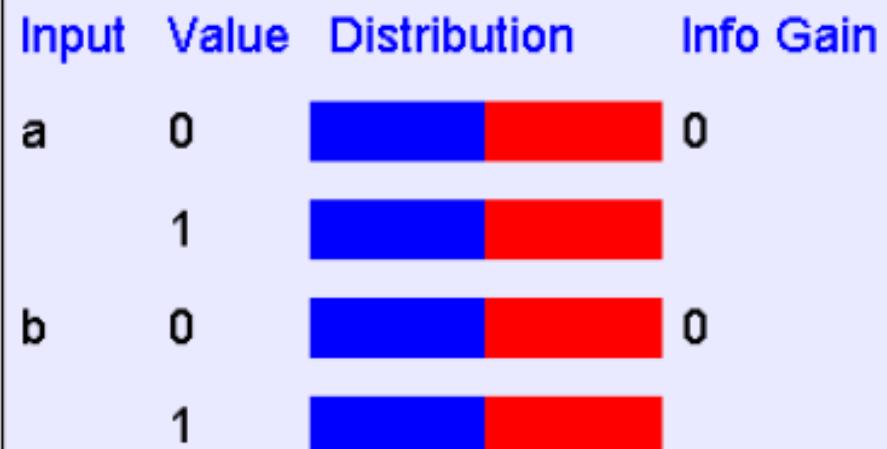
We are using a greedy heuristic.

$$y = a \text{ XOR } b$$

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

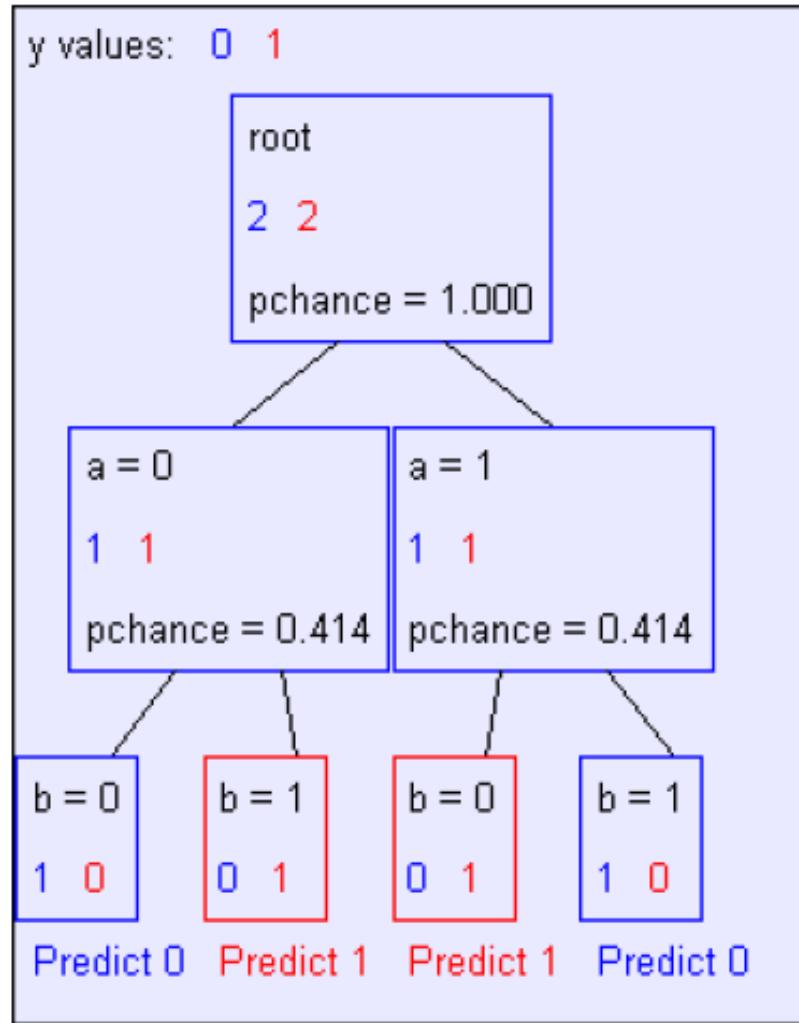
Information gains using the training set (4 records)

y values: 0 1



If we omit Base Case 3 for stopping:

The resulting decision tree:



| | | MODEL PREDICTIONS | |
|--------------|----------|-------------------|----------|
| | | Negative | Positive |
| GROUND TRUTH | Negative | TN | FP |
| | Positive | FN | TP |

$$y = a \text{ XOR } b$$

| 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 (e.g. χ^2).
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”
 - If all input values are the same, return a leaf node that says “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_X)
 - 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*:
- BASE CASE 1
- BASE CASE 2
- 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.

Averaging several decision trees

- Averaging reduces variance while keeping bias the same:
- If Y_1, \dots, Y_M are M uncorrelated random variables, each with mean μ and variance σ^2 , then

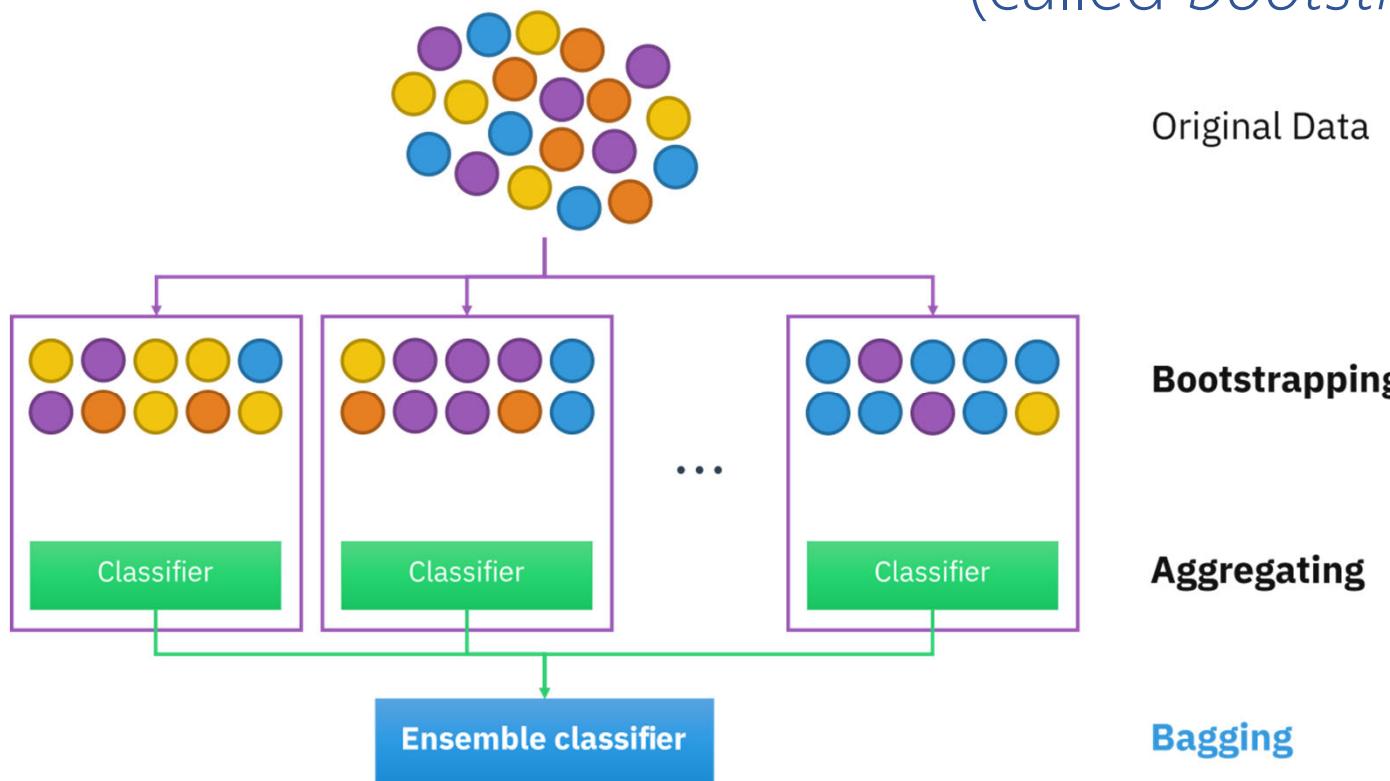
$$\mathbb{E} \left[\frac{1}{M} \sum_{k=1}^M Y_k \right] = \mu \quad \text{and} \quad \text{Var} \left[\frac{1}{M} \sum_{k=1}^M Y_k \right] = \frac{1}{M^2} \text{Var} \left[\sum_{k=1}^M Y_k \right] = \frac{1}{M^2} [M\sigma^2] = \frac{\sigma^2}{M}$$

using $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2 \times \text{Covar}(X, Y)$

Creating ensemble of DTs

1. Train M models using M samples of size N created by sampling with replacement, from N total data points.

(called *bootstrap replicates* of the data)



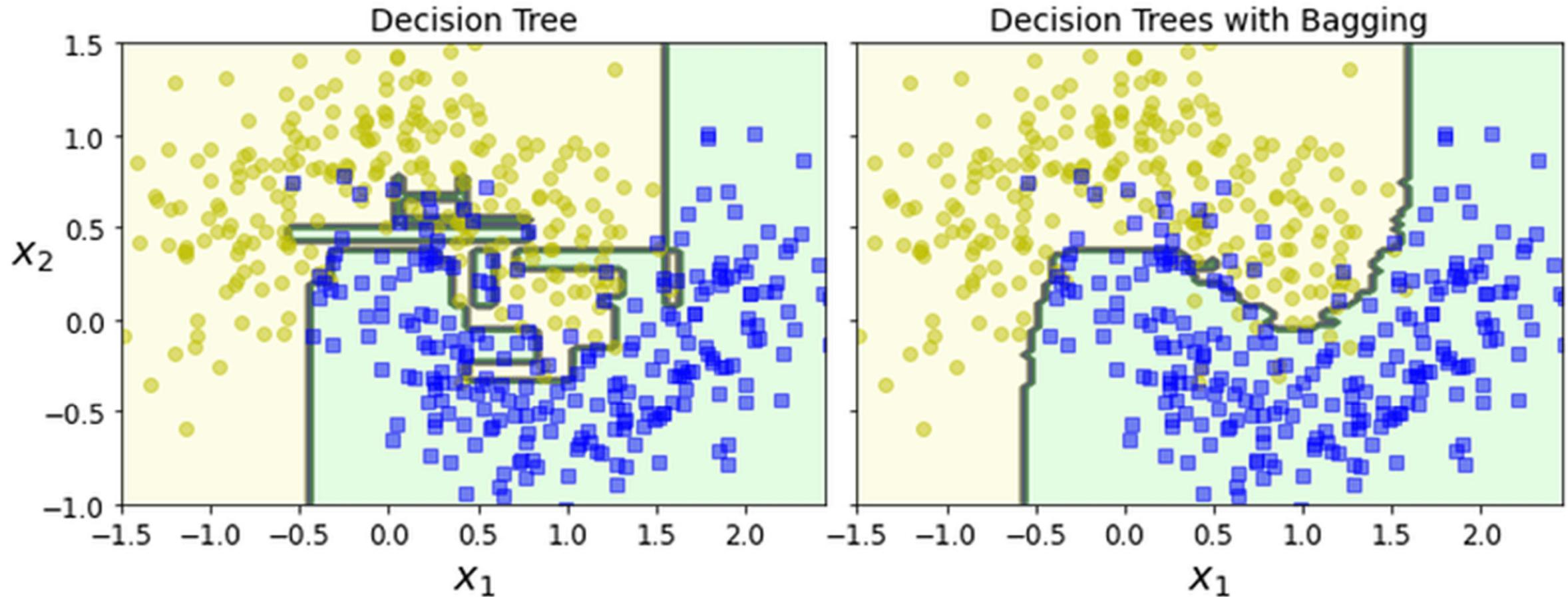
These trees will
be weakly
correlated.

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



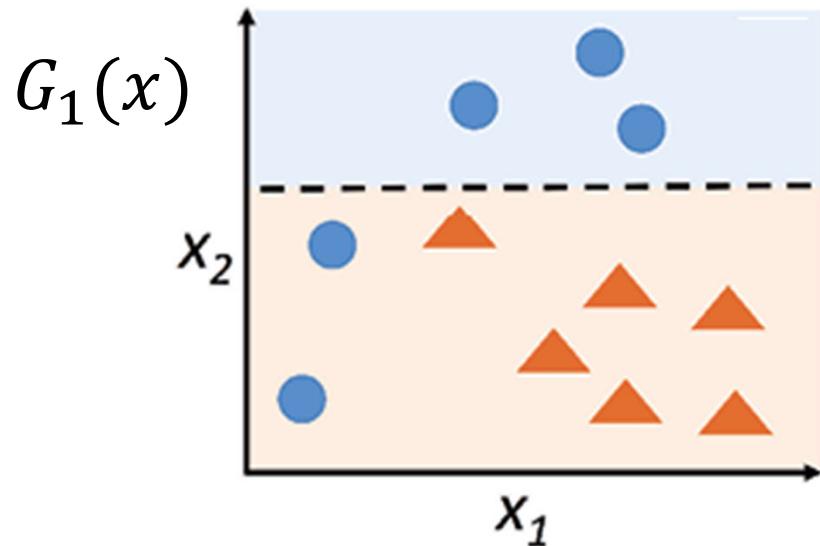
Other forms of ensemble modelling

- 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

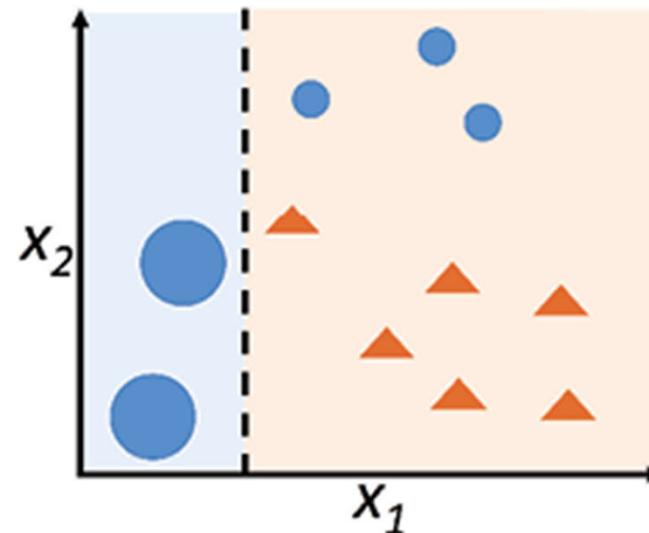
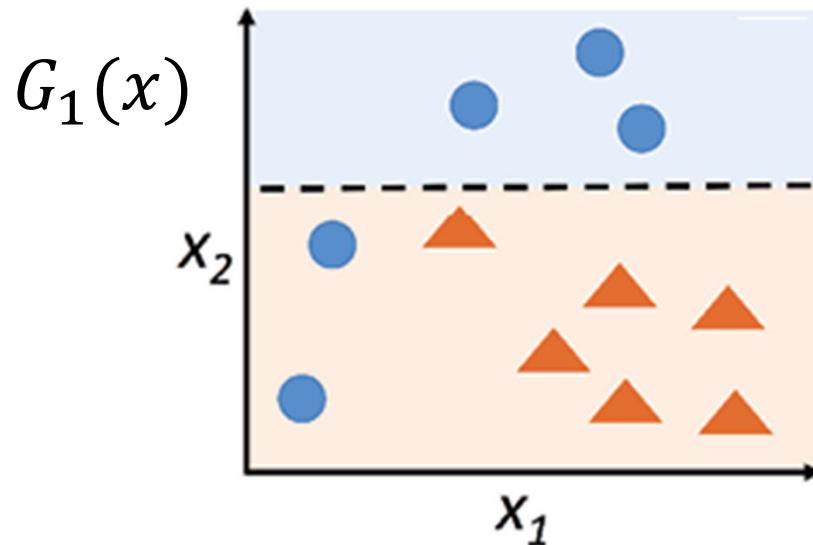
- Too hard to jointly optimize $\{G_m(x)\}$ and $\{\alpha_m\}$ in
$$y = \sum_{m=1}^M \alpha_m G_m(x).$$
- Instead, let's use a greedy approximation wherein we sequentially train the next $G_m(x)$ conditioned on all previous learned base models $\{G_{m-1}, G_{m-2}, \dots, G_1\}$, and their corresponding 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)



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

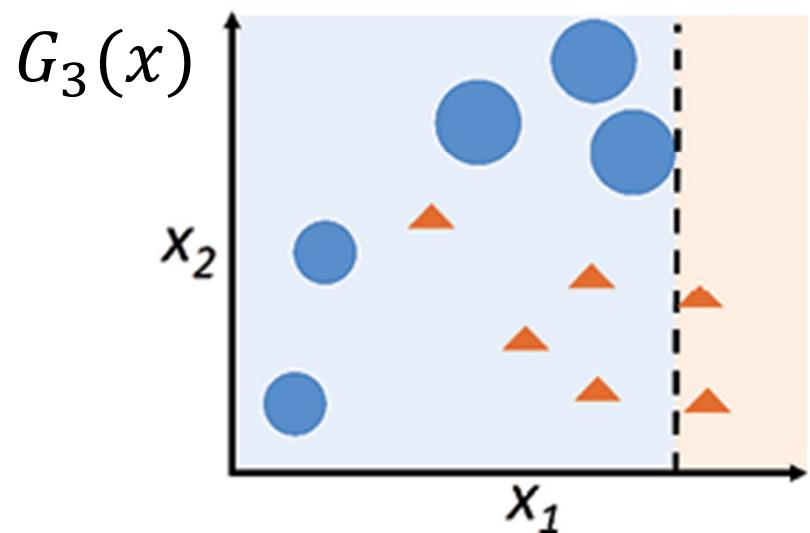
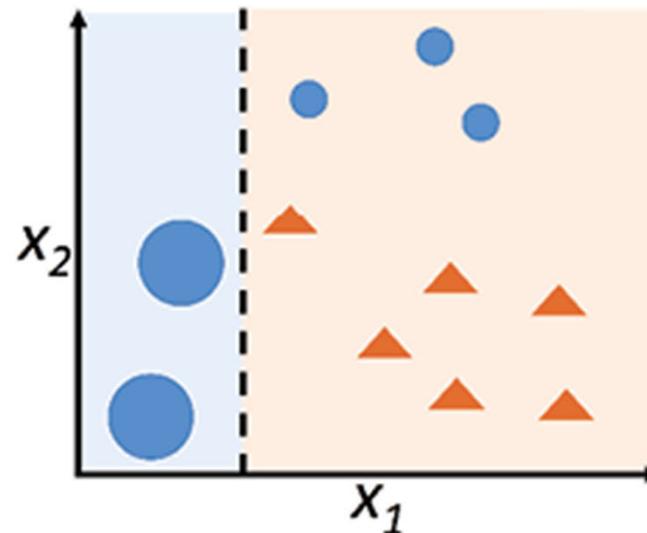
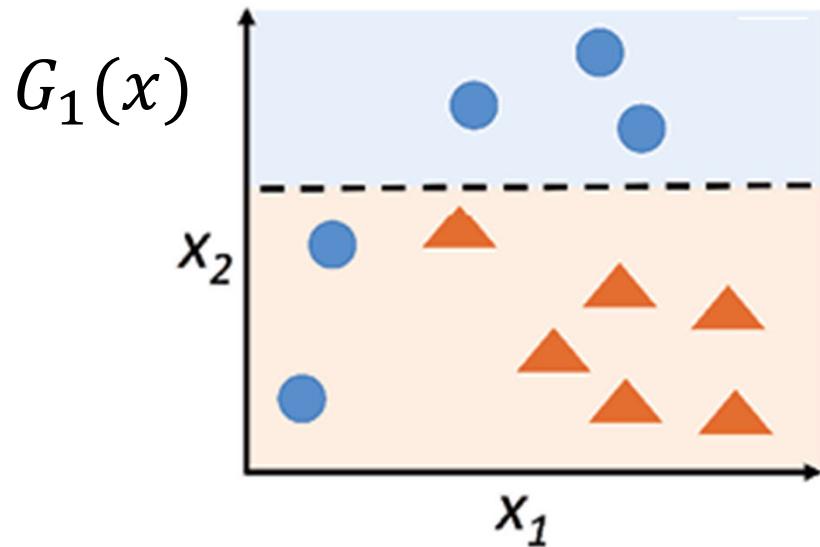
Boosting at an intuitive level (decision stump)



$$G_2(x)$$

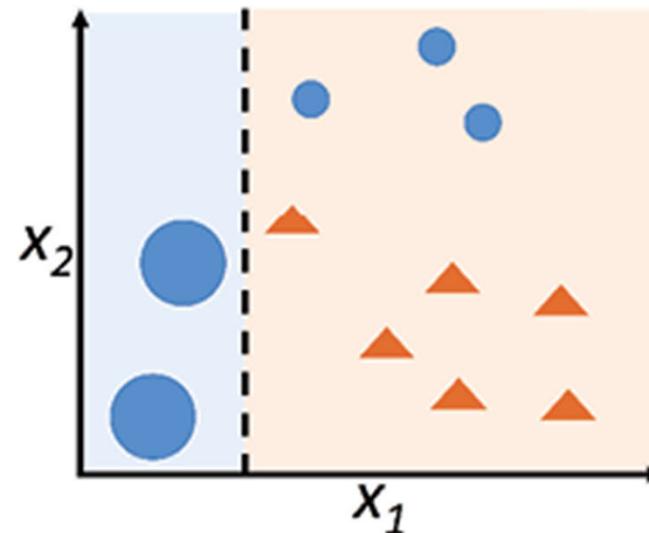
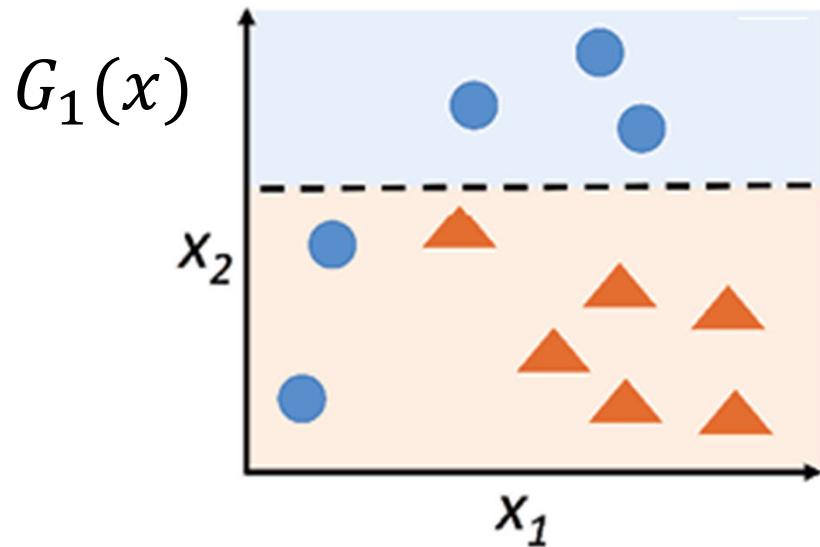
$$G(x) = \sum_{m=1}^M \alpha_m G_m(x)$$

Boosting at an intuitive level (decision stump)

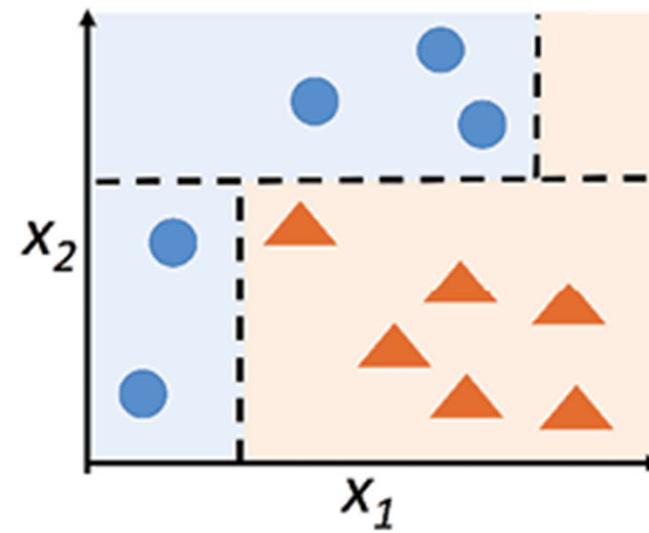
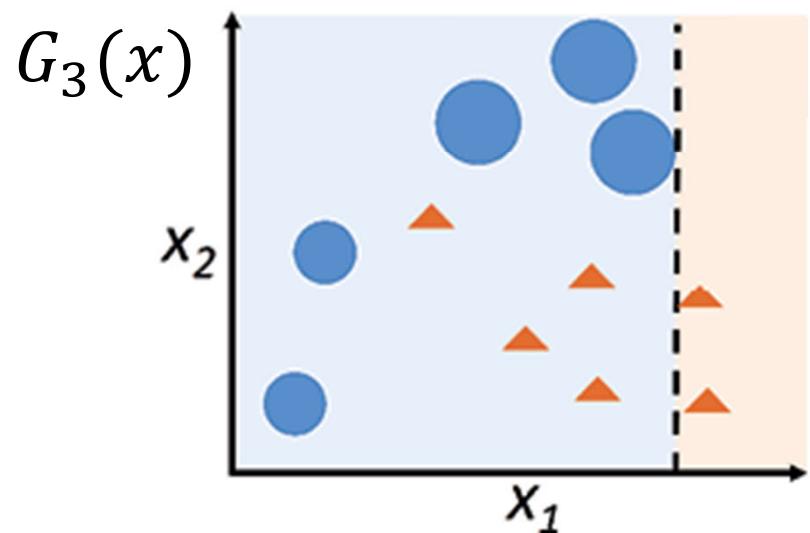


$$G(x) = \sum_{m=1}^M \alpha_m G_m(x)$$

Boosting at an intuitive level (decision stump)



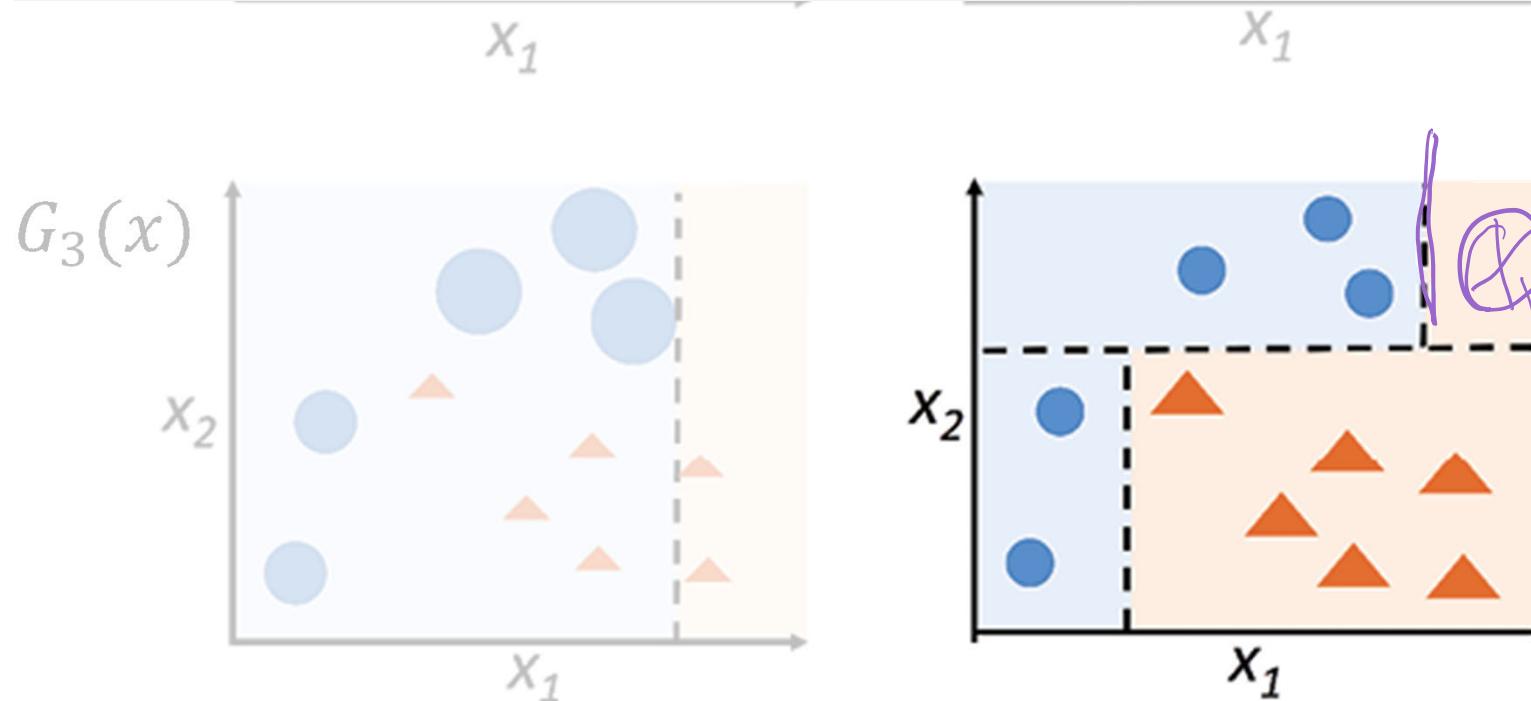
$G_2(x)$



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

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)