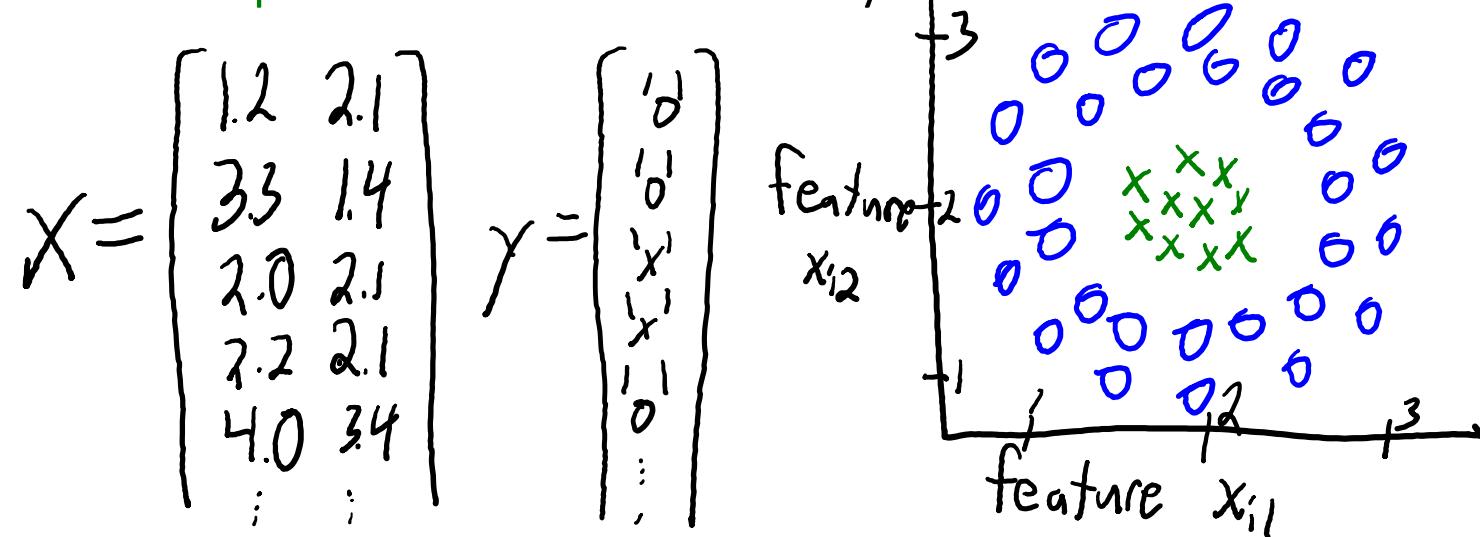


Which score function should a decision tree use?

- Shouldn't we just use accuracy score?
 - For leafs: yes, just maximize accuracy.
 - For internal nodes: not necessarily.
- Maybe no simple rule like ($egg > 0.5$) improves accuracy.
 - But this doesn't necessarily mean we should stop!

Example Where Accuracy Fails

- Consider a dataset with 2 features and 2 classes ('x' and 'o').
 - Because there are 2 features, we can draw 'X' as a scatterplot.
 - Colours and shapes denote the class labels 'y'.



- A decision stump would divide space by a horizontal or vertical line.
 - Testing whether $x_{i1} > t$ or whether $x_{i2} > t$.
- On this dataset no horizontal/vertical line improves accuracy.
 - Baseline is 'o', but need to get many 'o' wrong to get one 'x' right.

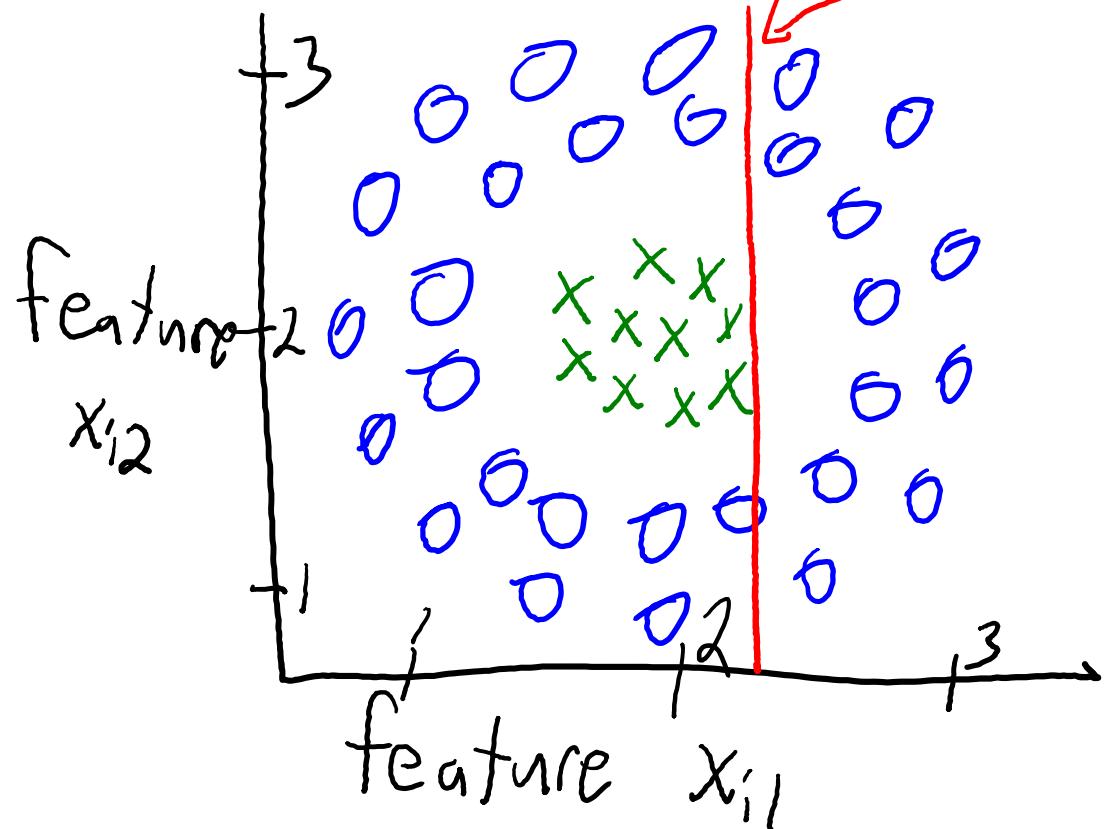
Which score function should a decision tree used?

- Most common score in practice is “**information gain**”.
 - “Choose split that decreases **entropy** of labels the most”.

$$\text{information gain} = \underbrace{\text{entropy}(y)}_{\substack{\text{entropy of labels} \\ \text{before split}}} - \underbrace{\frac{n_{\text{yes}}}{n} \text{entropy}(y_{\text{yes}}) + \frac{n_{\text{no}}}{n} \text{entropy}(y_{\text{no}})}_{\substack{\text{number of examples} \\ \text{satisfying rule}}} \underbrace{\text{entropy of labels for} \\ \text{examples satisfying rule.}}_{\text{example satisfying rule.}}$$

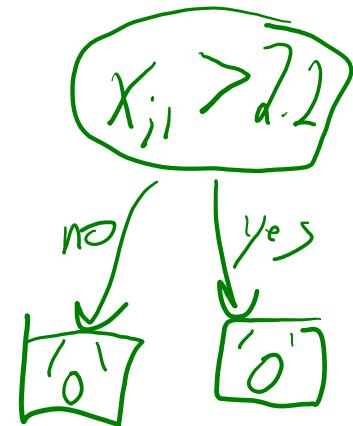
- Information gain for baseline rule (“do nothing”) is 0.
 - Infogain is large if labels are “more predictable” (“less random”) in next layer.
- Even if it does not increase classification accuracy at one depth, we hope that it makes classification easier at the next depth.

Example Where Accuracy Fails

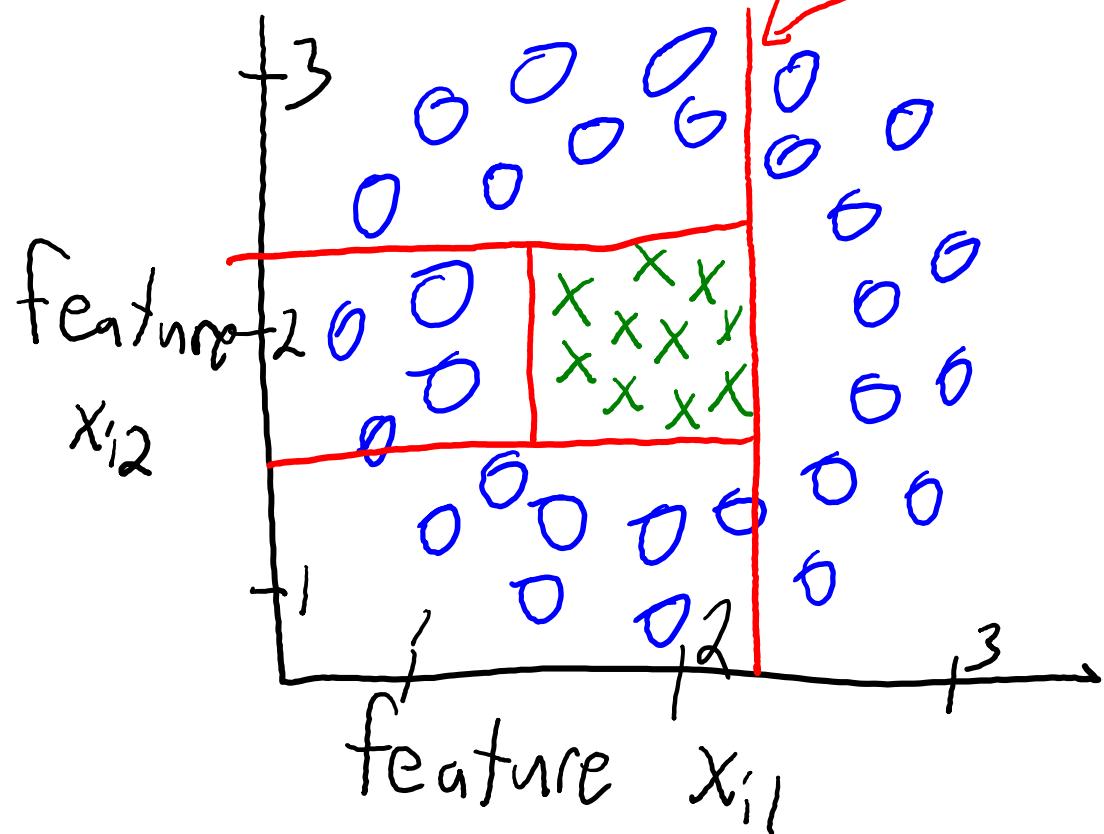


This split makes labels less random.
(Everything on the right is a '0')

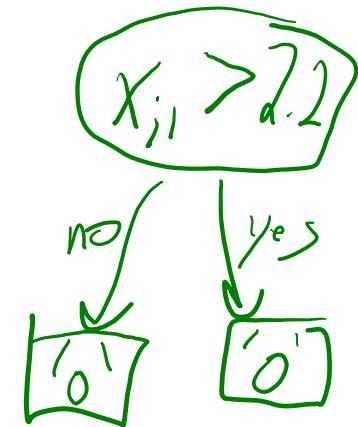
It did not improve accuracy.
(still classifies everything
as '0')



Example Where Accuracy Fails



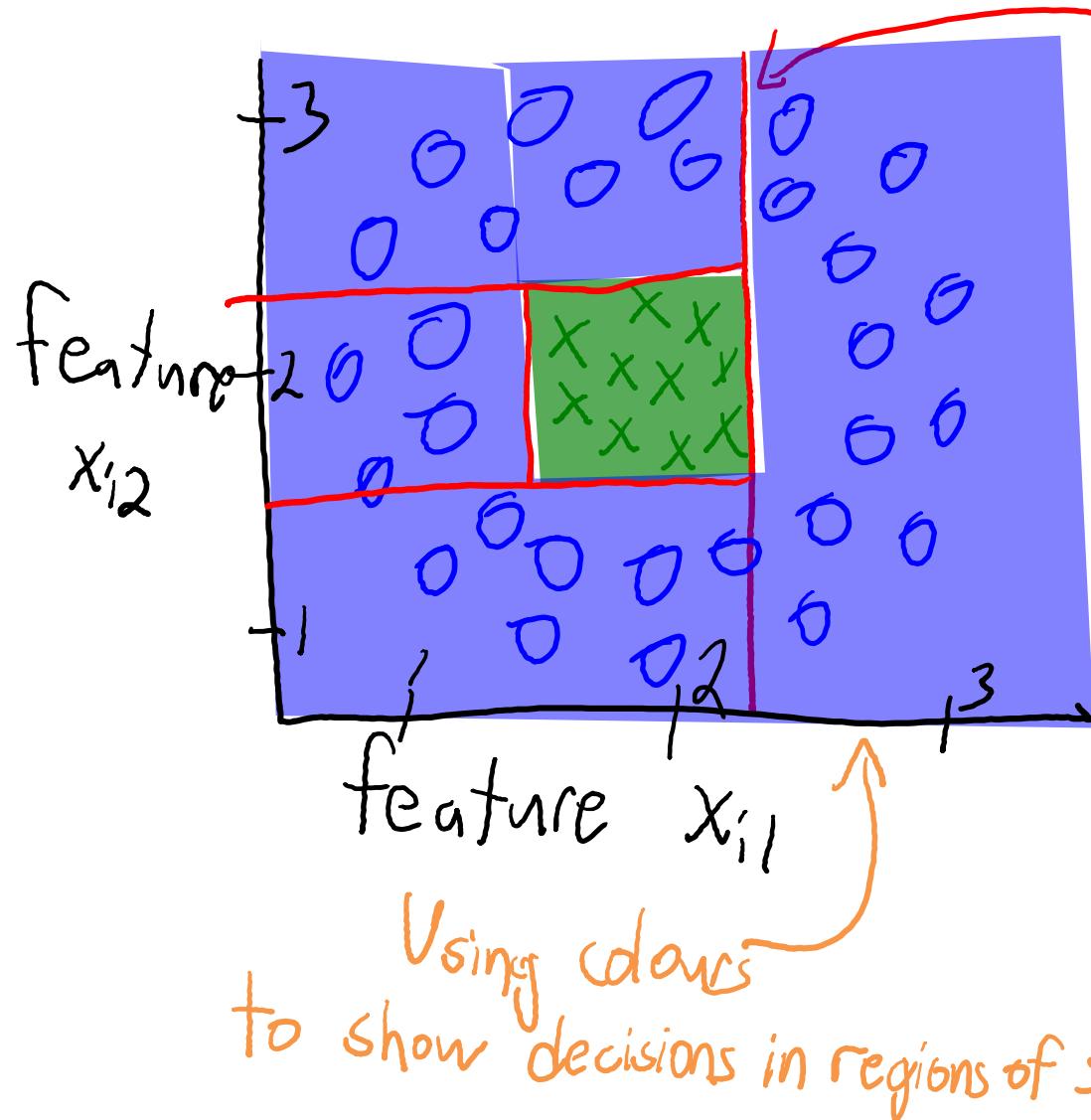
This split makes labels less random.
(Everything on the right is a 'o')



It did not improve accuracy.
(still classifies everything
as 'o')

But three more splits maximizing
infogain lead to perfect accuracy.

Example Where Accuracy Fails



This split makes labels less random.
(Everything on the right is a 'o')

↓
It did not improve accuracy.
(still classifies everything as 'o')

↓
But three more splits maximizing
infogain lead to perfect accuracy.

Entropy Function

Input: vector 'y' of length 'n' with numbers $\{1, 2, \dots, k\}$

counts = zeros(k)

for i in 1:n

 counts[y[i]] += 1

entropy = 0

for c in 1:k

 prob = counts[c]/n

 entropy -= prob * log(prob)

return entropy

Other Considerations for Food Allergy Example

- What types of **preprocessing** might we do?
 - **Data cleaning:** check for and fix missing/unreasonable values.
 - **Summary statistics:**
 - Can help identify “unclean” data.
 - Correlation might reveal an obvious dependence (“sick” \Leftrightarrow “peanuts”).
 - **Data transformations:**
 - Convert everything to same scale? (e.g., grams)
 - Add foods from day before? (maybe “sick” depends on multiple days)
 - Add date? (maybe what makes you “sick” changes over time).
 - **Data visualization:** look at a scatterplot of each feature and the label.
 - Maybe the visualization will show something weird in the features.
 - Maybe the pattern is really obvious!
- What you do might depend on how much data you have:
 - Very little data:
 - Represent food by common allergic ingredients (lactose, gluten, etc.)?
 - Lots of data:
 - Use more fine-grained features (bread from bakery vs. hamburger bun)?

Julia Decision Stump Code (not $O(n \log n)$ yet)

Input: feature matrix X and label vector y

$(n, d) = \text{size}(X)$

$\text{minError} = \sum(y \neq \text{mode}(y))$ compute error if you don't split (user-defined function "mode")
 $\text{minRule} = []$

for $j = 1:d$

 for $i = 1:n$

$t = X[i, j]$

$y_{\text{above}} = \text{mode}(y[X[:, j] > t])$

$y_{\text{below}} = \text{mode}(y[X[:, j] \leq t])$

$y_{\text{hat}} = \text{fill}(y_{\text{above}}, n)$

$y_{\text{hat}}[X[:, j] \leq t] = y_{\text{below}}$

$\text{error} = \sum(y_{\text{hat}} \neq y)$

 if $\text{error} < \text{minError}$

$\text{minError} = \text{error}$

$\text{minRule} = [j \ t]$

 for each feature ' j '

 for each example ' i '

 set threshold to feature ' j ' in example ' i '!

 find mode of label vector when feature ' j ' is above threshold

 find mode of label vector when feature ' j ' is below threshold.

classify all examples based on threshold

 count the number of errors.

store this rule if it has the lowest error so far.

Going from $O(n^2d)$ to $O(nd \log n)$ for Numerical Features

- Do we have to compute score from scratch?
 - As an example, assume we eat integer number of eggs:
 - So the rules ($\text{egg} > 1$) and ($\text{egg} > 2$) have same decisions, except when ($\text{egg} == 2$)).
- We can actually compute the best rule involving ‘egg’ in $O(n \log n)$:
 - Sort the examples based on ‘egg’, and use these positions to re-arrange ‘y’.
 - Go through the sorted values in order, updating the counts of #sick and #not-sick that both satisfy and don’t satisfy the rules.
 - With these counts, it’s easy to compute the classification accuracy (see bonus slide).
- Sorting costs $O(n \log n)$ per feature.
- Total cost of updating counts is $O(n)$ per feature.
- Total cost is reduced from $O(n^2d)$ to $O(nd \log n)$.
- This is a good runtime:
 - $O(nd)$ is the size of data, same as runtime up to a log factor.
 - We can apply this algorithm to huge datasets.

How do we fit stumps in $O(nd \log n)$?

- Let's say we're trying to find the best rule involving milk:

Egg	Milk	...	Sick?
0	0.7		1
1	0.7		1
0	0		0
1	0.6		1
1	0		0
2	0.6		1
0	1		1
2	0		1
0	0.3		0
1	0.6		0
2	0		1

First grab the milk column and sort it (using the sort positions to re-arrange the sick column). This step costs $O(n \log n)$ due to sorting.

Milk	Sick?
0	0
0	0
0	0
0	0
0.3	0
0.6	1
0.6	1
0.6	0
0.7	1
0.7	1
1	1

Now, we'll go through the milk values in order, keeping track of #sick and #not sick that are above/below the current value. E.g., #sick above 0.3 is 5.

With these counts, accuracy score is (sum of most common label above and below)/n.

How do we fit stumps in $O(nd \log n)$?

Milk	Sick?
0	0
0	0
0	0
0	0
0.3	0
0.6	1
0.6	1
0.6	0
0.7	1
0.7	1
1	1

Start with the baseline rule () which is always “satisfied”:

If satisfied, #sick=5 and #not-sick=6.

If not satisfied, #sick=0 and #not-sick=0.

This gives accuracy of $(6+0)/n = 6/11$.

Next try the rule ($\text{milk} > 0$), and update the counts based on these 4 rows:

If satisfied, #sick=5 and #not-sick=2.

If not satisfied, #sick=0 and #not-sick=4.

This gives accuracy of $(5+4)/n = 9/11$, which is better.

Next try the rule ($\text{milk} > 0.3$), and update the counts based on this 1 row:

If satisfied, #sick=5 and #not-sick=1.

If not satisfied, #sick=0 and #not-sick=5.

This gives accuracy of $(5+5)/n = 10/11$, which is better.

(and keep going until you get to the end...)

How do we fit stumps in $O(nd \log n)$?

Milk	Sick?
0	0
0	0
0	0
0	0
0.3	0
0.6	1
0.6	1
0.6	0
0.7	1
0.7	1
1	1

Notice that for each row, updating the counts only costs $O(1)$. Since there are $O(n)$ rows, total cost of updating counts is $O(n)$.

Instead of 2 labels (sick vs. not-sick), consider the case of ‘ k ’ labels:

- Updating the counts still costs $O(n)$, since each row has one label.
- But computing the ‘max’ across the labels costs $O(k)$, so cost is $O(kn)$.

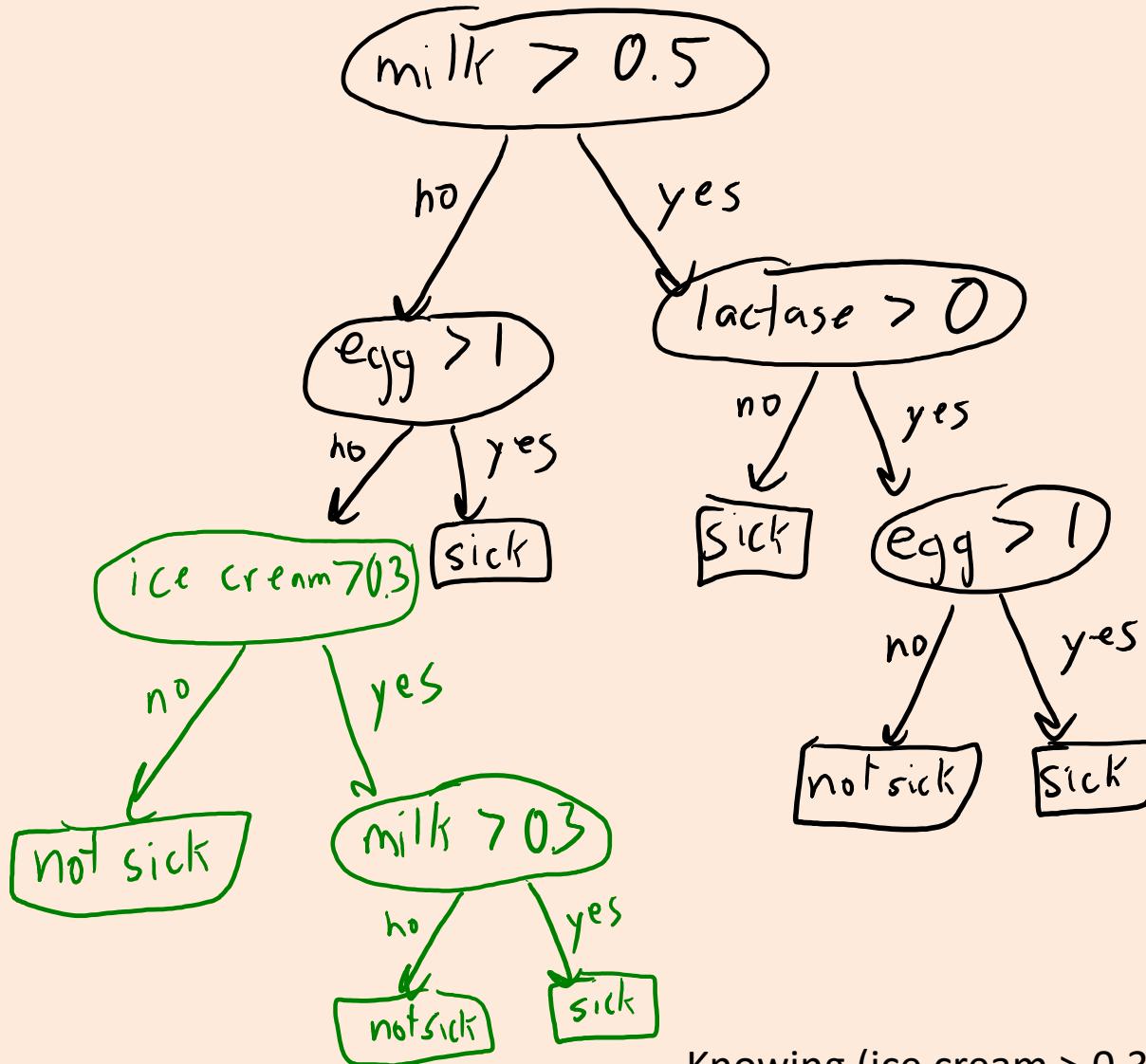
With ‘ k ’ labels, you can decrease cost using a “max-heap” data structure:

- Cost of getting max is $O(1)$, cost of updating heap for a row is $O(\log k)$.
- But $k \leq n$ (each row has only one label).
- So cost is in $O(\log n)$ for one row.

Since the above shows we can find best rule in one column in $O(n \log n)$, total cost to find best rule across all ‘ d ’ columns is $O(nd \log n)$.

Can decision trees re-visit a feature?

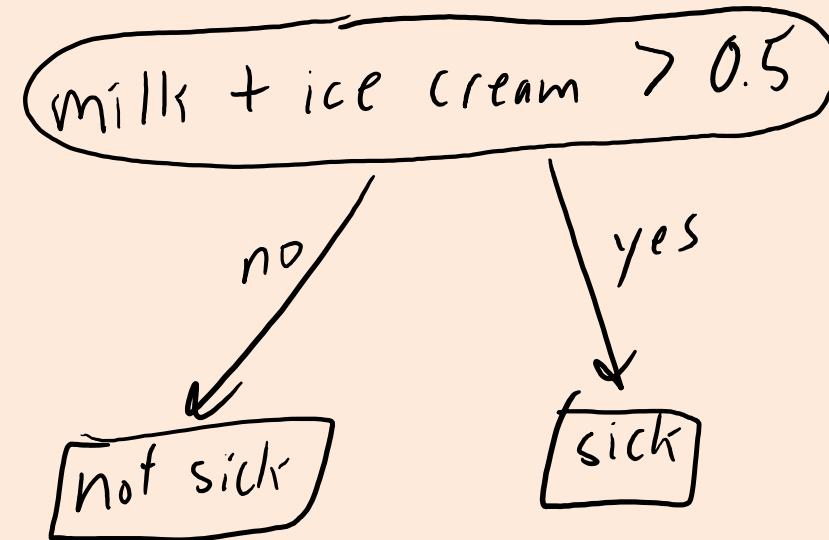
- Yes.



Knowing ($\text{ice cream} > 0.3$) makes small milk quantities relevant.

Can decision trees have more complicated rules?

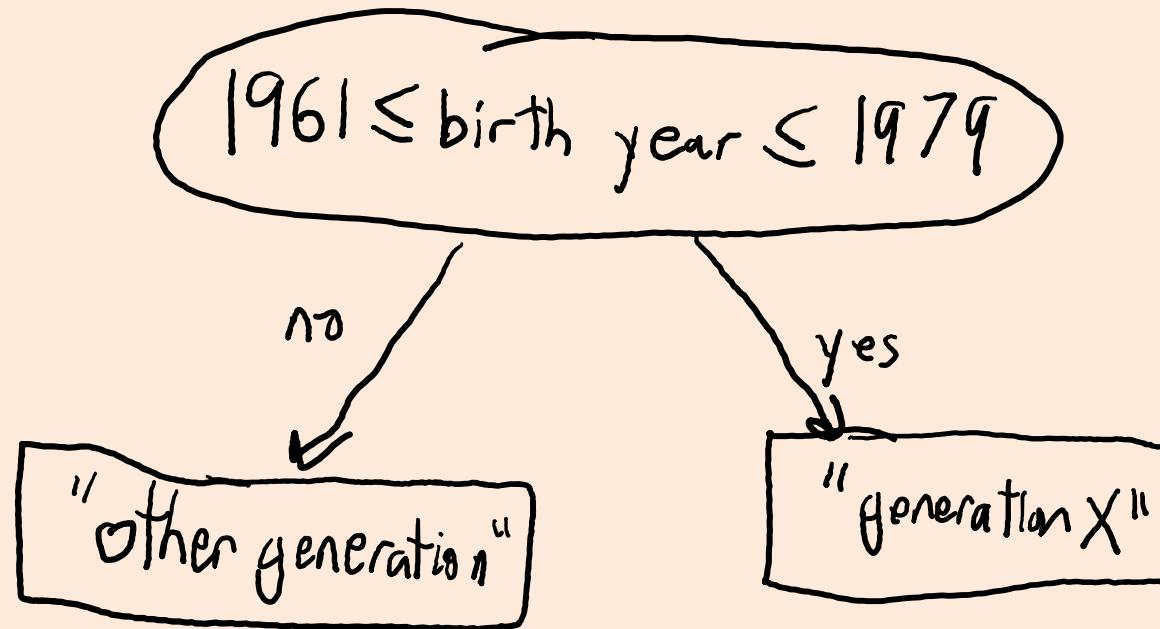
- Yes!
- Rules that depend on **more than one feature**:



- But now searching for the best rule can get expensive.

Can decision trees have more complicated rules?

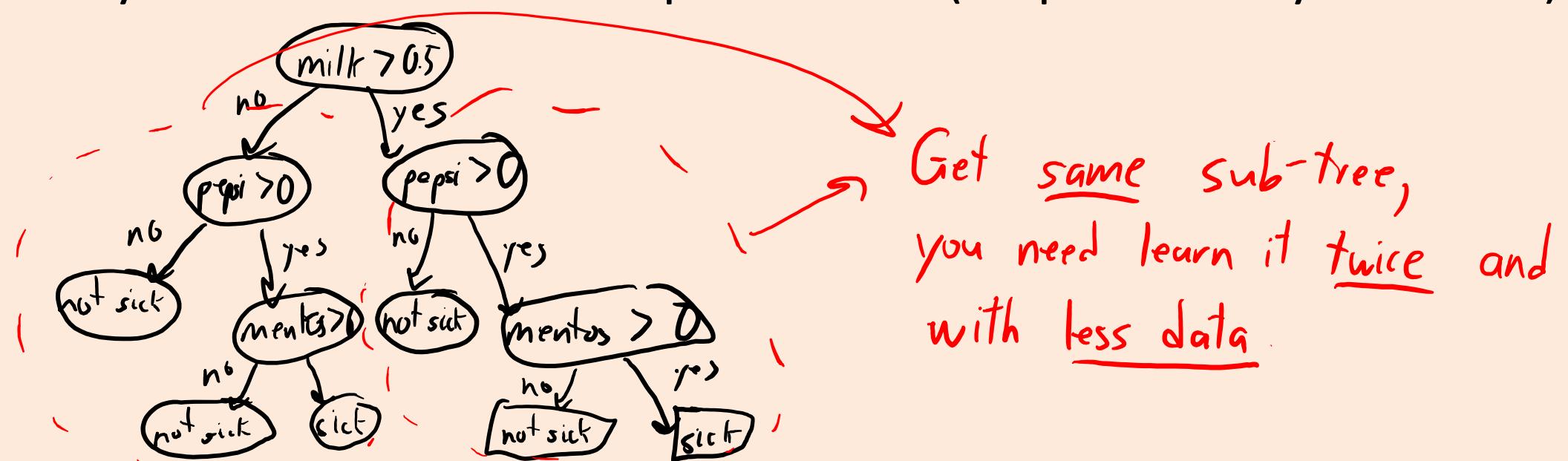
- Yes!
- Rules that depend on **more than one threshold**:



- **Very Simple Classification Rules Perform Well on Most Commonly Used Datasets**
 - Consider decision stumps based on multiple splits of 1 attribute.
 - Showed that this gives comparable performance to more-fancy methods on many datasets.

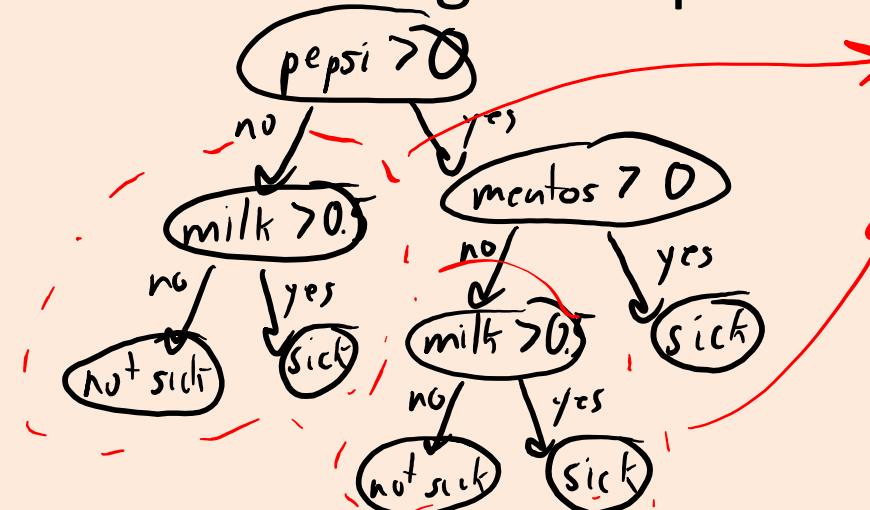
Does being greedy actually hurt?

- Can't you just go deeper to correct greedy decisions?
 - Yes, but you need to “re-discover” rules with less data.
- Consider that you are allergic to milk (and drink this often), and also get sick when you (rarely) combine diet coke with mentos.
- Greedy method should first split on milk (helps accuracy the most):



Does being greedy actually hurt?

- Can't you just go deeper to correct greedy decisions?
 - Yes, but you need to “re-discover” rules with less data.
- Consider that you are allergic to milk (and drink this often), and also get sick when you (rarely) combine diet coke with mentos.
- Greedy method should first split on milk (helps accuracy the most).
- Non-greedy method could get simpler tree (split on milk later):



Still has repeated structure,
but you should have more data
to estimate those splits.

Decision Trees with Probabilistic Predictions

- Often, we'll have multiple 'y' values at each leaf node.
- In these cases, we might **return probabilities** instead of a label.
- E.g., if in the leaf node we have 5 "sick" examples and 1 "not sick":
 - Return $p(y = \text{"sick"} | x_i) = 5/6$ and $p(y = \text{"not sick"} | x_i) = 1/6$.
- In general, a natural estimate of the probabilities at the leaf nodes:
 - Let ' n_k ' be the number of examples that arrive to leaf node 'k'.
 - Let ' n_{kc} ' be the number of times ($y == c$) in the examples at leaf node 'k'.
 - Maximum likelihood estimate for this leaf is $p(y = c | x_i) = n_{kc}/n_k$.

Alternative Stopping Rules

- There are more complicated rules for deciding when **not** to split.
- Rules based on **minimum sample size**.
 - Don't split any nodes where the number of examples is less than some 'm'.
 - Don't split any nodes that create children with less than 'm' examples.
 - These types of rules try to make sure that you have enough data to justify decisions.
- Alternately, you can use a **validation set** (see next lecture):
 - Don't split the node if it decreases an approximation of test accuracy.