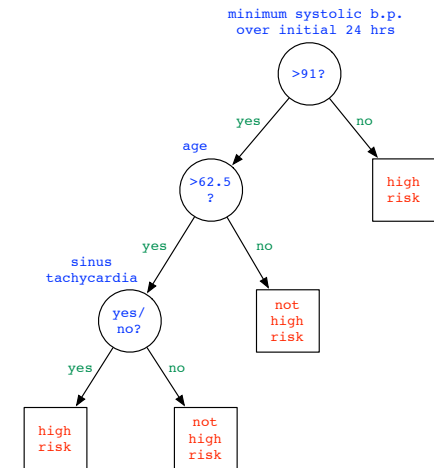# Decision trees

### CSE 250B

# Decision trees

Study at UCSD Medical Center, late 1970s.
Goal: identify patients at risk of dying within 30 days after heart attack.
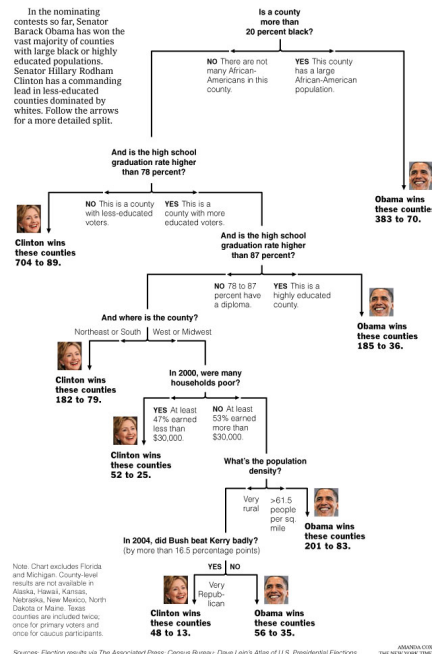
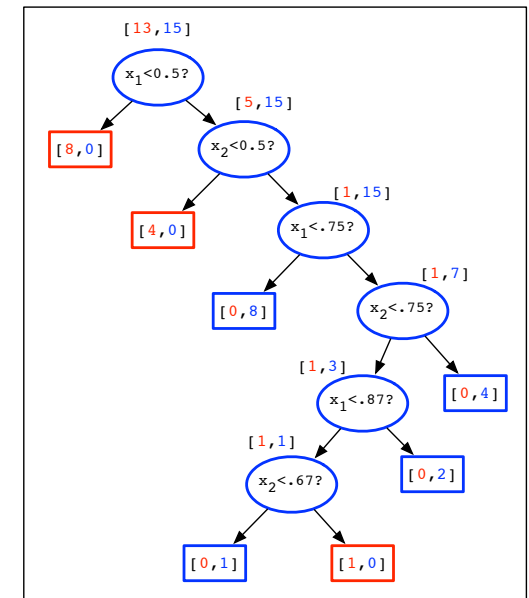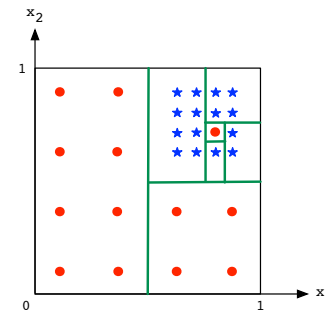Data set: 215 patients, 37 (=20%) died. 19 relevant variables.

minimum systolic b.p.
over initial 24 hrs

>91?

yes ——— no

age

>62.5
?

high
risk

yes ——— no

sinus
tachycardia

not
high
risk

yes/
no?

yes ——— no

high
risk

not
high
risk

## Decision Tree: The Obama-Clinton Divide

In the nominating contests so far, Senator Barack Obama has won the vast majority of counties with large black or highly educated populations. Senator Hillary Rodham Clinton has a commanding lead in less-educated counties dominated by whites. Follow the arrows for a more detailed split.

Is a county more than 20 percent black?

NO There are not many African-Americans in this county.

YES This county has a large African-American population.

And is the high school graduation rate higher than 78 percent?

NO This is a county with less-educated voters.

YES This is a county with more educated voters.

Clinton wins these counties 704 to 89.

And is the high school graduation rate higher than 87 percent?

Obama wins these counties 383 to 70.

NO 78 to 87 percent have a diploma.

YES This is a highly educated county.

And where is the county?

Obama wins these counties 185 to 36.

Northeast or South | West or Midwest

Clinton wins these counties 182 to 79.

In 2000, were many households poor?

YES At least 47% earned less than $30,000.

NO At least 53% earned more than $30,000.

Clinton wins these counties 52 to 25.

What's the population density?

Very rural | >61.5 people per sq. mile

In 2004, did Bush beat Kerry badly? (by more than 16.5 percentage points)

Obama wins these counties 201 to 83.

YES | NO

Very Repub-lican

Clinton wins these counties 48 to 13.

Obama wins these counties 56 to 35.

Note: Chart excludes Florida and Michigan. County-level results are not available in Alaska, Hawaii, Kansas, Nebraska, New Mexico, North Dakota or Maine. Texas counties are included twice, once for primary voters and once for caucus participants.

Sources: Election results via The Associated Press; Census Bureau; Dave Leip's Atlas of U.S. Presidential Elections

AMANDA COX / THE NEW YORK TIMES

# Example: building a decision tree

[13,15]

$x_1 < 0.5?$

[8,0]

[5,15]

$x_2 < 0.5?$

[4,0]

[1,15]

$x_1 < .75?$

[0,8]

[1,7]
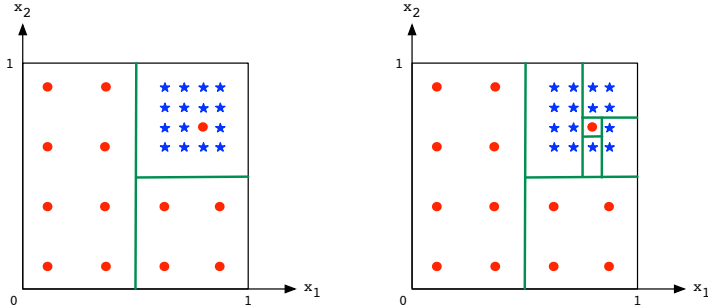
$x_2 < .75?$

[1,3]

$x_1 < .87?$

[0,4]

[1,1]

$x_2 < .67?$

[0,2]

[0,1]

[1,0]

# Overfitting?

Go back a few steps...



The final partition does better on the training data, but is more complex. That one point might have been an outlier anyway.

We have probably ended up **overfitting** the data.

# Decision tree issues

A very expressive family of classifiers:

- Can accommodate any type of data: real, Boolean, categorical, ...
- Can accommodate any number of classes
- Can fit any data set
- Statistically consistent

But this also means that there is serious danger of overfitting.

# Building a decision tree

Greedy algorithm: build tree top-down.

- Start with a single node containing all data points
- Repeat:
  - Look at all current leaves and all possible splits
  - Choose the split that most decreases the uncertainty in prediction

We need a measure of **uncertainty in prediction**.

# Uncertainty in prediction

Say there are two labels:

$$+ \quad p \text{ fraction of the points}$$
$$- \quad 1 - p \text{ fraction of the points}$$

What uncertainty score should we give to this?

❶ Misclassification rate

$$\min\{p, 1-p\}$$

❷ Gini index

$$2p(1-p)$$

❸ Entropy

$$p \log \frac{1}{p} + (1-p) \log \frac{1}{1-p}$$

# Uncertainty: $k$ classes

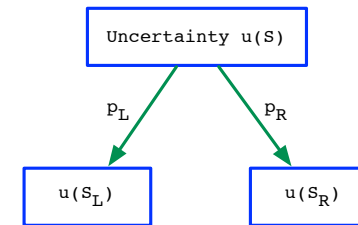Suppose there are $k$ classes, with probabilities $p_1, p_2, \ldots, p_k$.

|  | $k = 2$ | General $k$ |
|---|---|---|
| Misclassification rate | $\min\{p, 1-p\}$ | $1 - \max_i p_i = 1 - \|p\|_\infty$ |
| Gini index | $2p(1-p)$ | $\sum_{i \neq j} p_i p_j = 1 - \|p\|^2$ |
| Entropy | $p \log \dfrac{1}{p} + (1-p) \log \dfrac{1}{1-p}$ | $\sum_i p_i \log \dfrac{1}{p_i}$ |

# Benefit of a split

Let $u(S)$ be the uncertainty score for a set of labeled points $S$.

Consider a particular split:



Of the points in $S$:
- $p_L$ fraction go to $S_L$
- $p_R$ fraction go to $S_R$

Benefit of split = reduction in uncertainty:

$$\left( u(S) - \underbrace{(p_L\, u(S_L) + p_R\, u(S_R))}_{\text{expected uncertainty after split}} \right) \times |S|$$

# Benefit of a split: example



Initial uncertainty (Gini):

$$2 \times \frac{13}{28} \times \frac{15}{28}$$

[13,15]

$x_1 < 0.25?$

[4,0]   [9,15]

$$p_L u_L + p_R u_R = \frac{4}{28} \cdot 0 + \frac{24}{28} \cdot 2 \cdot \frac{9}{24} \cdot \frac{15}{24} = \frac{45}{112}$$

[13,15]

$x_1 < 0.5?$

[8,0]   [5,15]

$$p_L u_L + p_R u_R = \frac{8}{28} \cdot 0 + \frac{20}{28} \cdot 2 \cdot \frac{5}{20} \cdot \frac{15}{20} = \frac{30}{112}$$

# Building a decision tree

- Start with a single node containing all data points
- Repeat:
  - Look at all current leaves and all possible splits
  - Choose the split with the greatest benefit

When to stop?
- When each leaf is pure?
- When the tree is already pretty big?
- When each leaf has uncertainty below some threshold?

Common strategy: keep going until leaves are pure.
Then, shorten the tree by **pruning**, to correct for overfitting.

# What is overfitting?

Data comes from an unknown, underlying distribution $D$ on $\mathcal{X} \times \mathcal{Y}$. All we ever see are samples from $D$.

For a data set $(x_1, y_1), \ldots, (x_n, y_n)$, the **training error** of a classifier $h : \mathcal{X} \to \mathcal{Y}$ is
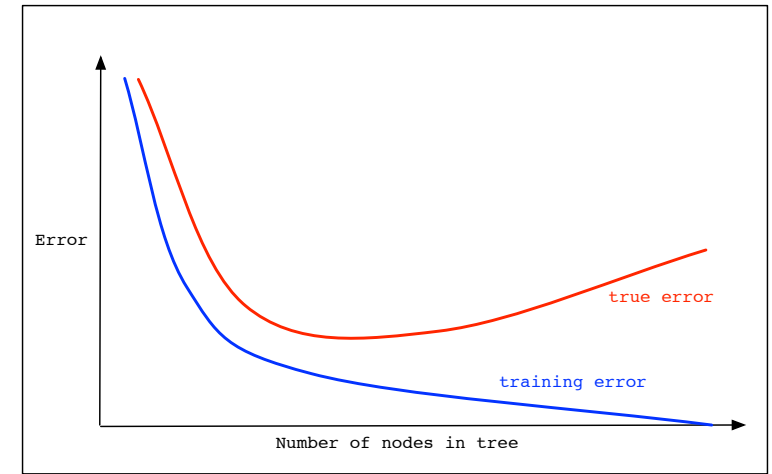
$$\widehat{\mathrm{err}}(h) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{1}(h(x_i) \neq y_i).$$

What we really care about is the **true error** of $h$ with respect to $D$:

$$\mathrm{err}(h) = \Pr_{(x,y) \sim D}(h(x) \neq y).$$

How are these two quantities related?

# Overfitting: picture



As we make our tree more and more complicated:
- training error keeps going down
- but, at some point, true error starts increasing!

# Overfitting: perspectives

- The true underlying distribution $D$ is the one whose structure we would like to capture.
- The training data reflects the structure of $D$, so it helps us.
- But it also has chance structure of its own – we must avoid modeling this.



Another perspective: it is absurd to fit a line to a point.
More generally, it is not good to use a model that is so complex that there isn't enough data to reliably estimate its parameters.

# Decision tree construction and pruning

❶ Split the training set into two parts
  - A smaller training set $S$
  - A validation set $V$ (surrogate test set, model of reality)
❷ Build a full decision tree using $S$
❸ Then prune using $V$
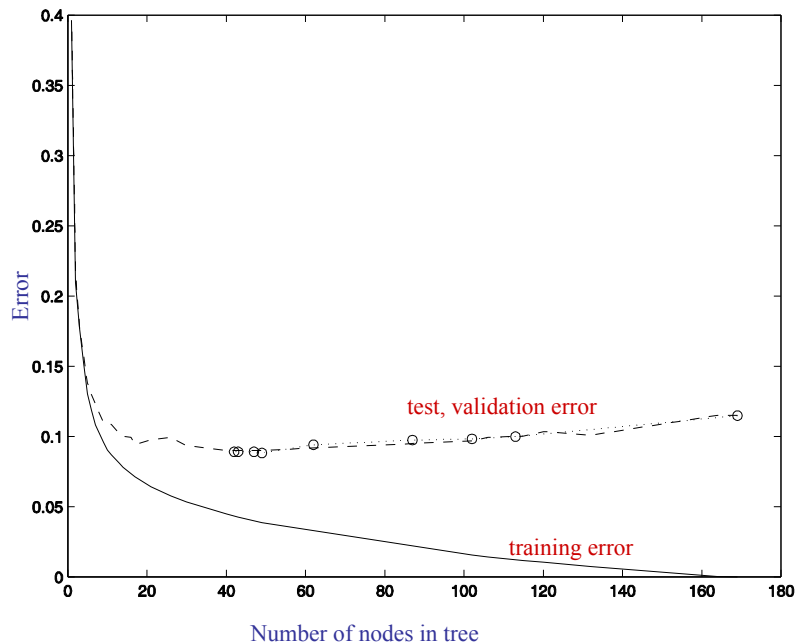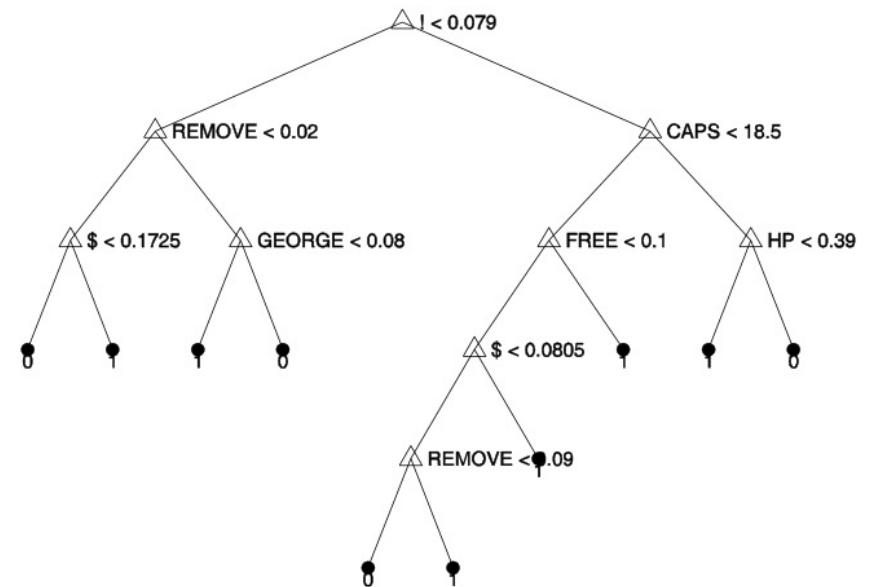  Use dynamic programming to find the pruning that does best on $V$

Of course, $V$ can have chance structure too — but its chance structure is unlikely to coincide with that of $S$.

# SPAMbase data set

- 4601 email messages, labeled SPAM or NOT SPAM.
- 39.4% are SPAM.
- Each email is represented by 57 features.
  - 48 check for specific words, e.g. FREE
  - 6 check for specific characters, e.g. !
  - 3 others, e.g. longest run of capitals

Randomly divide into three parts:

  *50% training data*
  *25% validation set*
  *25% test set*





## How accurate is the validation error?

For any classifier $h$ and underlying distribution $D$ on $\mathcal{X} \times \mathcal{Y}$:

**true error**    $\mathrm{err}(h) = \Pr_{(x,y)\sim D}(h(x) \neq y)$

**error on set $A$**    $\mathrm{err}(h, A) = \frac{1}{|A|} \sum_{(x,y)\in A} \mathbf{1}(h(x) \neq y)$

Suppose $A$ is chosen i.i.d. (independent, identically distributed) from $D$. Then (over the random choice of $A$),

$$\mathbb{E}[\mathrm{err}(h, A)] = \mathrm{err}(h)$$

and the standard deviation of $\mathrm{err}(h, A)$ is roughly $1/\sqrt{|A|}$.

Caution! In this scenario:
- Set $A$ is used to assess a single, prespecified classifier $h$.
- If $h$ was created using $A$ as a training set, the result does not apply. In such situations, $\mathrm{err}(h, A)$ could be a very poor estimate of $\mathrm{err}(h)$.