*Naïve Bayes & Logistic Regression,*
*See class website:*
>    **Mitchell's Chapter (required)**
>    **Ng & Jordan '02 (optional)**
*Gradient ascent and extensions:*
>    **Koller & Friedman Chapter 1.4**
*Decision Trees: many possible refs., e.g.,*
>    **Mitchell, Chapter 3**

# Logistic Regression (Continued) Generative v. Discriminative Decision Trees

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

February 1st, 2006

# Announcements

- Recitations stay on Thursdays
  - ☐ 5-6:30pm in Wean 5409
  - ☐ This week: Naïve Bayes & Logistic Regression

- **Extension** for the first homework:
  - ☐ **Due Wed. Feb 8th** beginning of class
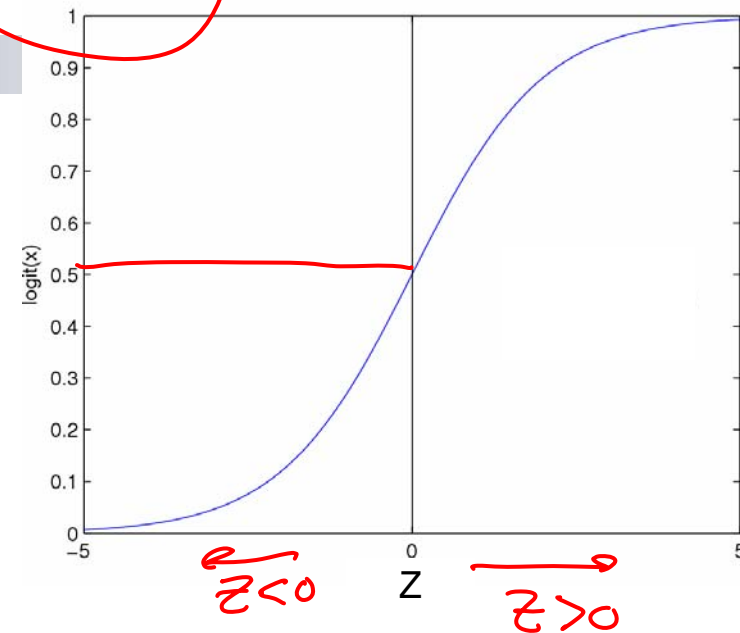  - ☐ Mitchell's chapter is most useful reading

# Logistic Regression
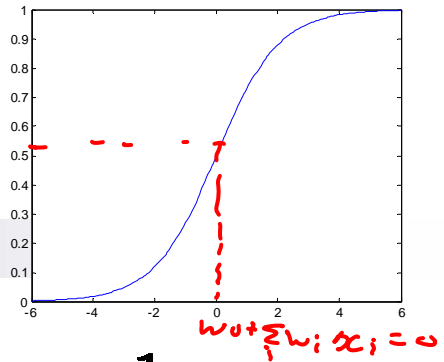
- ## Learn P(Y|**X**) directly!

  - ☐ Assume a particular functional form

  - ☐ Sigmoid applied to a linear function of the data:

  *regression part*

$$P(Y = 1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^{n} w_i X_i)}$$
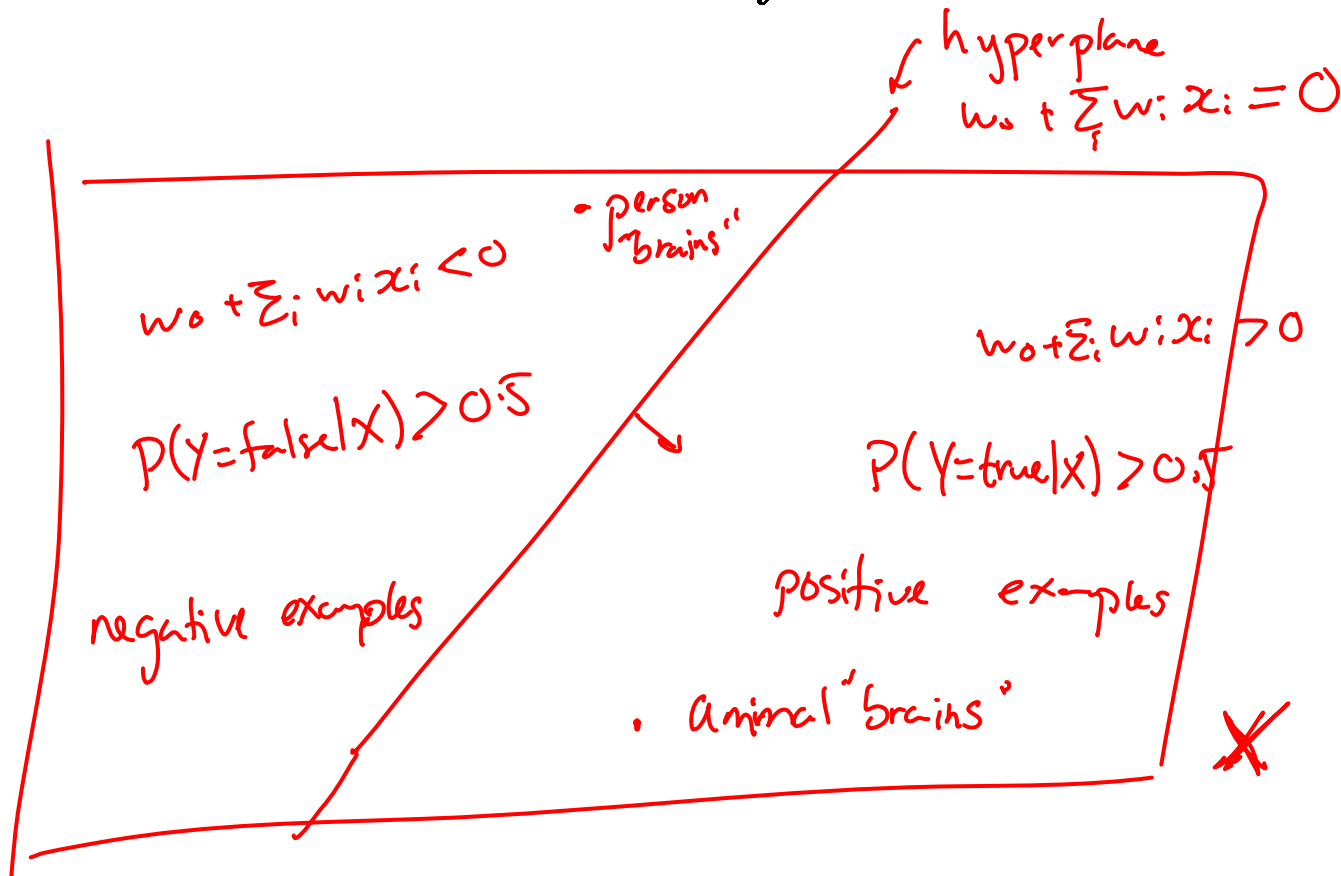


$z < 0$     z     $z > 0$

$P(Y=0|X) = 1 - P(Y=1|X)$

# Logistic Regression – a Linear classifier



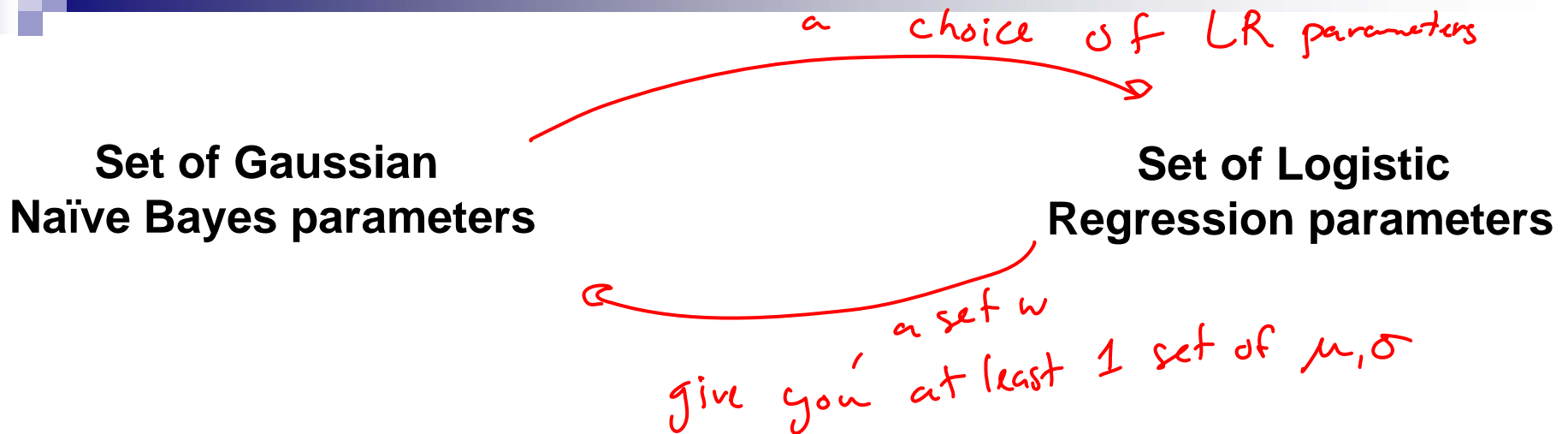$$g\left(w_0 + \sum_i w_i x_i\right) = \frac{1}{1 + e^{-\left(w_0 + \sum_i w_i x_i\right)}}$$

$w_0 + \sum_i w_i x_i = 0$

hyperplane
$w_0 + \sum_i w_i x_i = 0$

$w_0 + \sum_i w_i x_i < 0$

• person "brains"

$P(Y = false | X) > 0.5$

$w_0 + \sum_i w_i x_i > 0$

$P(Y = true | X) > 0.5$

negative examples

positive examples

• animal "brains"

# Logistic regression v. Naïve Bayes

- Consider learning f: X → Y, where
  - ☐ X is a vector of real-valued features, < X1 … Xn >
  - ☐ Y is boolean
- Could use a Gaussian Naïve Bayes classifier
  - ☐ assume all $X_i$ are conditionally independent given Y
  - ☐ model $P(X_i \mid Y = y_k)$ as Gaussian $N(\mu_{ik}, \sigma_i)$
  - ☐ model $P(Y)$ as Bernoulli$(\theta, 1-\theta)$

- What does that imply about the form of P(Y|X)?

$$P(Y = 1 \mid X = <X_1, ... X_n>) = \frac{1}{1 + exp(w_0 + \sum_i w_i X_i)}$$

**Cool!!!!**

# Gaussian Naïve Bayes v. Logistic Regression

*a choice of LR parameters*

**Set of Gaussian Naïve Bayes parameters**

**Set of Logistic Regression parameters**

*give you a set w at least 1 set of $\mu, \sigma$*

- Representation equivalence
  - **But only in a special case!!!** (GNB with class-independent variances)
- But what's the difference???
- **LR makes no assumptions about** $P(\mathbf{X}|Y)$ **in learning**!!!
- **Loss function!!!**
  - Optimize different functions → Obtain different solutions

# Logistic regression more generally

$w_i^T X = w_{0i} + \sum_j w_{ji} X_j$

- Logistic regression in more general case, where $Y \in \{Y_1 \dots Y_R\}$ : learn $R\text{-}1$ sets of weights

for $k<R$

$\sum_i P(y_i|X) = 1 = \sum_{i=1}^{R-1} e^{w_i^T X} + 1$

$1 + \sum_{j=1}^{R-1} e^{w_j^T X}$

$$P(Y = y_k|X) = \frac{\exp(w_{k0} + \sum_{i=1}^{n} w_{ki} X_i)}{1 + \sum_{j=1}^{R-1} \exp(w_{j0} + \sum_{i=1}^{n} w_{ji} X_i)}$$

for $k=R$ (normalization, so no weights for this class)

$$P(Y = y_R|X) = \frac{1}{1 + \sum_{j=1}^{R-1} \exp(w_{j0} + \sum_{i=1}^{n} w_{ji} X_i)}$$

categorical to numbers, e.g. $X_i = true \to +1$ / false $\to -1$

**Features can be discrete or continuous!**

# Logistic regression with more than 2 classes – an example

- $Y \in \{Y_1 \ldots Y_R\}$ : learn $R\text{-}1$ sets of weights

for $k < R$

$$P(Y = y_k | X) = \frac{\exp(w_{k0} + \sum_{i=1}^{n} w_{ki} X_i)}{1 + \sum_{j=1}^{R-1} \exp(w_{j0} + \sum_{i=1}^{n} w_{ji} X_i)}$$

for $k = R$ (normalization,
so no weights for this class)

$$P(Y = y_R | X) = \frac{1}{1 + \sum_{j=1}^{R-1} \exp(w_{j0} + \sum_{i=1}^{n} w_{ji} X_i)}$$

**Features can be discrete or continuous!**

$$\sum_{y_i} P(Y = y_i | X) = 1$$

$$Y = \{S, P, F\}$$

$$P(Y = S | X) \propto \frac{e^{w_{0S} + \sum_i w_{isX_i}}}{Z}$$

$$P(Y = P | X) \propto \frac{e^{w_{0P} + \sum_i w_{ip} X_i}}{Z}$$

$$P(Y = F | X) \propto \frac{1}{Z} \quad (\text{normalization})$$

$$P(Y = S|X) + P(Y = P|X) + P(Y = F|X) = 1$$

$$\frac{e^{w_{0s} + \sum_i w_{is}X_i}}{Z} + \frac{e^{w_{0p} + \sum_i w_{ip}X_i}}{Z} + \frac{1}{Z} = 1$$

$$Z = 1 + e^{w_{0s} + \sum_i w_{is}X_i} + e^{w_{0p} \sum w_{ip}X_i}$$

# Loss functions: Likelihood v. Conditional Likelihood

*chain rule of prob.:*

$$P(X, y \mid w) = P(y \mid x, w) \cdot P(X \mid w)$$

- Generative (Naïve Bayes) Loss function:

**Data likelihood**

$$\ln P(\mathcal{D} \mid \mathbf{w}) = \sum_{j=1}^{N} \ln P(\mathbf{x}^j, y^j \mid \mathbf{w})$$

$$= \sum_{j=1}^{N} \ln P(y^j \mid \mathbf{x}^j, \mathbf{w}) + \sum_{j=1}^{N} \ln P(\mathbf{x}^j \mid \mathbf{w})$$

*P(y|x)*      *P(x)*

- Discriminative models cannot compute P($\mathbf{x}^j$|$\mathbf{w}$)!
- But, discriminative (logistic regression) loss function:

**Conditional Data Likelihood**

$$\ln P(\mathcal{D}_Y \mid \mathcal{D}_{\mathbf{X}}, \mathbf{w}) = \sum_{j=1}^{N} \ln P(y^j \mid \mathbf{x}^j, \mathbf{w})$$

  - Doesn't waste effort learning P(X) – focuses on P(Y|$\mathbf{X}$) all that matters for classification

# Expressing Conditional Log Likelihood

$j^{th}$ example

$$l(\mathbf{w}) \equiv \sum_j \ln P(y^j | \mathbf{x}^j, \mathbf{w})$$

$$y^j = \begin{cases} 0 \\ 1 \end{cases}$$

$$P(Y = 0 | \mathbf{X}, \mathbf{w}) = \frac{1}{1 + exp(w_0 + \sum_i w_i X_i)}$$

$$P(Y = 1 | \mathbf{X}, \mathbf{w}) = \frac{exp(w_0 + \sum_i w_i X_i)}{1 + exp(w_0 + \sum_i w_i X_i)}$$
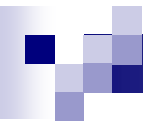
if $y_j = 0$

if $y_j = 1$

$$l(\mathbf{w}) = \sum_j y^j \ln P(y^j = 1 | \mathbf{x}^j, \mathbf{w}) + (1 - y^j) \ln P(y^j = 0 | \mathbf{x}^j, \mathbf{w})$$

$$\ln P(y^j = 1 | x^j, w)$$

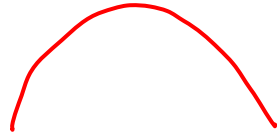$$= w_0 + \sum_i w_i x_i - \ln\left(1 + e^{w_0 + \sum_i w_i x_i}\right)$$

$$\ln P(y_j = 0 | x^j w) = \ln 1 - \ln\left(1 + e^{w_0 + \sum_i w_i x_i}\right)$$

# Maximizing Conditional Log Likelihood

$$P(Y = 0|X, W) = \frac{1}{1 + exp(w_0 + \sum_i w_i X_i)}$$

$$P(Y = 1|X, W) = \frac{exp(w_0 + \sum_i w_i X_i)}{1 + exp(w_0 + \sum_i w_i X_i)}$$

$$l(\mathbf{w}) \equiv \ln \prod_j P(y^j|\mathbf{x}^j, \mathbf{w})$$

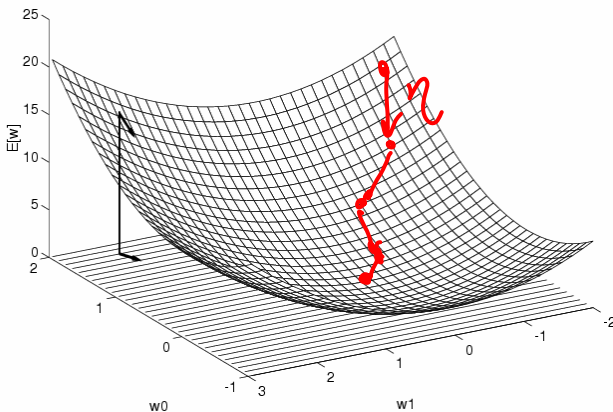$$= \sum_j y^j(w_0 + \sum_i^n w_i x_i^j) - \ln(1 + exp(w_0 + \sum_i^n w_i x_i^j))$$

Good news: $l(\mathbf{w})$ is concave function of $\mathbf{w} \rightarrow$ no locally optimal solutions

Bad news: no closed-form solution to maximize $l(\mathbf{w})$

Good news: concave functions easy to optimize

# Optimizing concave function – Gradient ascent

- Conditional likelihood for Logistic Regression is concave
  $\rightarrow$ Find optimum with gradient ascent



**Gradient:** $\qquad \nabla_{\mathbf{w}} l(\mathbf{w}) = [\frac{\partial l(\mathbf{w})}{\partial w_0}, \ldots, \frac{\partial l(\mathbf{w})}{\partial w_n}]'$

**Learning rate, $\eta > 0$**

**Update rule:** $\qquad \triangle \mathbf{w} = \eta \nabla_{\mathbf{w}} l(\mathbf{w})$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \frac{\partial l(\mathbf{w})}{\partial w_i}$$

local minima:

local      global

- Gradient ascent is simplest of optimization approaches

  with line search.

  - e.g., Conjugate gradient ascent much better (see reading)

# Maximize Conditional Log Likelihood: Gradient ascent

$\frac{\partial}{\partial w} e^{f(w)} = f'(w) e^{f(w)}$

$\frac{\partial}{\partial w} \ln f(w) = \frac{f'(w)}{f(w)}$

Derivative $\frac{\partial \ell}{\partial w_0}$ is slightly different !!

$$l(\mathbf{w}) = \sum_j y^j (w_0 + \sum_i^n w_i x_i^j) - \ln(1 + exp(w_0 + \sum_{i=1}^n w_i x_i^j))$$

$i = 1:n$

$\frac{\partial \ell}{\partial w_i} = \sum_j y^j x_i^j - \frac{x_i^j e^{w_0 + \sum_i w_i x_i^j}}{1 + e^{w_0 + \sum_i w_i x_i^j}}$

$= \sum_j x_i^j \left( y^j - \frac{e^{w_0 + \sum_i w_i x_i^j}}{1 + e^{w_0 + \sum_i w_i x_i^j}} \right)$

$= \sum_j x_i^j \left[ y^j - P(Y=1 \mid x^j, w) \right]$

$\frac{\partial \ell}{\partial w_0} = \sum_j [y^j - P(Y=1 \mid x^j, w)]$

---

Gradient ascent algorithm: iterate until change $< \varepsilon$

For all $i_{=1:n}$     $w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 \mid \mathbf{x}^j, \mathbf{w}^{(t)})]$
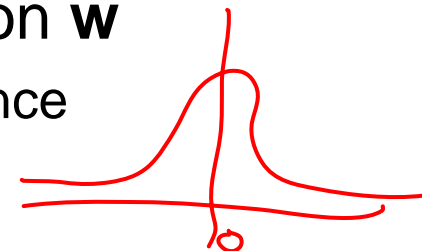
repeat     $i = 0$   slightly different

# That's all M(C)LE. How about MAP?

$$p(\mathbf{w} \mid Y, \mathbf{X}) \quad \propto \quad P(Y \mid \mathbf{X}, \mathbf{w})p(\mathbf{w})$$

- One common approach is to define priors on **w**
  - Normal distribution, zero mean, identity covariance
  - "Pushes" parameters towards zero
- Corresponds to ***Regularization***
  - Helps avoid very large weights and overfitting
  - Explore this in your homework
  - More on this later in the semester

- MAP estimate

$$\mathbf{w}^* = \arg\max_{\mathbf{w}} \ln \left[ p(\mathbf{w}) \prod_{j=1}^{N} P(y^j \mid \mathbf{x}^j, \mathbf{w}) \right]$$

prior

likelihood

# M(C)AP as Regularization

$w_i$ indep $w_j$

$w_i \sim N(0, \kappa^2)$

$$\ln \left[ p(\mathbf{w}) \prod_{j=1}^{N} P(y^j \mid \mathbf{x}^j, \mathbf{w}) \right]$$

$$p(\mathbf{w}) = \prod_{i} \frac{1}{\kappa\sqrt{2\pi}} \; e^{\frac{-w_i^2}{2\kappa^2}}$$

$= \ln p(w) + \ln \prod_{j=1}^{N} P(y^j | x_j \, w)$

$= \sum_i \ln \frac{1}{\kappa \sqrt{2\pi}} - \frac{w_i^2}{2\kappa^2} + \ell(w)$

maximize  $\ell(w)$  $- \sum_i \frac{w_i^2}{2\kappa^2}$
w

fit data

small parameters
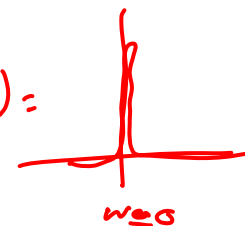
pick K large

can overfit $P(w)$:

$w=0$

pick K small

very biased

$P(w)$:

$w=0$

**Penalizes high weights, also applicable in linear regression (see homework)**

# Gradient of M(C)AP

$$\frac{\partial}{\partial w_i} \ln \left[ p(\mathbf{w}) \prod_{j=1}^{N} P(y^j \mid \mathbf{x}^j, \mathbf{w}) \right] \qquad p(\mathbf{w}) = \prod_i \frac{1}{\kappa\sqrt{2\pi}} \; e^{\frac{-w_i^2}{2\kappa^2}}$$

$$= \frac{\partial}{\partial w_i} \ln p(w) \quad + \quad \frac{\partial}{\partial w_i} \ln \underbrace{\prod_{j=1}^{N} P(y^j \mid x^j, w)}_{\ell(w)}$$

$$\hookrightarrow \frac{\partial}{\partial w_i} \left( ? - \frac{w_i^2}{2\kappa^2} \right)$$

$$\underbrace{- \frac{w_i}{\kappa^2}}_{} + \frac{\partial}{\partial w_i} \ell(w)$$

# MLE vs MAP

$$P(Y=1|x) = \frac{e^{w_0 + \sum w_i x_i}}{1 + e^{\cdots}}$$

$$P(Y=0|x) = \frac{1}{1 + e^{\cdots}}$$

- **Maximum conditional likelihood estimate**

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \ln \left[ \prod_{j=1}^{N} P(y^j \mid \mathbf{x}^j, \mathbf{w}) \right]$$

$$w_i \leftarrow w_i + \eta \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 \mid \mathbf{x}^j, \mathbf{w})]$$   *standard*

- **Maximum conditional a posteriori estimate**

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \ln \left[ p(\mathbf{w}) \prod_{j=1}^{N} P(y^j \mid \mathbf{x}^j, \mathbf{w}) \right]$$

$$\lambda = \frac{1}{\kappa^2}$$

$$w_i \leftarrow w_i + \eta \left\{ -\lambda w_i + \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 \mid \mathbf{x}^j, \mathbf{w})] \right\}$$

# Naïve Bayes vs Logistic Regression

Consider Y boolean, $X_i$ continuous, $X=<X_1 \ldots X_n>$

Number of parameters:

- NB: $4n + 1$
- LR: $n+1$

Estimation method:

- NB parameter estimates are uncoupled
- LR parameter estimates are coupled
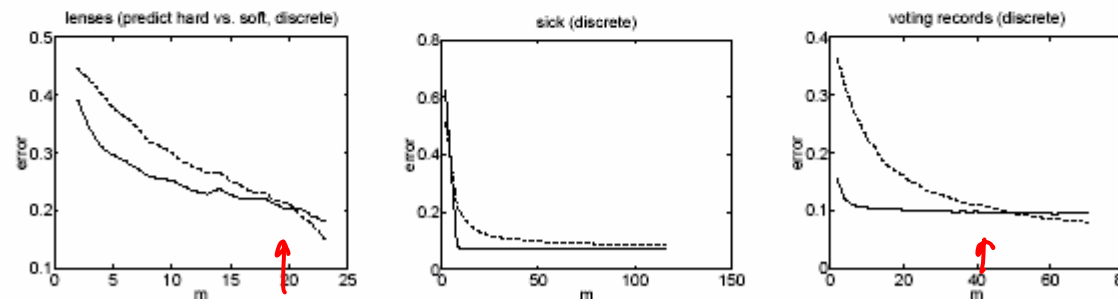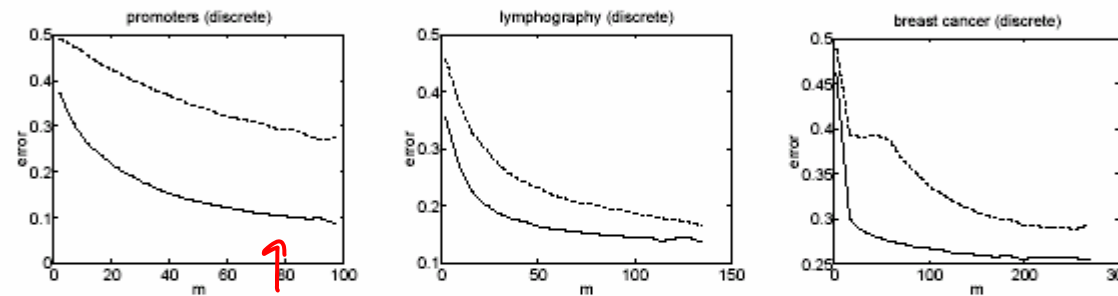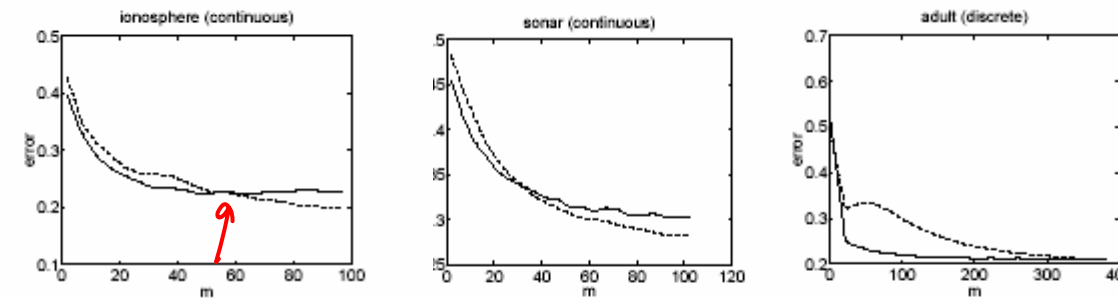
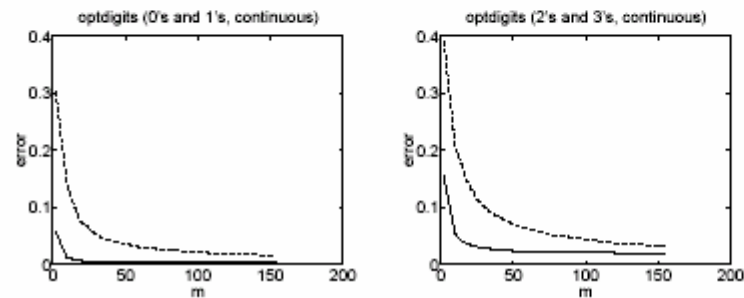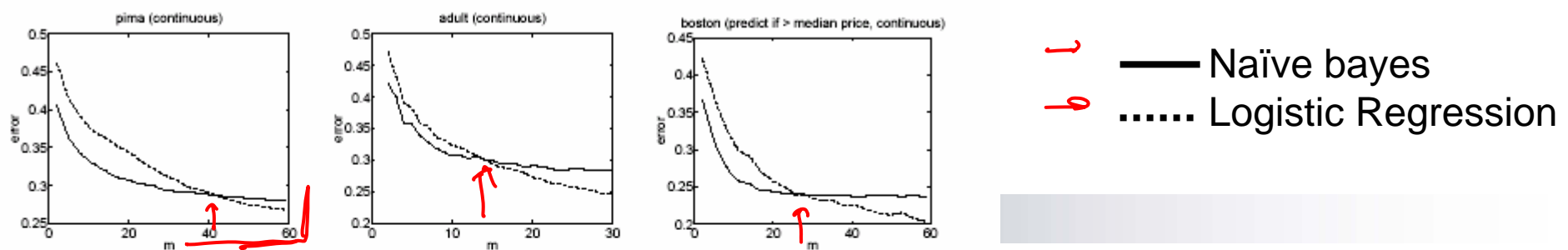# G. Naïve Bayes vs. Logistic Regression 1

- Generative and Discriminative classifiers

- Asymptotic comparison (# training examples → infinity)
  - when model correct    *indep assumptions*
    - GNB, LR produce identical classifiers

  - when model incorrect
    - LR is less biased – does not assume conditional independence
      - **therefore LR expected to outperform GNB**

# G. Naïve Bayes vs. Logistic Regression 2

- Generative and Discriminative classifiers


- Non-asymptotic analysis
  - convergence rate of parameter estimates, $n$ = # of attributes in X
    - Size of training data to get close to infinite data solution
    - GNB needs $O(\log n)$ samples
    - LR needs $O(n)$ samples


  - **GNB converges more quickly to its (perhaps less helpful) asymptotic estimates**

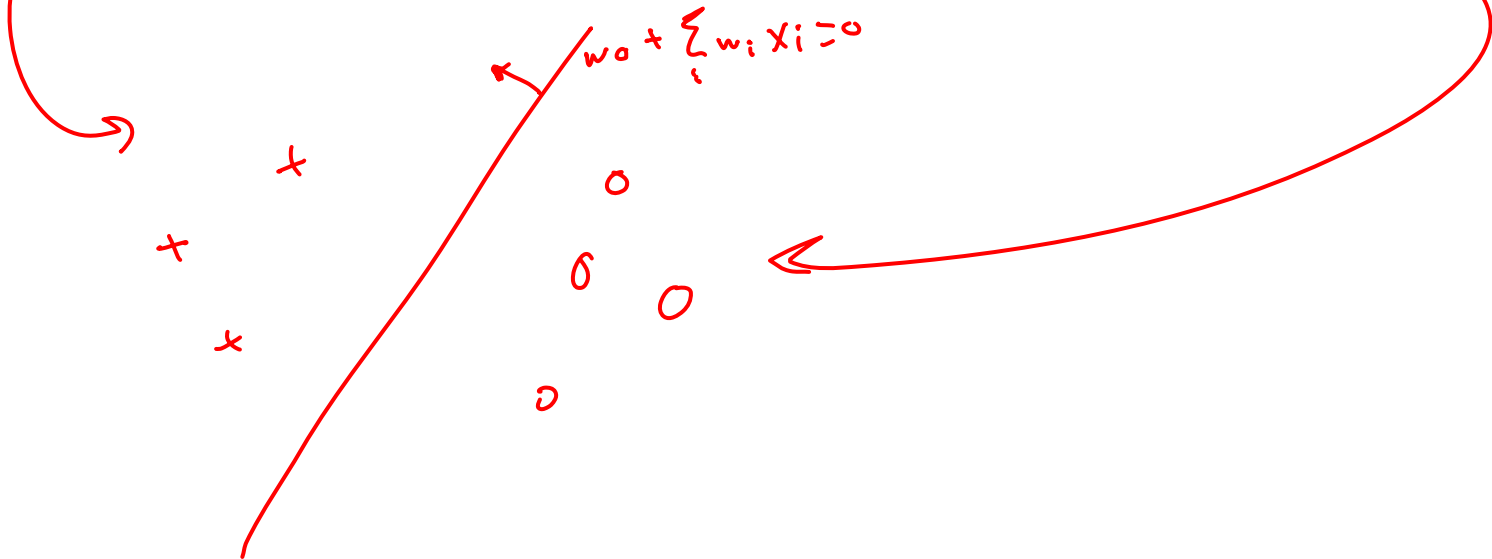Some experiments from UCI data sets

Figure 1: Results of 15 experiments on datasets from the UCI Machine Learning repository. Plots are of generalization error vs. $m$ (averaged over 1000 random train/test splits). Dashed line is logistic regression; solid line is naïve Bayes.

# What you should know about Logistic Regression (LR)

- **Gaussian Naïve Bayes with class-independent variances representationally equivalent to LR**
  - ☐ Solution differs because of objective (loss) function
- **In general, NB and LR make different assumptions**
  - ☐ NB: Features independent given class $\rightarrow$ assumption on P($\mathbf{X}$|Y)
  - ☐ LR: Functional form of P(Y|$\mathbf{X}$), no assumption on P($\mathbf{X}$|Y)
- **LR is a linear classifier**
  - ☐ decision rule is a hyperplane
- **LR optimized by conditional likelihood**
  - ☐ no closed-form solution
  - ☐ concave $\rightarrow$ global optimum with gradient ascent
  - ☐ Maximum conditional a posteriori corresponds to regularization
- **Convergence rates**
  - ☐ GNB (usually) needs more data
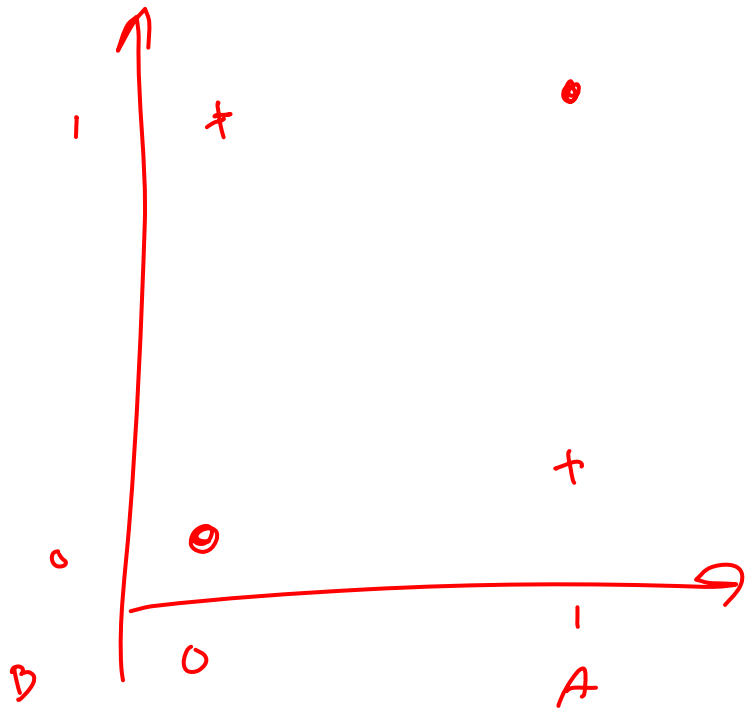  - ☐ LR (usually) gets to better solutions in the limit

# Linear separability

- A dataset is **linearly separable** iff ∃ a **separating hyperplane**:
  - ∃ **w**, such that:
    - $w_0 + \sum_i w_i x_i > 0$; if **x**=$\{x_1, \dots, x_n\}$ is a positive example
    - $w_0 + \sum_i w_i x_i < 0$; if **x**=$\{x_1, \dots, x_n\}$ is a negative example

$$w_0 + \sum_i w_i x_i = 0$$

# Not linearly separable data

- Some datasets are **not linearly separable!**
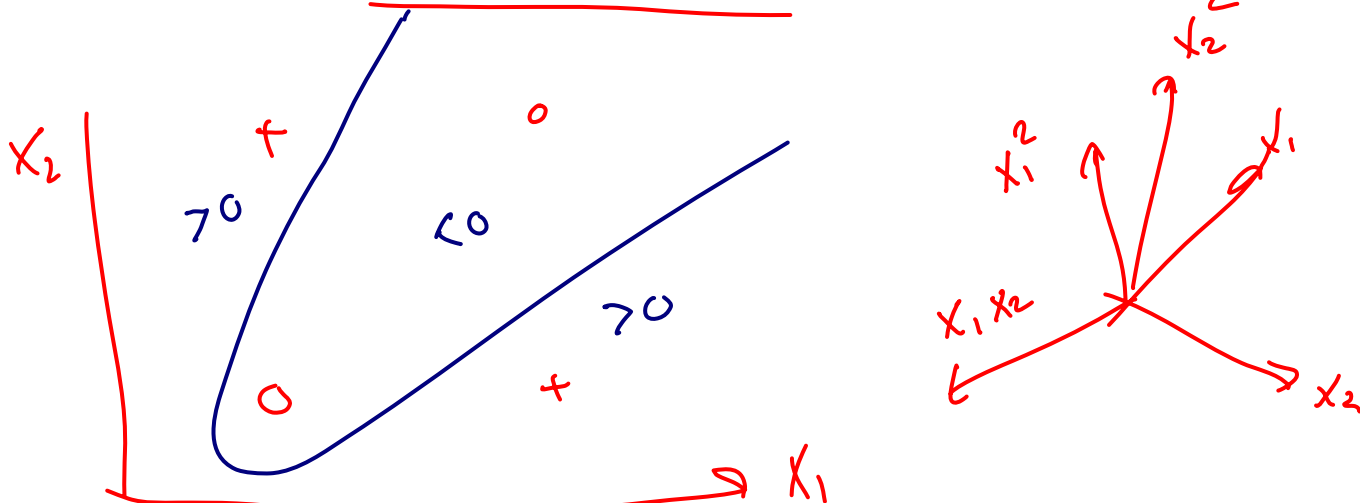


A xor B

no hyperplane!

# Addressing non-linearly separable data – Option 1, non-linear features

- Choose non-linear features, e.g.,
  - Typical linear features: $w_0 + \sum_i w_i x_i$
  - Example of non-linear features:
    - Degree 2 polynomials, $w_0 + \sum_i w_i x_i + \sum_{ij} w_{ij} x_i x_j$
- Classifier $h_{\mathbf{w}}(\mathbf{x})$ still linear in parameters $\mathbf{w}$
  - Usually easy to learn (closed-form or convex/concave optimization)
  - Data is linearly separable in higher dimensional spaces
  - More discussion later this semester

# Addressing non-linearly separable data – Option 2, non-linear classifier

- Choose a classifier $h_w(x)$ that is non-linear in parameters $w$, e.g.,
    - Decision trees, neural networks, nearest neighbor,…
- More general than linear classifiers
- But, can often be harder to learn (non-convex/concave optimization required)
- But, but, often very useful
- (BTW. Later this semester, we'll see that these options are not that different)
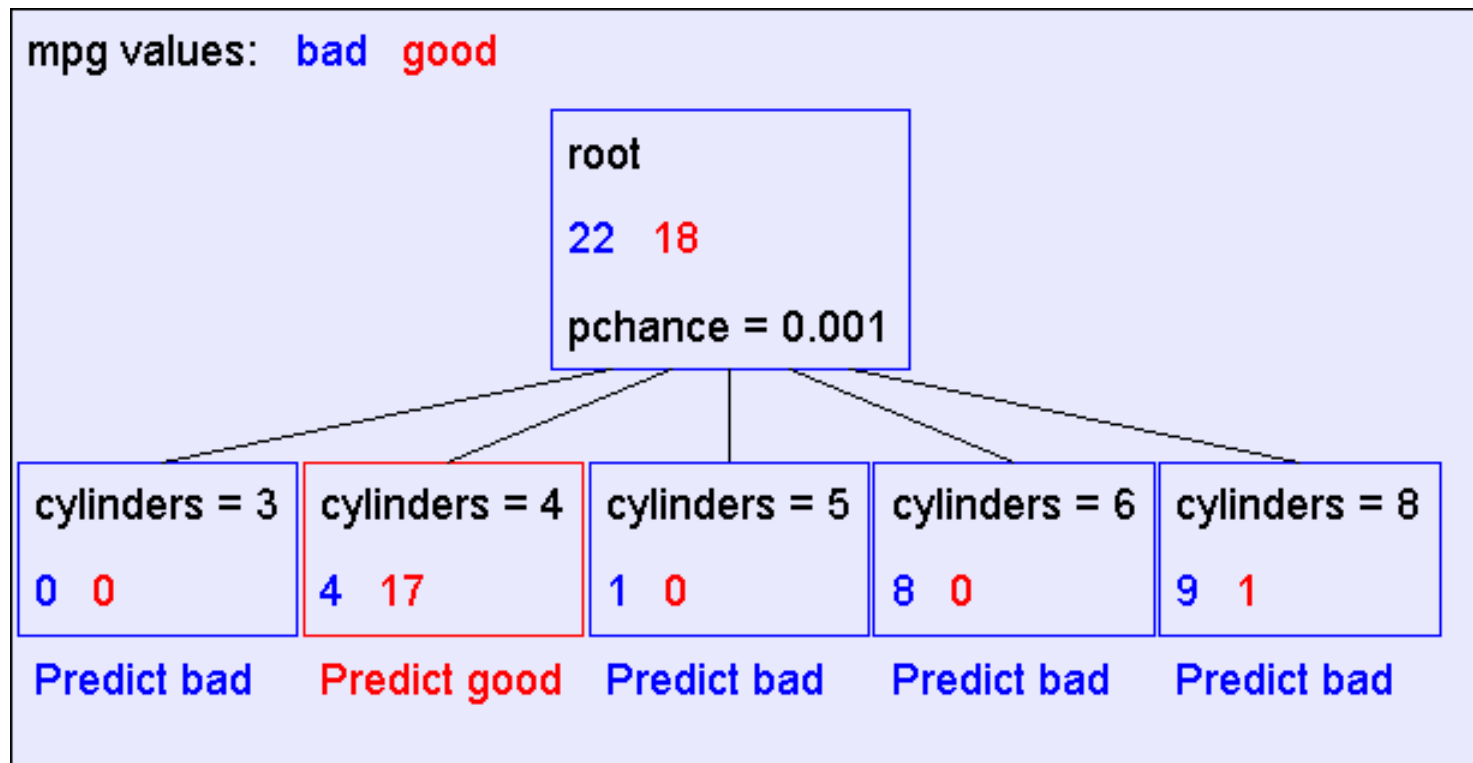
# A small dataset: Miles Per Gallon

Suppose we want to predict MPG

| 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 |

40 Records

From the UCI repository (thanks to Ross Quinlan)

# A Decision Stump

# Recursion Step

mpg values:  bad  good

root

22  18

pchance = 0.001

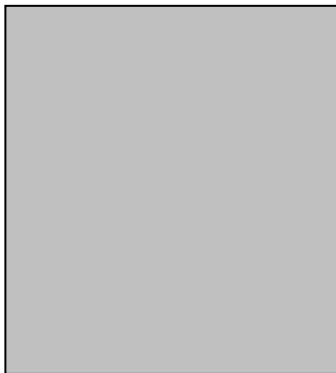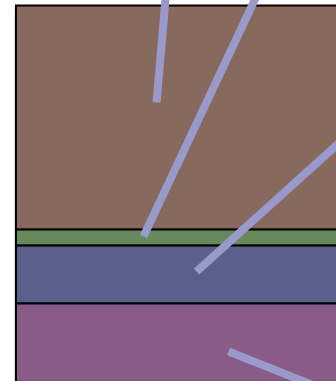| cylinders = 3 | cylinders = 4 | cylinders = 5 | cylinders = 6 | cylinders = 8 |
|---|---|---|---|---|
| 0  0 | 4  17 | 1  0 | 8  0 | 9  1 |
| Predict bad | Predict good | Predict bad | Predict bad | Predict bad |

Take the Original Dataset..

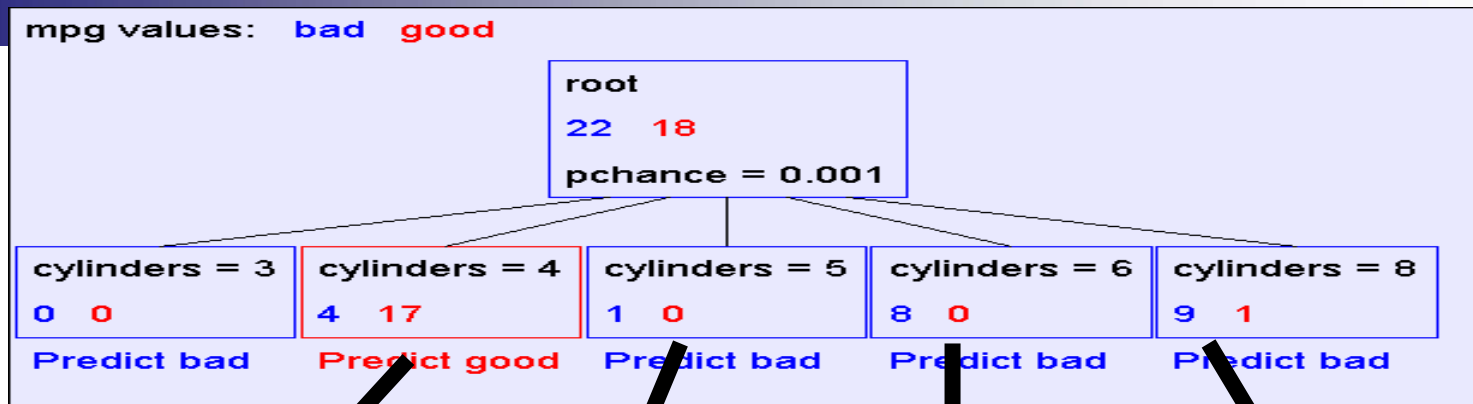And partition it according to the value of the attribute we split on

Records in which cylinders = 4
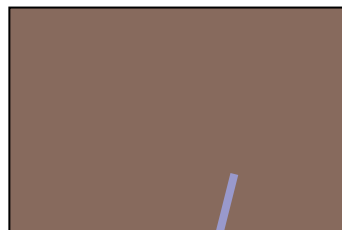
Records in which cylinders = 5

Records in which cylinders = 6

Records in which cylinders = 8

# Recursion Step

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 |

Build tree from These records..

Build tree from These records..

Build tree from These records..

Build tree from These records..

Records in which cylinders = 4

Records in which cylinders = 5

Records in which cylinders = 6

Records in which cylinders = 8

# Second level of tree



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 = america   maker = asia   maker = europe   horsepower = low   horsepower = medium   horsepower = high
0  10             2  5           2  2             0  0               0  1                   9  0
Predict good      Predict good   Predict bad      Predict bad        Predict good           Predict bad
```
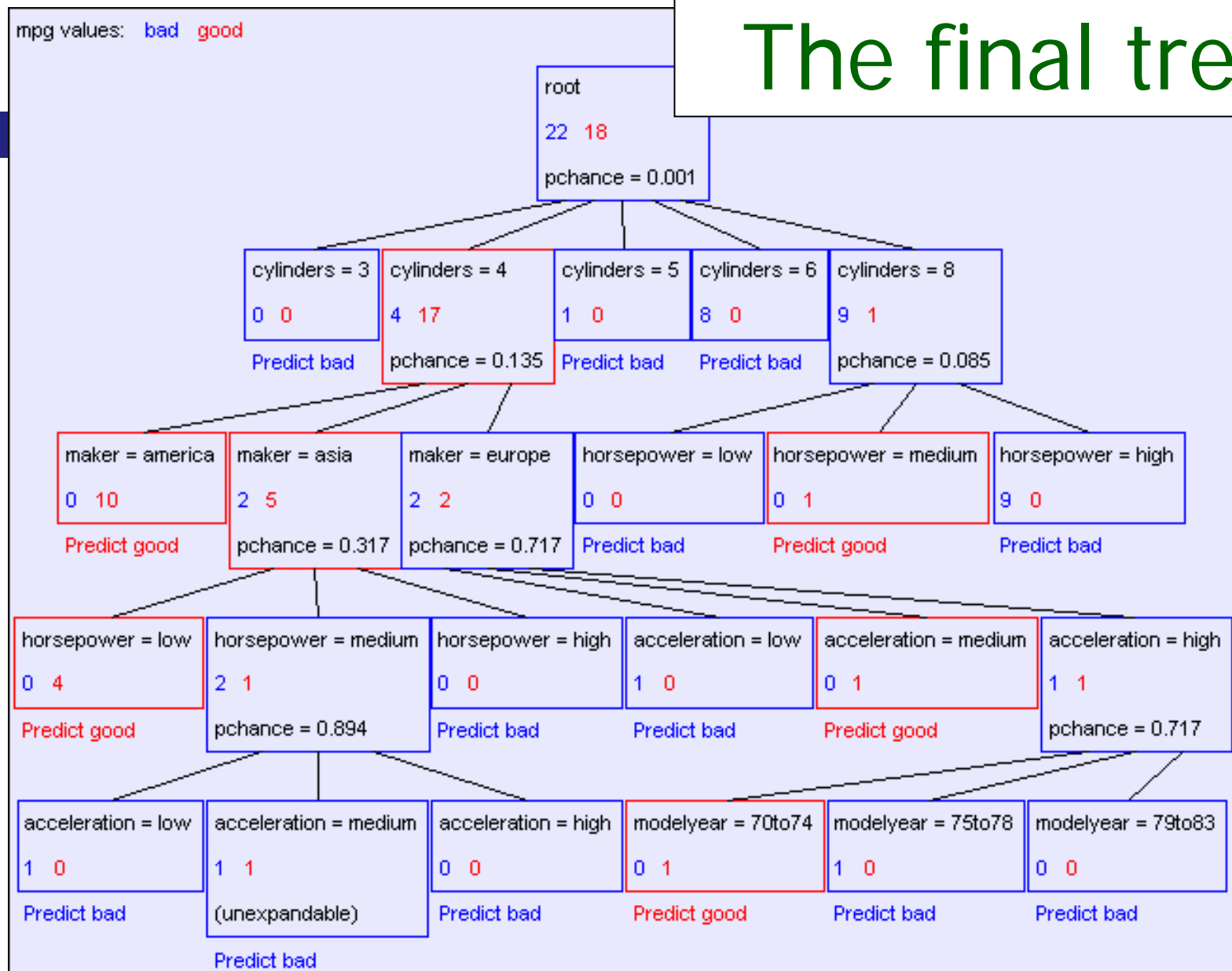
Recursively build a tree from the seven records in which there are four cylinders and the maker was based in Asia

(Similar recursion in the other cases)
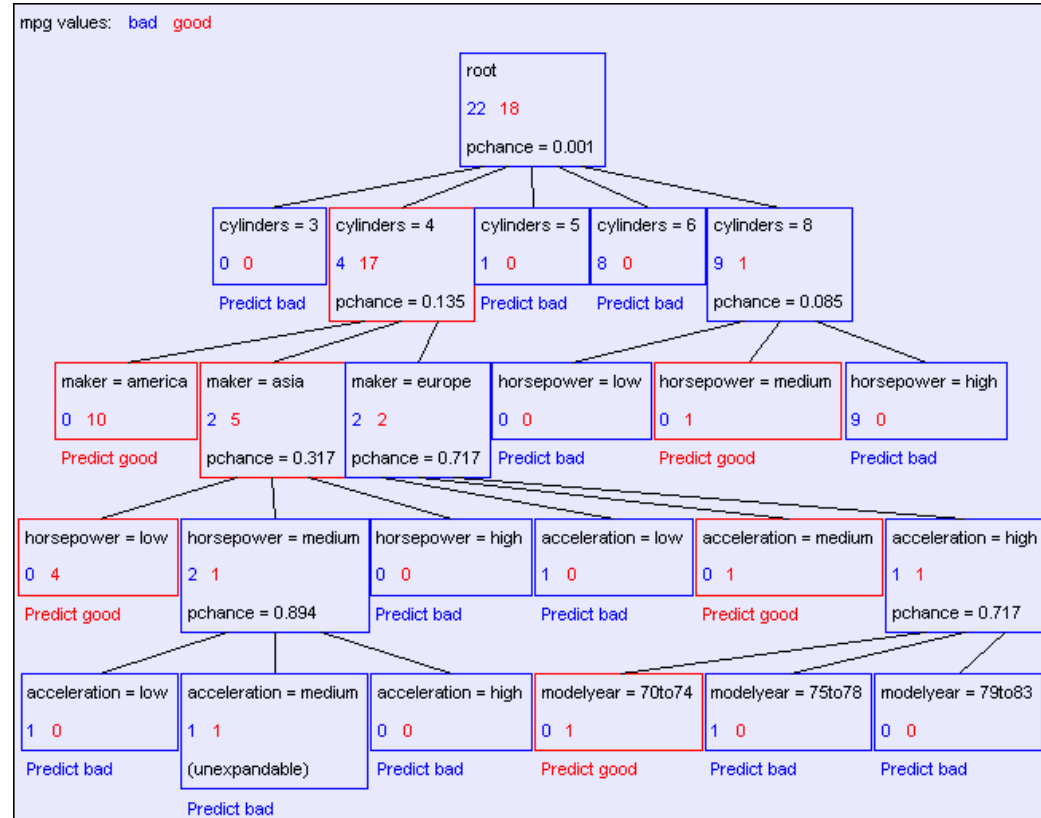
The final tree

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
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

# Classification of a new example

- Classifying a test example – traverse tree and report leaf label

# Are all decision trees equal?

- Many trees can represent the same concept

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

# Learning decision trees 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

# Choosing a good attribute

| X₁ | X₂ | 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 |

# Measuring uncertainty

- Good split if we are more certain about classification after split
  - Deterministic good (all true or all false)
  - Uniform distribution bad

| P(Y=A) = 1/2 | P(Y=B) = 1/4 | P(Y=C) = 1/8 | P(Y=D) = 1/8 |
|---|---|---|---|

| P(Y=A) = 1/4 | P(Y=B) = 1/4 | P(Y=C) = 1/4 | P(Y=D) = 1/4 |
|---|---|---|---|

# Entropy

Entropy $H(X)$ of a random variable $Y$

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

**More uncertainty, more entropy!**

*Information Theory interpretation:* $H(Y)$ is the expected number of bits needed to encode a randomly drawn value of $Y$ (under most efficient code)

# Andrew Moore's Entropy in a nutshell



Low Entropy

High Entropy

# Andrew Moore's Entropy in a nutshell



Low Entropy

..the values (locations of soup) sampled entirely from within the soup bowl

High Entropy

..the values (locations of soup) unpredictable... almost uniformly sampled throughout our dining room

# Information gain

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

- **Advantage of attribute – decrease in uncertainty**
  - □ Entropy of Y before you split

  - □ Entropy after split
    - Weight by probability of following each branch, i.e., normalized number of records

$$H(Y \mid X) = - \sum_{j=1}^{v} P(X = x_j) \sum_{i=1}^{k} P(Y = y_i \mid X = x_j) \log_2 P(Y = y_i \mid X = x_j)$$

- **Information gain is difference** $\quad IG(X) = H(Y) - H(Y \mid X)$

# Learning decision trees

- Start from empty decision tree
- Split on **next best attribute (feature)**
  - Use, for example, information gain to select attribute
  - Split on $\arg\max_i IG(X_i) = \arg\max_i H(Y) - H(Y \mid X_i)$
- Recurse

# Information Gain Example

wealth values: poor rich

gender  Female  14423  1769  ▮▮▮▮▮  H( wealth | gender = Female ) = 0.497654

      Male  22732  9918  ▮▮▮▮▮  H( wealth | gender = Male ) = 0.885847

H(wealth) = 0.793844   H(wealth|gender) = 0.757154

IG(wealth|gender) = 0.0366896

Suppose we want to predict MPG

# Look at all the information gains…



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 | | |
| maker | america | | 0.0437265 |
| | asia | | |

# A Decision Stump

Base Case One

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
nance = 0.717

Don't split a node if all matching records have the same output value

- horsepower = low
  0 0
  Predict bad
- horsepower = medium
  0 1
  Predict good
- horsepower = high
  9 0
  Predict bad

medium

- 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

1 0
Predict bad

1 1
(unexpandable)
Predict bad

medium

- acceleration = high
  0 0
  Predict bad
- modelyear = 70to74
  0 1
  Predict good
- modelyear = 75to78
  1 0
  Predict bad
- modelyear = 79to83
  0 0
  Predict bad

Base Case Two

Don't split a node if none of the attributes can create multiple non-empty children

# Base Case Two: No attributes can distinguish

# Base Cases

- Base Case One: If all records in current data subset have the same output then <span style="color:red">don't recurse</span>

- Base Case Two: If all records have exactly the same set of input attributes then <span style="color:red">don't recurse</span>

# Base Cases: An idea

- 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

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

y = a XOR b

The information gains:

The resulting decision tree:



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 | | |



y values: 0  1

root

2  2

Predict 0

# If we omit Base Case 3:

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

y = a XOR b

The resulting decision tree:

# Basic Decision Tree Building Summarized

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 and return a non-leaf node with $n_X$ children.

  - The $i$'th child should be built by calling

    BuildTree($DS_i$, *Output*)

    Where $DS_i$ built consists of all those records in DataSet for which X = $i$th distinct value of X.

# Real-Valued inputs

- What should we do if some of the inputs are real-valued?

| mpg | cylinders | displacemen | horsepower | weight | acceleration | modelyear | maker |
|---|---|---|---|---|---|---|---|
| | | | | | | | |
| good | 4 | 97 | 75 | 2265 | 18.2 | 77 | asia |
| bad | 6 | 199 | 90 | 2648 | 15 | 70 | america |
| bad | 4 | 121 | 110 | 2600 | 12.8 | 77 | europe |
| bad | 8 | 350 | 175 | 4100 | 13 | 73 | america |
| bad | 6 | 198 | 95 | 3102 | 16.5 | 74 | america |
| bad | 4 | 108 | 94 | 2379 | 16.5 | 73 | asia |
| bad | 4 | 113 | 95 | 2228 | 14 | 71 | asia |
| bad | 8 | 302 | 139 | 3570 | 12.8 | 78 | america |
| : | : | : | : | : | : | : | : |
| : | : | : | : | : | : | : | : |
| : | : | : | : | : | : | : | : |
| good | 4 | 120 | 79 | 2625 | 18.6 | 82 | america |
| bad | 8 | 455 | 225 | 4425 | 10 | 70 | america |
| good | 4 | 107 | 86 | 2464 | 15.5 | 76 | europe |
| bad | 5 | 131 | 103 | 2830 | 15.9 | 78 | europe |
| | | | | | | | |

Infinite number of possible split values!!!

Finite dataset, only finite number of relevant splits!

Idea One: Branch on each possible real value

# "One branch for each numeric value" idea:



mpg values: bad good

root
22 18
pchance = 0.222

| modelyear = 70 | modelyear = 71 | modelyear = 72 | modelyear = 73 | modelyear = 74 | modelyear = 75 | modelyear = 76 | modelyear = 77 | modelyear = 78 | modelyear = 79 | modelyear = 80 | modelyear = 81 | modelyear = 82 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 4 0 | 2 1 | 1 0 | 6 1 | 1 2 | 0 0 | 3 1 | 1 3 | 3 0 | 1 1 | 0 0 | 0 5 | 0 4 |
| Predict bad | Predict bad | Predict bad | Predict bad | Predict good | Predict bad | Predict bad | Predict good | Predict bad | Predict bad | Predict bad | Predict good | Predict good |

Hopeless: with such high branching factor will shatter the dataset and overfit

# Threshold splits

- Binary tree, split on attribute X
    - One branch: $X < t$
    - Other branch: $X \geq t$

# Choosing threshold split

- Binary tree, split on attribute X
  - One branch: X < t
  - Other branch: X $\geq$ t

- Search through possible values of *t*
  - Seems hard!!!

- But only finite number of *t*'s are important
  - Sort data according to X into $\{x_1,\ldots,x_m\}$
  - Consider split points of the form $x_i + (x_{i+1} - x_i)/2$

# A better idea: thresholded splits

- Suppose X is real valued

- Define *IG(Y|X:t)* as *H(Y) - H(Y|X:t)*

- Define *H(Y|X:t)* =
  $$H(Y|X < t) \, P(X < t) + H(Y|X >= t) \, P(X >= t)$$

  - *IG(Y|X:t)* is the information gain for predicting Y if all you know is whether X is greater than or less than *t*

- Then define $IG^*(Y|X) = max_t \, IG(Y|X:t)$

- For each real-valued attribute, use *IG*(Y|X)* for assessing its suitability as a split
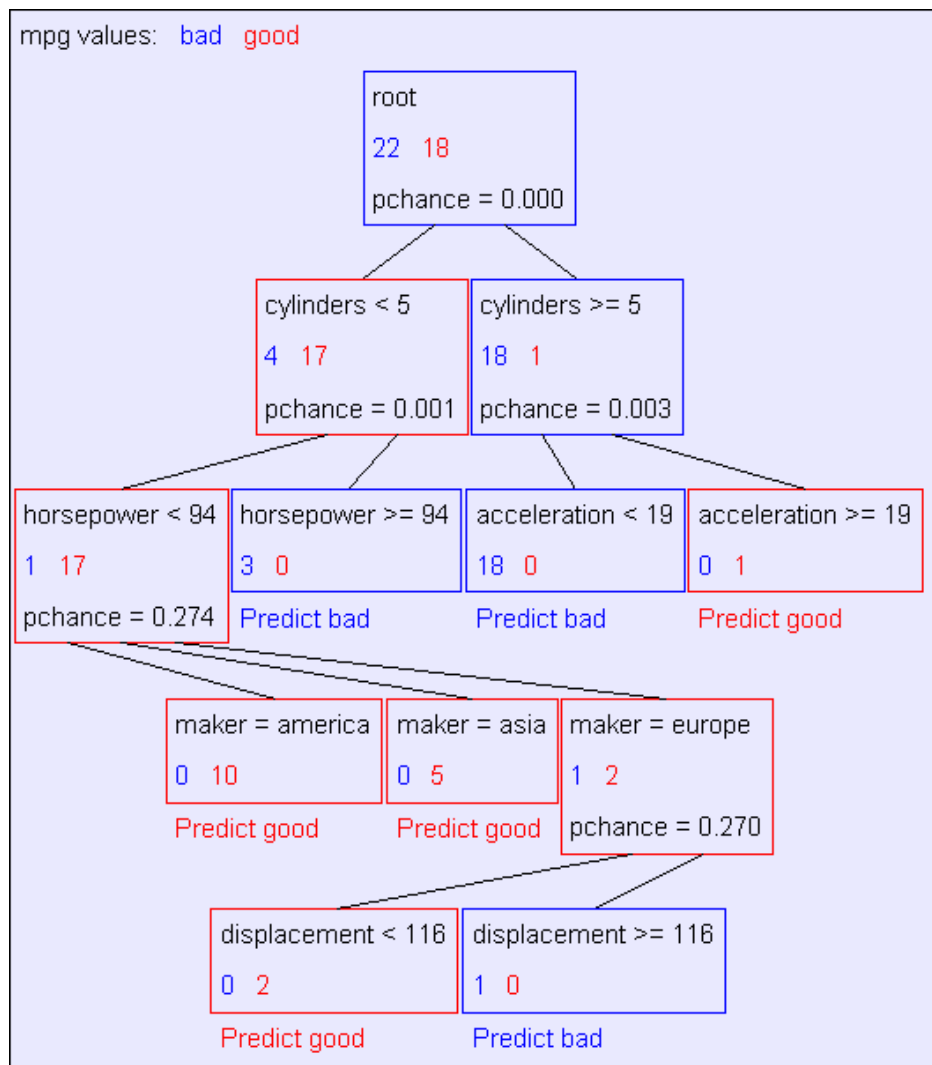
# Example with MPG

# Example tree using reals

# What you need to know about decision trees

- Decision trees are one of the most popular data mining tools
  - Easy to understand
  - Easy to implement
  - Easy to use
  - Computationally cheap (to solve heuristically)
- Information gain to select attributes (ID3, C4.5,…)
- Presented for classification, can be used for regression and density estimation too
- It's possible to get in trouble with overfitting (more next lecture)

# Acknowledgements

- Some of the material in the presentation is courtesy of Tom Mitchell, and of Andrew Moore, from his excellent collection of ML tutorials:

  - http://www.cs.cmu.edu/~awm/tutorials