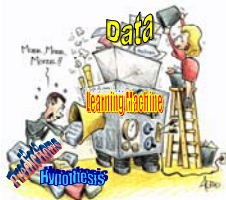# Machine Learning

## 10-701/15-781, Spring 2010

## Introduction to ML

## and

## Decision Trees

**Eric Xing**

**Lecture 1, January 11, 2010**

**Reading:** **Mitchell: Chap 1,3**

© Eric Xing @ CMU, 2006-2010

1

---

# Class Registration

- IF YOU ARE ON THE WAITING LIST: This class is now fully subscribed. You may want to consider the following options:

  - Take the class when it is offered again in the Fall semester;

  - Come to the first several lectures and see how the course develops. We will admit as many students from the waitlist as we can, once we see how many registered students drop the course during the first two weeks.
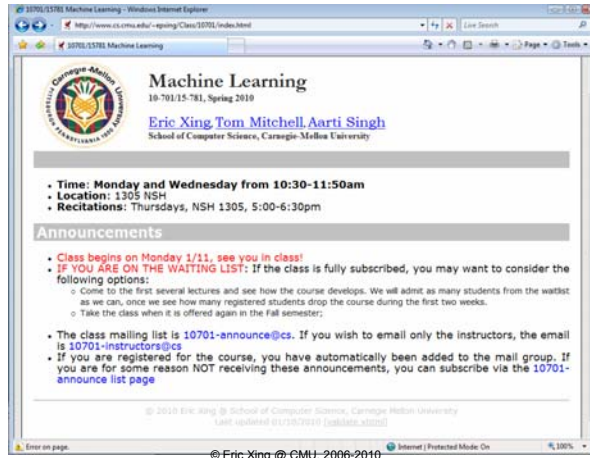
© Eric Xing @ CMU, 2006-2010

2

# Machine Learning 10-701/15-781

- Class webpage:
  - http://www.cs.cmu.edu/~epxing/Class/10701/

---

# Logistics

- Text book
  - Chris Bishop,  **Pattern Recognition and Machine Learning (required)**
  - Tom Mitchell,  **Machine Learning**
  - David Mackay,  **Information Theory, Inference, and Learning Algorithms**

- Mailing Lists:
  - To contact the instructors: 10701-instr@cs.cmu.edu
  - Class announcements list: 10701-announce@cs.cmu.edu.

- TA:
  - Amr Ahmed, GHC 6605, Office hours: TBA
  - Field Cady, Office hours: TBA
  - Ni Lao, NSH 4622, Office hours: TBA

- Class Assistant:
  - Michelle Martin, GHC 8001, x8-5527
  - Sharon Cavlovich, GHC, x8-5196
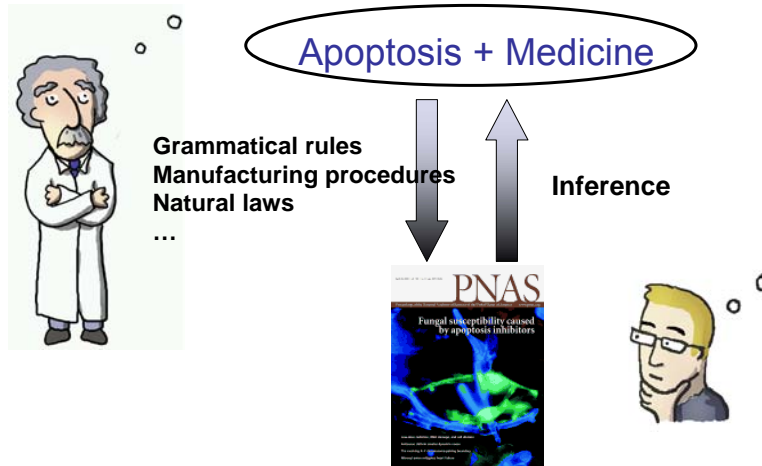
# Logistics

- 5 homework assignments: 30% of grade
    - Theory exercises
    - Implementation exercises

- Final project: 20% of grade
    - Applying machine learning to your research area
        - NLP, IR, Computational biology, vision, robotics …
    - Theoretical and/or algorithmic work
        - a more efficient approximate inference algorithm
        - a new sampling scheme for a non-trivial model …
    - 3-stage reports

- Two exams: 25% of grade each
    - Theory exercises and/or analysis: dates will be set by next week (no "ticket already booked", "I am in a conference", etc. excuse …)

- Policies …

5
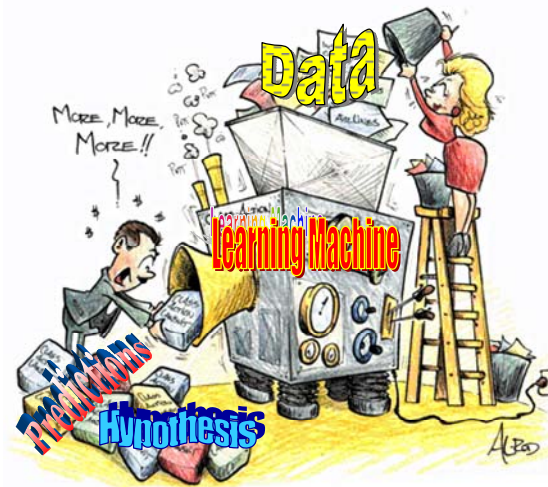
---

# What is Learning

**Learning is about seeking a predictive and/or executable understanding of natural/artificial subjects, phenomena, or activities from …**



Apoptosis + Medicine

**Grammatical rules
Manufacturing procedures
Natural laws
…**

**Inference**

PNAS

6

# Machine Learning

7

---

# Machine Learning

- Study of algorithms that
  - improve their <u>performance</u> P
  - at some <u>task</u> T
  - with <u>experience</u> E

## well-defined learning task: <P,T,E>

8

# Fetching a stapler from inside an office --- the Stanford STAIR robot



real time

9

---

# Machine Learning - Practice



Speech recognition

Information retrieval

Computer vision

Games

Pedigree

Evolution

Robotic control

Planning

10

# Natural language processing and speech recognition

- Now most pocket **Speech Recognizers** or **Translators** are running on some sort of learning device --- the more you play/use them, the smarter they become!
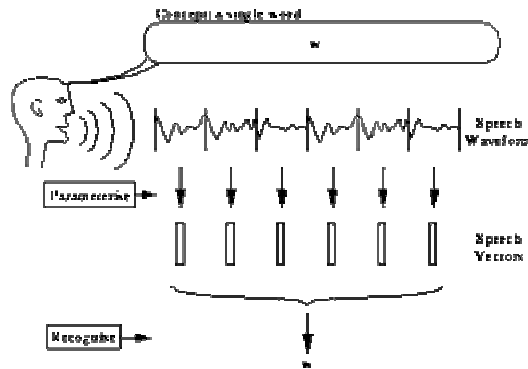


Fig. 1.2 Isolated Word Problem

11

# Object Recognition

- Behind a security camera, most likely there is a computer that is learning and/or checking!

12

# Robotic Control I

- Now cars can find their own ways!

13

# Robotic Control II

- The **best** helicopter pilot is now a computer!
  - it runs a program that learns how to fly and make acrobatic maneuvers by itself!
  - no taped instructions, joysticks, or things like that …

A. Ng 2005

14

# Text Mining

● **We want:**

● Reading, digesting, and categorizing a vast text database is too much for human!

15

---

# Understanding Brain Activities



Reading
a noun
(vs verb)

[Rustandi et al., 2005]

16

# Bioinformatics

Where is the gene?

17

---

# Evolution

ancestor

?

T years

$Q_h$

$Q_m$

A

C

18

# Paradigms of Machine Learning

- Supervised Learning
  - Given $D = \{\mathbf{X}_i, \mathbf{Y}_i\}$, learn $F(\;\cdot\;;\theta)$, s.t.: $\mathbf{Y}_i = F(\mathbf{X}_i)$ $D^{new} = \{\mathbf{X}_j\} \Rightarrow \{\mathbf{Y}_j\}$

- Unsupervised Learning
  - Given $D = \{\mathbf{X}_i\}$, learn $F(\;\cdot\;;\theta)$, s.t.: $\mathbf{Y}_i = F(\mathbf{X}_i)$ $D^{new} = \{\mathbf{X}_j\} \Rightarrow \{\mathbf{Y}_j\}$

- Reinforcement Learning
  - Given $D = \{\text{env, actions, rewards, simulator/trace/real game}\}$

    learn $\begin{array}{l}\text{policy}: e, r \to a \\ \text{utility}: a, e \to r\end{array}$, s.t. $\{\text{env, new real game}\} \Rightarrow a_1, a_2, a_3 \ldots$

- Active Learning
  - Given $D \sim G(\cdot)$, learn $F(\cdot)\,|\,D$ and $D^{new} \sim G'(\cdot)$ s.t. $D^{all} \Rightarrow G'(\cdot), \text{policy}, \{\mathbf{Y}_j\}$

19

---

# Machine Learning - Theory

For the learned $F(;\,\theta)$

- Consistency (value, pattern, …)
- Bias versus variance
- Sample complexity
- Learning rate
- Convergence
- Error bound
- Confidence
- Stability
- …

**PAC Learning Theory**

**(supervised concept learning)**

**# examples (m)**

**representational complexity (H)**

**error rate (ε)**

**failure probability (δ)**

$$m \geq \frac{1}{\epsilon}(\ln |H| + \ln(1/\delta))$$

20

# Machine Learning

Machine Learning seeks to develop theories and computer systems for

- **representing;**
- **classifying, clustering and recognizing;**
- **reasoning under uncertainty;**
- **predicting;**
- **and reacting to**
- **…**

**complex, real world information, based on the system's own experience with data, and (hopefully) under a explicit model or mathematical framework, that**

- **can be formally characterized and analyzed**
- **can take into account human prior knowledge**
- **can generalize and adapt across data and domains**
- **can operate automatically and autonomously**
- **and can be interpreted and perceived by human.**

21

# Growth of Machine Learning

- Machine learning already the preferred approach to
  - Speech recognition, Natural language processing
  - Computer vision
  - Medical outcomes analysis
  - Robot control
  - …

- This ML niche is growing (why?)

ML apps.

All software apps.

22

# Growth of Machine Learning

- Machine learning already the preferred approach to
  - Speech recognition, Natural language processing
  - Computer vision
  - Medical outcomes analysis
  - Robot control
  - …

  **ML apps.**

  **All software apps.**

- This ML niche is growing
  - Improved machine learning algorithms
  - Increased data capture, networking
  - Software too complex to write by hand
  - New sensors / IO devices
  - Demand for self-customization to user, environment

23

---

**Inference**
**Prediction**
**Decision-Making under uncertainty**
**…**

→ **Statistical Machine Learning**
→ **Function Approximation:** $F(\ |\theta)?$

24

# Classification

- sickle-cell anemia



$$f(x_n^1, x_n^2, \cdots, x_n^k) \quad f \quad f$$

---

# Function Approximation

- **Setting**:
  - Set of possible instances $X$
  - Unknown target function $f: X \rightarrow Y$
  - Set of function hypotheses $H = \{ h \mid h: X \rightarrow Y \}$

- **Given**:
  - Training examples $\{<x_i, y_i>\}$ of unknown target function $f$

- **Determine**:
  - Hypothesis $h \in H$ that best approximates $f$

# A Tax-Fraud detection problem:

- What F to use?
  - Hypothesis

- How to use?

**Query Data**

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

---

# Apply Model to Query Data

Start from the root of tree.

**Query Data**

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

**Refund**

Yes → **NO**

No → **MarSt**

Single, Divorced → **TaxInc**

Married → **NO**

< 80K → **NO**

> 80K → **YES**

# Apply Model to Test Data

**Query Data**

| Refund | Marital Status | Taxable Income | Cheat |
|--------|---------------|----------------|-------|
| No | Married | 80K | ? |

```
        Refund
     Yes /    \ No
       NO      MarSt
            Single, Divorced /   \ Married
                  TaxInc            NO
             < 80K /   \ > 80K
                NO      YES
```

© Eric Xing @ CMU, 2006-2010

29

---

# Apply Model to Test Data

**Query Data**

| Refund | Marital Status | Taxable Income | Cheat |
|--------|---------------|----------------|-------|
| No | Married | 80K | ? |

```
        Refund
     Yes /    \ No
       NO      MarSt
            Single, Divorced /   \ Married
                  TaxInc            NO
             < 80K /   \ > 80K
                NO      YES
```

© Eric Xing @ CMU, 2006-2010

30

# Apply Model to Test Data

### Query Data

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

```
        Refund
       /      \
    Yes        No
    /            \
   NO           MarSt
              /       \
   Single, Divorced   Married
        /                \
     TaxInc              NO
    /      \
  < 80K   > 80K
   /        \
  NO        YES
```

31

---

# Apply Model to Test Data

### Query Data

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

```
        Refund
       /      \
    Yes        No
    /            \
   NO           MarSt
              /       \
   Single, Divorced   Married
        /                \
     TaxInc              NO
    /      \
  < 80K   > 80K
   /        \
  NO        YES
```

32

# Apply Model to Test Data

**Query Data**

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

**Refund**

Yes — **NO**

No — **MarSt**

Single, Divorced — **TaxInc**

Married — **NO**

< 80K — **NO**

> 80K — **YES**

Assign Cheat to "No"

---

# A hypothesis for *TaxFraud*

- Input: a vector of attributes
  - X=[Refund,MarSt,TaxInc]
- Output:
  - Y= Cheating or Not
- *H* as a procedure:

**Refund**

Yes — **NO**

No — **MarSt**

Single, Divorced — **TaxInc**

Married — **NO**

< 80K — **NO**

> 80K — **YES**

- Each internal node: test one attribute $X_i$
- Each branch from a node: selects one value for $X_i$
- Each leaf node: predict Y

## A Tree to Predict C-Section Risk

- Learned from medical records of 1000 wonman

  Negative examples are C-sections

```
[833+,167-]  .83+ .17-
Fetal_Presentation = 1: [822+,116-]  .88+ .12-
| Previous_Csection = 0: [767+,81-]  .90+ .10-
| | Primiparous = 0: [399+,13-]  .97+ .03-
| | Primiparous = 1: [368+,68-]  .84+ .16-
| | | Fetal_Distress = 0: [334+,47-]  .88+ .12-
| | | | Birth_Weight < 3349: [201+,10.6-]  .95+
| | | | Birth_Weight >= 3349: [133+,36.4-]  .78+
| | | Fetal_Distress = 1: [34+,21-]  .62+ .38-
| Previous_Csection = 1: [55+,35-]  .61+ .39-
Fetal_Presentation = 2: [3+,29-]  .11+ .89-
Fetal_Presentation = 3: [8+,22-]  .27+ .73-
```

---

## Expressiveness

- Decision trees can express any function of the input attributes.
- E.g., for Boolean functions, truth table row → path to leaf:



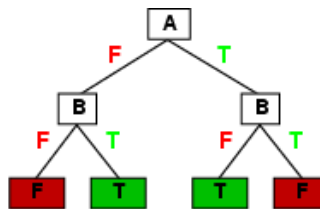- Trivially, there is a consistent decision tree for any training set with one path to leaf for each example (unless *f* nondeterministic in *x*) but it probably won't generalize to new examples

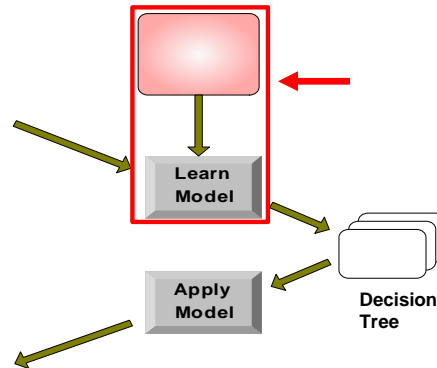- Prefer to find more compact decision trees

## Decision Tree Learning

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 1 | Yes | Large | 125K | No |
| 2 | No | Medium | 100K | No |
| 3 | No | Small | 70K | No |
| 4 | Yes | Medium | 120K | No |
| 5 | No | Large | 95K | Yes |
| 6 | No | Medium | 60K | No |
| 7 | Yes | Large | 220K | No |
| 8 | No | Small | 85K | Yes |
| 9 | No | Medium | 75K | No |
| 10 | No | Small | 90K | Yes |

**Learn Model**

**Apply Model**

**Decision Tree**

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 11 | No | Small | 55K | ? |
| 12 | Yes | Medium | 80K | ? |
| 13 | Yes | Large | 110K | ? |
| 14 | No | Small | 95K | ? |
| 15 | No | Large | 67K | ? |

---

## Example of a Decision Tree

*categorical*   *categorical*   *continuous*   *class*

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

**Training Data**

*Splitting Attributes*

**Refund**
Yes — **NO**
No — **MarSt**

Single, Divorced — **TaxInc**
Married — **NO**

< 80K — **NO**
> 80K — **YES**

**Model: Decision Tree**

# Another Example of Decision Tree

categorical  categorical  continuous  class

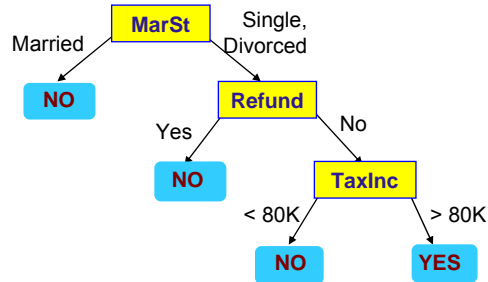| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

**Training Data**

**MarSt**

Married → **NO**

Single, Divorced → **Refund**

Refund: Yes → **NO**

Refund: No → **TaxInc**

TaxInc: < 80K → **NO**

TaxInc: > 80K → **YES**

**There could be more than one tree that fits the same data!**

39

---

# Top-Down Induction of DT

Main loop:

1. $A \leftarrow$ the "best" decision attribute for next *node*

2. Assign $A$ as decision attribute for *node*

3. For each value of $A$, create new descendant of *node*

4. Sort training examples to leaf nodes

5. If training examples perfectly classified, Then STOP, Else iterate over new leaf nodes

Which attribute is best?

[29+, 35−]  A1=?
  t / \ f
[21+, 5−]  [8+, 30−]

[29+, 35−]  A2=?
  t / \ f
[18+, 33−]  [11+, 2−]

40

# Tree Induction

- Greedy strategy.
  - Split the records based on an attribute test that optimizes certain criterion.

- Issues
  - Determine how to split the records
    - How to specify the attribute test condition?
    - How to determine the best split?
  - Determine when to stop splitting

41

---

# Tree Induction

- Greedy strategy.
  - Split the records based on an attribute test that optimizes certain criterion.

- Issues
  - Determine how to split the records
    - How to specify the attribute test condition?
    - How to determine the best split?
  - Determine when to stop splitting

42

# How to Specify Test Condition?

- Depends on attribute types
  - Nominal
  - Ordinal
  - Continuous

- Depends on number of ways to split
  - 2-way split
  - Multi-way split

# Splitting Based on Nominal Attributes

- Multi-way split: Use as many partitions as distinct values.

CarType
Family    Luxury
Sports

- Binary split: Divides values into two subsets.
                Need to find optimal partitioning.

{Sports, Luxury}    CarType    {Family}        OR        {Family, Luxury}    CarType    {Sports}

# Splitting Based on Ordinal Attributes

- Multi-way split: Use as many partitions as distinct values.

Size
Small — Medium — Large

- Binary split:  Divides values into two subsets.
                Need to find optimal partitioning.

{Small, Medium}  Size  {Large}        OR        {Medium, Large}  Size  {Small}

- What about this split?    {Small, Large}  Size  {Medium}

© Eric Xing @ CMU, 2006-2010                                                    45

---

# Splitting Based on Continuous Attributes

- Different ways of handling
  - Discretization to form an ordinal categorical attribute
    - Static – discretize once at the beginning
    - Dynamic – ranges can be found by equal interval bucketing, equal frequency bucketing (percentiles), or clustering.

  - Binary Decision: $(A < v)$ or $(A \geq v)$
    - consider all possible splits and finds the best cut
    - can be more compute intensive

© Eric Xing @ CMU, 2006-2010                                                    46

## Splitting Based on Continuous Attributes

Taxable
Income
> 80K?

## Tree Induction

- Greedy strategy.
  - Split the records based on an attribute test that optimizes certain criterion.

- Issues
  - Determine how to split the records
    - How to specify the attribute test condition?
    - How to determine the best split?
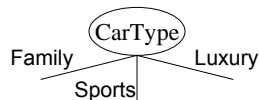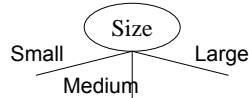  - Determine when to stop splitting

# How to determine the Best Split

- Idea: a good attribute splits the examples into subsets that are (ideally) "all positive" or "all negative"



**Homogeneous,**
**Low degree of impurity**

**Non-homogeneous,**
**High degree of impurity**

- Greedy approach:
  - Nodes with homogeneous class distribution are preferred
- Need a measure of node impurity:

---

# How to compare attribute?

- Entropy
  - Entropy *H(X)* of a random variable *X*

$$H(X) = -\sum_{i=1}^{N} P(x = i) \log_2 P(x = i)$$

  - *H(X)* is the expected number of bits needed to encode a randomly drawn value of *X* (under most efficient code)
  - Why?

  Information theory:
  Most efficient code assigns $-\log_2 P(X=i)$ bits to encode the message *X=I*,
  So, expected number of bits to code one random *X* is:

$$-\sum_{i=1}^{N} P(x = i) \log_2 P(x = i)$$

# How to compare attribute?

- Conditional Entropy
  - Specific conditional entropy $H(X|Y=v)$ of $X$ given $Y=v$ :

  $$H(X|y=j) = -\sum_{i=1}^{N} P(x=i|y=j) \log_2 P(x=i|y=j)$$

  - Conditional entropy $H(X|Y)$ of $X$ given $Y$ :

  $$H(X|Y) = -\sum_{j \in Val(y)} P(y=j) \log_2 H(X|y=j)$$

  - Mututal information (aka information gain) of $X$ and $Y$ :

  $$\begin{aligned} I(X;Y) &= H(X) - H(X|Y) = H(Y) - H(Y|X) \\ &= H(X) + H(Y) - H(X,Y) \end{aligned}$$

51

# Sample Entropy



- $S$ is a sample of training examples
- $p_+$ is the proportion of positive examples in $S$
- $p_-$ is the proportion of negative examples in $S$
- Entropy measure the impurity of $S$

$$H(S) \equiv -p_+ \log_2 p_+ - p_- \log_2 p_-$$

52

# Examples for computing Entropy

$$H(X) = -\sum_{i=1}^{N} P(x=i) \log_2 P(x=i)$$

| C1 | 0 |
|----|---|
| C2 | 6 |

P(C1) = 0/6 = 0    P(C2) = 6/6 = 1

Entropy = − 0 log 0 − 1 log 1 = − 0 − 0 = 0

| C1 | 1 |
|----|---|
| C2 | 5 |

P(C1) = 1/6    P(C2) = 5/6

Entropy = − (1/6) log$_2$ (1/6) − (5/6) log$_2$ (1/6) = 0.65

| C1 | 2 |
|----|---|
| C2 | 4 |

P(C1) = 2/6    P(C2) = 4/6

Entropy = − (2/6) log$_2$ (2/6) − (4/6) log$_2$ (4/6) = 0.92

53

---

# Information Gain

- Information Gain:

$$GAIN_{split} = Entropy(p) - \left( \sum_{i=1}^{k} \frac{n_i}{n} Entropy(i) \right)$$

Parent Node, p is split into k partitions; $n_i$ is number of records in partition i



[29+, 35−]  A1=?    [29+, 35−]  A2=?
t   f    t   f
[21+, 5−]  [8+, 30−]    [18+, 33−]  [11+, 2−]

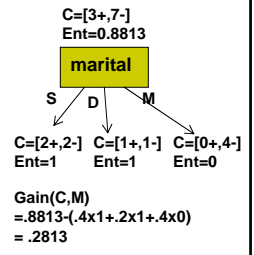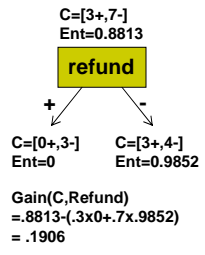Gain(S,A) = mutual information between A and target class variable over sample S

- Measures Reduction in Entropy achieved because of the split. Choose the split that achieves most reduction (maximizes GAIN)
- Used in ID3 and C4.5
- Disadvantage: Tends to prefer splits that result in large #of partitions, each being small but pure.

54

# Exercise

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|---------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

categorical  categorical  continuous  class

**Training Data**

C=[3+,7-]
Ent=0.8813

**refund**

+        -

C=[0+,3-]        C=[3+,4-]
Ent=0        Ent=0.9852

Gain(C,Refund)
=.8813-(.3x0+.7x.9852)
= .1906

C=[3+,7-]
Ent=0.8813

**marital**

S    D    M

C=[2+,2-]  C=[1+,1-]  C=[0+,4-]
Ent=1        Ent=1        Ent=0

Gain(C,M)
=.8813-(.4x1+.2x1+.4x0)
= .2813

**Which one should be at the root?**

▪**Choose the best classifier!**
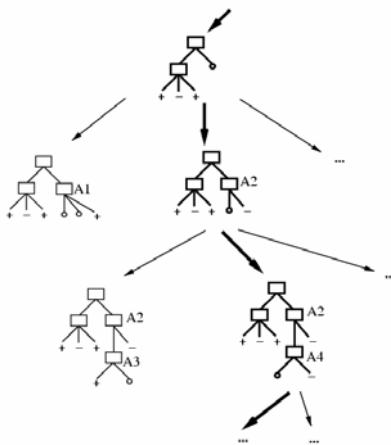
55

---

# Stopping Criteria for Tree Induction

- Stop expanding a node when all the records belong to the same class

- Stop expanding a node when all the records have similar attribute values

- Early termination (to be discussed later)

56

# Decision Tree Based Classification

- Advantages:
  - Inexpensive to construct
  - Extremely fast at classifying unknown records
  - Easy to interpret for small-sized trees
  - Accuracy is comparable to other classification techniques for many simple data sets

- Example: C4.5
  - Simple depth-first construction.
  - Uses Information Gain
  - Sorts Continuous Attributes at each node.
  - Needs entire data to fit in memory.
  - Unsuitable for Large Datasets.
    - Needs out-of-core sorting.
  - You can download the software from:
    http://www.cse.unsw.edu.au/~quinlan/c4.5r8.tar.gz

57

---

# Which Tree Should We Output?



- ID3 performs heuristic search through space of decision trees

- It stops at smallest acceptable tree. Why?

> Occam's razor: prefer the simplest hypothesis that fits the data

58

# Practical Issues of DT

- Underfitting and Overfitting

- Missing Values

**Will be covered in recitation!**

59

---

# Summary: what you should know:

- Well posed function approximation problems:
  - Instance space, X
  - Sample of labeled training data { <$x_i$, $y_i$>}
  - Hypothesis space, H = { f: X→Y }

- Learning is a search/optimization problem over H
  - Various objective functions
    - minimize training error (0-1 loss)
    - among hypotheses that minimize training error, select smallest (?)

- Decision tree learning
  - Greedy top-down learning of decision trees (ID3, C4.5, ...)
  - Overfitting and tree/rule post-pruning
  - Extensions…

60

# Questions to think about (1)

- ID3 and C4.5 are heuristic algorithms that search through the space of decision trees. Why not just do an exhaustive search?

# Questions to think about (2)

- Consider target function f: <x1,x2> → y, where x1 and x2 are real-valued, y is boolean. What is the set of decision surfaces describable with decision trees that use each attribute at most once?

# Questions to think about (3)

- Why use Information Gain to select attributes in decision trees? What other criteria seem reasonable, and what are the tradeoffs in making this choice?

# Additional material:

# Hypothesis spaces

How many distinct decision trees with *n* Boolean attributes?

= number of Boolean functions

= number of distinct truth tables with $2^n$ rows = $2^{2^n}$

- E.g., with 6 Boolean attributes, there are 18,446,744,073,709,551,616 trees

# Notes on Overfitting

- Overfitting results in decision trees that are more complex than necessary

- Training error no longer provides a good estimate of how well the tree will perform on previously unseen records

- **Which Tree Should We Output?**
  - Occam's razor: prefer the simplest hypothesis that fits the data

# Occam's Razor

- Given two models of similar generalization errors, one should prefer the simpler model over the more complex model

- For complex models, there is a greater chance that it was fitted accidentally by errors in data

- Therefore, one should include model complexity when evaluating a model

# Minimum Description Length (MDL)



- Cost(Model,Data) = Cost(Data|Model) + Cost(Model)
  - Cost is the number of bits needed for encoding.
  - Search for the least costly model.
- Cost(Data|Model) encodes the misclassification errors.
- Cost(Model) uses node encoding (number of children) plus splitting condition encoding.

# How to Address Overfitting

- Pre-Pruning (Early Stopping Rule)
  - Stop the algorithm before it becomes a fully-grown tree
  - Typical stopping conditions for a node:
    - Stop if all instances belong to the same class
    - Stop if all the attribute values are the same
  - More restrictive conditions:
    - Stop if number of instances is less than some user-specified threshold
    - Stop if class distribution of instances are independent of the available features (e.g., using $\chi^2$ test)
    - Stop if expanding the current node does not improve impurity measures (e.g., Gini or information gain).

# How to Address Overfitting…

- Post-pruning
  - Grow decision tree to its entirety
  - Trim the nodes of the decision tree in a bottom-up fashion
  - If generalization error improves after trimming, replace sub-tree by a leaf node.
  - Class label of leaf node is determined from majority class of instances in the sub-tree
  - Can use MDL for post-pruning

# Handling Missing Attribute Values

- Missing values affect decision tree construction in three different ways:
  - Affects how impurity measures are computed
  - Affects how to distribute instance with missing value to child nodes
  - Affects how a test instance with missing value is classified

# Computing Impurity Measure

| Tid | Refund | Marital Status | Taxable Income | Class |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | ? | Single | 90K | Yes |

**Missing value**

**Before Splitting:**
Entropy(Parent)
= -0.3 log(0.3)-(0.7)log(0.7) = 0.8813

| | Class = Yes | Class = No |
|-----------|------|------|
| Refund=Yes | 0 | 3 |
| Refund=No | 2 | 4 |
| Refund=? | 1 | 0 |

**Split on Refund:**

Entropy(Refund=Yes) = 0

Entropy(Refund=No)
= -(2/6)log(2/6) – (4/6)log(4/6) = 0.9183

Entropy(Children)
= 0.3 (0) + 0.6 (0.9183) = 0.551

Gain = 0.9 × (0.8813 – 0.551) = 0.3303

# Distribute Instances

| Tid | Refund | Marital Status | Taxable Income | Class |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |

| Tid | Refund | Marital Status | Taxable Income | Class |
|-----|--------|----------------|----------------|-------|
| 10 | ? | Single | 90K | Yes |

**Refund**

Yes      No

| Class=Yes | 0 + 3/9 |
|-----------|---------|
| Class=No | 3 |

| Class=Yes | 2 + 6/9 |
|-----------|---------|
| Class=No | 4 |

**Probability that Refund=Yes is 3/9**

**Probability that Refund=No is 6/9**

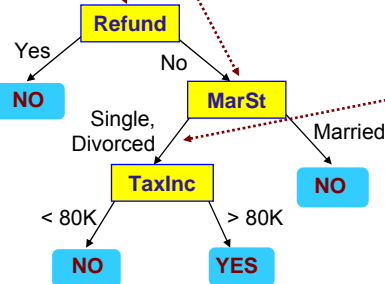**Assign record to the left child with weight = 3/9 and to the right child with weight = 6/9**

**Refund**

Yes      No

| Class=Yes | 0 |
|-----------|---|
| Class=No | 3 |

| Cheat=Yes | 2 |
|-----------|---|
| Cheat=No | 4 |

© Eric Xing @ CMU, 2006-2010

73

---

# Classify Instances

**New record:**

| Tid | Refund | Marital Status | Taxable Income | Class |
|-----|--------|----------------|----------------|-------|
| 11 | No | ? | 85K | ? |

|  | Married | Single | Divorced | Total |
|--|---------|--------|----------|-------|
| Class=No | 3 | 1 | 0 | 4 |
| Class=Yes | 6/9 | 1 | 1 | 2.67 |
| Total | 3.67 | 2 | 1 | 6.67 |

**Refund**

Yes    No

NO

**MarSt**

Single, Divorced     Married

**TaxInc**    NO

< 80K    > 80K

NO    YES

**Probability that Marital Status = Married is 3.67/6.67**

**Probability that Marital Status ={Single,Divorced} is 3/6.67**

© Eric Xing @ CMU, 2006-2010

74