# Machine Learning

**10-701/15-781, Spring 2008**

## Decision Trees

**Eric Xing**

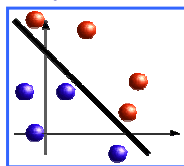**Lecture 6, February 4, 2008**

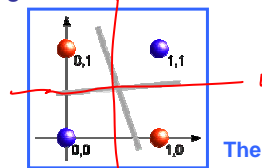**Reading: Chap. 1.6, CB & Chap 3, TM**

---

# Learning non-linear functions

**f: X → Y**

- X (vector of) continuous and/or discrete vars
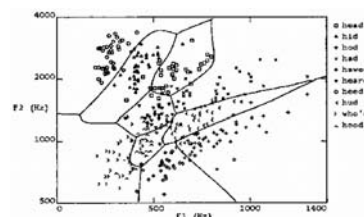- Y discrete vars

- Linear separator

- f might be non-linear function

  **The XOR gate**

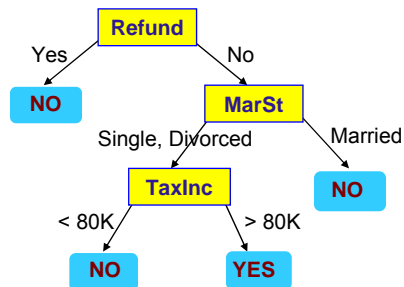  **Speech recognition**

# A hypothesis for *TaxFraud*

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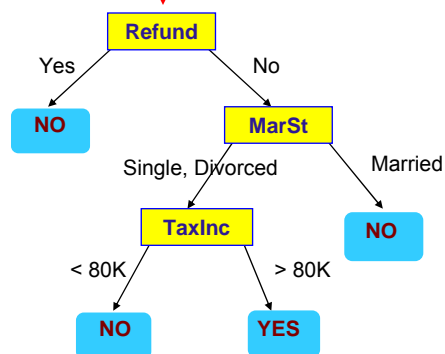  - X=[Refund,MarSt,TaxInc]
- Output:
  - Y= Cheating or Not
- *H* as a procedure:

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

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

---

# 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 /   \ Married
                   TaxInc        NO
            < 80K /    \ > 80K
                NO      YES
```

# 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



# 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

3

# Apply Model to Test Data

**Query Data**

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

- Refund
  - Yes → NO
  - No → MarSt
    - Single, Divorced → TaxInc
      - < 80K → NO
      - > 80K → YES
    - Married → NO

---

# Apply Model to Test Data

**Query Data**

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

- Refund
  - Yes → NO
  - No → MarSt
    - Single, Divorced → TaxInc
      - < 80K → NO
      - > 80K → YES
    - Married → NO

## 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

< 80K / > 80K

NO   YES

NO

Assign Cheat to "No"

---

## A Tree to Predict C-Section Risk

- Learned from medical records of 1000 wonman
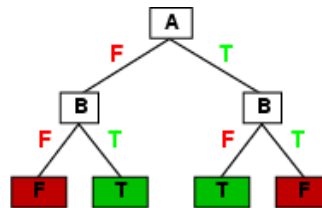
  Negative examples are C-sections

*Fe*

```
[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-
```

5

# Expressiveness

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

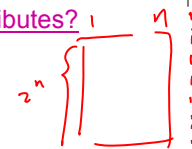| A | B | C | A xor B |
|---|---|---|---|
| F | F | T | F |
| F | T | T | T |
| T | F | T | T |
| T | T | T | F |
| ? | ? | F | ? |

- 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

# 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

# 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

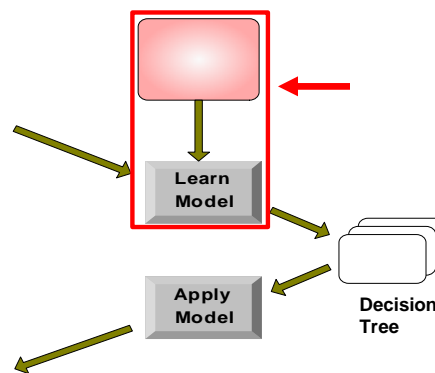How many purely conjunctive hypotheses (e.g., *Hungry* $\land \lnot$*Rain*)?

- Each attribute can be in (positive), in (negative), or out
  $\Rightarrow 3^n$ distinct conjunctive hypotheses
- More expressive hypothesis space
  - increases chance that target function can be expressed
  - increases number of hypotheses consistent with training set
    $\Rightarrow$ may get worse predictions

# 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 |

| 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 | ? |

Learn Model

Apply Model

Decision Tree

# 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        Married

TaxInc          NO

< 80K          > 80K

NO          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          Single, Divorced

Married

NO          Refund

Yes          No

NO          TaxInc
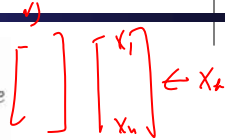
< 80K          > 80K

NO          YES

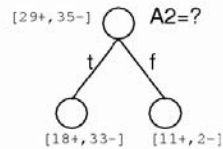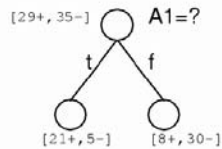**There could be more than one tree that fits the same data!**

## 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−]

---

## 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

# 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
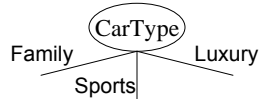
# How to Specify Test Condition?

- Depends on attribute types
  - Nominal
  - Ordinal
  - Continuous

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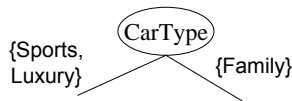  - 2-way split
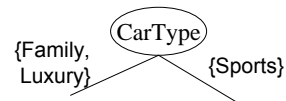  - Multi-way split

# Splitting Based on Nominal Attributes

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

CarType
Family — Sports — Luxury

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

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

---
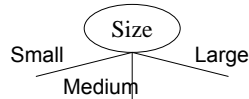
# Splitting Based on Ordinal Attributes

- 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}

# 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

---

# Splitting Based on Continuous Attributes

# 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* :

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

## 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*

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

## 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

# 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

- 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 number of partitions, each being small but pure.

```
[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

---

# Splitting Based on INFO...

- Gain Ratio:

$$GainRATIO_{split} = \frac{GAIN_{Split}}{SplitINFO} \qquad SplitINFO = -\sum_{i=1}^{k} \frac{n_i}{n} \log \frac{n_i}{n}$$

Parent Node, p is split into k partitions

n_i is the number of records in partition i

- Adjusts Information Gain by the entropy of the partitioning (SplitINFO). Higher entropy partitioning (large number of small partitions) is penalized!
- Used in C4.5
- Designed to overcome the disadvantage of Information Gain

# Exercise

*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**

---

# 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

# 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)

# 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

# Practical Issues of Classification

- Underfitting and Overfitting

- Missing Values

- Costs of Classification
  -- Later lectures

---

# Underfitting and Overfitting



$$E_t(h_0^*) \subseteq E_t(h')$$
$$h' \in H.$$

$$\exists \, h'' \in H$$
$$E_{Test}(h'') \subset E_{Test}(h^*)$$

**Underfitting**: when model is too simple, both training and test errors are large

# 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
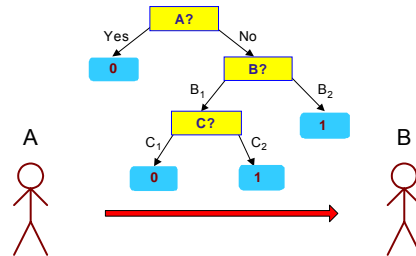
# Occam's Razor

- 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)



| X | y |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0 |
| $X_3$ | 0 |
| $X_4$ | 1 |
| … | … |
| $X_n$ | 1 |

| X | y |
|---|---|
| $X_1$ | ? |
| $X_2$ | ? |
| $X_3$ | ? |
| $X_4$ | ? |
| … | … |
| $X_n$ | ? |

- 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

*[Handwritten annotations:]*

$N-1 \cdot 1$

Leave -one -out - CV

K-fold - CV

Tr | hold out

$\frac{}{T}$  $E(h_1)$

$\rightarrow E_T^1(h_2)$

$E_T(h_3)$

---

# 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

22

# 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 $= 0.9 \times (0.8813 - 0.551) = 0.3303$

---

# 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 |

**Refund**

Yes / No

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

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

| 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**

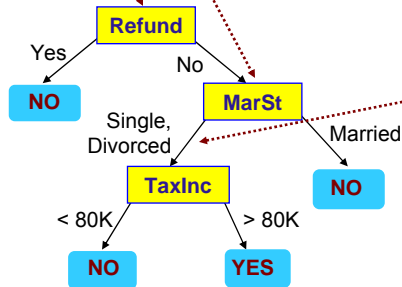**Assign record to the left child with weight = 3/9 and to the right child with weight = 6/9**

# Classify Instances

**New record:**

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

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

|            | 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  |

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

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