

FUNDAMENTALS OF MACHINE LEARNING

DECISION TREES

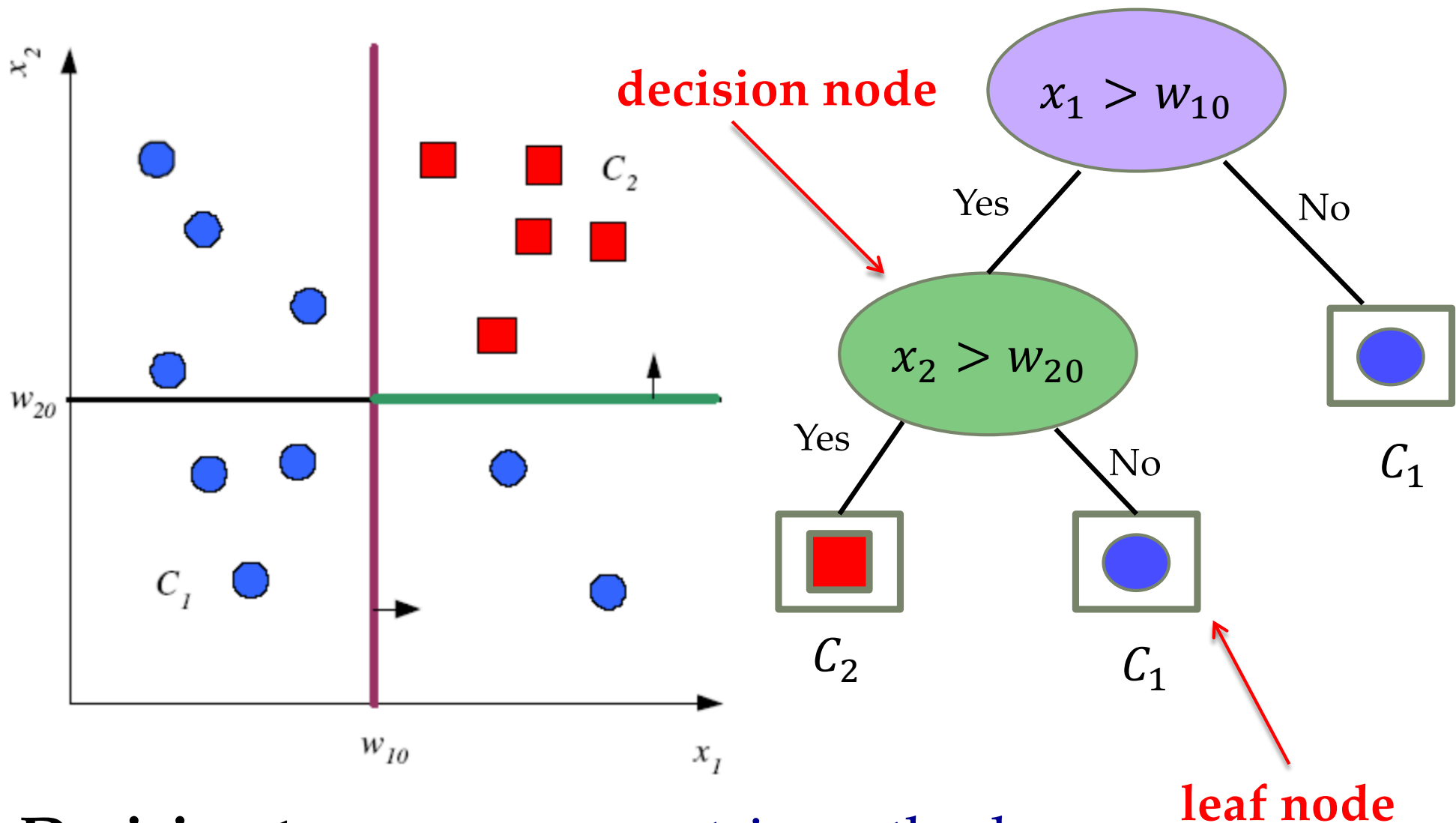
CSCI3320

Prof. John C.S. Lui, CSE Department, CUHK
Introduction to Machine Learning

Decision Tree

- An efficient *non-parametric method* for **supervised learning**, which is often used for classification and regression
- **Decision tree** is a *hierarchical data structure* which relies on **divide-and-conquer strategy**
- We will explore various learning algorithms to build the tree from a given labeled training sample
- Also will learn how to **convert a tree to a set of rules**
- **Parametric method**: define a model over whole input space, learn its parameters from all training data
- **Nonparametric method**: divide input space into regions (clusters) defined by some **distance measure**. For each input, the corresponding local model computed from the training data in that region is used.

Tree Uses Nodes and Leaves



Decision tree: nonparametric method

Divide and Conquer

□ Internal decision nodes

- **Univariate:** Uses a single attribute (or single dimension), x_i

 - **Numeric** x_i : Binary split : $x_i > w_m$

 - **Discrete** x_i : **n -way split** for n possible values (**example**)

- **Multivariate:** Uses all attributes, x

□ Leaves

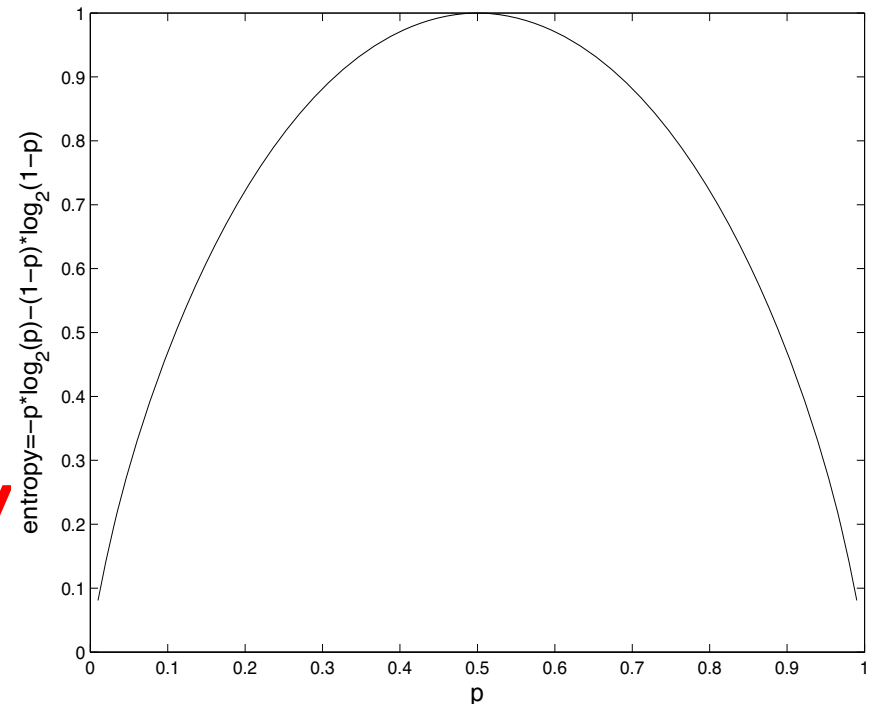
- **Classification:** Class labels, or proportions

- **Regression:** Numeric; r average, or local fit

- Learning is **greedy**: **find the ``best'' split recursively**

Univariate Classification Trees (ID3,CART,C4.5)

- For *classification tree*, the goodness of a split is measured by the **impurity measure**.
- $N = \#$ of instances in the root, $N_m = \#$ of instances at node m , $N_m^i = \#$ of instances of N_m belongs to class C_i
- The probability of class C_i is
$$\hat{P}(C_i | \mathbf{x}, m) \equiv p_m^i = \frac{N_m^i}{N_m}$$
- Node m is **pure** if p_m^i is 0 or 1
- Measure of impurity is **entropy**



$$\mathcal{I}_m = - \sum_{i=1}^K p_m^i \log_2 p_m^i$$

Explain entropy

Univariate Classification Trees (ID3,CART,C4.5)

- There are many possible *measures*.
- For a two class problem, we have $p^1 \equiv p$ and $p^2 = 1 - p$,
- A **nonnegative function** to measure impurity of a split
$$\phi(p, 1 - p)$$
- The function needs to satisfy:
 - $\phi(1/2, 1/2) \geq \phi(p, 1 - p)$, for any $p \in [0, 1]$.
 - $\phi(0, 1) = \phi(1, 0) = 0$.
 - $\phi(p, 1 - p)$ is increasing in p on $[0, 1/2]$ and decreasing in p on $[1/2, 1]$.
- **Possible functions**
 - **Entropy:** $\phi(p, 1 - p) = -p \log_2 p - (1 - p) \log_2 (1 - p)$
 - **Gini Index:** $\phi(p, 1 - p) = 2p(1 - p)$
 - **Misclassification error:** $\phi(p, 1 - p) = 1 - \max(p, 1 - p)$

Best Split for univariate classification tree

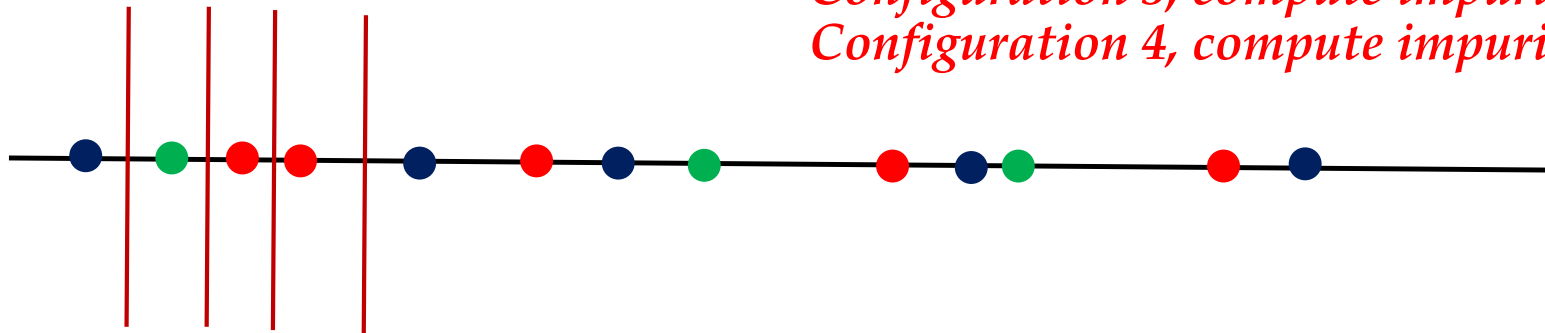
- If node m is pure, generate a leaf and stop, otherwise split and continue recursively
- Define impurity **after split**: N_{mj} of N_m take branch j . N_{mj}^i belong to C_i

$$\hat{P}(C_i | \mathbf{x}, m, j) \equiv p_{mj}^i = \frac{N_{mj}^i}{N_{mj}} \quad \mathcal{I}'_m = - \sum_{j=1}^n \frac{N_{mj}}{N_m} \sum_{i=1}^K p_{mj}^i \log_2 p_{mj}^i \quad \text{explain}$$

- For discrete attributes, there are n outcomes
- Find the variable and split that minimize impurity (among all variables -- and split positions for numeric variables)

How to do the splitting?

- Consider doing a splitting on feature x_i
- We have the following training instances



Configuration 1, compute impurity
Configuration 2, compute impurity
Configuration 3, compute impurity
Configuration 4, compute impurity....

- Compute impurity for different configurations using previous formula
- Choose the configuration which has the minimum
- Note: There are many ways to do splitting: exhaustive, binary search...etc

GenerateTree(\mathcal{X})

If NodeEntropy(\mathcal{X}) < θ_I /* equation 9.3 */ (*or entropy*)

Create leaf labelled by majority class in \mathcal{X}

Return

$i \leftarrow \text{SplitAttribute}(\mathcal{X})$

For each branch of \mathbf{x}_i

Find \mathcal{X}_i falling in branch

GenerateTree(\mathcal{X}_i)

complexity
parameter

SplitAttribute(\mathcal{X})

MinEnt \leftarrow MAX

For all attributes $i = 1, \dots, d$

If \mathbf{x}_i is discrete with n values

Split \mathcal{X} into $\mathcal{X}_1, \dots, \mathcal{X}_n$ by \mathbf{x}_i

$e \leftarrow \text{SplitEntropy}(\mathcal{X}_1, \dots, \mathcal{X}_n)$ /* equation 9.8 */ (*or \mathcal{I}'_m*)

If $e < \text{MinEnt}$ MinEnt $\leftarrow e$; bestf $\leftarrow i$

Else /* \mathbf{x}_i is numeric */

For all possible splits

Split \mathcal{X} into $\mathcal{X}_1, \mathcal{X}_2$ on \mathbf{x}_i

$e \leftarrow \text{SplitEntropy}(\mathcal{X}_1, \mathcal{X}_2)$

If $e < \text{MinEnt}$ MinEnt $\leftarrow e$; bestf $\leftarrow i$

Return bestf

Why is
this a
heuristic?

Univariate **Regression** Trees

- Construction is similar to classification tree, except replace entropy by some regression measures
- Let \mathcal{X}_m be the set of instances which reached node m . We also define **error (or variance)** at node m as:

$$b_m(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x} \in \mathcal{X}_m: \mathbf{x} \text{ reaches node } m \\ 0 & \text{otherwise} \end{cases} \quad \text{where } N_m = |\mathcal{X}_m| = \sum_t b_m(\mathbf{x}^t)$$
$$E_m = \frac{1}{N_m} \sum_t (r^t - g_m)^2 b_m(\mathbf{x}^t) \quad g_m = \frac{\sum_t b_m(\mathbf{x}^t) r^t}{\sum_t b_m(\mathbf{x}^t)}$$

- If at node m , the error is acceptable (or $E_m < \theta_r$), then a leaf node is created and it stores the g_m value
- If the error is not acceptable, data reaching node m is split further to reduce the sum of errors in the branches

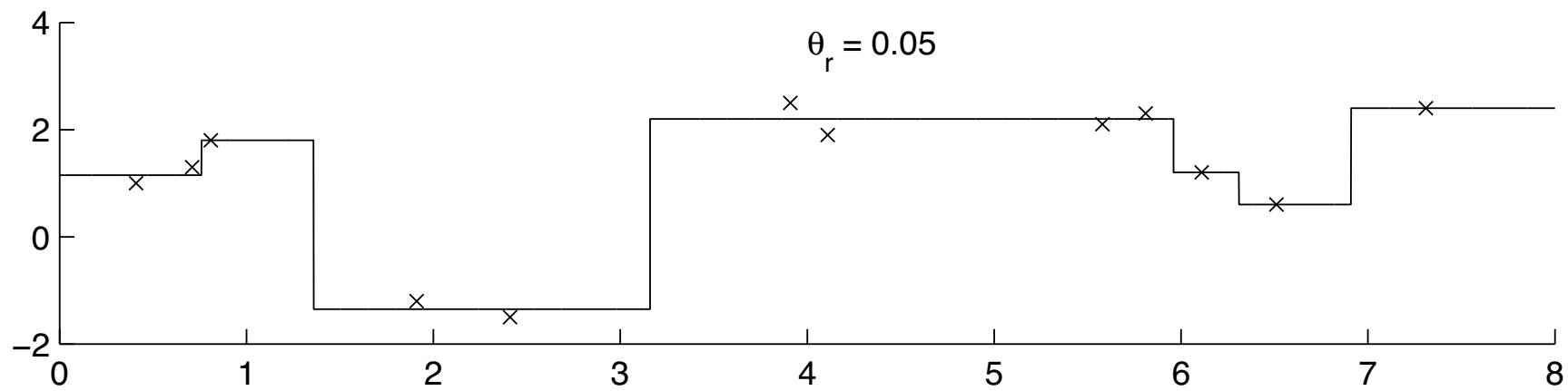
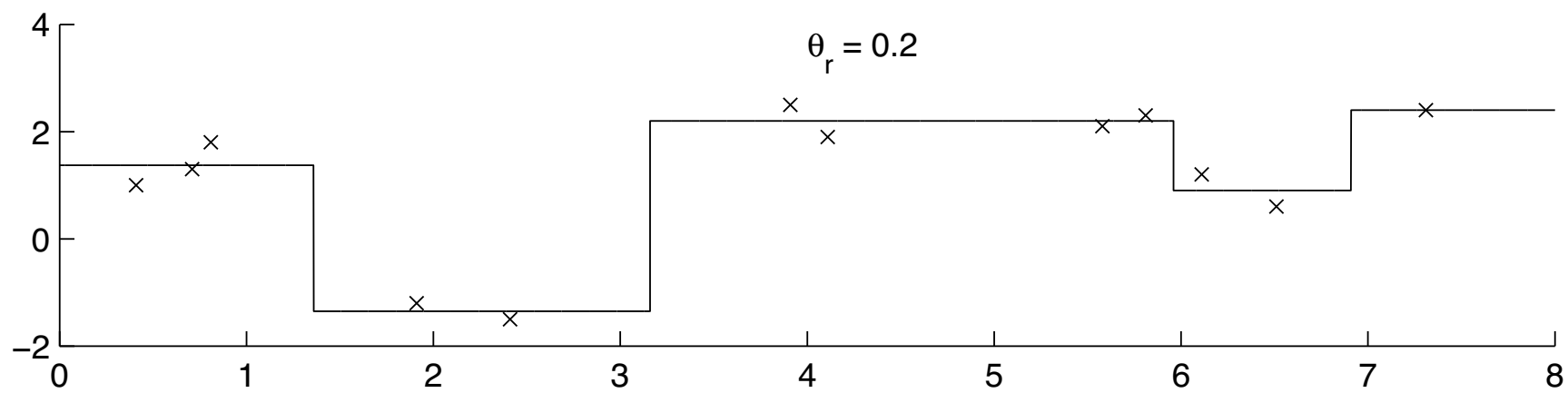
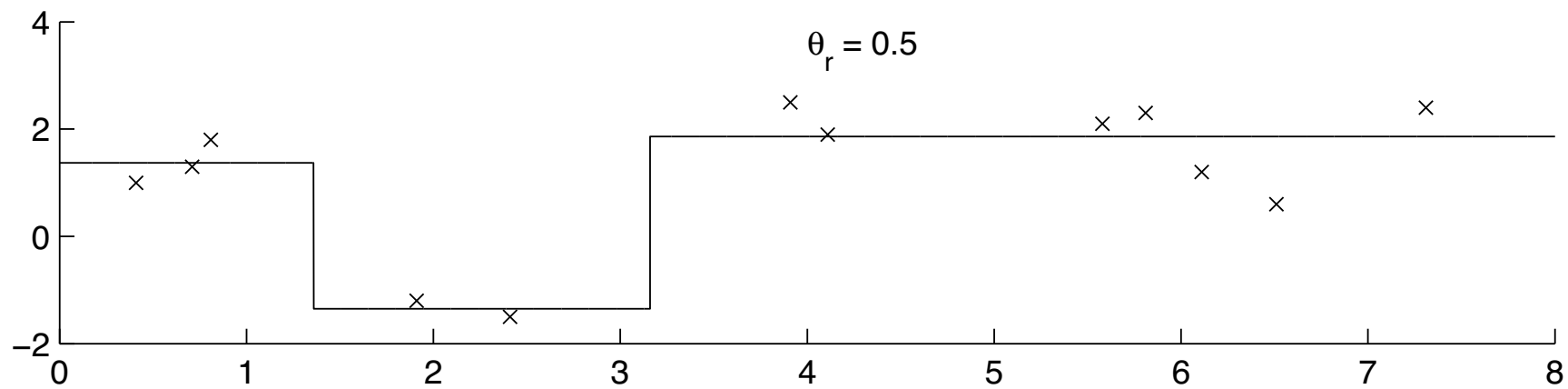
Univariate Regression Trees

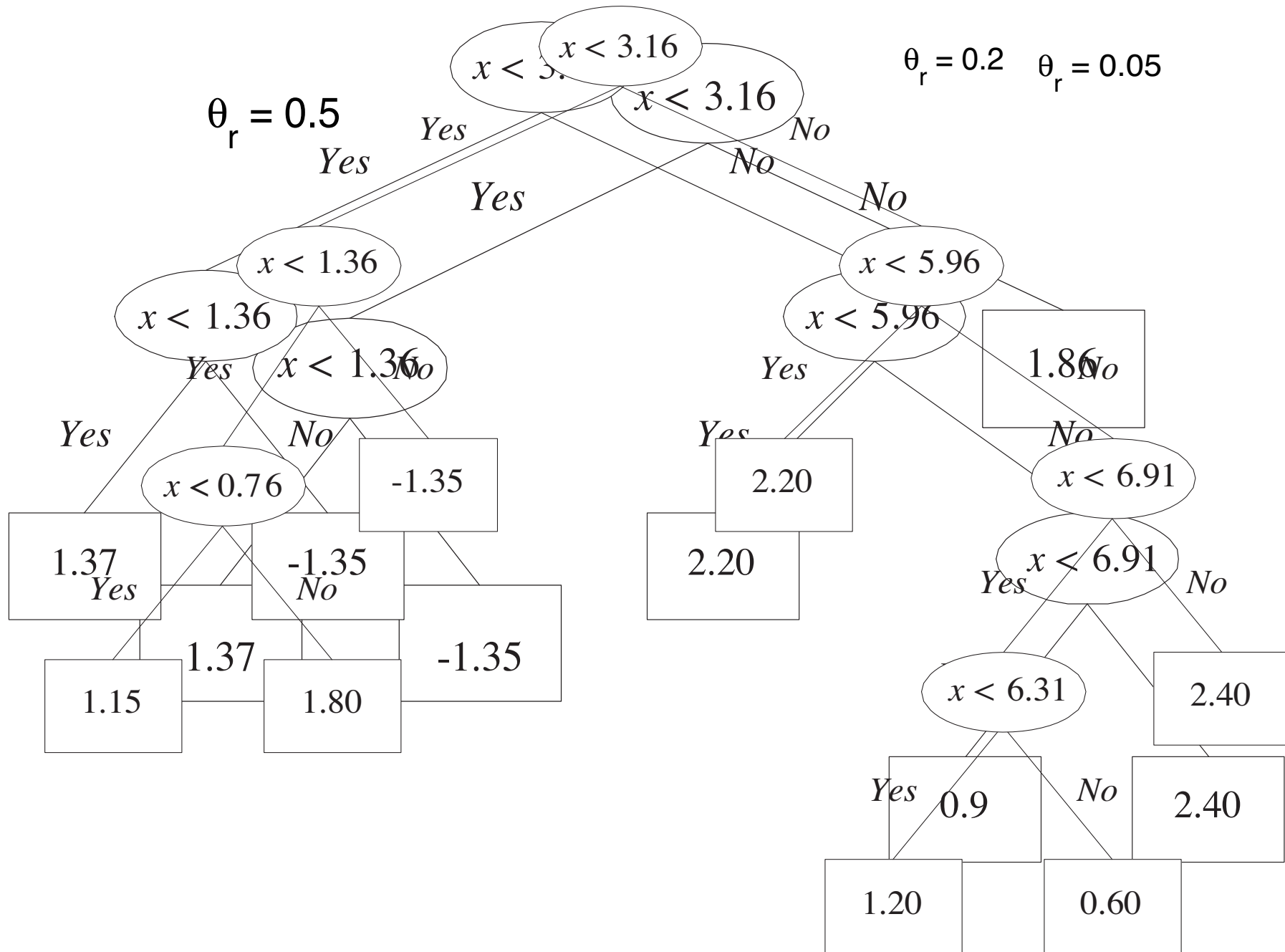
- To split at a node, we look for the attribute (and split threshold for a numeric attribute) that minimizes the error. Then continue recursively.
- Let \mathcal{X}_{mj} be the subset of \mathcal{X}_m taking branch j : $\cup_{j=1}^n \mathcal{X}_{mj} = \mathcal{X}_m$
- Define error as:

$$b_{mj}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x} \in \mathcal{X}_{mj}: \mathbf{x} \text{ reaches node } m \text{ and takes branch } j \\ 0 & \text{otherwise} \end{cases}$$

$$g_{mj} = \frac{\sum_t b_{mj}(\mathbf{x}^t) r^t}{\sum_t b_{mj}(\mathbf{x}^t)} \quad E'_m = \frac{1}{N_m} \sum_j \sum_t (r^t - g_{mj})^2 b_{mj}(\mathbf{x}^t)$$

- Find lowest E'_m . For tree construction, replace entropy calculation with E'_m in the alg, and θ_I by θ_r



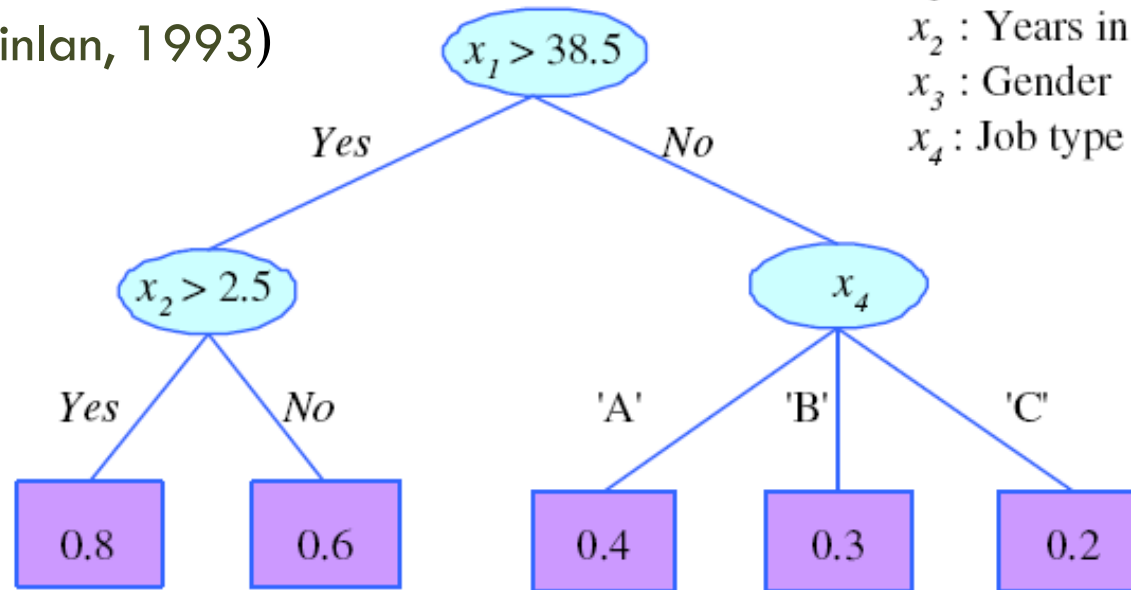


Pruning Trees

- Remove sub-trees for better generalization (decrease variance)
 - ▣ **Pre-pruning:** Early stopping
 - ▣ **Post-pruning:** Grow the whole tree then prune subtrees that overfit on the pruning set
- Pre-pruning is faster, post-pruning is more accurate (requires a separate pruning set)

Rule Extraction from Trees

C4.5Rules
(Quinlan, 1993)



x_1 : Age
 x_2 : Years in job
 x_3 : Gender
 x_4 : Job type

Use decision tree
As **feature extraction**.
In the figure, we
Don't need x_3

- R1: IF (age > 38.5) AND (years-in-job > 2.5) THEN $y = 0.8$
R2: IF (age > 38.5) AND (years-in-job ≤ 2.5) THEN $y = 0.6$
R3: IF (age ≤ 38.5) AND (job-type = 'A') THEN $y = 0.4$
R4: IF (age ≤ 38.5) AND (job-type = 'B') THEN $y = 0.3$
R5: IF (age ≤ 38.5) AND (job-type = 'C') THEN $y = 0.2$

Learning Rules from Data **Directly**

- Rule induction is similar to tree induction but
 - ▣ tree induction is breadth-first,
 - ▣ rule induction is **depth-first; one rule at a time**
- Rule set contains rules; rules are **conjunctions of conditions** on discrete or numeric *attributes*
- Rule is said to **cover** an example if all terms of the rule evaluate to true for that example
- Once a rule is added to the rule base, all training samples covered by the rule are removed, the process continues
- **Sequential covering:** Generate rules one at a time until all positive examples are covered
- IREP (Fürnkranz and Widmer, 1994), **Ripper** (Cohen, 1995)


```

Ripper(Pos, Neg, k)
  RuleSet  $\leftarrow$  LearnRuleSet(Pos, Neg)
  For  $k$  times
    RuleSet  $\leftarrow$  OptimizeRuleSet(RuleSet, Pos, Neg)
LearnRuleSet(Pos, Neg)
  RuleSet  $\leftarrow \emptyset$ 
  DL  $\leftarrow$  DescLen(RuleSet, Pos, Neg)
  Repeat
    Rule  $\leftarrow$  LearnRule(Pos, Neg)
    Add Rule to RuleSet
    DL'  $\leftarrow$  DescLen(RuleSet, Pos, Neg)
    If DL' > DL + 64
      PruneRuleSet(RuleSet, Pos, Neg)
      Return RuleSet
    If DL' < DL DL  $\leftarrow$  DL'
    Delete instances covered by Rule from Pos and Neg
  Until Pos =  $\emptyset$ 
  Return RuleSet

```

```

PruneRuleSet(RuleSet,Pos,Neg)
  For each Rule  $\in$  RuleSet in reverse order
    DL  $\leftarrow$  DescLen(RuleSet,Pos,Neg)
    DL'  $\leftarrow$  DescLen(RuleSet-Rule,Pos,Neg)
    IF DL' < DL Delete Rule from RuleSet
  Return RuleSet

OptimizeRuleSet(RuleSet,Pos,Neg)
  For each Rule  $\in$  RuleSet
    DL0  $\leftarrow$  DescLen(RuleSet,Pos,Neg)
    DL1  $\leftarrow$  DescLen(RuleSet-Rule+
      ReplaceRule(RuleSet,Pos,Neg),Pos,Neg)
    DL2  $\leftarrow$  DescLen(RuleSet-Rule+
      ReviseRule(RuleSet,Rule,Pos,Neg),Pos,Neg)
    If DL1 = min(DL0,DL1,DL2)
      Delete Rule from RuleSet and
      add ReplaceRule(RuleSet,Pos,Neg)
    Else If DL2 = min(DL0,DL1,DL2)
      Delete Rule from RuleSet and
      add ReviseRule(RuleSet,Rule,Pos,Neg)
  Return RuleSet

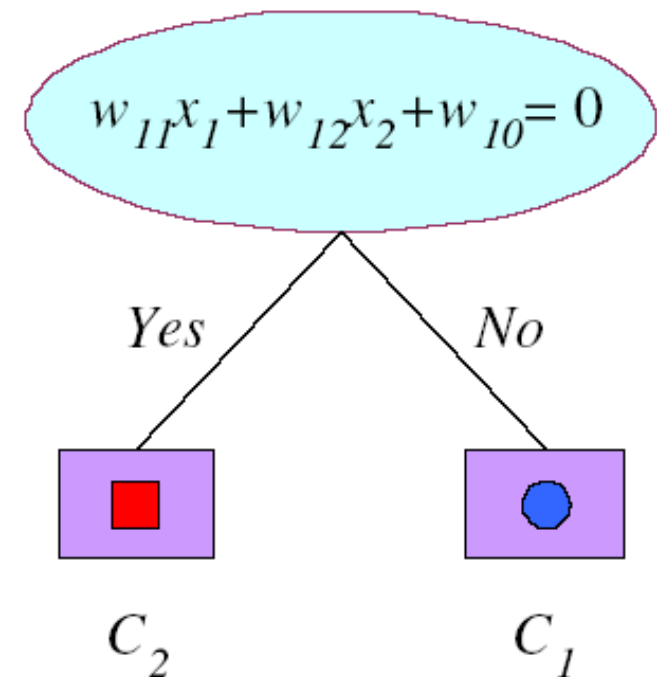
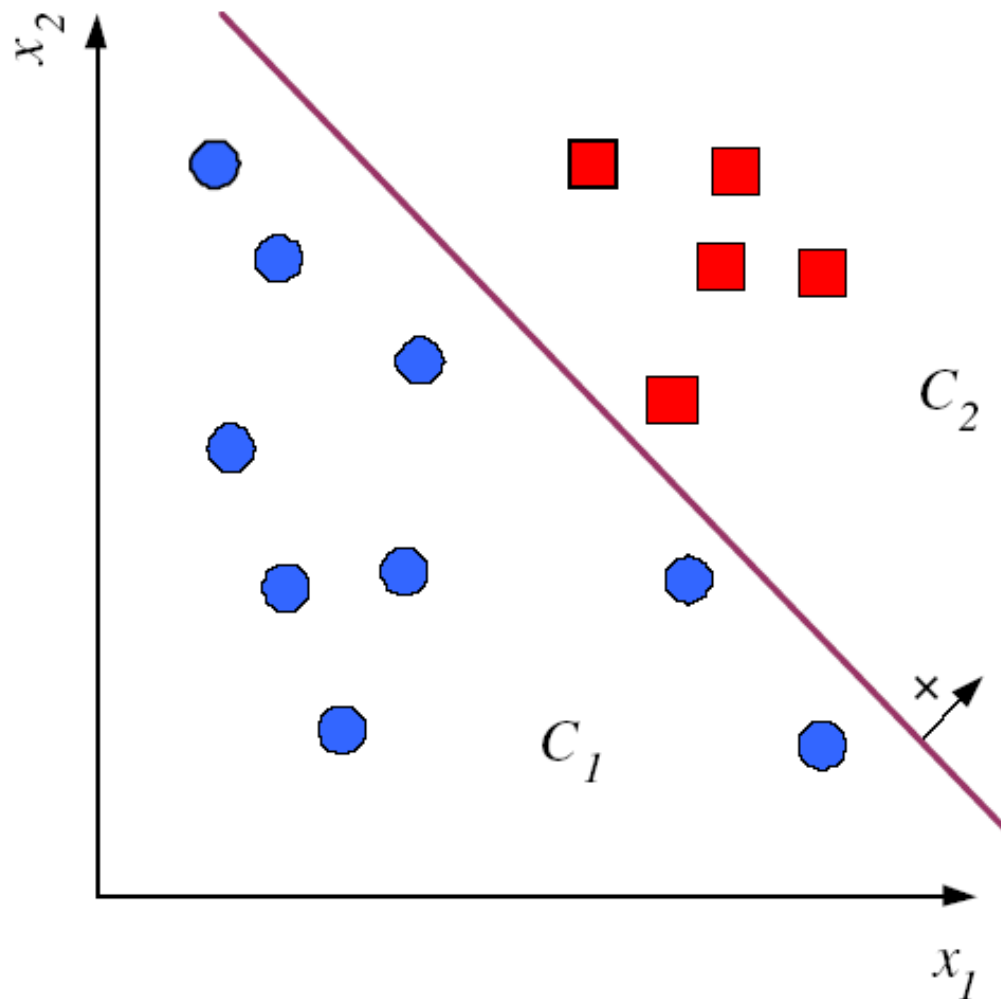
```

Multivariate Trees

- Unlike univariate tree, at a decision node, **ALL** input dimensions can be used
- When all inputs are numeric, a binary linear multivariate node is defined as: $f_m(\mathbf{x}) : \mathbf{w}_m^T \mathbf{x} + w_{m0} > 0$
- The above equation defines a *hyperplane* (see figure)
- For discrete attributes should be represented by 0/1 dummy numeric variables
- Keep dividing until we have leaf nodes which are defined by a *polyhedra* in the input space

Multivariate Trees

Linear multivariate hyperplane



Multivariate Trees

- Instead of linear multivariate discriminant, we can use nonlinear multivariate discriminant at a node.

- Example: *quadratic multivariate discriminant*

$$f_m(\mathbf{x}) : \mathbf{x}^T \mathbf{W}_m \mathbf{x} + \mathbf{w}_m^T \mathbf{x} + w_{m0} > 0$$

- Example: *sphere node*

$$f_m(\mathbf{x}) : \|\mathbf{x} - \mathbf{c}_m\| \leq \alpha_m$$

- Many other discriminant methods, but note that these are all “heuristics”