# Decision Trees

- Decision tree
  - A flow-chart-like tree structure
  - Internal node denotes a test on an attribute
  - Branch represents an outcome of the test
  - Leaf nodes represent class labels or class distribution

# Catching tax-evasion

| 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 |

Tax-return data for year 2011
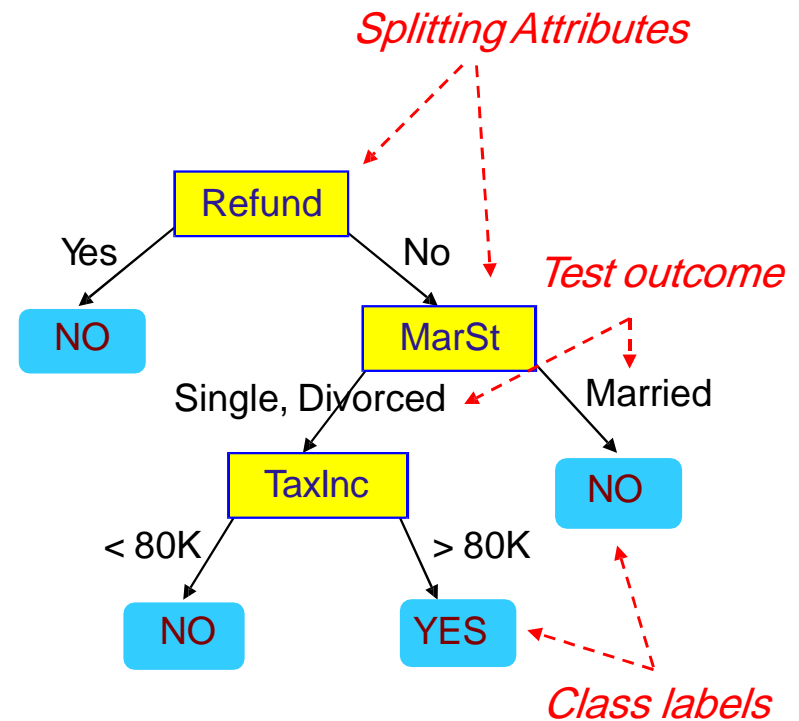
A new tax return for 2012
Is this a cheating tax return?

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

An instance of the classification problem: learn a method for discriminating between records of different classes (cheaters vs non-cheaters)

# 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

*Test outcome*

MarSt
Single, Divorced → TaxInc
Married → NO

TaxInc
< 80K → NO
> 80K → YES

*Class labels*

Model: Decision Tree

# Another Example of Decision Tree



| 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

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

# Decision Tree Classification Task

| 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 |

Training Set

Tree Induction algorithm

Induction

**Learn Model**

**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 | ? |

Test Set

**Apply Model**

Deduction

# Apply Model to Test Data

Start from the root of tree.

**Test Data**

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

# Apply Model to Test Data

**Test Data**

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

# Apply Model to Test Data

**Test Data**

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

# Apply Model to Test Data

**Test Data**

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

# Apply Model to Test Data

**Test Data**

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

**Refund**

Yes → **NO**

No → **MarSt**

**MarSt**

Single, Divorced → **TaxInc**

Married → **NO**

**TaxInc**

< 80K → **NO**

> 80K → **YES**

# Apply Model to Test Data

**Test Data**

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



Assign Cheat to "No"

# Decision Tree Classification Task

| 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 |

Training Set

Induction

Tree Induction algorithm

Learn Model

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 | ? |

Test Set

Apply Model

Deduction

# Decision Tree Induction

- Many Algorithms:
  - Hunt's Algorithm (one of the earliest)
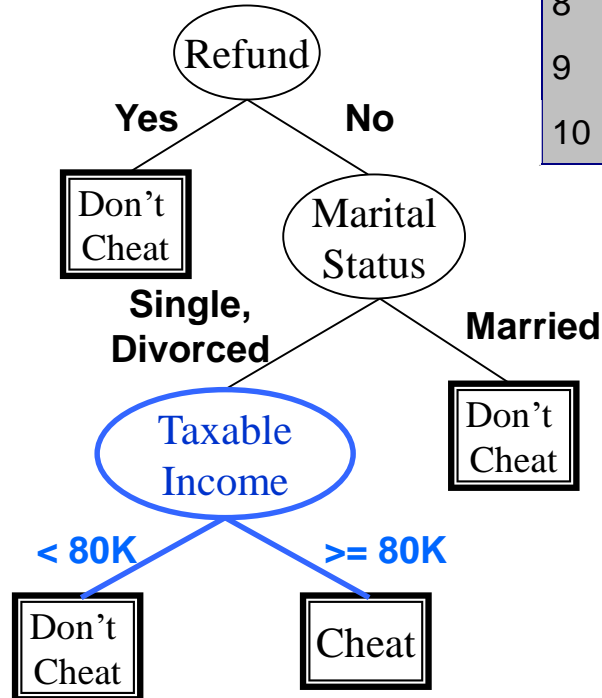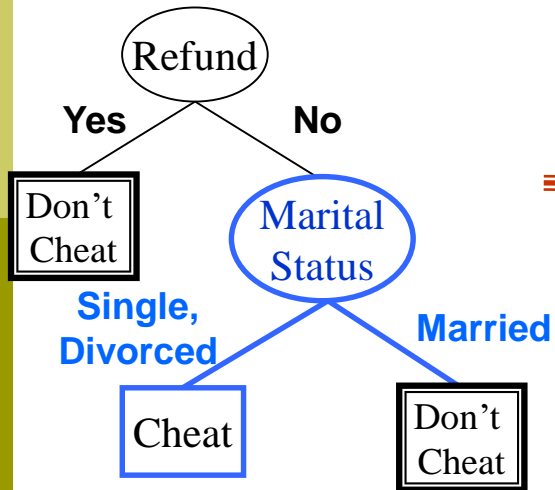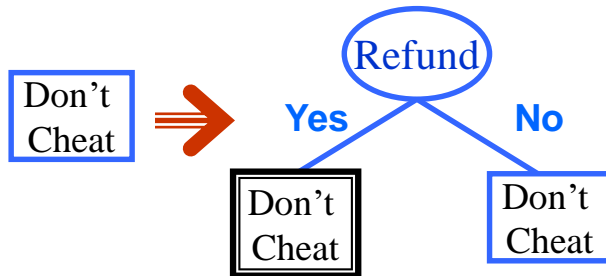  - CART
  - ID3, C4.5
  - SLIQ,SPRINT

# General Structure of Hunt's Algorithm

- Let $D_t$ be the set of training records that reach a node t
- General Procedure:
  - If $D_t$ contains records that belong the same class $y_t$, then t is a leaf node labeled as $y_t$
  - If $D_t$ is an empty set, then t is a leaf node labeled by the default class, $y_d$
  - If $D_t$ contains records that belong to more than one class, use an attribute test to split the data into smaller subsets. Recursively apply the procedure to each subset.

| 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 |

$D_t$

?

# Hunt's Algorithm

| 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** |

# Constructing decision-trees (pseudocode)

GenDecTree(Sample S, Features F)

1. If stopping_condition(S,F) = true then
   a. leaf = createNode()
   b. leaf.label = Classify(S)
   c. return leaf
2. root = createNode()
3. root.test_condition = findBestSplit(S,F)
4. $V = \{v| v$ a possible outcome of root.test_condition$\}$
5. for *each* value v∈V:
   a. $S_v := \{s|$ root.test_condition(s) = v and s ∈ S$\}$;
   b. child = GenDecTree($S_v$, F);
   c. Add child as a descent of root and label the edge (root→child) as v
6. return root

# 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?
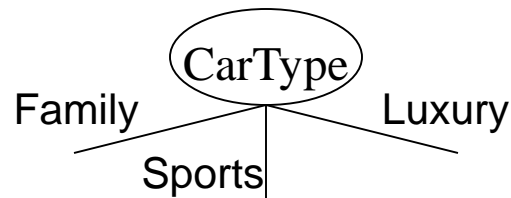  - Determine when to stop splitting

# 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
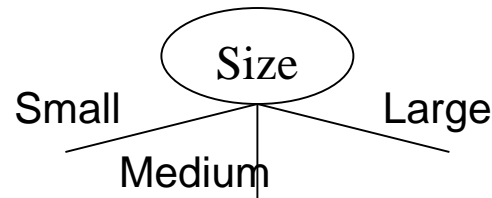    - 2-way split
    - Multi-way split

# Splitting Based on Nominal Attributes

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

```
        (CarType)
   Family  |  Luxury
        Sports
```

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

```
{Sports,   (CarType)                        {Family,   (CarType)
 Luxury}       {Family}          OR          Luxury}       {Sports}
```

# Splitting Based on Ordinal Attributes

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

```
        ( Size )
Small   /   |   \   Large
           Medium
```

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

```
{Small,    ( Size )
Medium}   /       \   {Large}
```

OR

```
{Medium,   ( Size )
Large}    /       \   {Small}
```

- **What about this split?**

```
{Small,    ( Size )
Large}    /       \   {Medium}
```
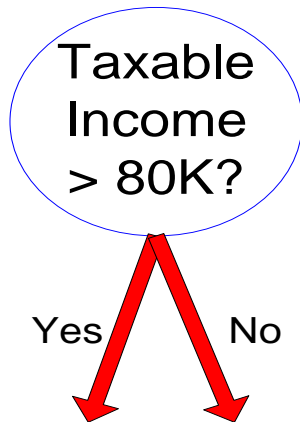
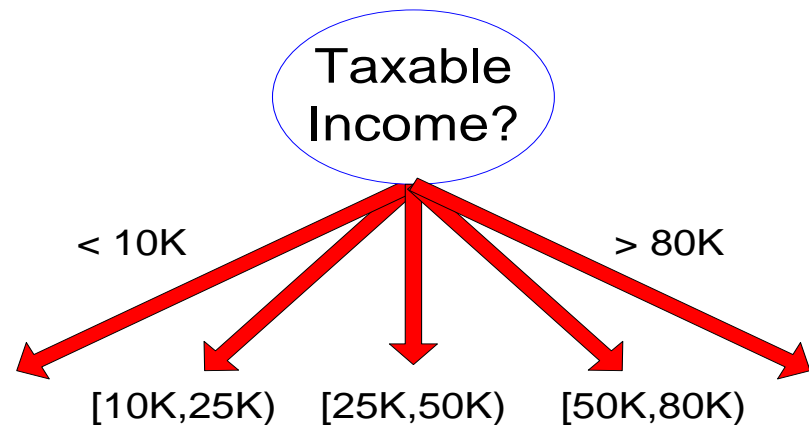# Splitting Based on Continuous Attributes

- ◻ 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



(i) Binary split

(ii) Multi-way split

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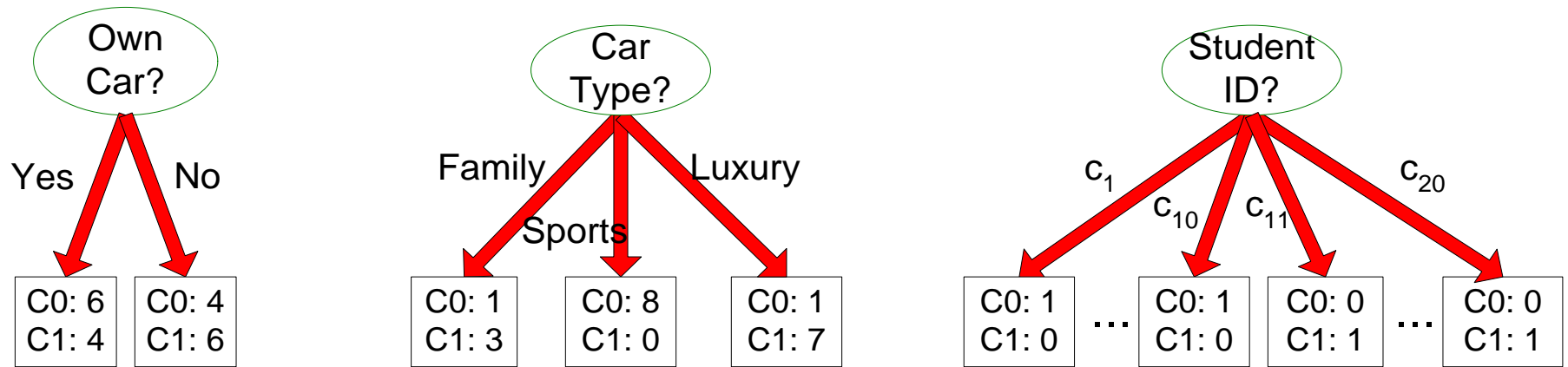
- Greedy strategy.
  - Split the records based on an attribute test that optimizes certain criterion.

- 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

**Before Splitting: 10 records of class 0,**
**10 records of class 1**



**Which test condition is the best?**

# How to determine the Best Split

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

|       |
|-------|
| C0: 5 |
| C1: 5 |

**Non-homogeneous,**

**High degree of impurity**

|       |
|-------|
| C0: 9 |
| C1: 1 |

**Homogeneous,**

**Low degree of impurity**

# Measures of Node Impurity

- Gini Index

- Entropy

- Misclassification error

# How to Find the Best Split

**Before Splitting:**

| C0 | **N00** |
|----|---------|
| C1 | **N01** |

→ **M0**

**A?**

Yes                                    No

| Node N1 |

| Node N2 |

| C0 | **N10** |
|----|---------|
| C1 | **N11** |

| C0 | **N20** |
|----|---------|
| C1 | **N21** |

↓ **M1**

↓ **M2**

**M12**

**B?**

Yes                                    No

| Node N3 |

| Node N4 |

| C0 | **N30** |
|----|---------|
| C1 | **N31** |

| C0 | **N40** |
|----|---------|
| C1 | **N41** |

↓ **M3**

↓ **M4**

**M34**

**Gain = M0 – M12 vs M0 – M34**

# Measure of Impurity: GINI

- Gini Index for a given node t :

$$GINI(t) = 1 - \sum_j [p(j \mid t)]^2$$

(NOTE: $p(j/t)$ is the relative frequency of class j at node t).

- Maximum ($1 - 1/n_c$) when records are equally distributed among all classes, implying least interesting information
- Minimum (0.0) when all records belong to one class, implying most interesting information

| C1 | 0 |
|----|---|
| C2 | 6 |
| Gini=0.000 ||

| C1 | 1 |
|----|---|
| C2 | 5 |
| Gini=0.278 ||

| C1 | 2 |
|----|---|
| C2 | 4 |
| Gini=0.444 ||

| C1 | 3 |
|----|---|
| C2 | 3 |
| Gini=0.500 ||

# Examples for computing GINI

$$GINI(t) = 1 - \sum_j [p(j\,|\,t)]^2$$

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

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

Gini = 1 – P(C1)² – P(C2)² = 1 – 0 – 1 = 0

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

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

Gini = 1 – (1/6)² – (5/6)² = 0.278

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

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

Gini = 1 – (2/6)² – (4/6)² = 0.444

# Splitting Based on GINI

- Used in CART, SLIQ, SPRINT.
- When a node p is split into k partitions (children), the quality of split is computed as,

$$GINI_{split} = \sum_{i=1}^{k} \frac{n_i}{n} GINI(i)$$

where,   $n_i$ = number of records at child i,

   n  = number of records at node p.

# Binary Attributes: Computing GINI Index

- Splits into two partitions
- Effect of Weighing partitions:
  - Larger and Purer Partitions are sought for.

|  | Parent |
|---|---|
| C1 | 6 |
| C2 | 6 |
| **Gini = 0.500** ||

B?

Yes        No

Node N1       Node N2

**Gini(N1)**
$= 1 - (5/6)^2 - (2/6)^2$
$= 0.194$

**Gini(N2)**
$= 1 - (1/6)^2 - (4/6)^2$
$= 0.528$

|  | N1 | N2 |
|---|---|---|
| C1 | 5 | 1 |
| C2 | 2 | 4 |
| **Gini=0.333** |||

**Gini(Children)**
$= 7/12 * 0.194 +$
$\quad 5/12 * 0.528$
$= 0.333$

# Categorical Attributes: Computing Gini Index

- For each distinct value, gather counts for each class in the dataset
- Use the count matrix to make decisions

Multi-way split

| CarType | | |
|---|---|---|
| **Family** | **Sports** | **Luxury** |
| 1 | 2 | 1 |
| 4 | 1 | 1 |
| **0.393** | | |

(First table row labels C1, C2, Gini)

| | CarType | | |
|---|---|---|---|
| | **Family** | **Sports** | **Luxury** |
| **C1** | 1 | 2 | 1 |
| **C2** | 4 | 1 | 1 |
| **Gini** | 0.393 | | |

Two-way split
(find best partition of values)

| | CarType | |
|---|---|---|
| | **{Sports, Luxury}** | **{Family}** |
| **C1** | 3 | 1 |
| **C2** | 2 | 4 |
| **Gini** | 0.400 | |

| | CarType | |
|---|---|---|
| | **{Sports}** | **{Family, Luxury}** |
| **C1** | 2 | 2 |
| **C2** | 1 | 5 |
| **Gini** | 0.419 | |

# Continuous Attributes: Computing Gini Index

- Use Binary Decisions based on one value
- Several Choices for the splitting value
  - Number of possible splitting values
    = Number of distinct values
- Each splitting value has a count matrix associated with it
  - Class counts in each of the partitions, A < v and A ≥ v
- Simple method to choose best v
  - For each v, scan the database to gather count matrix and compute its Gini index
  - Computationally Inefficient! Repetition of work.

| 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** |

Taxable Income > 80K?

Yes    No

# Continuous Attributes: Computing Gini Index...

- For efficient computation: for each attribute,
  - Sort the attribute on values
  - Linearly scan these values, each time updating the count matrix and computing gini index
  - Choose the split position that has the least gini index

| Cheat | No | | No | | No | | Yes | | Yes | | Yes | | No | | No | | No | | No | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Taxable Income** | | | | | | | | | | | | | | | | | | | | |
| Sorted Values → | 60 | | 70 | | 75 | | 85 | | 90 | | 95 | | 100 | | 120 | | 125 | | 220 | |
| Split Positions → | 55 | | 65 | | 72 | | 80 | | 87 | | 92 | | 97 | | 110 | | 122 | | 172 | | 230 | |
| | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > |
| **Yes** | 0 | 3 | 0 | 3 | 0 | 3 | 0 | 3 | 1 | 2 | 2 | 1 | 3 | 0 | 3 | 0 | 3 | 0 | 3 | 0 |
| **No** | 0 | 7 | 1 | 6 | 2 | 5 | 3 | 4 | 3 | 4 | 3 | 4 | 3 | 4 | 4 | 3 | 5 | 2 | 6 | 1 | 7 | 0 |
| **Gini** | 0.420 | | 0.400 | | 0.375 | | 0.343 | | 0.417 | | 0.400 | | *0.300* | | 0.343 | | 0.375 | | 0.400 | | 0.420 | |

# Alternative Splitting Criteria based on INFO

- Entropy at a given node t:

$$Entropy(t) = -\sum_j p(j \mid t) \log p(j \mid t)$$

(NOTE: $p(j \mid t)$ is the relative frequency of class j at node t).

- Measures homogeneity of a node.
  - Maximum ($\log n_c$) when records are equally distributed among all classes implying least information
  - Minimum (0.0) when all records belong to one class, implying most information

- Entropy based computations are similar to the GINI index computations

# Examples for computing Entropy

$$Entropy(t) = -\sum_j p(j \mid t) \log_2 p(j \mid t)$$

| 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

# Splitting Based on INFO...

- 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.

# Splitting Based on INFO...

- Gain Ratio:

$$GainRATIO_{split} = \frac{GAIN_{Split}}{SplitINFO}$$

$$SplitINFO = -\sum_{i=1}^{k} \frac{n_i}{n} \log \frac{n_i}{n}$$

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

# Splitting Criteria based on Classification Error

- Classification error at a node t :

$$Error(t) = 1 - \max_i P(i \mid t)$$

- Measures misclassification error made by a node.
  - Maximum ($1 - 1/n_c$) when records are equally distributed among all classes, implying least interesting information
  - Minimum (0.0) when all records belong to one class, implying most interesting information

# Examples for Computing Error

$$Error(t) = 1 - \max_i P(i \mid t)$$

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

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

Error = 1 – max (0, 1) = 1 – 1 = 0

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

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

Error = 1 – max (1/6, 5/6) = 1 – 5/6 = 1/6

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

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

Error = 1 – max (2/6, 4/6) = 1 – 4/6 = 1/3

# Comparison among Splitting Criteria

**For a 2-class problem:**

# 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

# Other Issues

- Data Fragmentation
- Expressiveness

# Data Fragmentation

- Number of instances gets smaller as you traverse down the tree

- Number of instances at the leaf nodes could be too small to make any statistically significant decision

- You can introduce a lower bound on the number of items per leaf node in the stopping criterion.

# Expressiveness

- A classifier defines a function that discriminates between two (or more) classes.

- The expressiveness of a classifier is the class of functions that it can model, and the kind of data that it can separate

  - When we have discrete (or binary) values, we are interested in the class of boolean functions that can be modeled

  - If the data-points are real vectors we talk about the decision boundary that the classifier can model

# Decision Boundary



- Border line between two neighboring regions of different classes is known as decision boundary

- Decision boundary is parallel to axes because test condition involves a single attribute at-a-time

# Expressiveness

- Decision tree provides expressive representation for learning discrete-valued function
  - But they do not generalize well to certain types of Boolean functions
    - Example: parity function:
      - Class = 1 if there is an even number of Boolean attributes with truth value = True
      - Class = 0 if there is an odd number of Boolean attributes with truth value = True
    - For accurate modeling, must have a complete tree

- Less expressive for modeling continuous variables
  - Particularly when test condition involves only a single attribute at-a-time

# Oblique Decision Trees



$$x + y < 1$$

Class = +    Class = ●
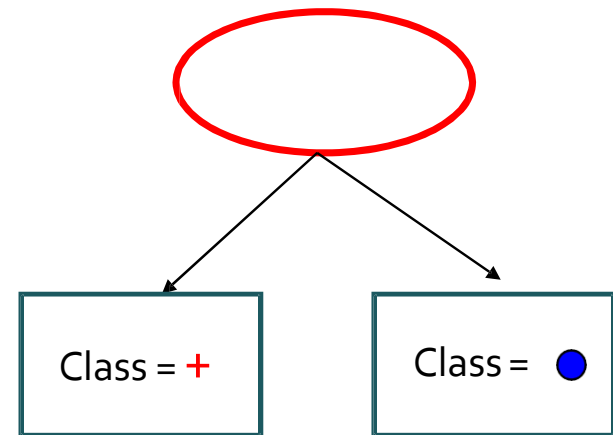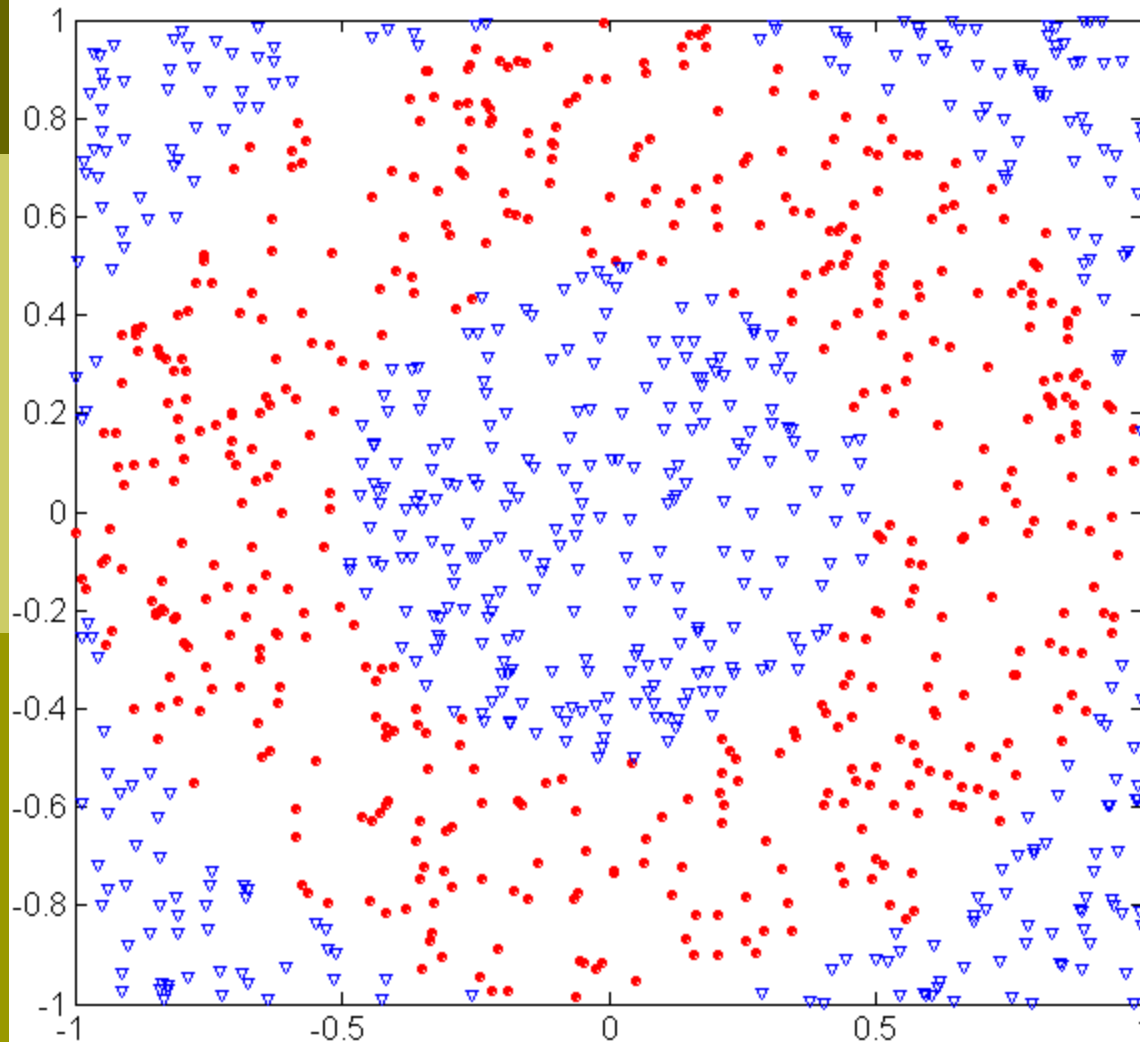
- Test condition may involve multiple attributes
- More expressive representation
- Finding optimal test condition is computationally expensive

# Practical Issues of Classification

- Underfitting and Overfitting

- Missing Values

- Costs of Classification

# Underfitting and Overfitting (Example)



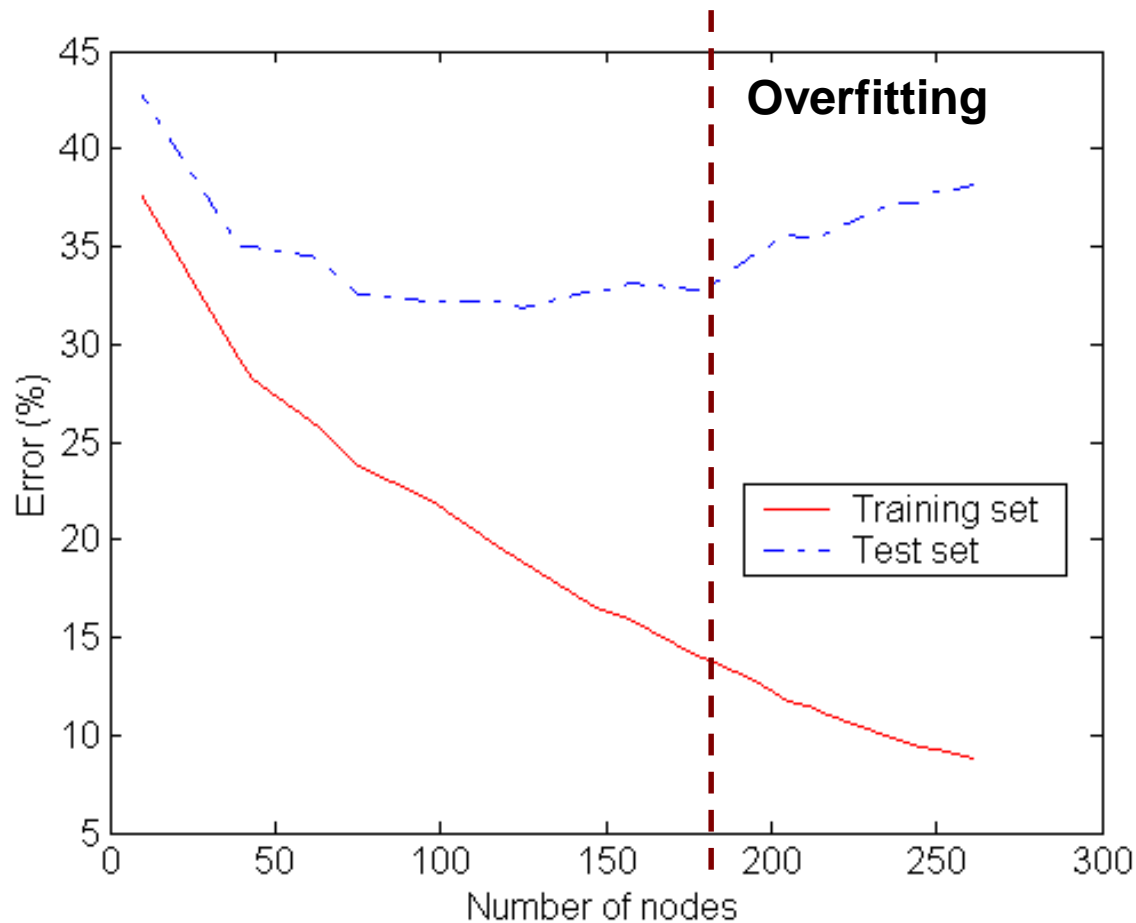**500 circular and 500 triangular data points.**

**Circular points:**

$0.5 \leq \text{sqrt}(x_1^2+x_2^2) \leq 1$

**Triangular points:**

$\text{sqrt}(x_1^2+x_2^2) > 0.5$ or

$\text{sqrt}(x_1^2+x_2^2) < 1$

# Underfitting and Overfitting



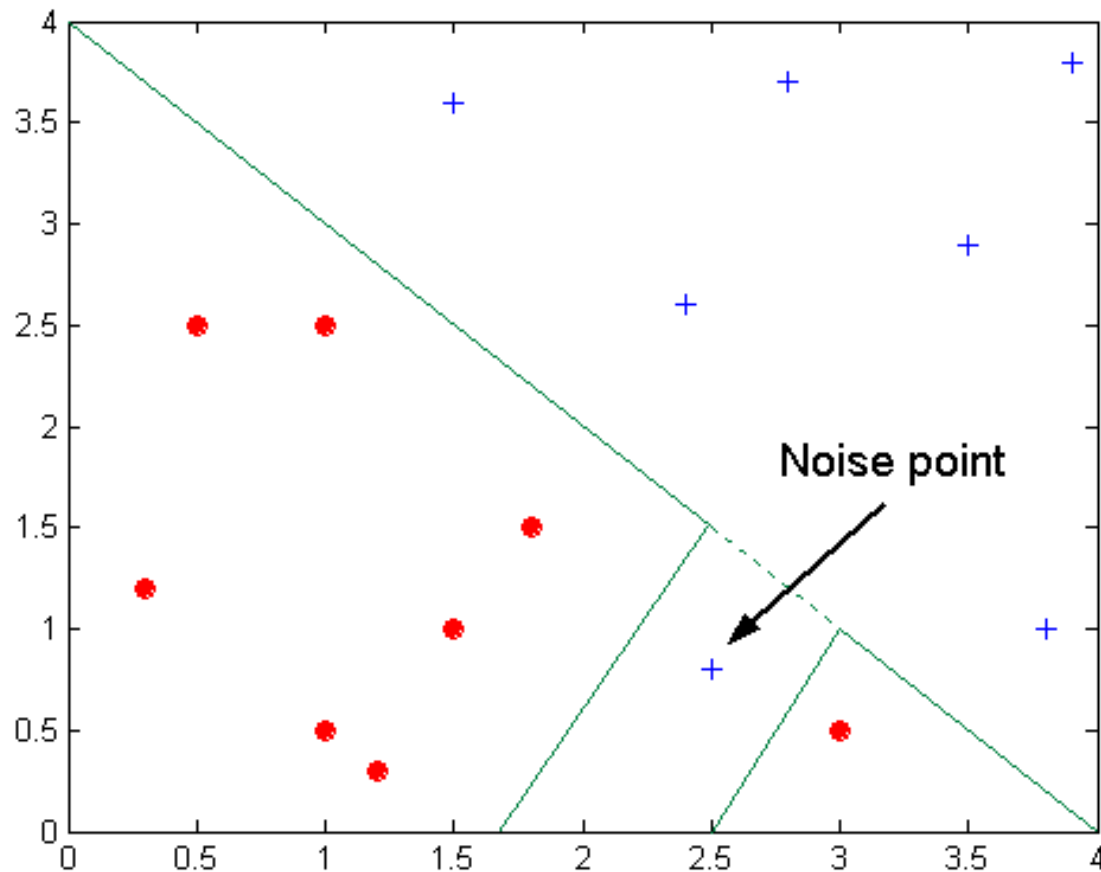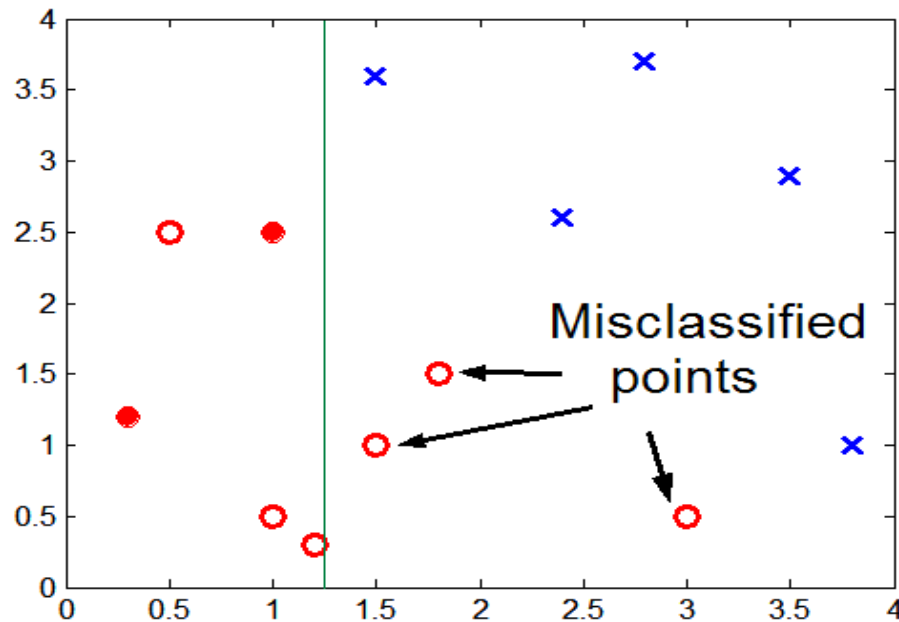**Underfitting**: when model is too simple, both training and test errors are large

# Overfitting due to Noise



**Decision boundary is distorted by noise point**

# Overfitting due to Insufficient Examples



**Lack of data points in the lower half of the diagram makes it difficult to predict correctly the class labels of that region**

**- Insufficient number of training records in the region causes the decision tree to predict the test examples using other training records that are irrelevant to the classification task**

# 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

- Need new ways for estimating errors

# Estimating Generalization Errors

- Re-substitution errors: error on training ($\Sigma$ e(t) )
- Generalization errors: error on testing ($\Sigma$ e'(t))
- Methods for estimating generalization errors:
  - Optimistic approach:  e'(t) = e(t)
  - Pessimistic approach：
    - For each leaf node: e'(t) = (e(t)+0.5)
    - Total errors: e'(T) = e(T) + N $\times$ 0.5 (N: number of leaf nodes)
    - For a tree with 30 leaf nodes and 10 errors on training (out of 1000 instances):
      Training error = 10/1000 = 1%

      Generalization error = (10 + 30$\times$0.5)/1000 = 2.5%
  - Reduced error pruning (REP):
    - uses validation data set to estimate generalization error

# 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

# Example of Post-Pruning

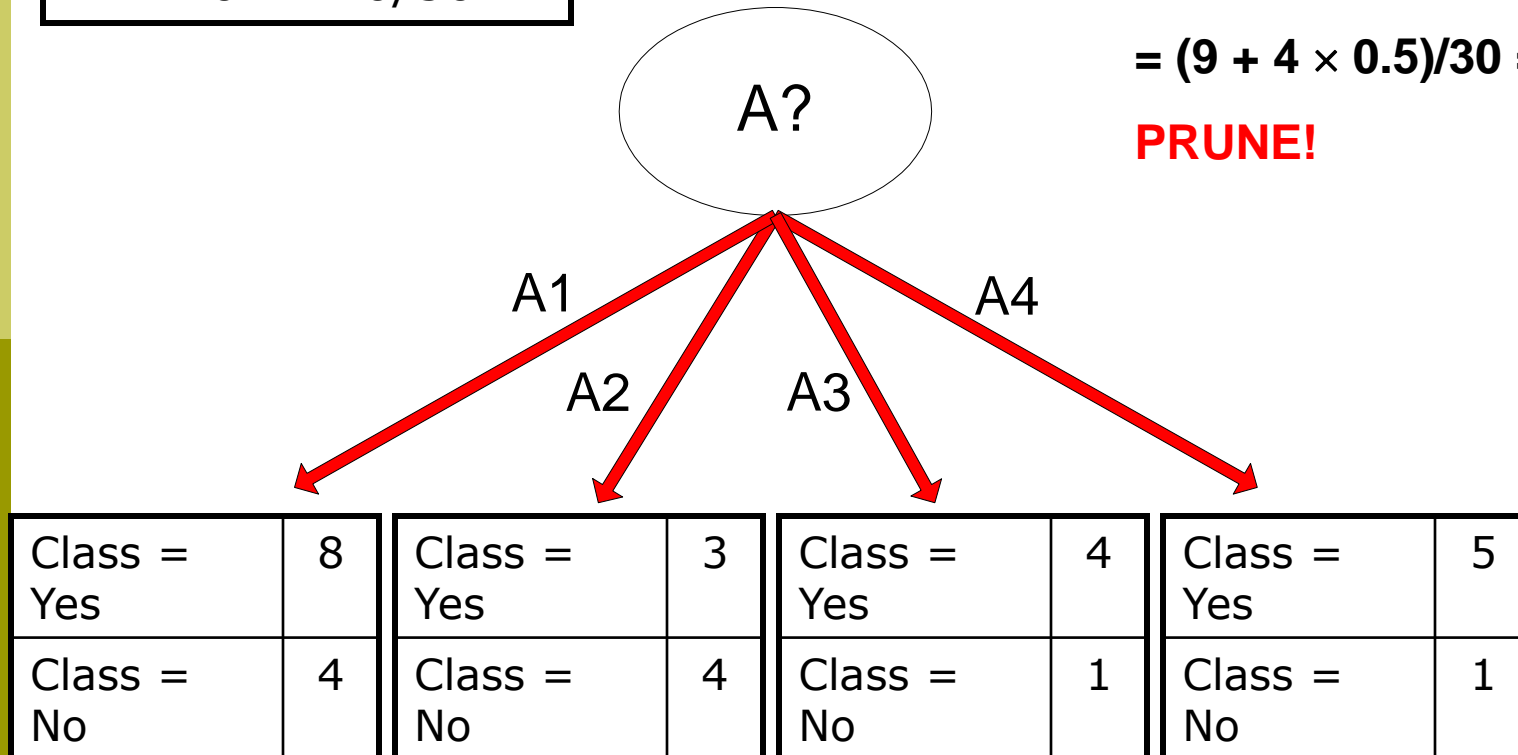| Class = Yes | 20 |
|---|---|
| Class = No | 10 |
| Error = 10/30 | |

**Training Error (Before splitting) = 10/30**

**Pessimistic error = (10 + 0.5)/30 = 10.5/30**

**Training Error (After splitting) = 9/30**

**Pessimistic error (After splitting)**

$$= (9 + 4 \times 0.5)/30 = 11/30$$

**PRUNE!**

A?

A1        A2        A3        A4

| Class = Yes | 8 | Class = Yes | 3 | Class = Yes | 4 | Class = Yes | 5 |
|---|---|---|---|---|---|---|---|
| Class = No | 4 | Class = No | 4 | Class = No | 1 | Class = No | 1 |

# 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)
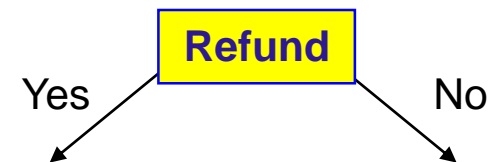= 0.3 (0) + 0.6 (0.9183) = 0.551

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** |

**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**

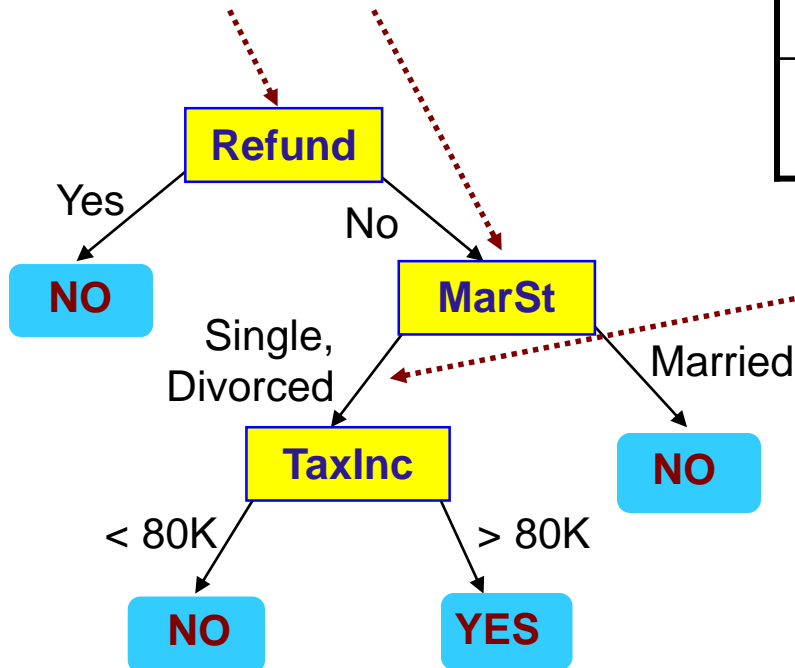**Assign record to the left child with weight = 3/9 and to the right child with weight = 6/9**

# Classify Instances

**New record:**

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

|  | 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 → TaxInc
    - < 80K → **NO**
    - > 80K → **YES**
  - Married → **NO**

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

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

# Regression trees

$$X^{(1)} \quad \dots \quad X^{(p)} \quad Y$$

| Gender | Rich? | Num. Children | # travel per yr. | Age |
|--------|-------|---------------|------------------|-----|
| F | No | 2 | 5 | 38 |
| M | No | 0 | 2 | 25 |
| M | Yes | 1 | 0 | 72 |
| : | : | : | : | : |

**Num Children?**

≥ 2                    < 2

**Gender?**

Female              Male

Predicted age=39      Predicted age=36

Average (fit a constant ) using training data at the leaves

28

# Scalable Decision Tree Induction Methods

- **SLIQ** (EDBT'96 — Mehta et al.)
  - Builds an index for each attribute and only class list and the current attribute list reside in memory
- **SPRINT** (VLDB'96 — J. Shafer et al.)
  - Constructs an attribute list data structure
- **PUBLIC** (VLDB'98 — Rastogi & Shim)
  - Integrates tree splitting and tree pruning: stop growing the tree earlier
- **RainForest** (VLDB'98 — Gehrke, Ramakrishnan & Ganti)
  - Builds an AVC-list (attribute, value, class label)
- **BOAT** (PODS'99 — Gehrke, Ganti, Ramakrishnan & Loh)
  - Uses bootstrapping to create several small samples

# Bagging: Bootstrap Aggregation

- Leo Breiman (1994)

- Take repeated bootstrap samples from training set $D$

- *Bootstrap sampling*: Given set $D$ containing $N$ training examples, create $D'$ by drawing $N$ examples at random with replacement from $D$.

- Bagging:

  – Create $k$ bootstrap samples $D_1 \ldots D_k$.

  – Train disGnct classifier on each $D_i$.

  – Classify new instance by majority vote / average.
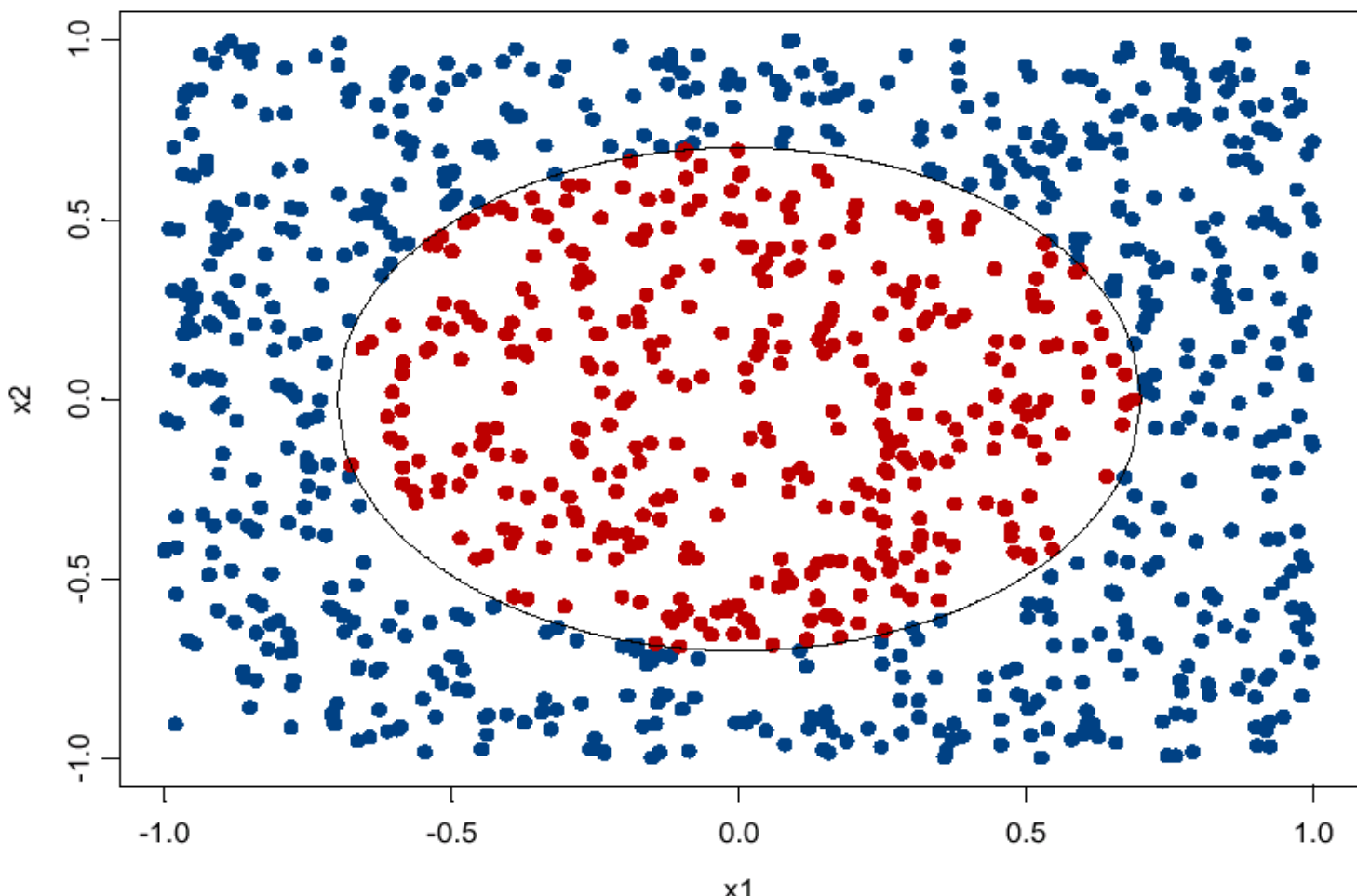
# General Idea

# Example of Bagging

- Sampling with replacement

Data ID

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Original Data | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| Bagging (Round 1) | 7 | 8 | 10 | 8 | 2 | 5 | 10 | 10 | 5 | 9 |
| Bagging (Round 2) | 1 | 4 | 9 | 1 | 2 | 3 | 2 | 7 | 3 | 2 |
| Bagging (Round 3) | 1 | 8 | 5 | 10 | 5 | 5 | 9 | 6 | 3 | 7 |

Training Data

- Build classifier on each bootstrap sample

- Each data point has probability $(1 - 1/n)^n$ of being selected as test data

- Training data = $1 - (1 - 1/n)^n$ of the original data

# Bagging Example

# CART decision boundary

# 100 bagged trees

# Random Forests

- Ensemble method specifically designed for decision tree classifiers

- Introduce two sources of randomness: "Bagging" and "Random input vectors"
  - Bagging method: each tree is grown using a bootstrap sample of training data
  - Random vector method: At each node, best split is chosen from a random sample of $m$ aIributes instead of all aIributes
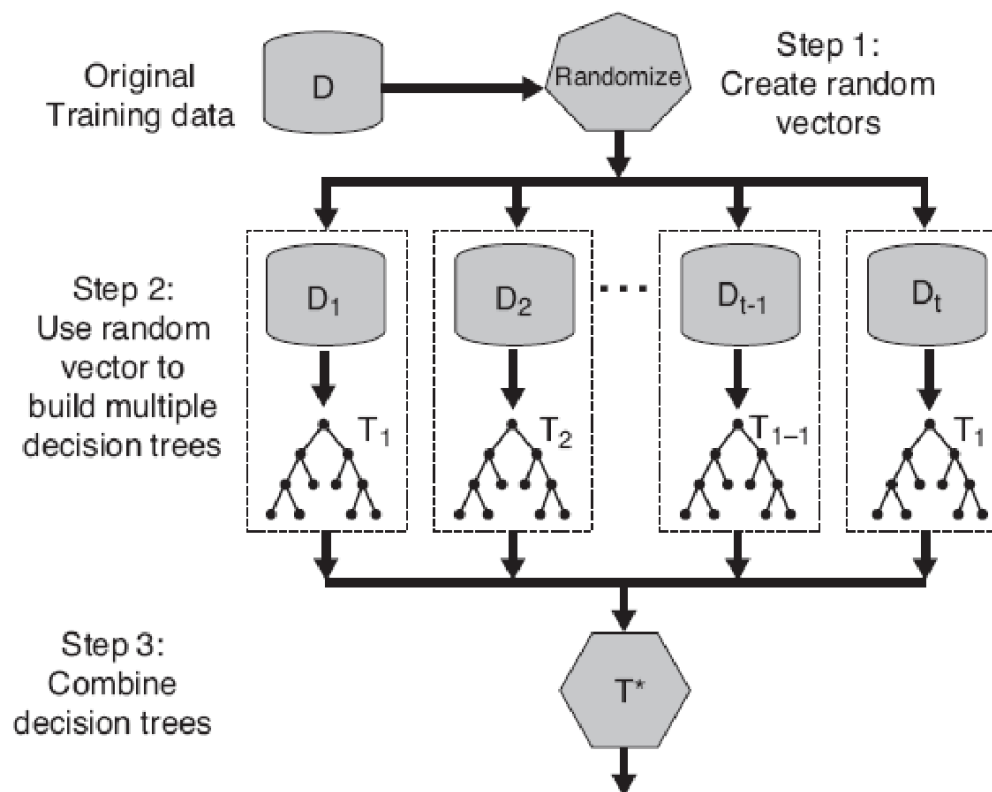
# Random Forests



Figure 5.40. Random forests.

# Random Forests Algorithm

1. For $b = 1$ to $B$:

    (a) Draw a bootstrap sample $\mathbf{Z}^*$ of size $N$ from the training data.

    (b) Grow a random-forest tree $T_b$ to the bootstrapped data, by recursively repeating the following steps for each terminal node of the tree, until the minimum node size $n_{min}$ is reached.

        i. Select $m$ variables at random from the $p$ variables.
        ii. Pick the best variable/split-point among the $m$.
        iii. Split the node into two daughter nodes.

2. Output the ensemble of trees $\{T_b\}_1^B$.
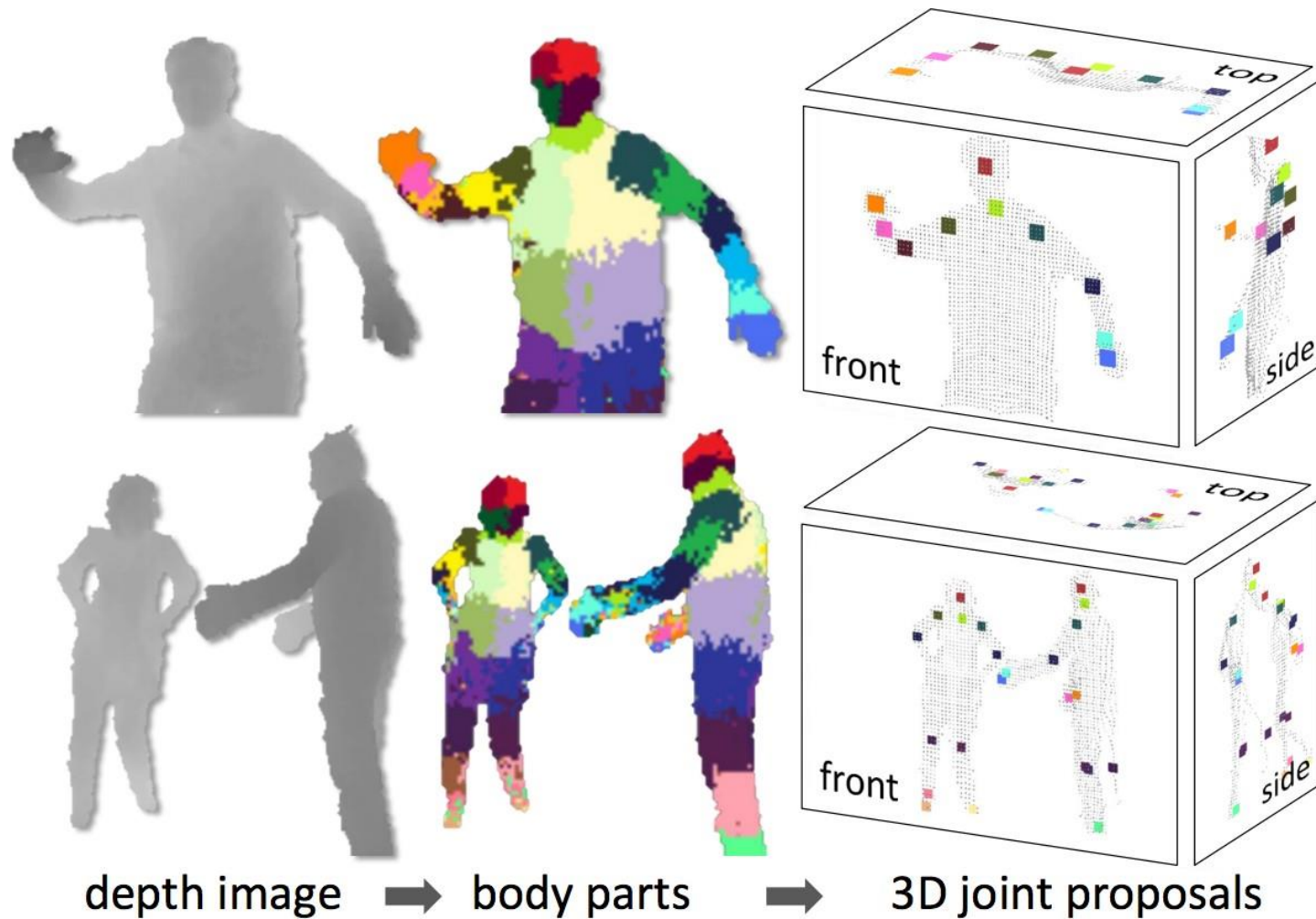
To make a prediction at a new point $x$:

*Regression:* $\hat{f}_{rf}^B(x) = \frac{1}{B} \sum_{b=1}^{B} T_b(x)$.

*Classification:* Let $\hat{C}_b(x)$ be the class prediction of the $b$th random-forest tree. Then $\hat{C}_{rf}^B(x) = majority\ vote\ \{\hat{C}_b(x)\}_1^B$.
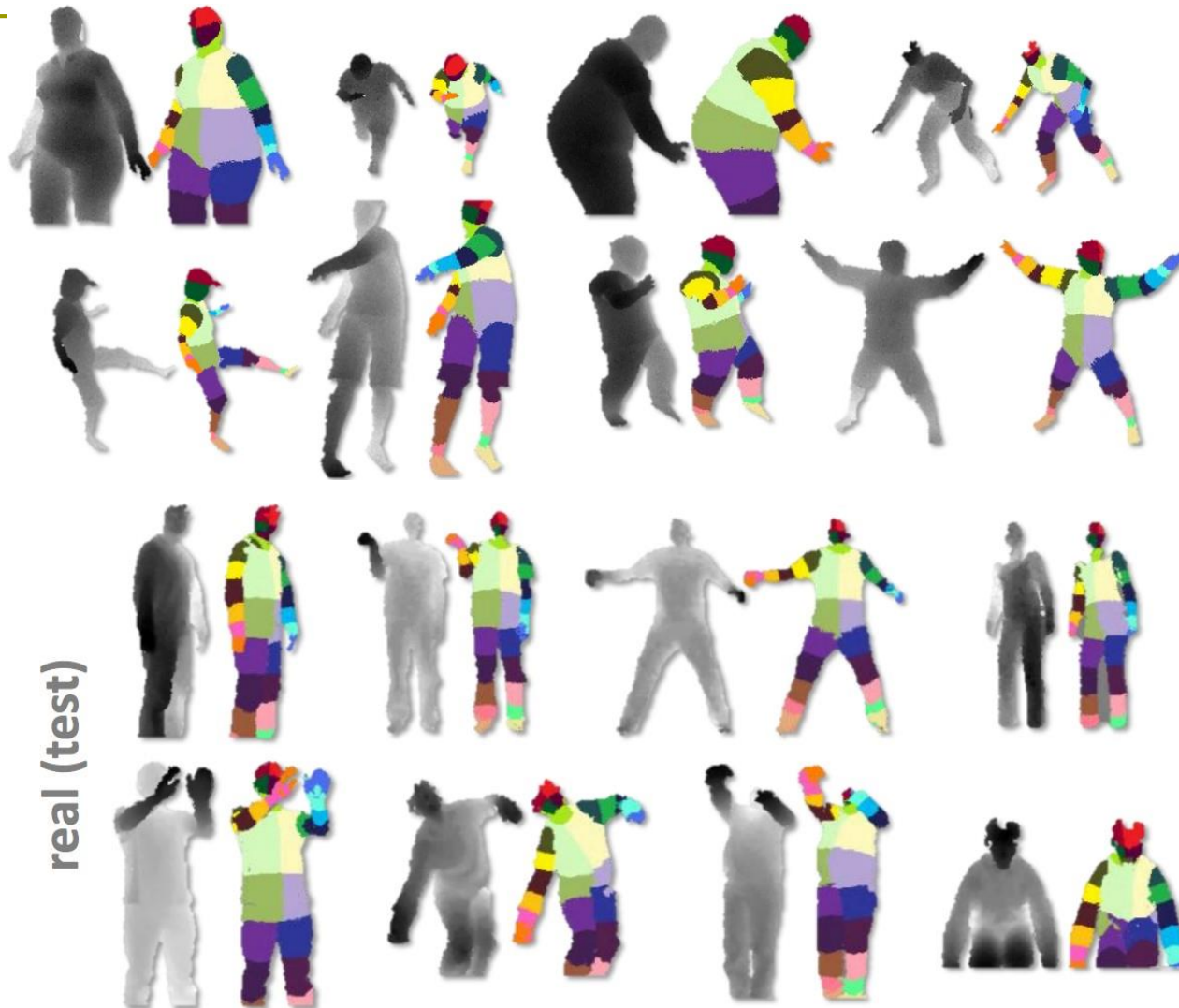
- Decision trees are in XBox



[J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, A. Blake. Real-Time Human Pose

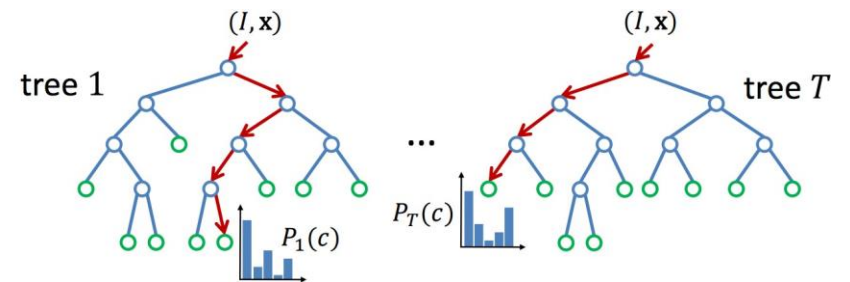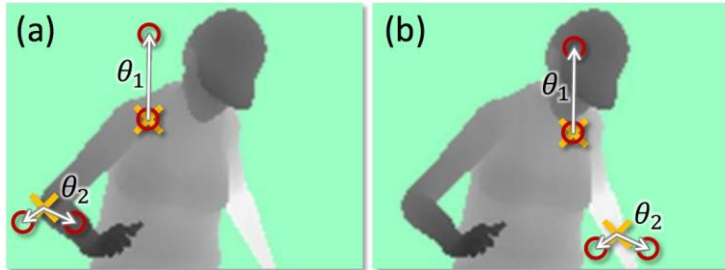# Decision trees are in XBox: Classifying body parts
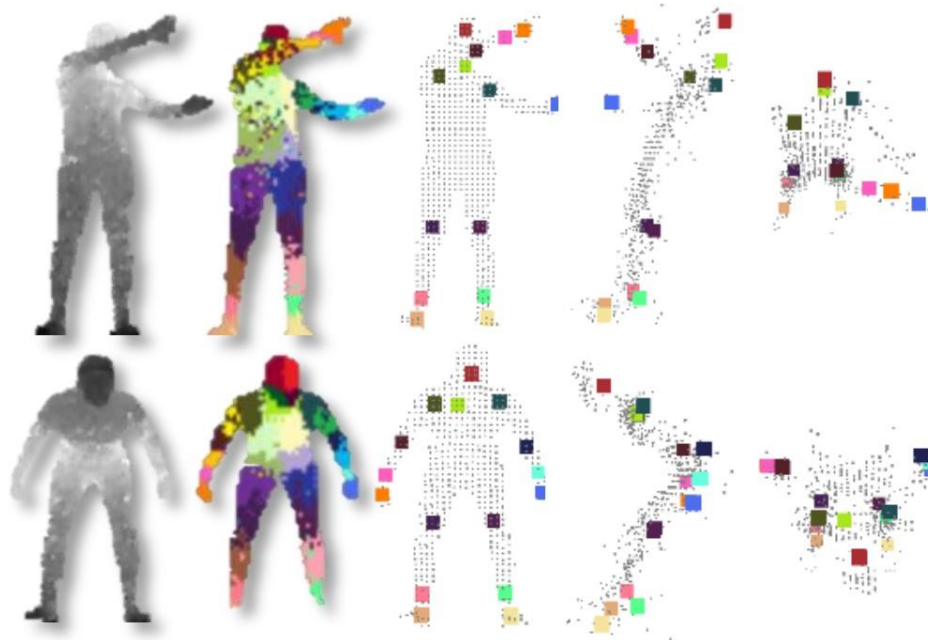


depth image ➡ body parts ➡ 3D joint proposals

- Trained on million(s) of examples

- Trained on million(s) of examples



- Results:

# Applications of Decision Trees

- Can express any Boolean function, but most useful when function depends critically on few attributes

- Bad on: parity, majority functions; also not well-suited to continuous attributes

- Practical Applications:

    - Flight simulator: 20 state variables; 90K examples based on expert pilot's actions; auto-pilot tree
    - Yahoo Ranking Challenge
    - Random Forests: Microsoft Kinect Pose Estimation