



# Data Mining: Classification2 (Chapter 4&5)

郑子彬 副教授

中山大学 数据科学与计算机学院

<http://dm16.github.io>

<http://www.inpluslab.com>

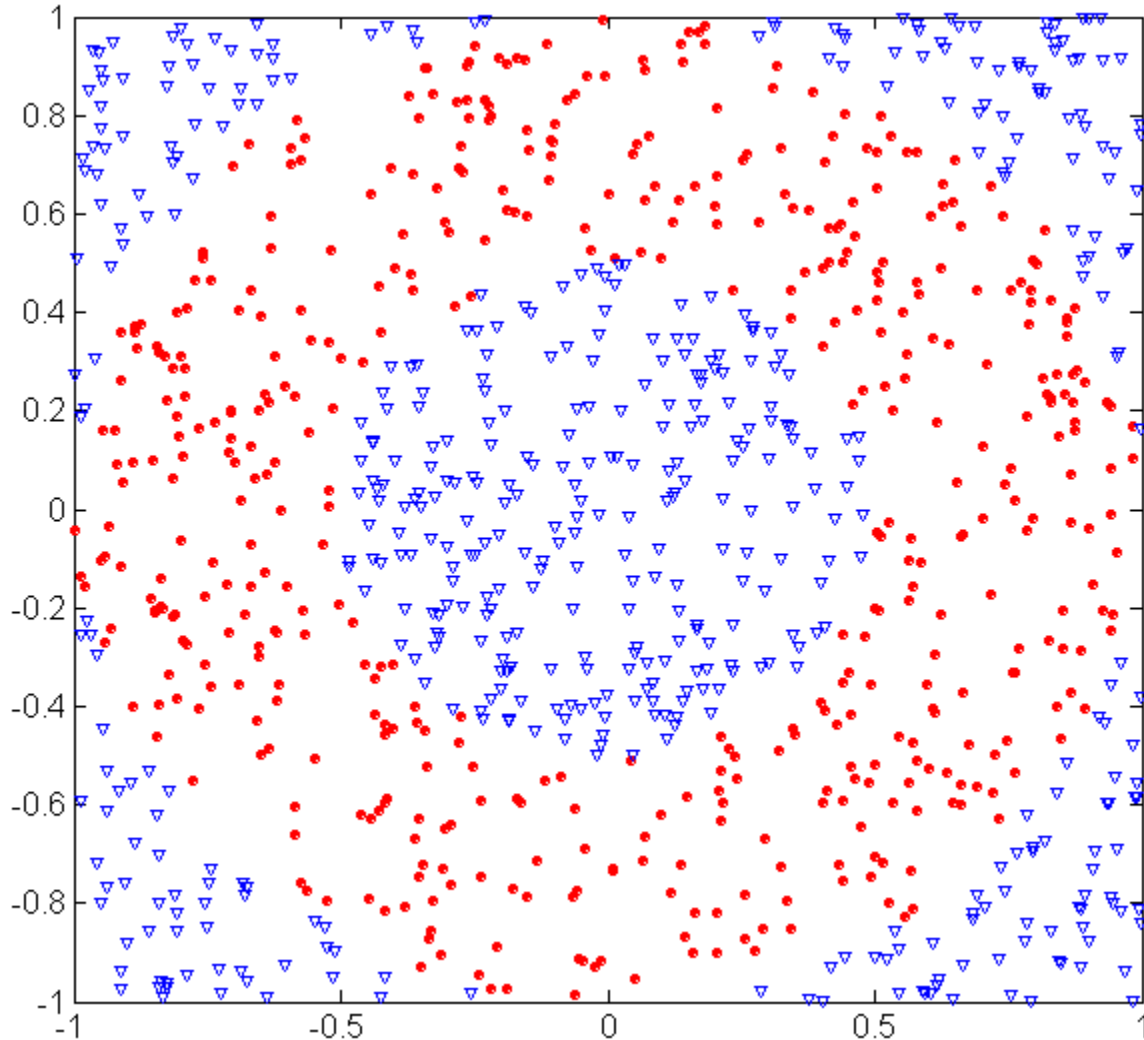


# Practical Issues of Classification



- Underfitting and Overfitting
- Missing Values
- Costs of Classification

# Underfitting and Overfitting (Example)



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

**Circular points:**

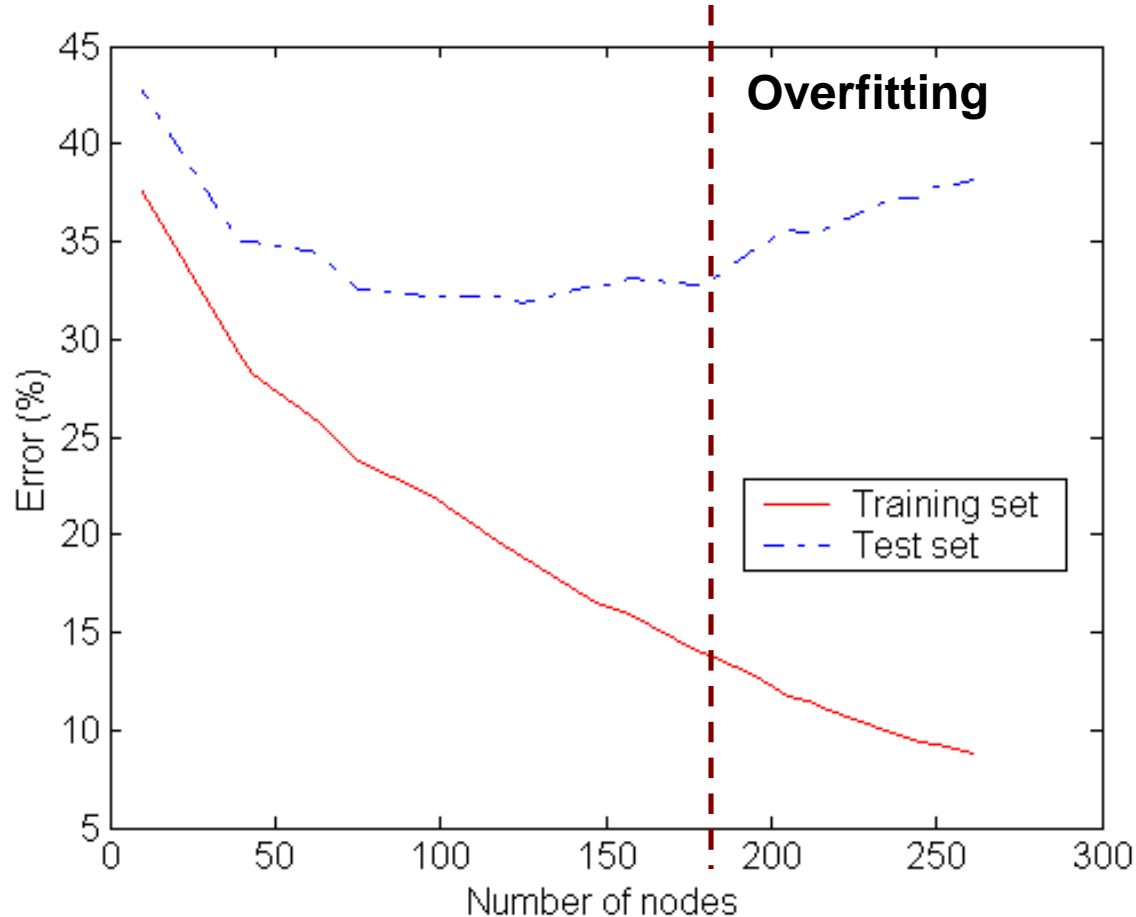
$$0.5 \leq \sqrt{x_1^2 + x_2^2} \leq 1$$

**Triangular points:**

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

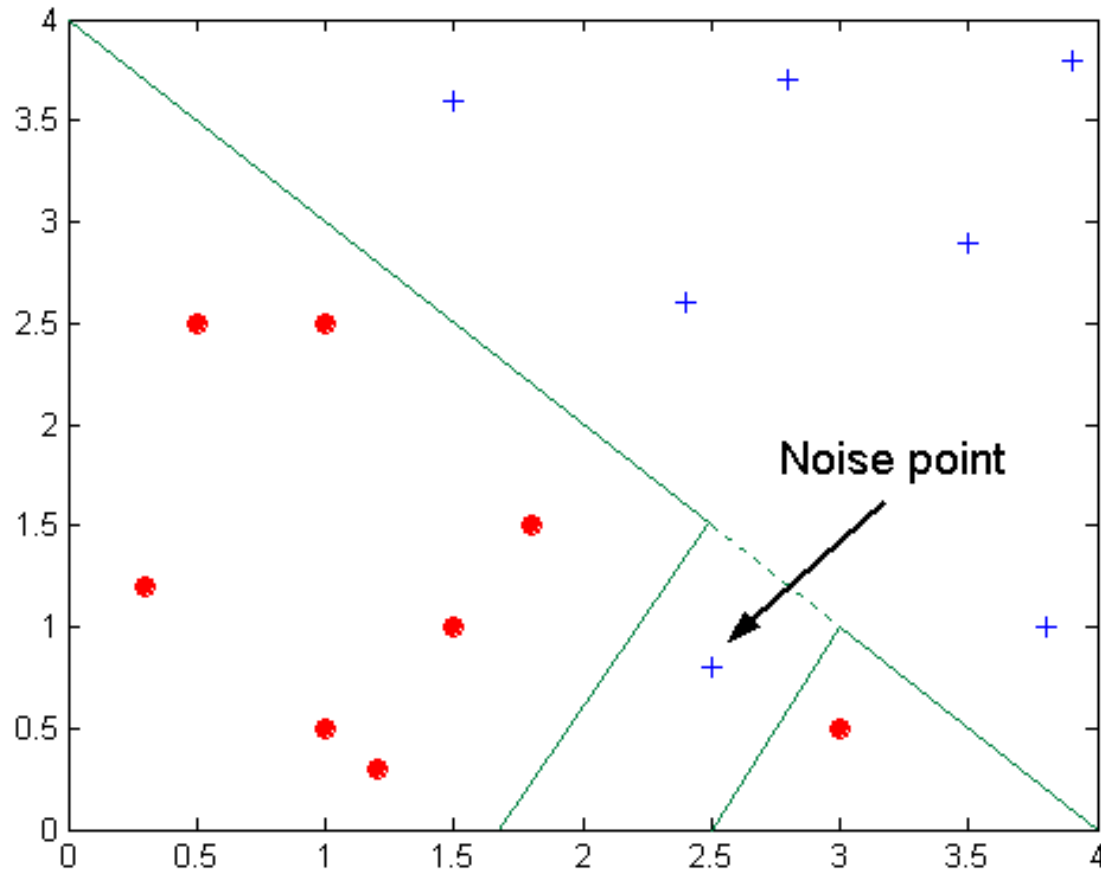
$$\sqrt{x_1^2 + x_2^2} < 1$$

# Underfitting and Overfitting



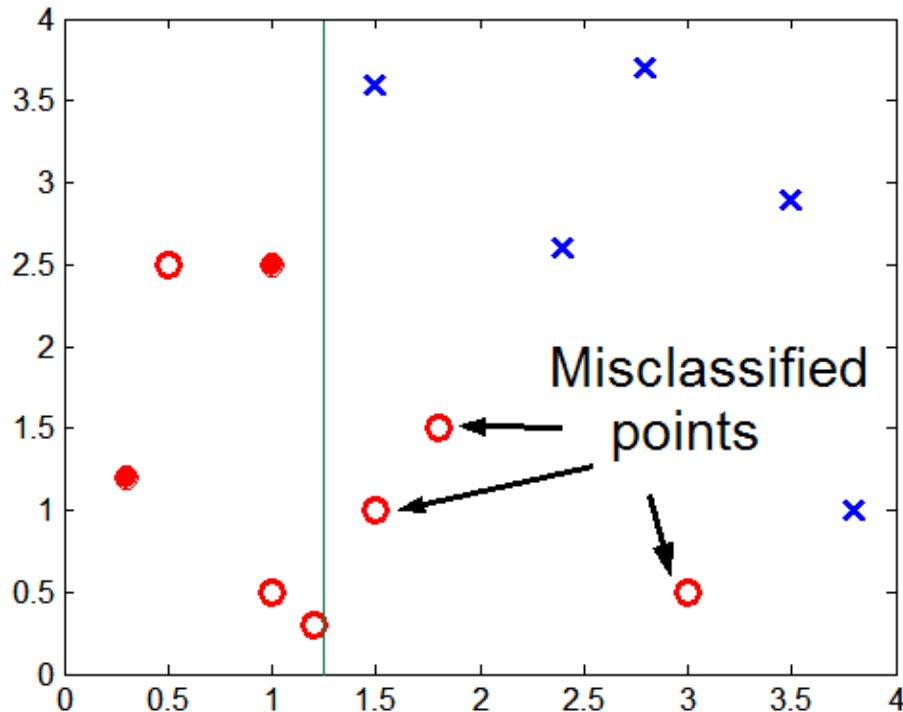
**Underfitting:** when model is too simple, both training and test errors are large

# 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**

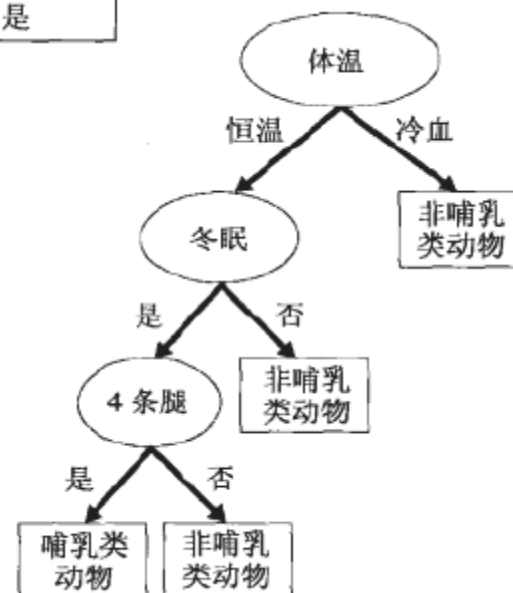
- Insufficient number of training records in the region causes the decision tree to predict the test examples using other training records that are irrelevant to the classification task**

# Example



- 人、大象、海豚被误分类
- 训练误差为0，检验误差为30%

名称	体温	胎生	4 条腿	冬眠	类标号
蜥蜴	冷血	否	是	是	否
虹鳟	冷血	是	否	否	否
鹰	恒温	否	否	否	否
弱夜鹰	恒温	否	否	是	否
鸭嘴兽	恒温	否	是	是	是



# Multiple Comparison Procedure



- 预测10个股票交易日的升降，每日的正确率50%  
正确的预测8次及以上的概率

$$\frac{C_{10}^8 + C_{10}^9 + C_{10}^{10}}{2^{10}} = 0.0547$$

- 50个股票分析家中选择一位

$$1 - (1 - 0.0547)^{50} = 0.9399$$

没人能准确预测 $\geq 8$ 次的概率

- 决策树选择变量
- 赌球预测的例子



# 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



- **Training error:** error on training ( $\sum e(t)$ )
- **Generalization errors:** error on testing ( $\sum e'(t)$ )
- Methods for estimating generalization errors:
  - **Optimistic approach:**  $e'(t) = e(t)$

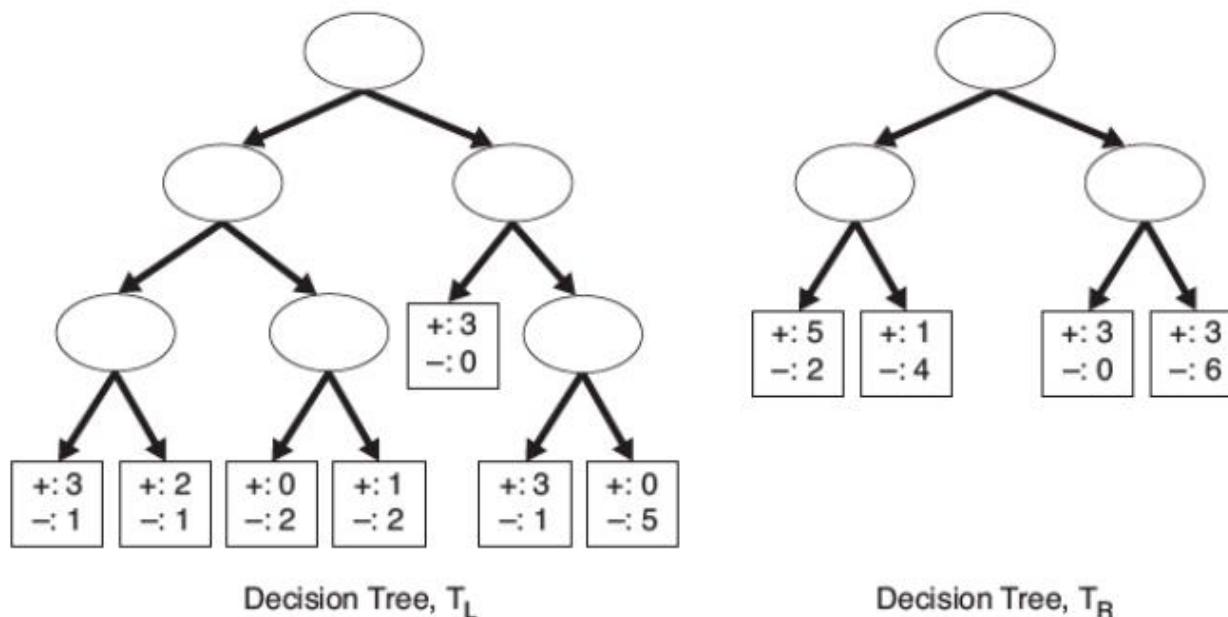


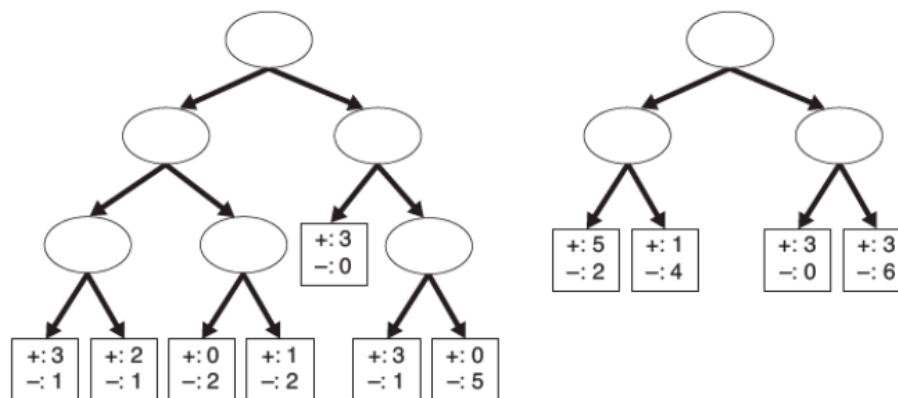
Figure 4.27. Example of two decision trees generated from the same training data.

# Estimating Generalization Errors



- **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 =  $(10 + 30 \times 0.5)/1000 = 2.5\%$



$$e_g(T_L) = \frac{4 + 7 \times 0.5}{24} = \frac{7.5}{24} = 0.3125$$

$$e_g(T_R) = \frac{6 + 4 \times 0.5}{24} = \frac{8}{24} = 0.3333$$

# Occam's Razor (奥卡姆剃刀)

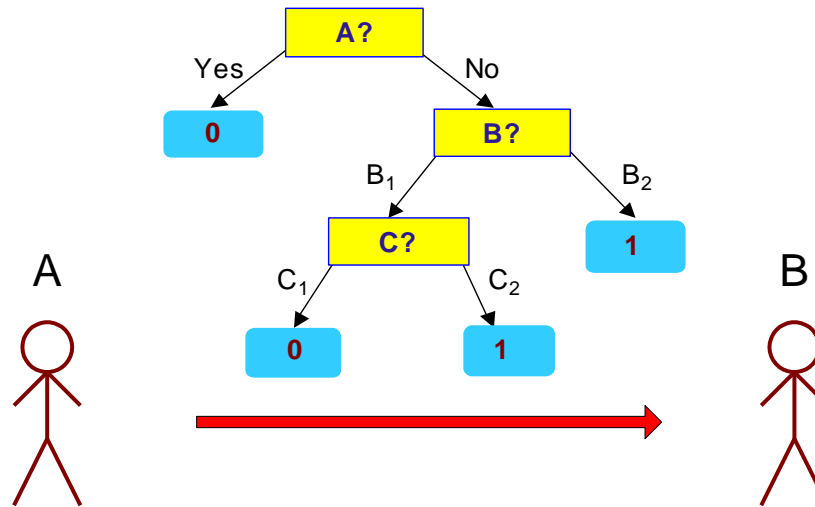


- Given two models of similar generalization errors, one should prefer the simpler model over the more complex model
- For complex models, there is a greater chance that it was fitted accidentally by errors in data
- Therefore, one should include model complexity when evaluating a model

# Minimum Description Length (MDL)



X	y
X <sub>1</sub>	1
X <sub>2</sub>	0
X <sub>3</sub>	0
X <sub>4</sub>	1
...	...
X <sub>n</sub>	1



X	y
X <sub>1</sub>	?
X <sub>2</sub>	?
X <sub>3</sub>	?
X <sub>4</sub>	?
...	...
X <sub>n</sub>	?

- $\text{Cost}(\text{Model}, \text{Data}) = \text{Cost}(\text{Data} | \text{Model}) + \text{Cost}(\text{Model})$ 
  - Cost is the number of bits needed for encoding.
  - Search for the least costly model.
- $\text{Cost}(\text{Data} | \text{Model})$  encodes the misclassification errors.
- $\text{Cost}(\text{Model})$  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).
  - Drawback:
    - ◆ Hard to determine the threshold
    - ◆ May have better gain in the following split

# 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
  - More accurate
  - Drawback:
    - ◆ Computational complexity

# Example of Post-Pruning



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

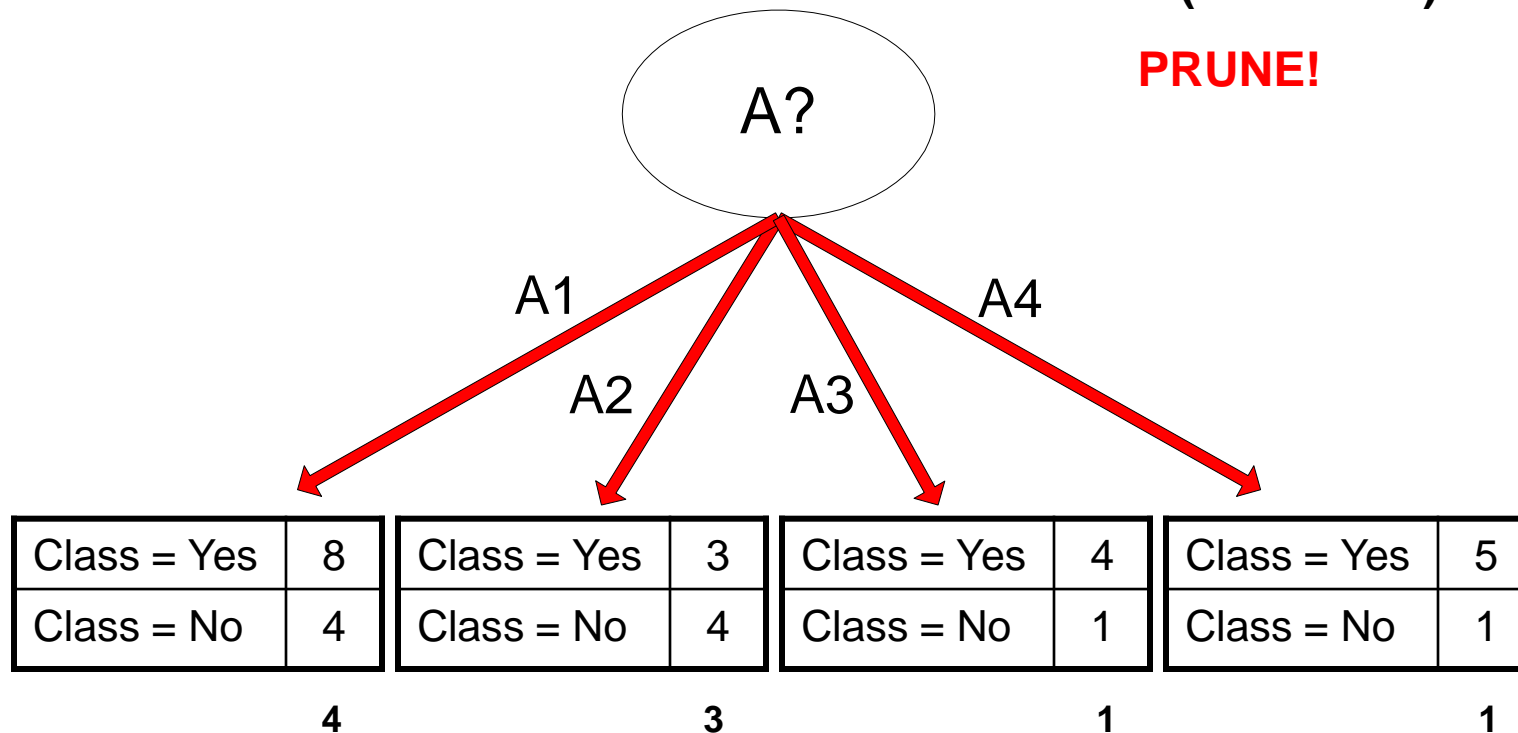
Training Error (Before splitting) = 10/30

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!**





# Cross Validation



- Divide the dataset into training set and testing set
  - Large training set: testing may not accurate
  - Large testing set: model may not accurate
- 2-fold cross validation
- K-fold cross validation
- Leave-one-out approach ( $k=N$ )
  - Computation intensive
  - Large variance on the testing

How to address?

# 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)

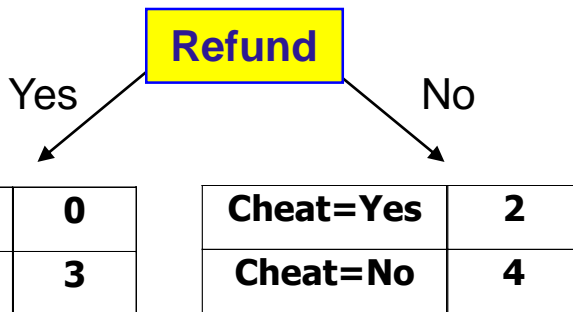
$$= 0.3 (0) + 0.6 (0.9183) = 0.551$$

$$\text{Gain} = 0.9 \times (0.8813 - 0.551) = 0.3303$$

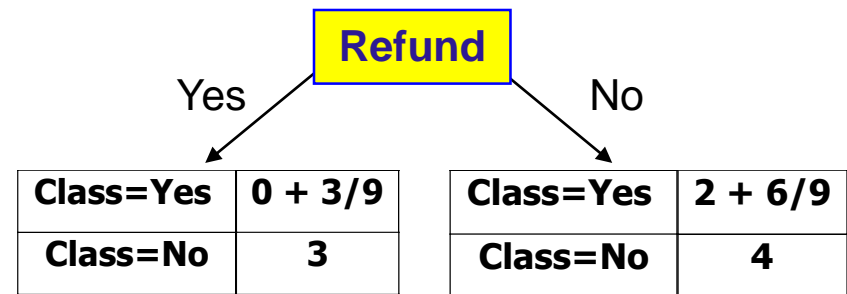
# Distribute Instances



<i>Tid</i>	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



<i>Tid</i>	Refund	Marital Status	Taxable Income	Class
10	?	Single	90K	Yes



Probability that Refund=Yes is 3/9

Probability that Refund=No is 6/9

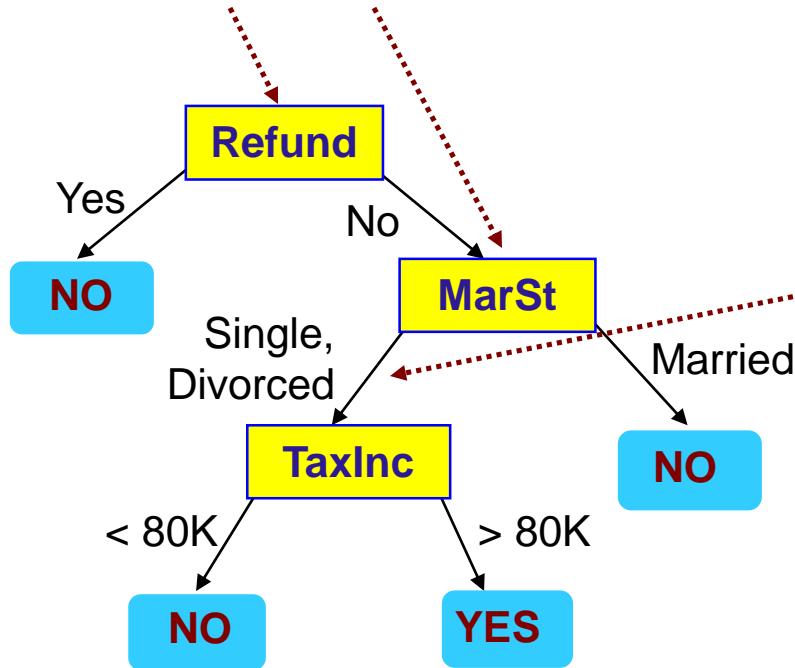
Assign record to the left child with weight = 3/9 and to the right child with weight = 6/9

# 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

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

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

# Metrics for Performance Evaluation



- Focus on the predictive capability of a model
  - Rather than how fast it takes to classify or build models, scalability, etc.
- Confusion Matrix:

ACTUAL CLASS	PREDICTED CLASS		
		Class=Yes	Class=No
	Class=Yes	a	b
	Class=No	c	d

a: TP (true positive)  
b: FN (false negative)  
c: FP (false positive)  
d: TN (true negative)

# Metrics for Performance Evaluation...



ACTUAL CLASS	PREDICTED CLASS	
	Class=Yes	Class=No
Class=Yes	a (TP)	b (FN)
	c (FP)	d (TN)

- Most widely-used metric:

$$\text{Accuracy} = \frac{a + d}{a + b + c + d} = \frac{TP + TN}{TP + TN + FP + FN}$$

# Limitation of Accuracy



- What is the limitation of accuracy?
- Example: Paper acceptance prediction model
  - Author name, Author number, hometown
  - University
  - Title length
  - Number of Reference
  - Number of Figure and table



# Limitation of Accuracy



- Consider a 2-class problem
  - Number of Class 0 examples = 9990
  - Number of Class 1 examples = 10
- If model predicts everything to be class 0, accuracy is  $9990/10000 = 99.9\%$ 
  - Accuracy is misleading because model does not detect any class 1 example

# Cost Matrix



	PREDICTED CLASS		
	$C(i j)$	Class=Yes	Class=No
	Class=Yes	$C(\text{Yes} \text{Yes})$	$C(\text{No} \text{Yes})$
	Class=No	$C(\text{Yes} \text{No})$	$C(\text{No} \text{No})$

$C(i|j)$ : Cost of misclassifying class  $j$  example as class  $i$

# Computing Cost of Classification



Cost Matrix	PREDICTED CLASS		
	C(i j)	+	-
	+	-1	100
	-	1	0

Model $M_1$	PREDICTED CLASS		
ACTUAL CLASS		+	-
	+	150	40
	-	60	250

Accuracy = 80%

Cost = 3910

Model $M_2$	PREDICTED CLASS		
ACTUAL CLASS		+	-
	+	250	45
	-	5	200

Accuracy = 90%

Cost = 4255

# Cost-Sensitive Measures



$$\text{Precision (p)} = \frac{a}{a + c}$$

准确率：所有的预测为正类的结果里面，有多少比例是准确的？  
搜索返回的结果页面中，有多少比例是相关的？

$$\text{Recall (r)} = \frac{a}{a + b}$$

召回率：所有实际中的正类，有多少比例被预测准确？  
[查全率] 所有相关的页面中，有多少比例被返回给用户？

相互制约，搜索引擎严格一些：准确率高，查全率低；  
搜索引擎宽松一些：准确率低，查全率高。

$$\text{F - measure (F)} = \frac{2rp}{r + p} = \frac{2a}{2a + b + c}$$

	PREDICTED CLASS	
	Class=Yes	Class=No
	Class=Yes	Class=No
ACTUAL CLASS	a	b
	c	d

# ROC (Receiver Operating Characteristic)

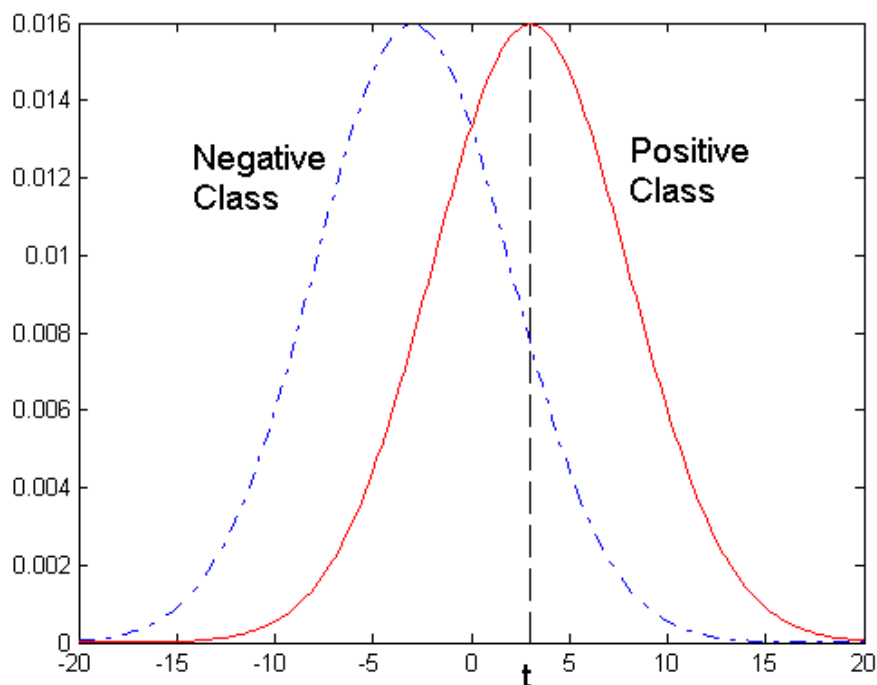


- Developed in 1950s for signal detection theory to analyze noisy signals
  - Characterize the trade-off between positive hits and false alarms
- ROC curve plots TPR (on the y-axis) against FPR (on the x-axis)
  - TP rate,  $TPR = TP / (TP + FN)$  所有实际正类里面，有多少比例被预测准确？等同于Recall
  - FP rate,  $FPR = FP / (FP + TN)$  所有的实际的负类里面，有多少比例被预测成正样本？误报，虚惊
- Performance of each classifier represented as a point on the ROC curve

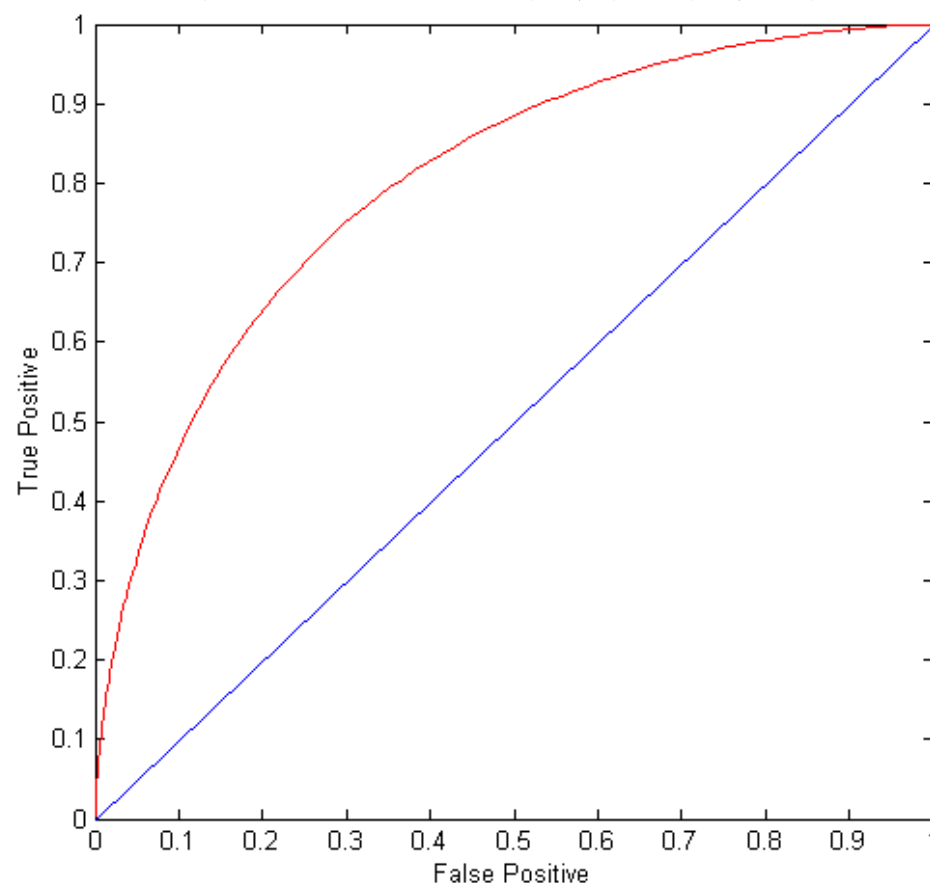
# ROC Curve



- 1-dimensional data set containing 2 classes (positive and negative)
- any points located at  $x > t$  is classified as positive



击中率：被预测为正的正样本数/正样本总数



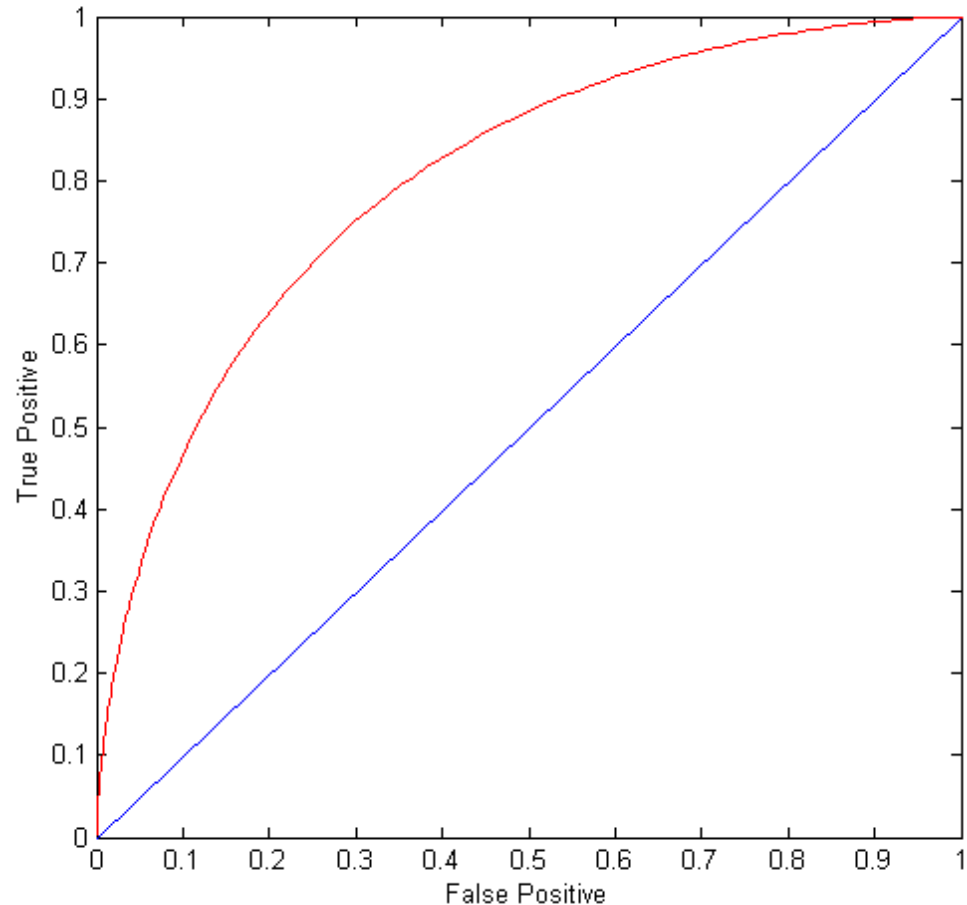
虚惊率：被预测为正的负样本数/负样本总数

# ROC Curve

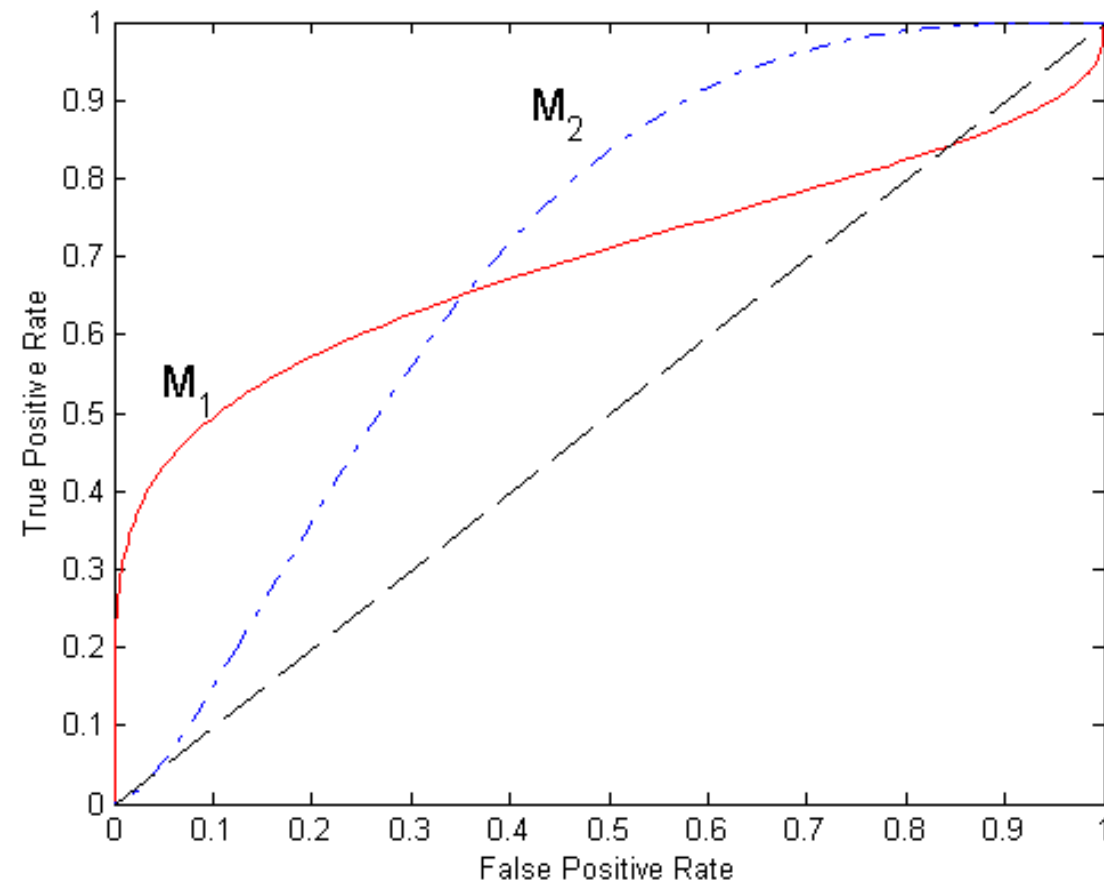


(TP,FP):

- (0,0): declare everything to be negative class
- (1,1): declare everything to be positive class
- (1,0): ideal
- Diagonal line:
  - Random guessing



# Using ROC for Model Comparison



- No model consistently outperform the other
  - $M_1$  is better for small FPR
  - $M_2$  is better for large FPR
- Area Under the ROC curve (**AUC**)
  - Ideal:
    - Area = 1
  - Random guess:
    - Area = 0.5



# How to Construct an ROC curve



Instance	$P(+ A)$	True Class
1	0.95	+
2	0.93	+
3	0.87	-
4	0.85	-
5	0.85	-
6	0.85	+
7	0.76	-
8	0.53	+
9	0.43	-
10	0.25	+

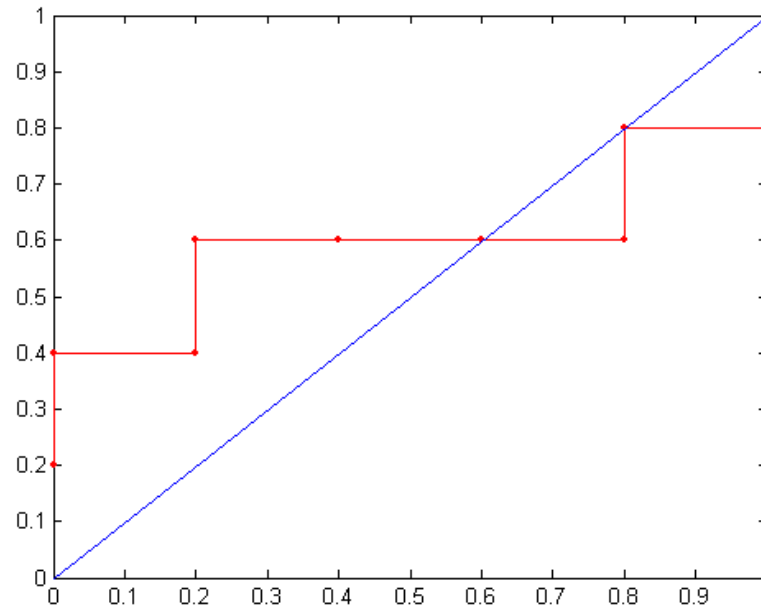
- Use classifier that produces posterior probability for each test instance  $P(+|A)$
- Sort the instances according to  $P(+|A)$  in decreasing order
- Apply threshold at each unique value of  $P(+|A)$
- Count the number of TP, FP, TN, FN at each threshold
- TP rate,  $TPR = TP/(TP+FN)$
- FP rate,  $FPR = FP/(FP + TN)$

# How to construct an ROC curve



Class	+	-	+	-	-	-	+	-	+	+	
Threshold $\geq$	0.25	0.43	0.53	0.76	0.85	0.85	0.85	0.87	0.93	0.95	1.00
TP	5	4	4	3	3	3	3	2	2	1	0
FP	5	5	4	4	3	2	1	1	0	0	0
TN	0	0	1	1	2	3	4	4	5	5	5
FN	0	1	1	2	2	2	2	3	3	4	5
TPR	1	0.8	0.8	0.6	0.6	0.6	0.6	0.4	0.4	0.2	0
FPR	1	1	0.8	0.8	0.6	0.4	0.2	0.2	0	0	0

ROC Curve:





# Chapter 5

# Rule-Based Classifier



- Classify records by using a collection of “if...then...” rules
- Rule:  $(Condition) \rightarrow y$ 
  - where
    - ◆ *Condition* is a conjunctions of attributes
    - ◆  $y$  is the class label
  - *LHS*: rule antecedent or condition
  - *RHS*: rule consequent
  - Examples of classification rules:
    - ◆  $(\text{Blood Type}=\text{Warm}) \wedge (\text{Lay Eggs}=\text{Yes}) \rightarrow \text{Birds}$
    - ◆  $(\text{Taxable Income} < 50\text{K}) \wedge (\text{Refund}=\text{Yes}) \rightarrow \text{Evade}=\text{No}$

# Rule-based Classifier (Example)



Name	Blood Type	Give Birth	Can Fly	Live in Water	Class
human	warm	yes	no	no	mammals
python	cold	no	no	no	reptiles
salmon	cold	no	no	yes	fishes
whale	warm	yes	no	yes	mammals
frog	cold	no	no	sometimes	amphibians
komodo	cold	no	no	no	reptiles
bat	warm	yes	yes	no	mammals
pigeon	warm	no	yes	no	birds
cat	warm	yes	no	no	mammals
leopard shark	cold	yes	no	yes	fishes
turtle	cold	no	no	sometimes	reptiles
penguin	warm	no	no	sometimes	birds
porcupine	warm	yes	no	no	mammals
eel	cold	no	no	yes	fishes
salamander	cold	no	no	sometimes	amphibians
gila monster	cold	no	no	no	reptiles
platypus	warm	no	no	no	mammals
owl	warm	no	yes	no	birds
dolphin	warm	yes	no	yes	mammals
eagle	warm	no	yes	no	birds

R1: (Give Birth = no)  $\wedge$  (Can Fly = yes)  $\rightarrow$  Birds

R2: (Give Birth = no)  $\wedge$  (Live in Water = yes)  $\rightarrow$  Fishes

R3: (Give Birth = yes)  $\wedge$  (Blood Type = warm)  $\rightarrow$  Mammals

R4: (Give Birth = no)  $\wedge$  (Can Fly = no)  $\rightarrow$  Reptiles

R5: (Live in Water = sometimes)  $\rightarrow$  Amphibians

# Application of Rule-Based Classifier



- A rule  $r$  **covers** an instance  $x$  if the attributes of the instance satisfy the condition of the rule

R1: (Give Birth = no)  $\wedge$  (Can Fly = yes)  $\rightarrow$  Birds

R2: (Give Birth = no)  $\wedge$  (Live in Water = yes)  $\rightarrow$  Fishes

R3: (Give Birth = yes)  $\wedge$  (Blood Type = warm)  $\rightarrow$  Mammals

R4: (Give Birth = no)  $\wedge$  (Can Fly = no)  $\rightarrow$  Reptiles

R5: (Live in Water = sometimes)  $\rightarrow$  Amphibians

Name	Blood Type	Give Birth	Can Fly	Live in Water	Class
hawk	warm	no	yes	no	?
grizzly bear	warm	yes	no	no	?

The rule R1 covers a hawk  $\Rightarrow$  Bird

The rule R3 covers the grizzly bear  $\Rightarrow$  Mammal

# Rule Coverage and Accuracy



- Coverage of a rule:
  - Fraction of records that satisfy the antecedent of a rule
- Accuracy of a rule:
  - Fraction of records that satisfy both the antecedent and consequent of a rule

<i>Tid</i>	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	No	Single	90K	Yes

(Status=Single) → No

Coverage = 40%, Accuracy = 50%

# How does Rule-based Classifier Work?



R1: (Give Birth = no)  $\wedge$  (Can Fly = yes)  $\rightarrow$  Birds

R2: (Give Birth = no)  $\wedge$  (Live in Water = yes)  $\rightarrow$  Fishes

R3: (Give Birth = yes)  $\wedge$  (Blood Type = warm)  $\rightarrow$  Mammals

R4: (Give Birth = no)  $\wedge$  (Can Fly = no)  $\rightarrow$  Reptiles

R5: (Live in Water = sometimes)  $\rightarrow$  Amphibians

Name	Blood Type	Give Birth	Can Fly	Live in Water	Class
lemur	warm	yes	no	no	?
turtle	cold	no	no	sometimes	?
dogfish shark	cold	yes	no	yes	?

A lemur triggers rule R3, so it is classified as a mammal

A turtle triggers both R4 and R5

A dogfish shark triggers none of the rules

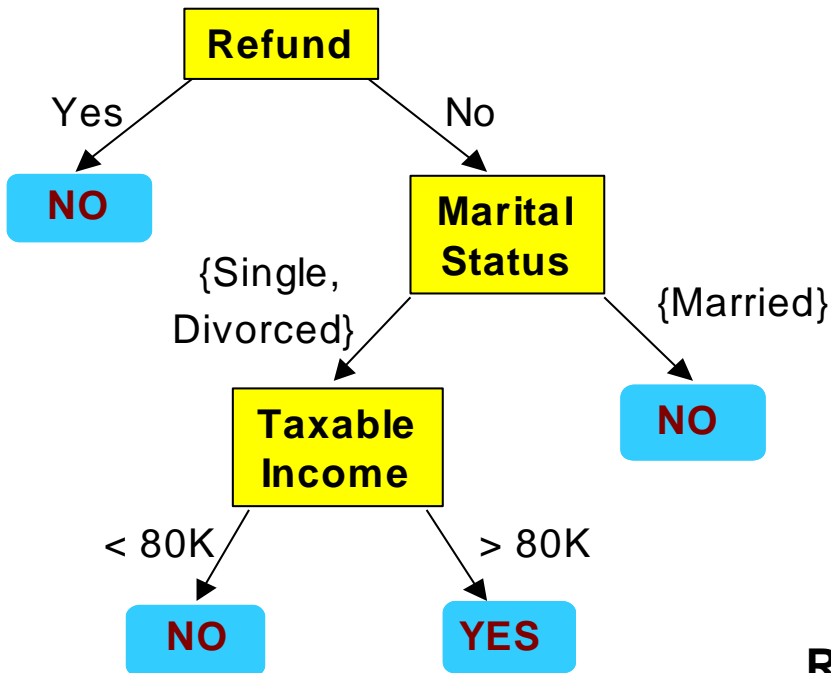


# Characteristics of Rule-Based Classifier



- Mutually exclusive rules （互斥规则）
  - Classifier contains mutually exclusive rules if the rules are independent of each other
  - Every record is covered by at most one rule
- Exhaustive rules （穷举规则）
  - Classifier has exhaustive coverage if it accounts for every possible combination of attribute values
  - Each record is covered by at least one rule

# From Decision Trees To Rules



## Classification Rules

(Refund=Yes) ==> No

(Refund=No, Marital Status={Single, Divorced}, Taxable Income<80K) ==> No

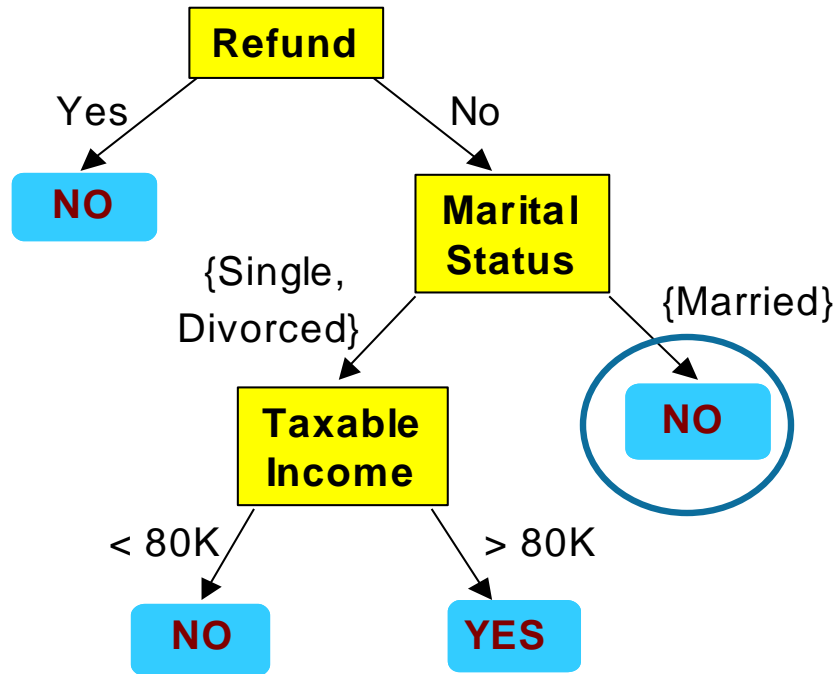
(Refund=No, Marital Status={Single, Divorced}, Taxable Income>80K) ==> Yes

(Refund=No, Marital Status={Married}) ==> No

Rules are mutually exclusive and exhaustive

Rule set contains as much information as the tree

# Rules Can Be Simplified



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

Initial Rule:  $(\text{Refund}=\text{No}) \wedge (\text{Status}=\text{Married}) \rightarrow \text{No}$

Simplified Rule:  $(\text{Status}=\text{Married}) \rightarrow \text{No}$

# Effect of Rule Simplification



- Rules are no longer mutually exclusive
  - A record may trigger more than one rule
  - Solution?
    - ◆ Ordered rule set
    - ◆ Unordered rule set – use voting schemes
    - ◆ 基于规则的准确性加权
- Rules are no longer exhaustive
  - A record may not trigger any rules
  - Solution?
    - ◆ Use a default class

# Ordered Rule Set



- Rules are rank ordered according to their priority
  - An ordered rule set is known as a decision list
- When a test record is presented to the classifier
  - It is assigned to the class label of the highest ranked rule it has triggered
  - If none of the rules fired, it is assigned to the default class

R1: (Give Birth = no)  $\wedge$  (Can Fly = yes)  $\rightarrow$  Birds

R2: (Give Birth = no)  $\wedge$  (Live in Water = yes)  $\rightarrow$  Fishes

R3: (Give Birth = yes)  $\wedge$  (Blood Type = warm)  $\rightarrow$  Mammals

R4: (Give Birth = no)  $\wedge$  (Can Fly = no)  $\rightarrow$  Reptiles

R5: (Live in Water = sometimes)  $\rightarrow$  Amphibians

Name	Blood Type	Give Birth	Can Fly	Live in Water	Class
turtle	cold	no	no	sometimes	?

# Rule Ordering Schemes



- Rule-based ordering
  - Individual rules are ranked based on their quality
- Class-based ordering
  - Rules that belong to the same class appear together

## Rule-based Ordering

(Refund=Yes) ==> No

(Refund=No, Marital Status={Single,Divorced},  
Taxable Income<80K) ==> No

(Refund=No, Marital Status={Single,Divorced},  
Taxable Income>80K) ==> Yes

(Refund=No, Marital Status={Married}) ==> No

## Class-based Ordering

(Refund=Yes) ==> No

(Refund=No, Marital Status={Single,Divorced},  
Taxable Income<80K) ==> No

(Refund=No, Marital Status={Married}) ==> No

(Refund=No, Marital Status={Single,Divorced},  
Taxable Income>80K) ==> Yes

# Building Classification Rules



- Direct Method:
  - ◆ Extract rules directly from data
  - ◆ e.g.: RIPPER, CN2, Holte's 1R
  
- Indirect Method:
  - ◆ Extract rules from other classification models (e.g. decision trees, neural networks, etc).
  - ◆ e.g: C4.5rules

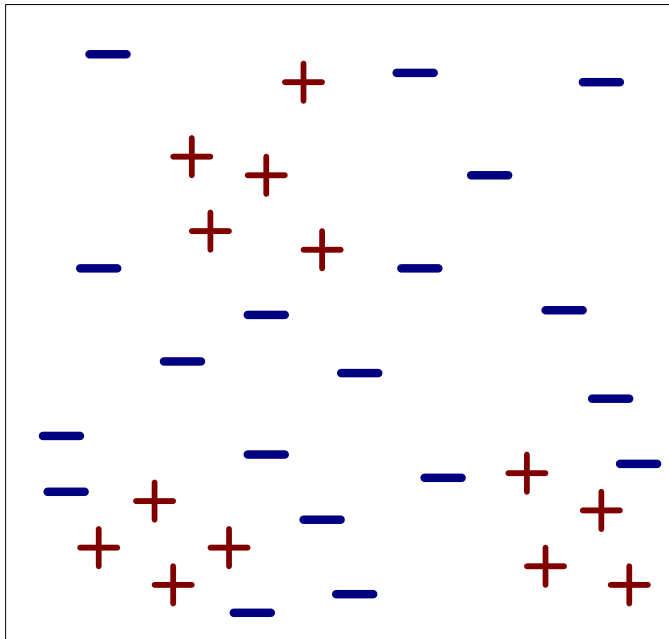
# Direct Method: Sequential Covering



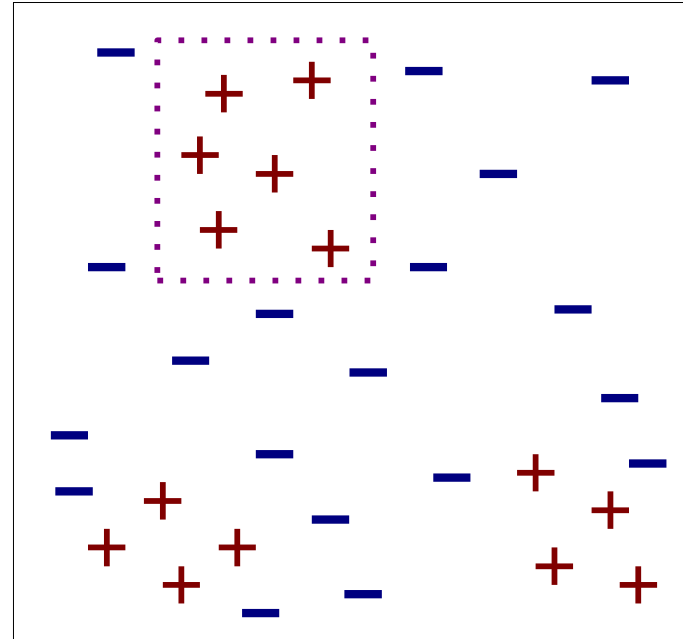
1. Start from an empty rule
2. Grow a rule using the Learn-One-Rule function
3. Remove training records covered by the rule
4. Repeat Step (2) and (3) until stopping criterion is met



# Example of Sequential Covering

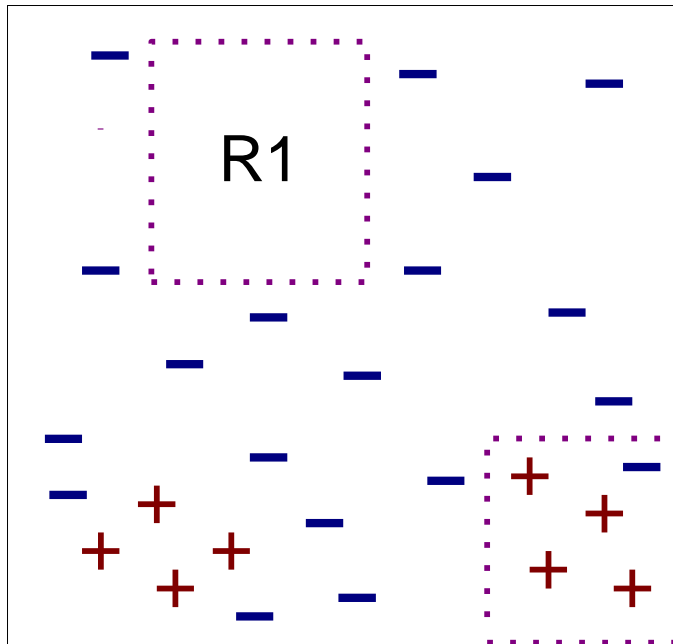


(i) Original Data

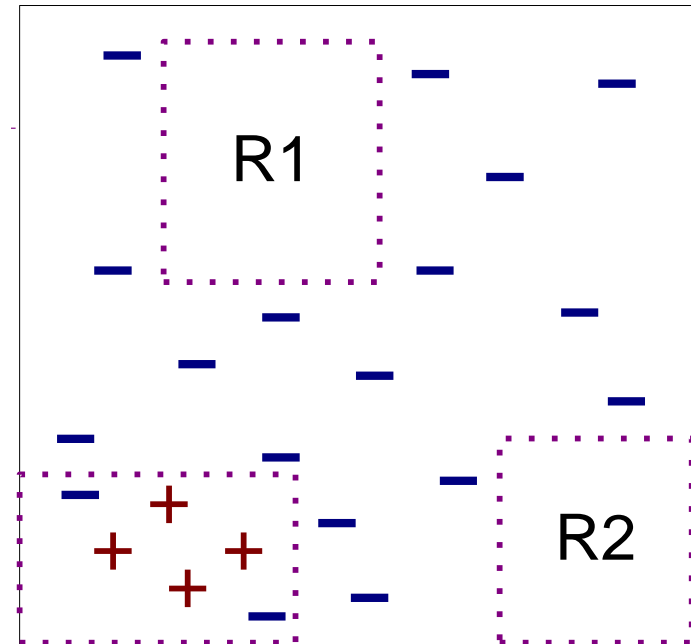


(ii) Step 1

# Example of Sequential Covering...



(iii) Step 2



(iv) Step 3

# Aspects of Sequential Covering

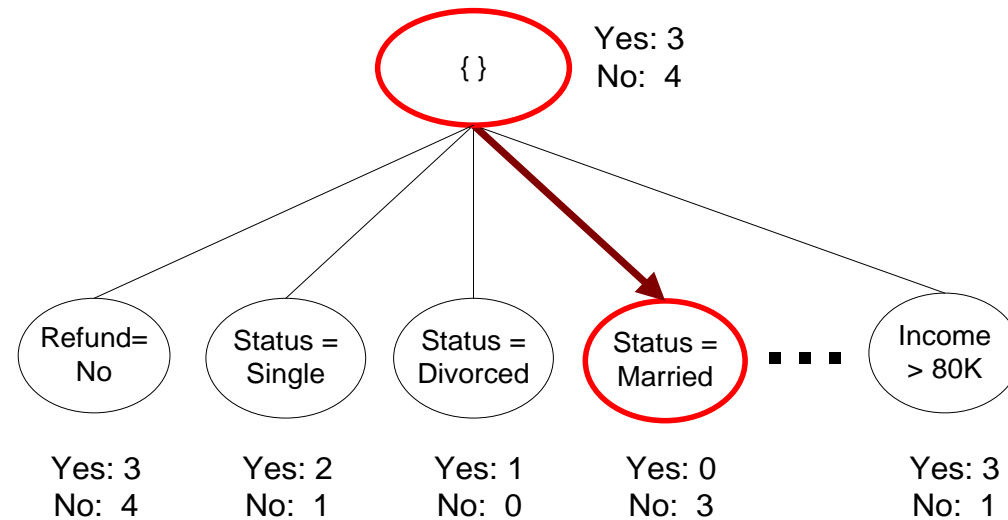


- Rule Growing
- Instance Elimination
- Rule Evaluation
- Stopping Criterion
- Rule Pruning

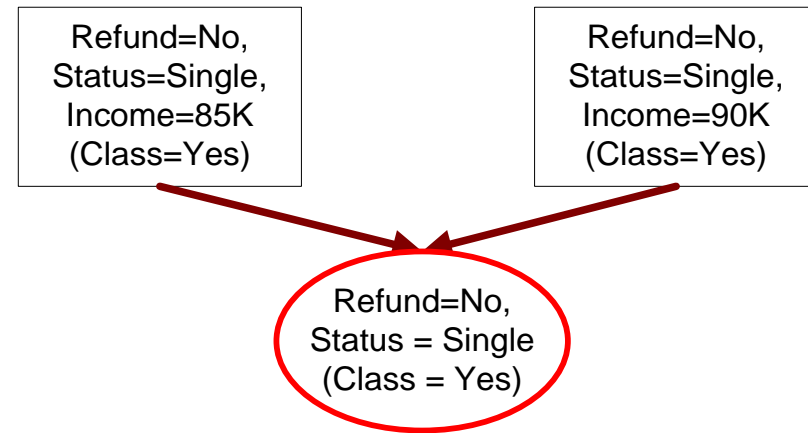
# Rule Growing



- Two common strategies



(a) General-to-specific



(b) Specific-to-general

# Rule Growing (Examples)



- CN2 Algorithm:

- Start from an empty conjunct:  $\{\}$
- Add conjuncts that minimizes the entropy measure:  $\{A\}, \{A,B\}, \dots$
- Determine the rule consequent by taking majority class of instances covered by the rule

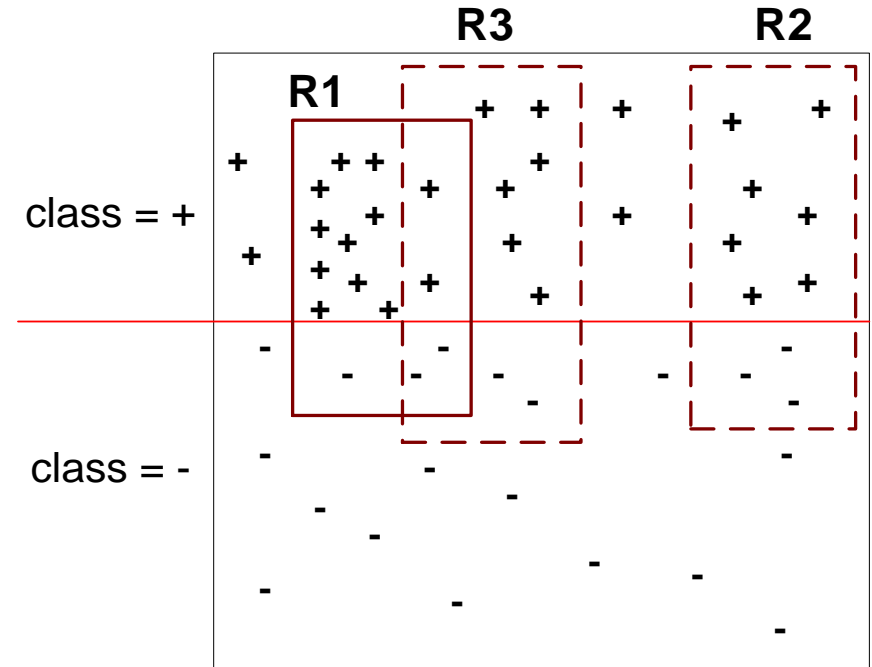
- RIPPER Algorithm:

- Start from an empty rule:  $\{\} \Rightarrow \text{class}$
  - Add conjuncts that maximizes FOIL's information gain measure:
    - ◆  $R0: \{\} \Rightarrow \text{class}$  (initial rule)
    - ◆  $R1: \{A\} \Rightarrow \text{class}$  (rule after adding conjunct)
    - ◆  $\text{Gain}(R0, R1) = p1 [ \log (p1/(p1+n1)) - \log (p0/(p0 + n0)) ]$
    - ◆ where  $t$ : number of positive instances covered by both  $R0$  and  $R1$ 
      - $p0$ : number of positive instances covered by  $R0$
      - $n0$ : number of negative instances covered by  $R0$
      - $p1$ : number of positive instances covered by  $R1$
      - $n1$ : number of negative instances covered by  $R1$
- $p1/(p1+n1)$  : accuracy越高越好**  
 **$P1$ : 个数越多越好**

# Instance Elimination



- Why do we need to eliminate instances?
  - Otherwise, the next rule is identical to previous rule
- Why do we remove positive instances?
  - Ensure that the next rule is different
- Why do we remove negative instances?
  - Prevent underestimating accuracy of rule
  - Compare rules R2 and R3 in the diagram



# Rule Evaluation



- Metrics:

- Accuracy  $= \frac{n_c}{n}$

- Laplace  $= \frac{n_c + 1}{n + k}$

- M-estimate  $= \frac{n_c + kp}{n + k}$

$n$  : Number of instances covered by rule

$n_c$  : Number of positive instances covered by rule

$k$  : Number of classes

$p$  : Prior probability

R1: 50 positive and 5 negative 50/55 vs 51/57

R2: 2 positive and 2 negative 2/2 vs 3/4

降低小coverage的数值，避免为0

**Laplace:** 每个类出现概率均等

**M-estimate:** 对不同的类区别对待，出现概率大的相应值也大

# Stopping Criterion and Rule Pruning



- Stopping criterion
  - Compute the gain
  - If gain is not significant, discard the new rule
- Rule Pruning
  - Similar to post-pruning of decision trees
  - Reduced Error Pruning:
    - ◆ Remove one of the conjuncts in the rule
    - ◆ Compare error rate on validation set before and after pruning
    - ◆ If error improves, prune the conjunct



# Summary of Direct Method



- Grow a single rule
- Remove Instances from rule
- Prune the rule (if necessary)
- Add rule to Current Rule Set
- Repeat

# Advantages of Rule-Based Classifiers



- As highly expressive as decision trees
- Easy to interpret
- Easy to generate
- Can classify new instances rapidly
- Performance comparable to decision trees

# 谢谢

<http://dm16.github.io>

<http://www.inpluslab.com>