

Decision Trees

Recap and Decision Trees for Regression

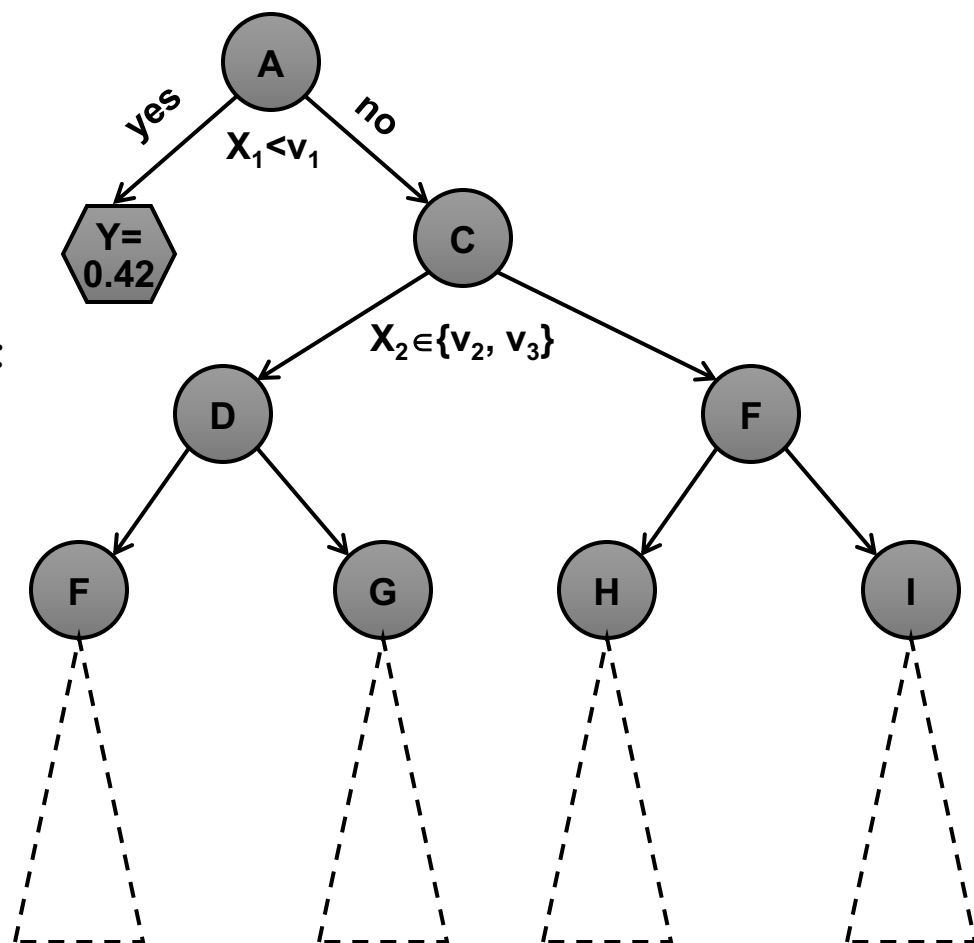
Decision Trees

► Input features:

- N features: X_1, X_2, \dots, X_N
- Each X_j has domain D_j
 - Categorical: $D_j = \{\text{red, blue}\}$
 - Numerical: $D_j = (0, 10)$
- Y is output variable with domain D_Y :
 - Categorical: Classification
 - Numerical: Regression

► Task:

- Given input data vector x_i predict y_i



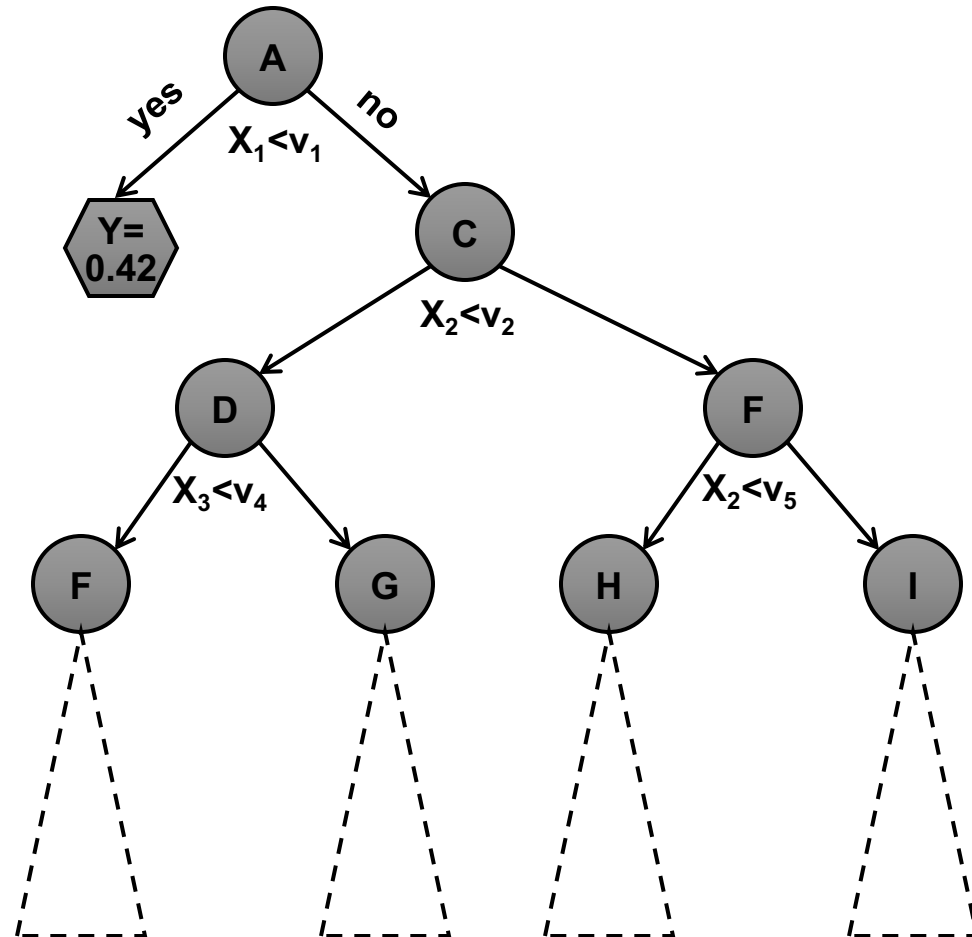
Decision Trees (I)

► Decision trees:

- Split the data at each internal node
- Each leaf node makes a prediction

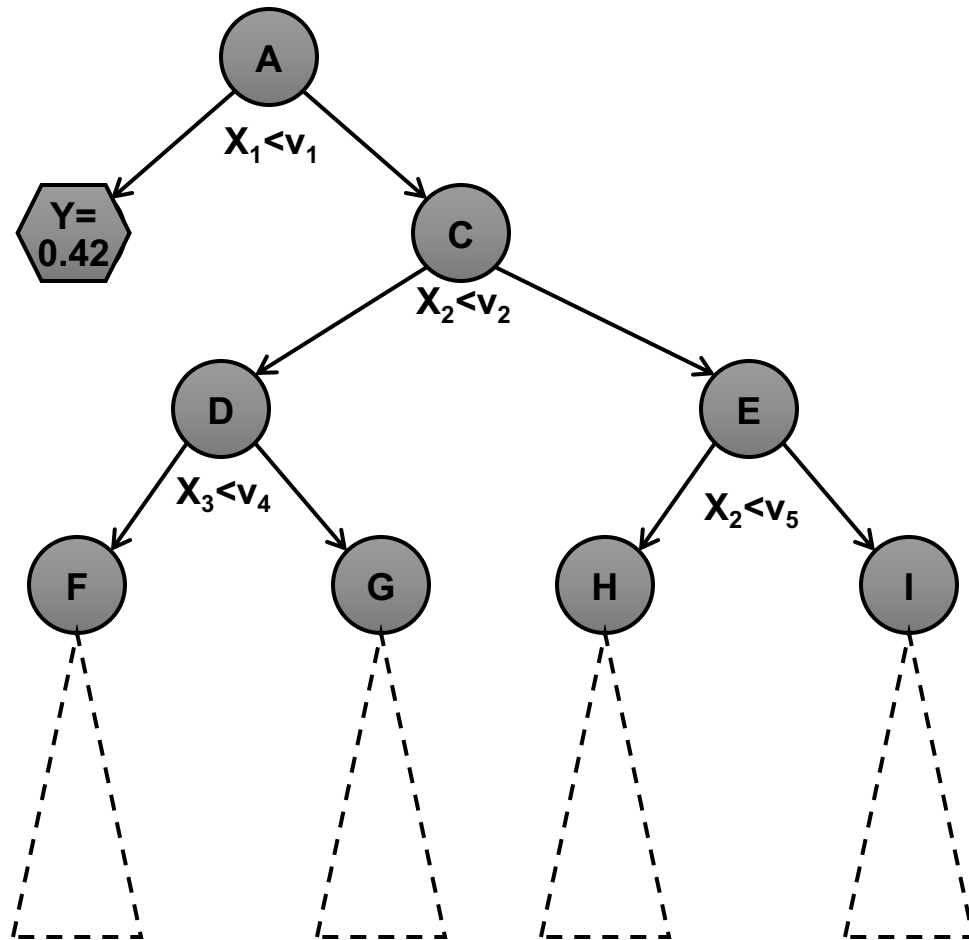
► Setting:

- Binary splits: $X_j < v$
- Numerical attrs.
- Regression



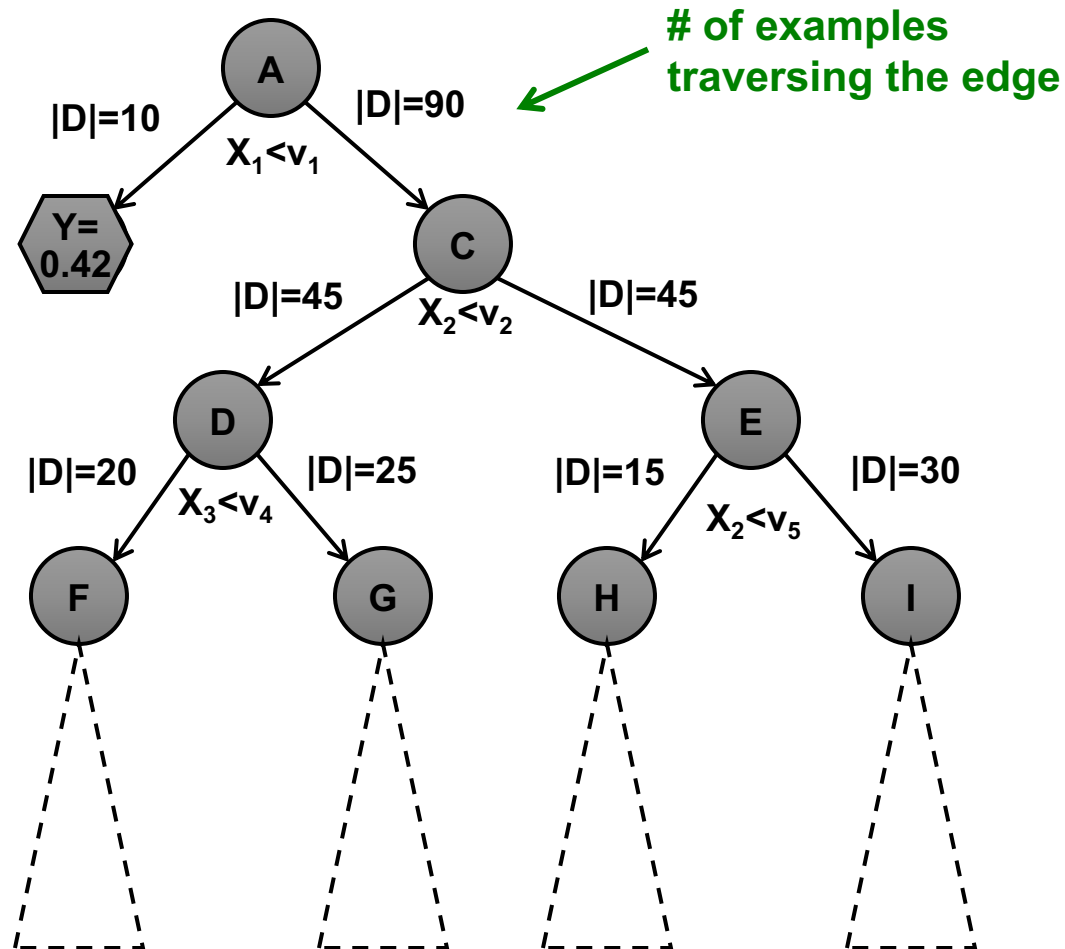
How to make predictions?

- ▶ **Input:** Example x_i
- ▶ **Output:** Predicted y_i'
- ▶ “Drop” x_i down the tree until it hits a leaf node
- ▶ Predict the value stored in the leaf that x_i hits



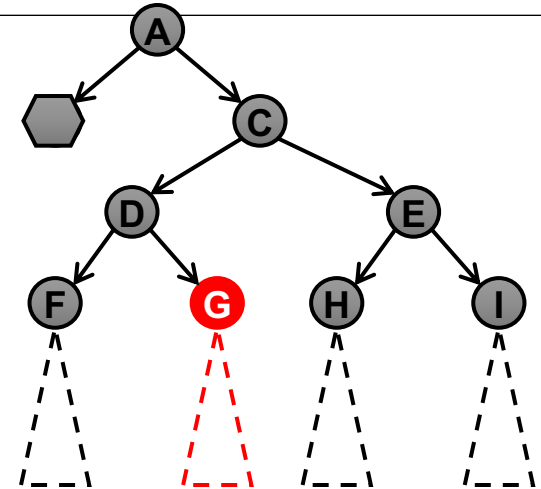
How to construct a tree?

- ▶ Training dataset D^* , $|D^*|=100$ examples

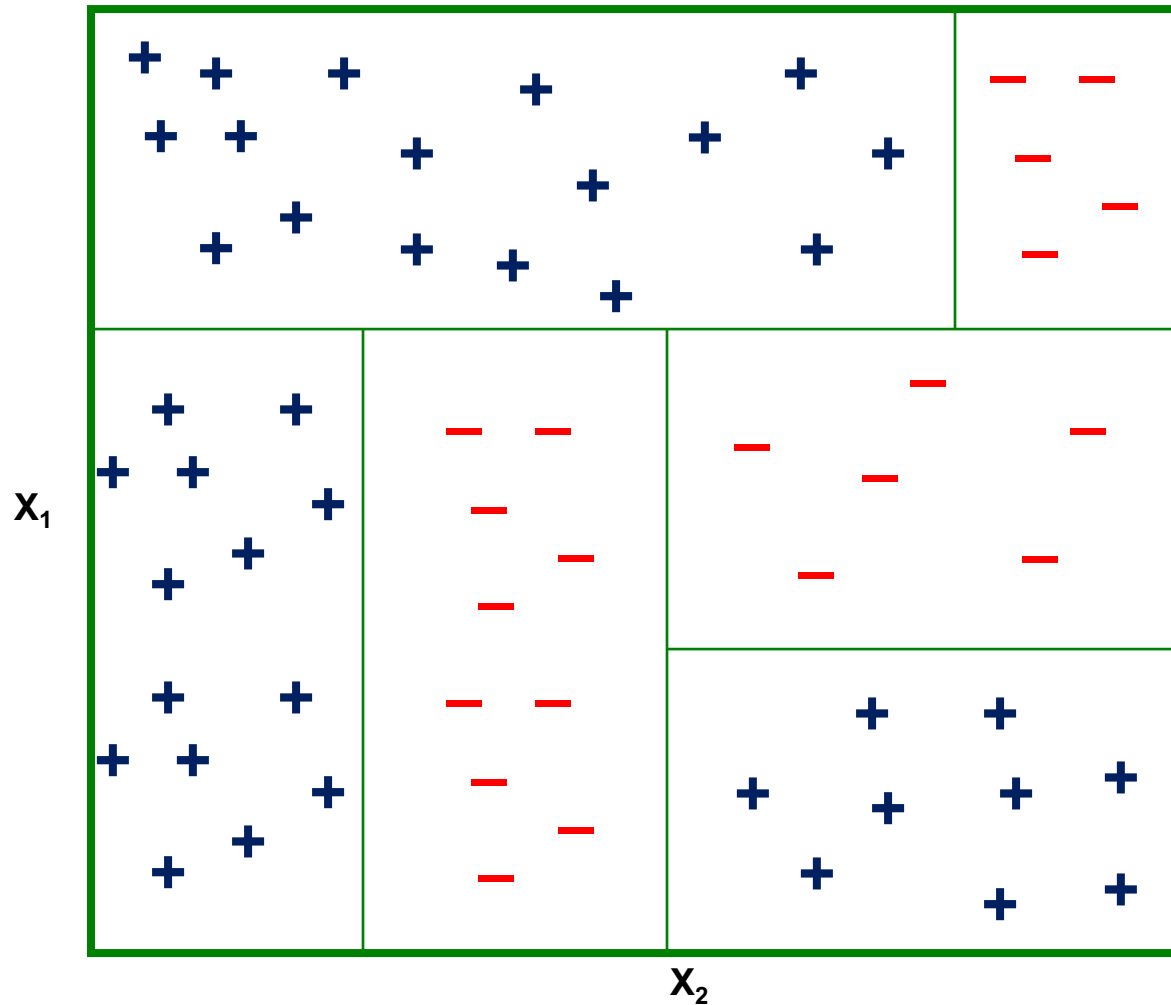


How to construct a tree?

- ▶ Imagine we are currently at some node G
 - ▶ Let D_G be the data reaches G
- ▶ There is a decision we have to make:
 - ▶ If so, which variable and which value do we use for a split?
 - ▶ If not, how do we make a prediction?
 - ▶ We need to build a “predictor node”



How to construct a tree?



How to construct a tree?

Algorithm 1 InMemoryBuildNode

Require: Node n , Data $D \subseteq D^*$

- 1: $(n \rightarrow \text{split}, D_L, D_R) = \text{FindBestSplit}(D)$
 - 2: if $\text{StoppingCriteria}(D_L)$ then
 - 3: $n \rightarrow \text{left_prediction} = \text{FindPrediction}(D_L)$
 - 4: else
 - 5: $\text{InMemoryBuildNode}(n \rightarrow \text{left}, D_L)$
 - 6: if $\text{StoppingCriteria}(D_R)$ then
 - 7: $n \rightarrow \text{right_prediction} = \text{FindPrediction}(D_R)$
 - 8: else
 - 9: $\text{InMemoryBuildNode}(n \rightarrow \text{right}, D_R)$
-

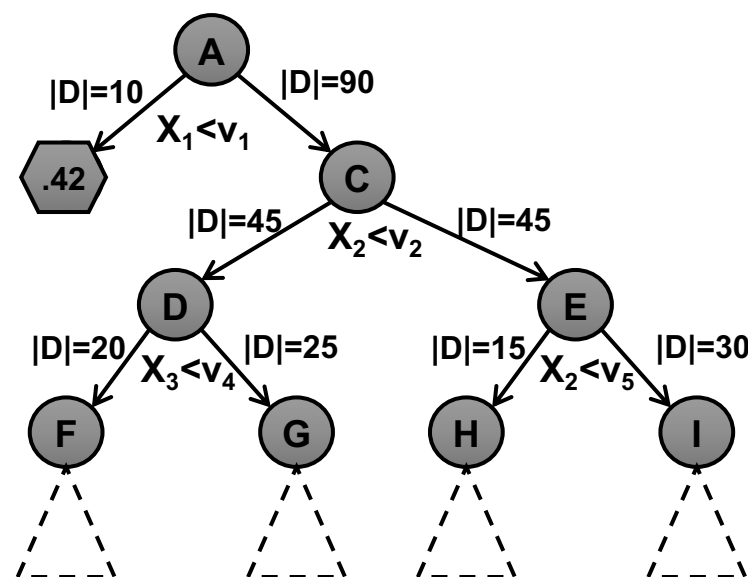
How to construct a tree?

- ▶ **How to split?** Pick attribute & value that optimizes some criterion

- ▶ **Classification:**
Information Gain

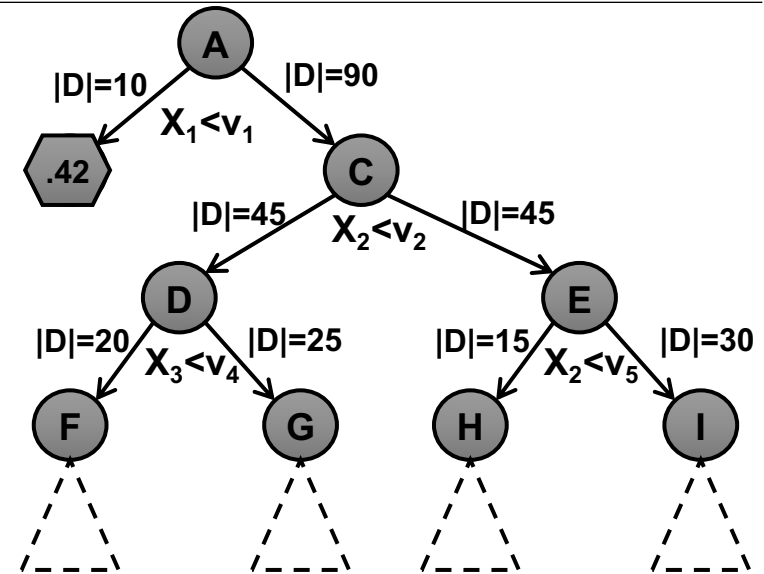
- ▶ $IG(Y|X) = H(Y) - H(Y|X)$

- ▶ Entropy: $H(Z) = -\sum_{j=1}^m p_j \log p_j$
- ▶ Conditional entropy: $H(W|Z) = -\sum_{j=1}^m P(Z = v_j) H(W|Z = v_j)$
 - ▶ Suppose Z takes m values ($v_1 \dots v_m$)
 - ▶ $H(W|Z=v)$... Entropy of W among the records in which Z has value v



How to construct a tree?

- ▶ **How to split?** Pick attribute & value that optimizes some criterion
- ▶ **Regression:**
 - ▶ Find split (X_i, v) that creates D, D_L, D_R : parent, left, right child datasets and maximizes:
 $|D| \cdot \text{Var}(D)$
 $- (|D_L| \cdot \text{Var}(D_L) + |D_R| \cdot \text{Var}(D_R))$
 - ▶ For ordered domains sort X_i and consider a split between each pair of adjacent values
 - ▶ For categorical X_i find best split based on subsets (Breiman's algorithm)



How to construct a tree?

► When to stop?

► 1) When the leaf is “pure”

- E.g., $\text{Var}(y_i) < \varepsilon$

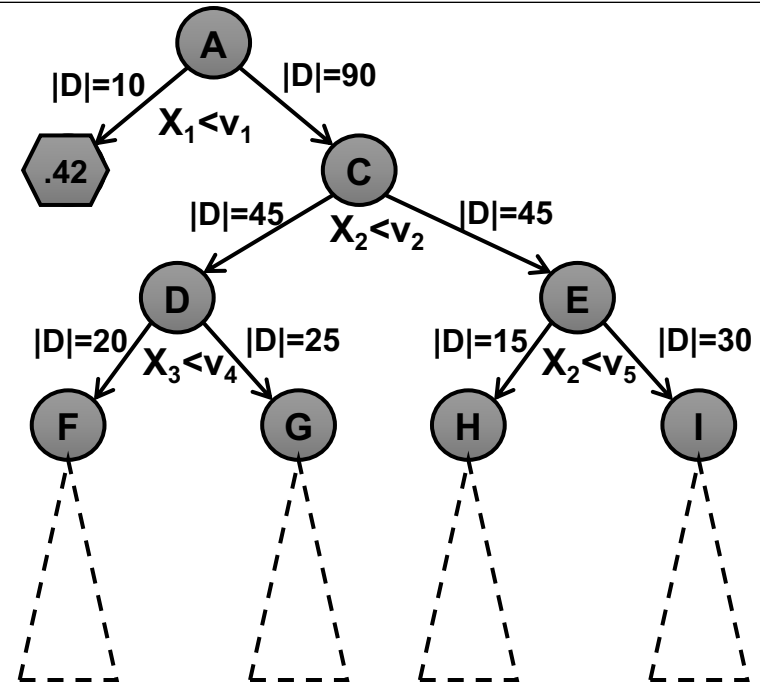
► 2) When # of examples in the leaf is too small

- E.g., $|D| \leq 10$

► How to predict?

► Predictor:

- **Regression:** Avg. y_i of the examples in the leaf
- **Classification:** Most common y_i in the leaf



End of Recap

Estimating Generalization Errors

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

Example of Post-Pruning

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

Training Error (Before splitting) = 10/30

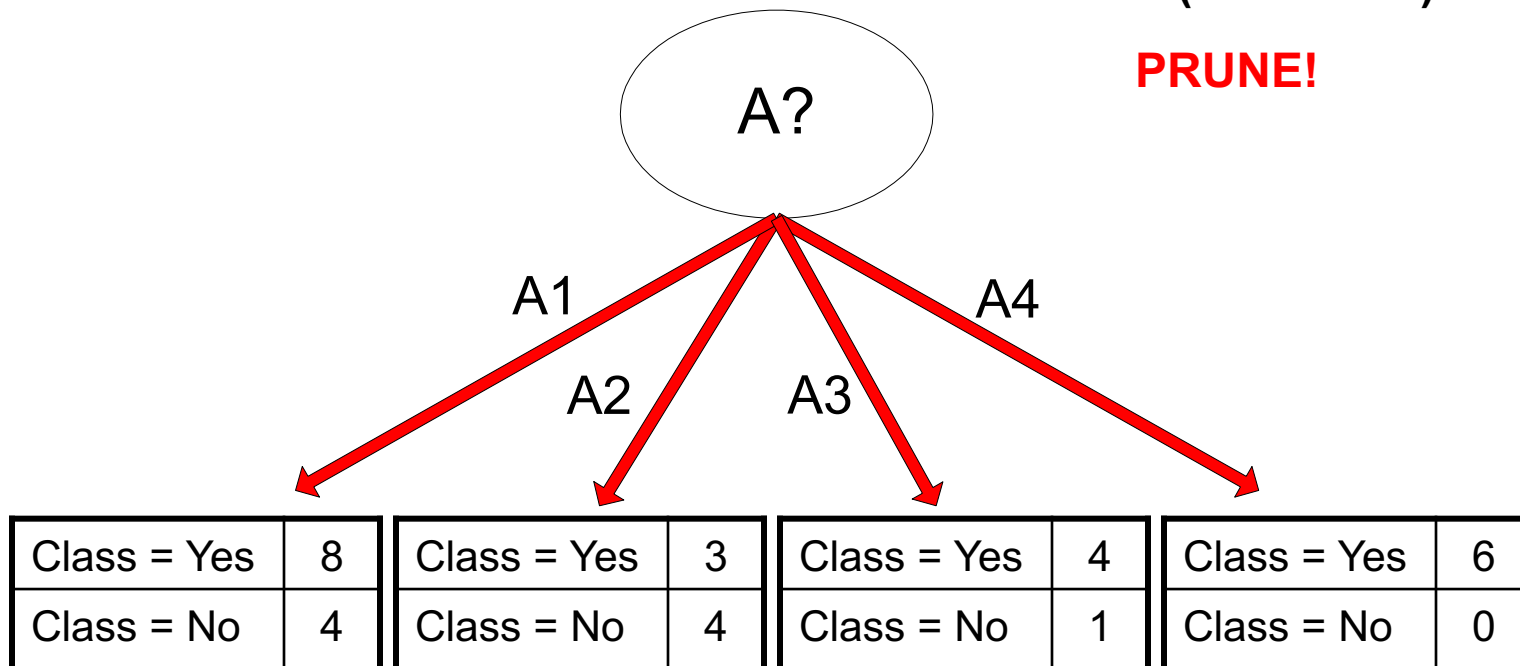
Pessimistic error = $(10 + 0.5)/30 = 10.5/30$

Training Error (After splitting) = 9/30

Pessimistic error (After splitting)

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

PRUNE!



Examples of Post-pruning

► Optimistic error?

Don't prune for both cases

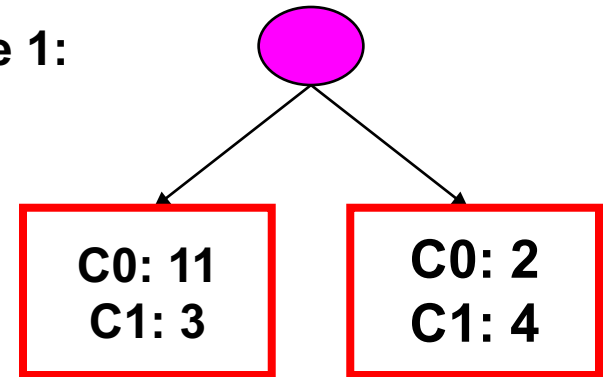
► Pessimistic error?

Don't prune case 1, prune case 2

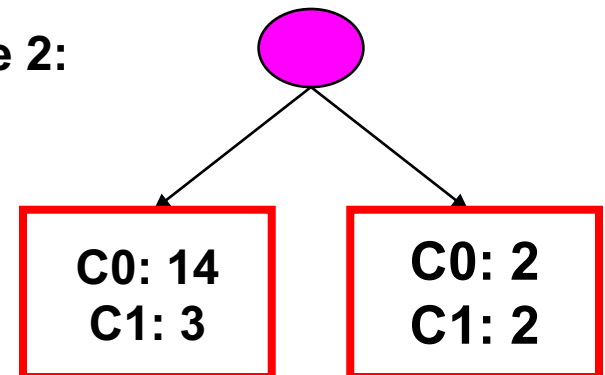
► Reduced error pruning?

Depends on validation set

Case 1:



Case 2:



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

<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	?	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)

$$= -(2/6) \log(2/6) - (4/6) \log(4/6) = 0.9183$$

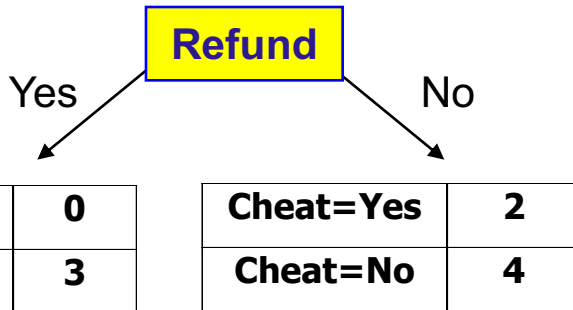
Entropy(Children)

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

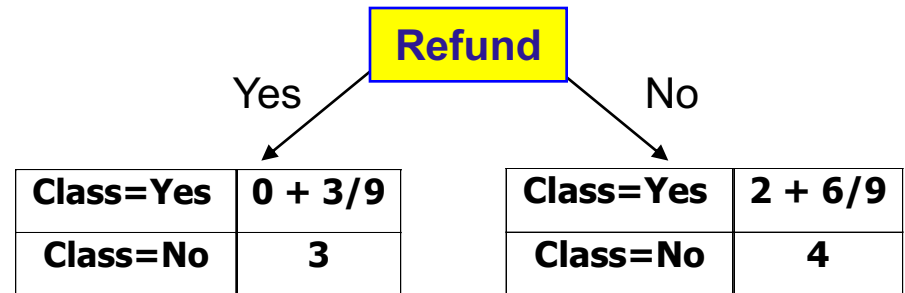
$$\text{Gain} = 0.9 \times (0.8813 - 0.551) = 0.3303$$

Distribute Instances

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



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



Probability that Refund=Yes is 3/9

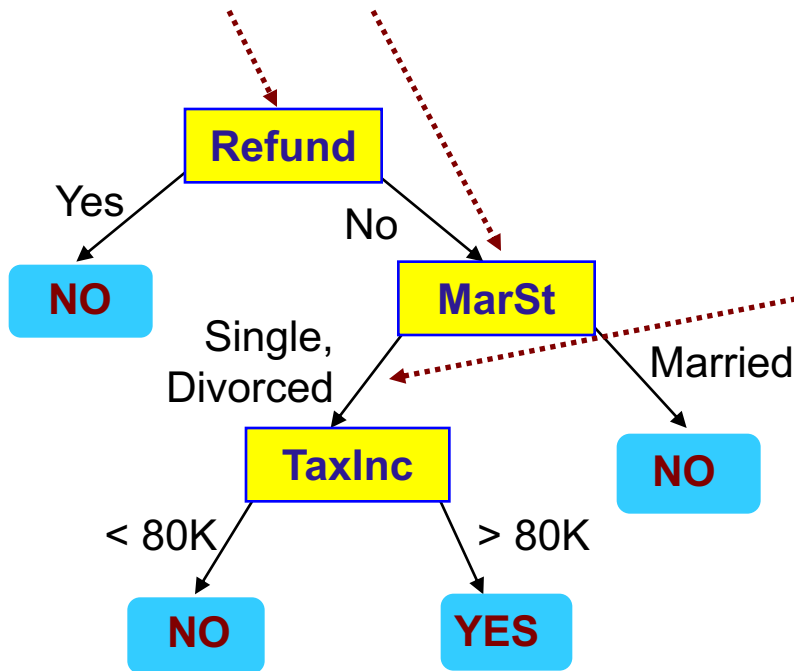
Probability that Refund=No is 6/9

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

Classify Instances

New record:

<i>Tid</i>	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 $3/6.67$

Other Issues

- ▶ Data Fragmentation
- ▶ Search Strategy
- ▶ Expressiveness
- ▶ Tree Replication

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

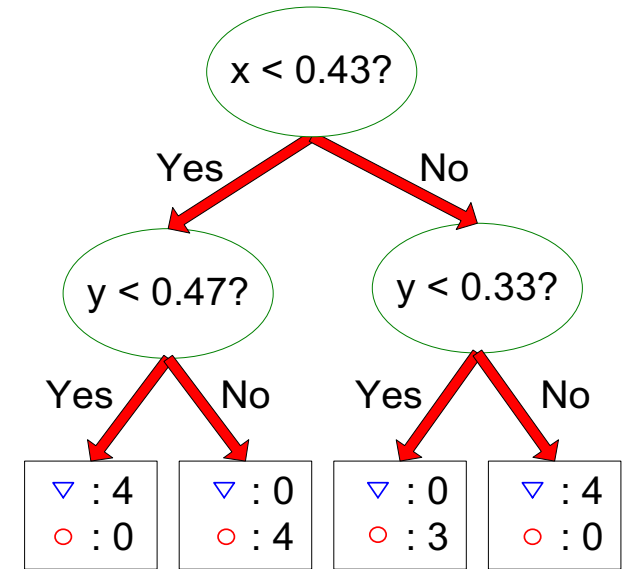
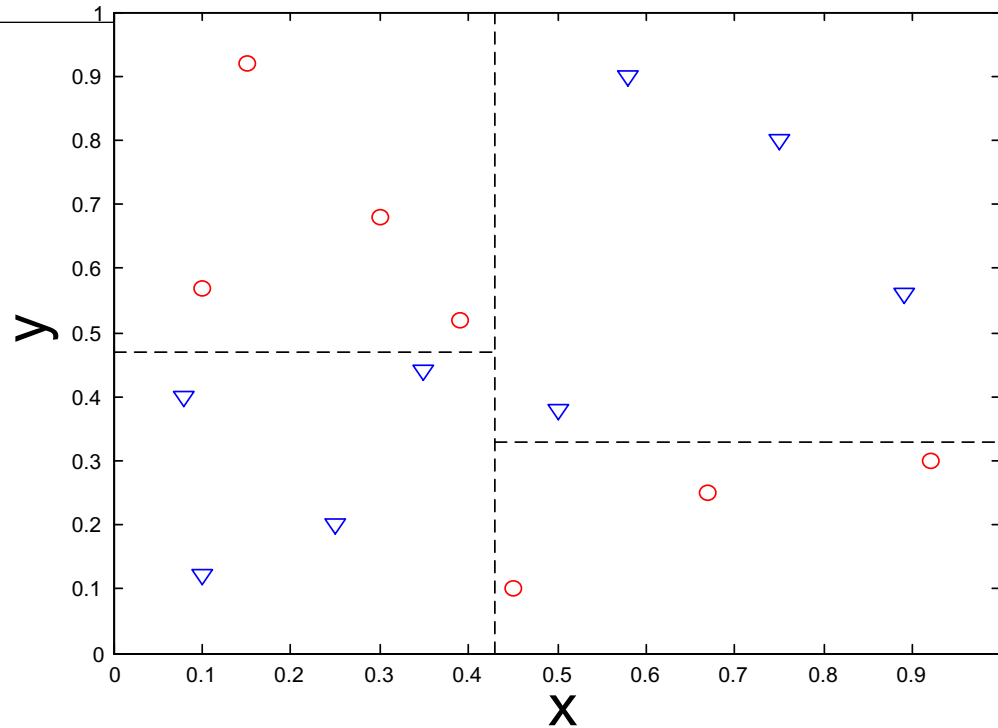
Search Strategy

- ▶ Finding an optimal decision tree is NP-hard
- ▶ The algorithm presented so far uses a greedy, top-down, recursive partitioning strategy to induce a reasonable solution
- ▶ Other strategies?
 - ▶ Bottom-up
 - ▶ Bi-directional

Expressiveness

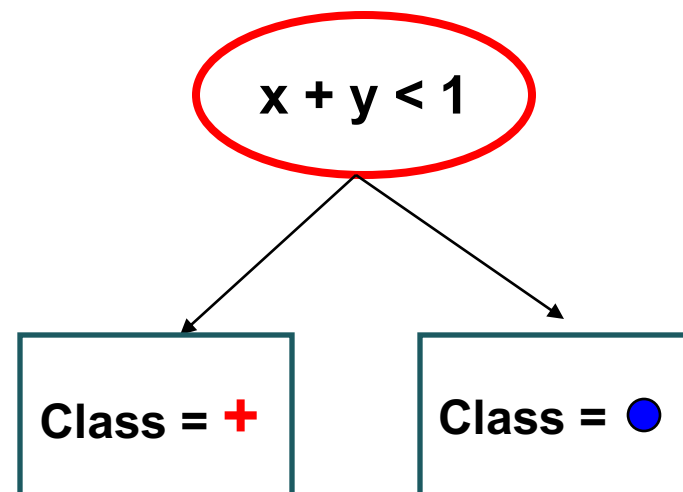
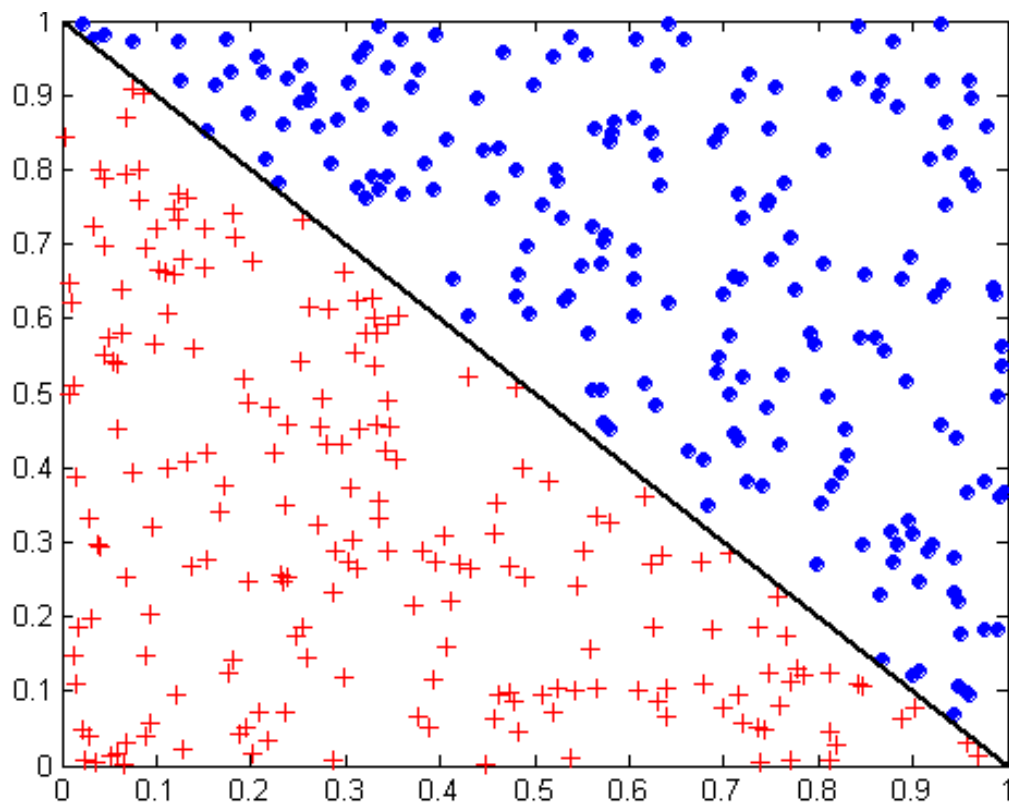
- ▶ Decision tree provides expressive representation for learning discrete-valued function
 - ▶ But they do not generalize well to certain types of Boolean functions
 - ▶ Example: parity function:
 - ▶ Class = 1 if there is an even number of Boolean attributes with truth value = True
 - ▶ Class = 0 if there is an odd number of Boolean attributes with truth value = True
 - ▶ For accurate modeling, must have a complete tree
- ▶ Not expressive enough for modeling continuous variables
 - ▶ Particularly when test condition involves only a single attribute at-a-time

Decision Boundary



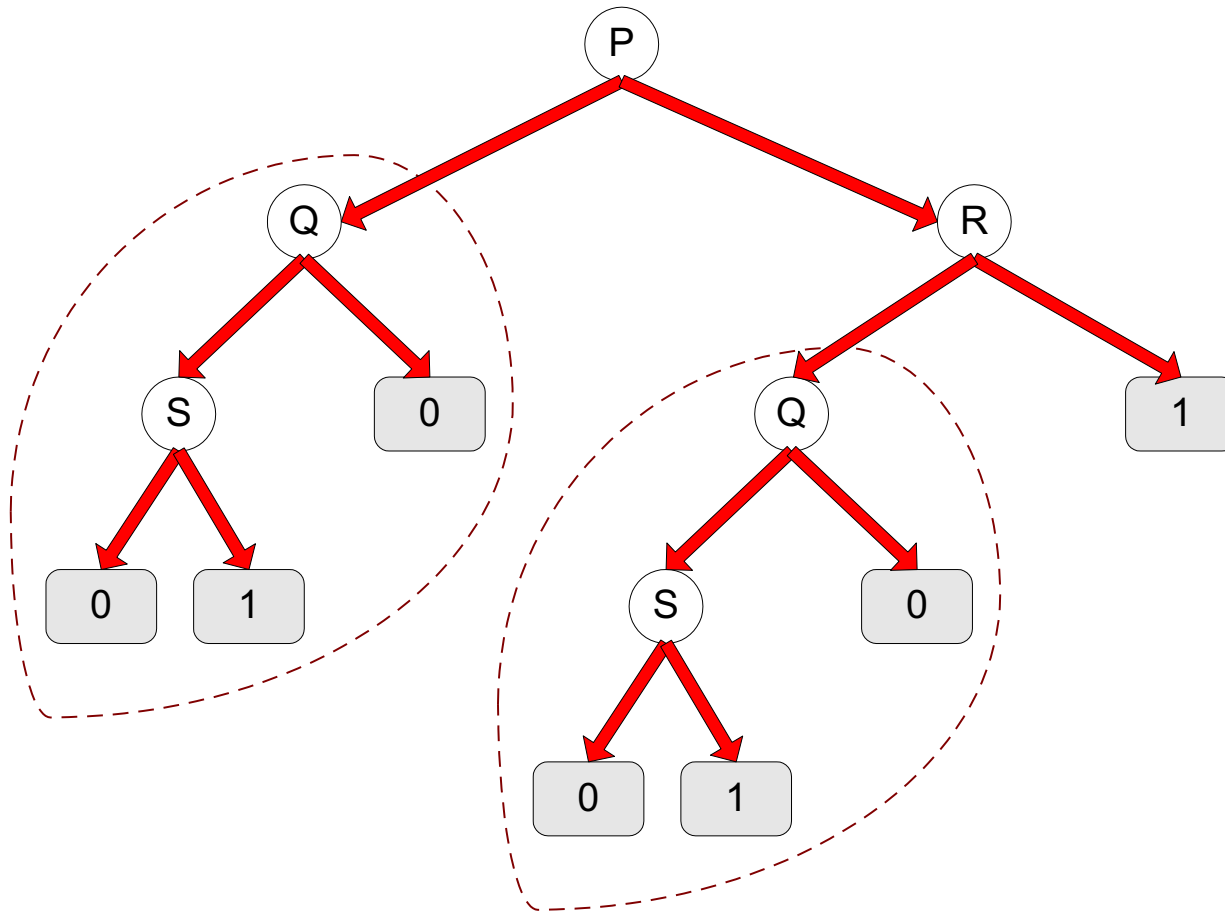
- Border line between two neighboring regions of different classes is known as decision boundary
- Decision boundary is parallel to axes because test condition involves a single attribute at-a-time

Oblique Decision Trees



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

Tree Replication



- Same subtree appears in multiple branches

Model Evaluation

- ▶ Metrics for Performance Evaluation
 - ▶ How to evaluate the performance of a model?
- ▶ Methods for Performance Evaluation
 - ▶ How to obtain reliable estimates?
- ▶ Methods for Model Comparison
 - ▶ How to compare the relative performance among competing models?

Model Evaluation

- ▶ **Metrics for Performance Evaluation**
 - ▶ How to evaluate the performance of a model?
- ▶ **Methods for Performance Evaluation**
 - ▶ How to obtain reliable estimates?
- ▶ **Methods for Model Comparison**
 - ▶ How to compare the relative performance among competing models?

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	Class=No
	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

- ▶ 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		
ACTUAL 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 vs Accuracy

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

Accuracy is proportional to cost if

1. $C(\text{Yes}|\text{No}) = C(\text{No}|\text{Yes}) = q$
2. $C(\text{Yes}|\text{Yes}) = C(\text{No}|\text{No}) = p$

$$N = a + b + c + d$$

$$\text{Accuracy} = (a + d)/N$$

Cost	PREDICTED CLASS		
ACTUAL CLASS		Class=Yes	Class=No
	Class=Yes	p	q
	Class=No	q	p

$$\text{Cost} = p(a + d) + q(b + c)$$

$$= p(a + d) + q(N - a - d)$$

$$= qN - (q - p)(a + d)$$

$$= N[q - (q - p) \times \text{Accuracy}]$$

Cost-Sensitive Measures

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

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

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

- Precision is biased towards $C(\text{Yes}|\text{Yes})$ & $C(\text{Yes}|\text{No})$
- Recall is biased towards $C(\text{Yes}|\text{Yes})$ & $C(\text{No}|\text{Yes})$
- F-measure is biased towards all except $C(\text{No}|\text{No})$

$$\text{Weighted Accuracy} = \frac{w_1 a + w_4 d}{w_1 a + w_2 b + w_3 c + w_4 d}$$

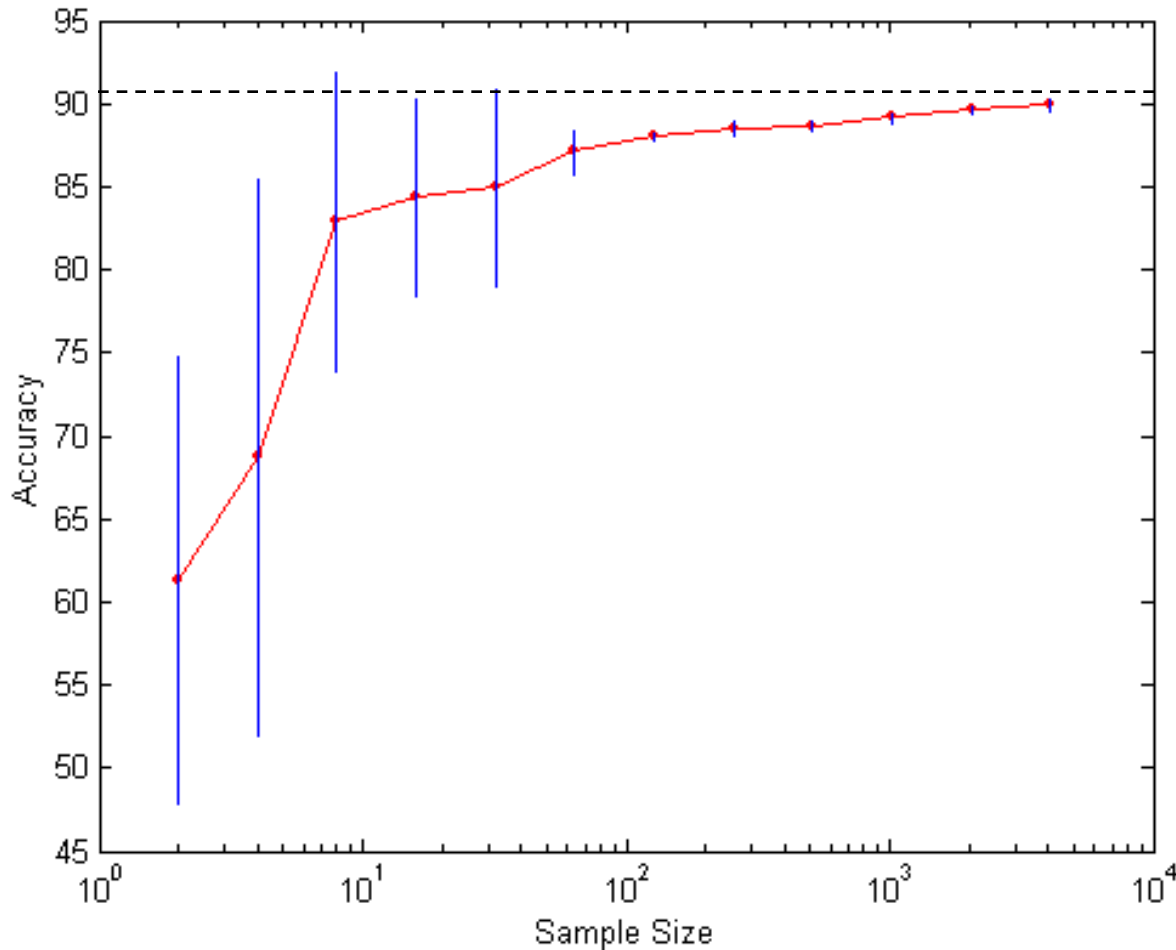
Model Evaluation

- ▶ Metrics for Performance Evaluation
 - ▶ How to evaluate the performance of a model?
- ▶ Methods for Performance Evaluation
 - ▶ How to obtain reliable estimates?
- ▶ Methods for Model Comparison
 - ▶ How to compare the relative performance among competing models?

Methods for Performance Evaluation

- ▶ How to obtain a reliable estimate of performance?
- ▶ Performance of a model may depend on other factors besides the learning algorithm:
 - ▶ Class distribution
 - ▶ Cost of misclassification
 - ▶ Size of training and test sets

Learning Curve



- Learning curve shows how accuracy changes with varying sample size
- Requires a sampling schedule for creating learning curve:
 - Arithmetic sampling (Langley, et al)
 - Geometric sampling (Provost et al)

Effect of small sample size:

- Bias in the estimate
- Variance of estimate

Methods of Estimation

- ▶ Holdout
 - ▶ Reserve $2/3$ for training and $1/3$ for testing
- ▶ Random subsampling
 - ▶ Repeated holdout
- ▶ Cross validation
 - ▶ Partition data into k disjoint subsets
 - ▶ k -fold: train on $k-1$ partitions, test on the remaining one
 - ▶ Leave-one-out: $k=n$
- ▶ Bootstrap
 - ▶ Sampling with replacement

Model Evaluation

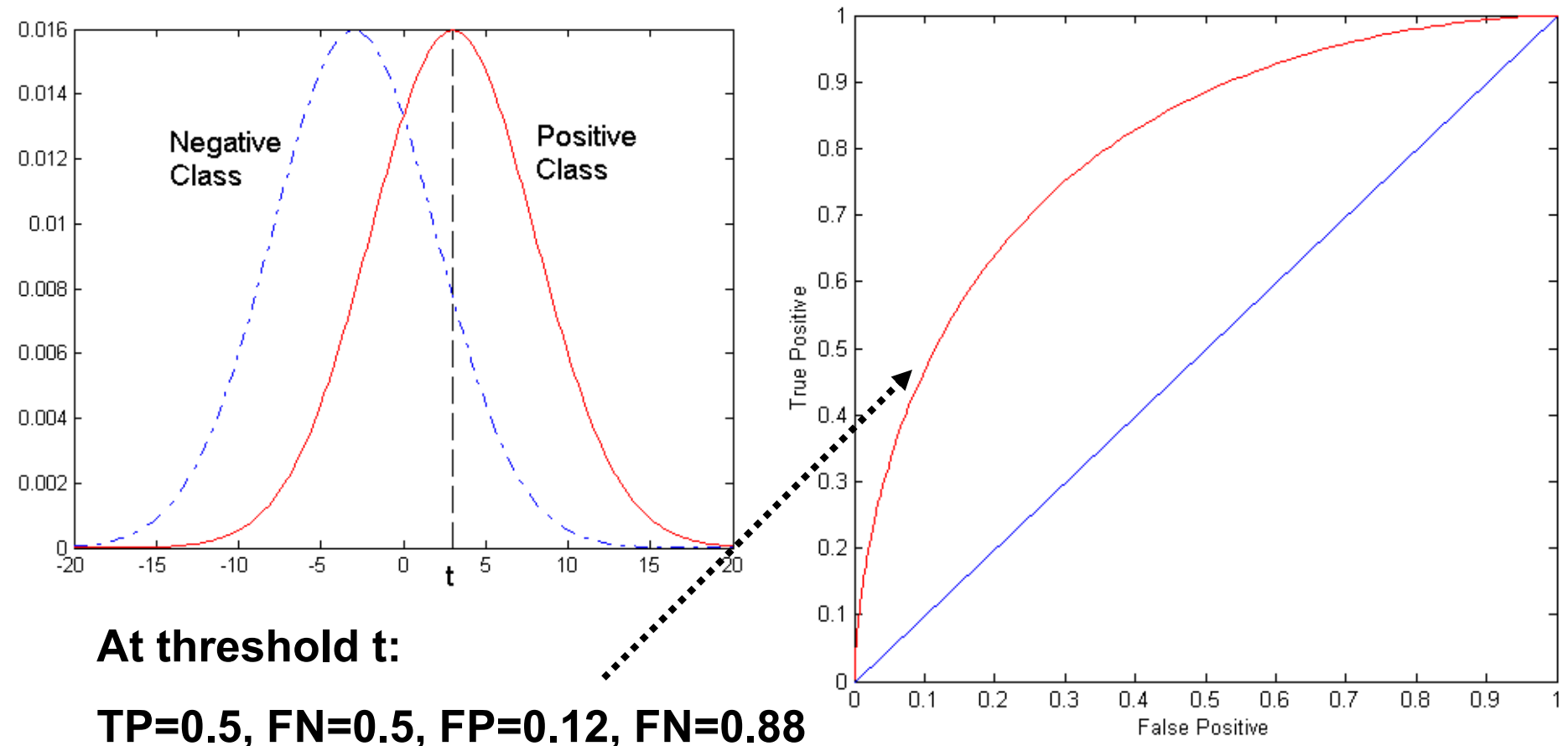
- ▶ Metrics for Performance Evaluation
 - ▶ How to evaluate the performance of a model?
- ▶ Methods for Performance Evaluation
 - ▶ How to obtain reliable estimates?
- ▶ **Methods for Model Comparison**
 - ▶ How to compare the relative performance among competing models?

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 TP (on the y-axis) against FP (on the x-axis)
- ▶ Performance of each classifier represented as a point on the ROC curve
 - ▶ changing the threshold of algorithm, sample distribution or cost matrix changes the location of the point

ROC Curve

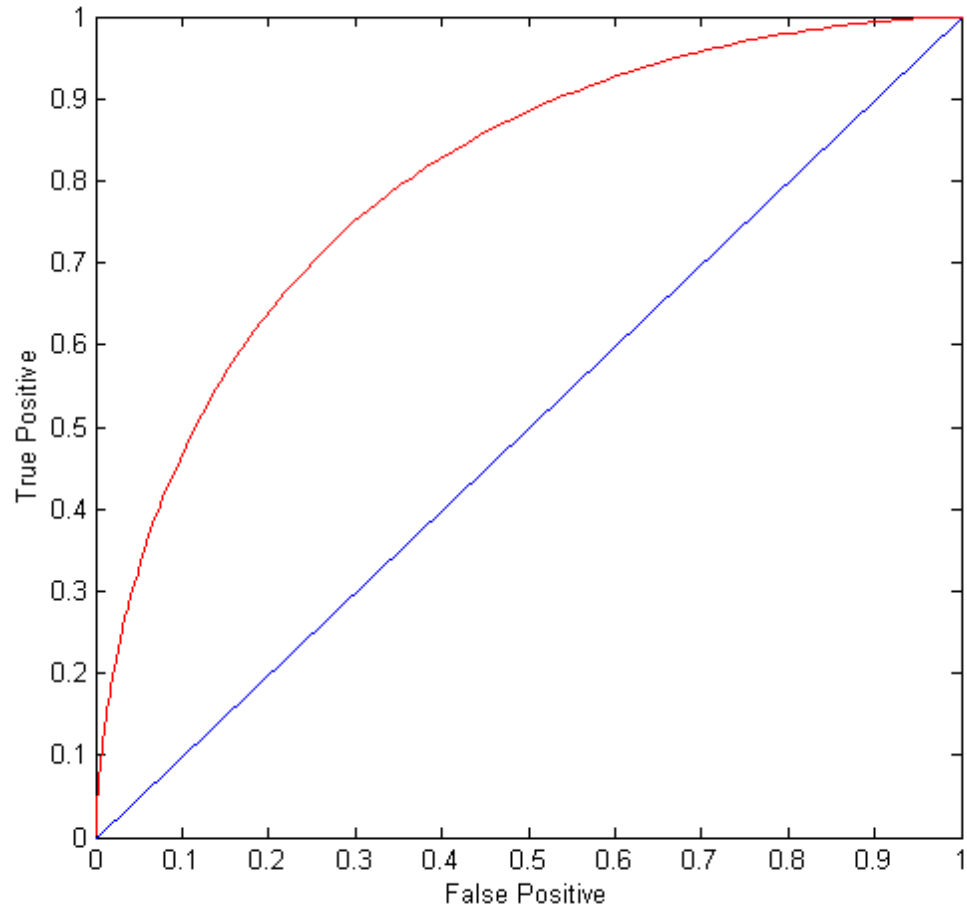
- 1-dimensional data set containing 2 classes (positive and negative)
- any points located at $x > t$ is classified as positive



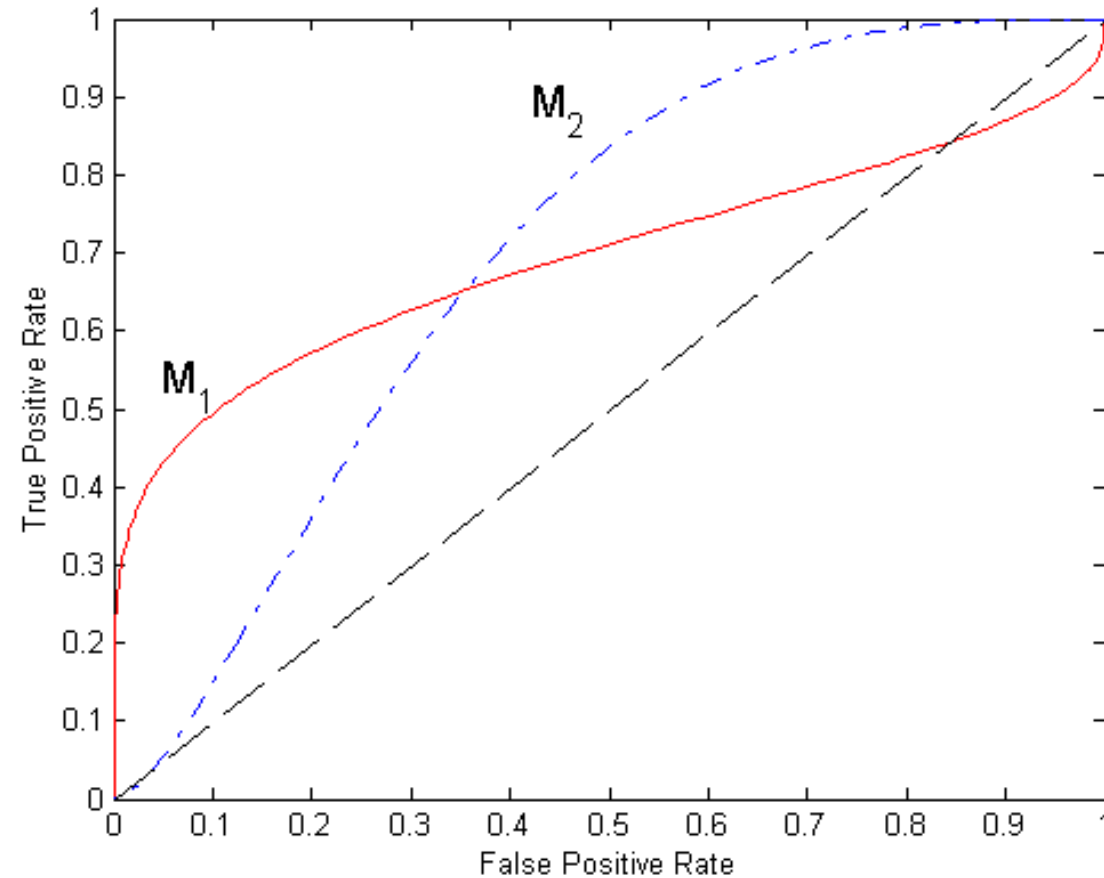
ROC Curve

(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
 - ▶ Below diagonal line:
 - ▶ prediction is opposite of the true class



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
 - Ideal:
 - Area = 1
 - Random guess:
 - Area = 0.5

How to Construct an ROC curve

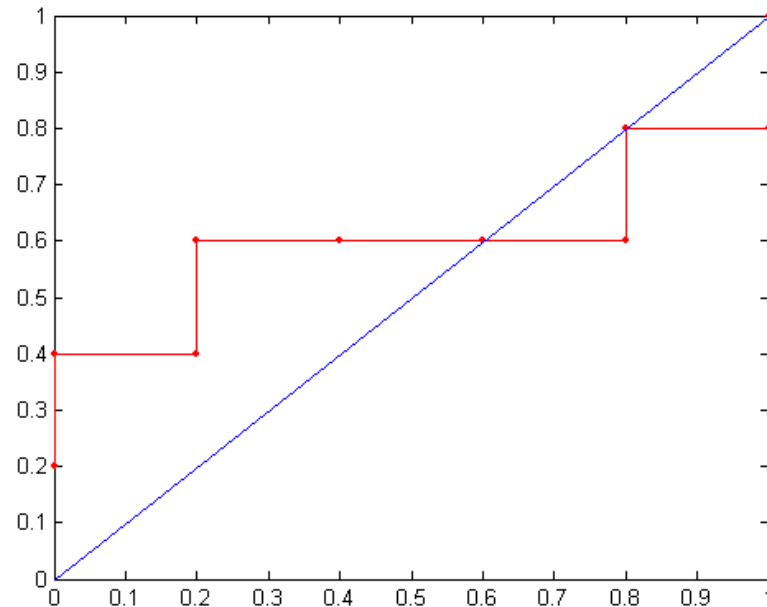
Instance	$P(+ A)$	True Class
1	0.95	+
2	0.93	+
3	0.87	-
4	0.85	-
5	0.85	-
6	0.85	+
7	0.76	-
8	0.53	+
9	0.43	-
10	0.25	+

- 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 a ROC curve

Class	+	-	+	-	-	-	+	-	+	+	
Threshold >=	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:



Test of Significance

- ▶ Given two models:
 - ▶ Model M1: accuracy = 85%, tested on 30 instances
 - ▶ Model M2: accuracy = 75%, tested on 5000 instances
- ▶ Can we say M1 is better than M2?
 - ▶ How much confidence can we place on accuracy of M1 and M2?
 - ▶ Can the difference in performance measure be explained as a result of random fluctuations in the test set?

Literature

- ▶ Chapter 4 (except 4.6) from the Tan et. al. Textbook.