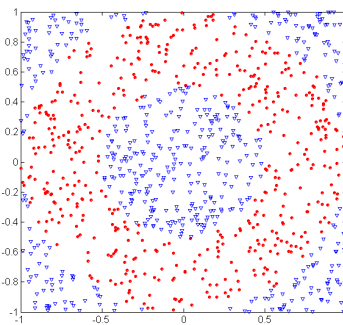


Underfitting and Overfitting (Example)

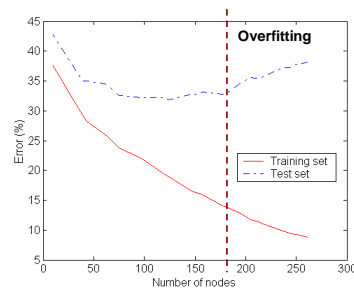


500 circular and 500 triangular data points.

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

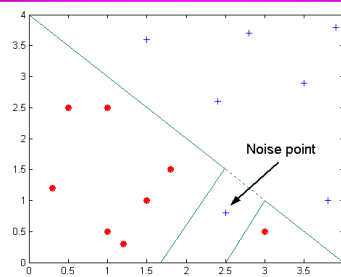
Triangular points:
 $\sqrt{x_1^2 + x_2^2} < 0.5$ or
 $\sqrt{x_1^2 + x_2^2} > 1$

Underfitting and Overfitting



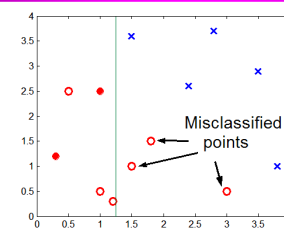
Underfitting: when model is too simple, both training and test errors are large

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

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 ($\sum e(t)$)
- **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

Occam's Razor

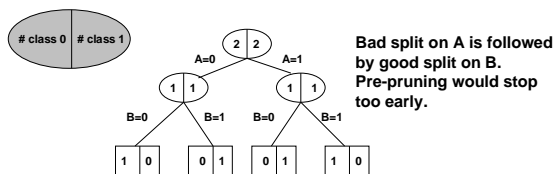
- Given two models of similar generalization errors, one should prefer the simpler model over the more complex model
- For complex models, there is a greater chance that it was fitted accidentally by errors in data
- Therefore, one should include model complexity when evaluating a model

How to Address Overfitting

- Pre-Pruning (Early Stopping Rule)**
 - Stop the algorithm before it becomes a fully-grown tree
 - Possible conditions:
 - Stop if number of instances is less than some user-specified threshold.
 - Stop if expanding the current node does not improve impurity measures (e.g., Gini or information gain) by at least some threshold.

Disadvantage of Pre-Pruning

- Since we use a hill-climbing search, looking only one step ahead, pre-pruning might stop too early.
- Extreme example: exclusive OR, A XOR B.



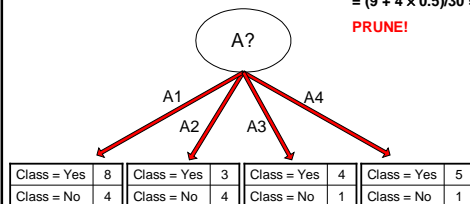
How to Address Overfitting...

- 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

Example of Post-Pruning

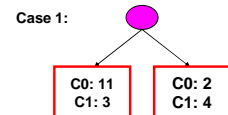
Class = Yes	20
Class = No	10
Error = 10/30	

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!



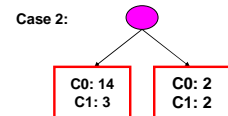
Examples of Post-pruning

- 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



Error based pruning in C4.5 (J48)

Suppose we observe 2 errors out of 7 in a leaf node.
The point estimate for the error rate in that leaf is $2/7=0.286$.

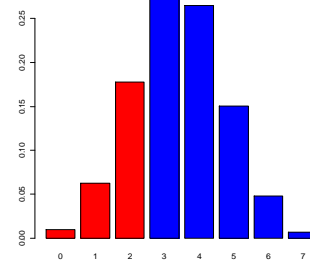
If the probability of error is $2/7$ then the probability of observing 2 errors or less is 0.68 (binomial distribution with $n=7$, and $p=2/7$).

As a pessimistic estimate of the error rate in this leaf node, we are going to find the value of p , such that the probability of 2 errors or less is relatively small, say 0.25. This turns out to be the case for $p = 0.4861$. Hence $[0, 0.4861]$ can be regarded as a 75% right-one-sided confidence interval for p .

For $p=0.659$ the probability of observing 2 errors or less is 0.05. Values higher than 0.659 give even smaller probabilities to observing 2 errors or less, and are therefore highly unlikely.

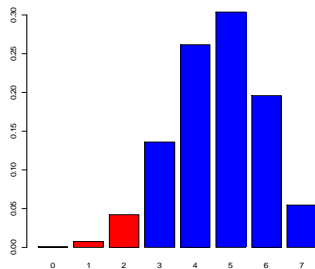
Error based pruning in C4.5 (J48)

Binomial distribution with $n=7$, and $p = 0.4861$. $P(X \leq 2) = 0.25$. This is the sum of the height of the red bars. p denotes the probability of error, and X the number of errors.

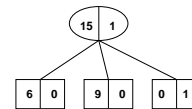


Error based pruning in C4.5 (J48)

Binomial distribution with $n=7$, and $p = 0.659$. $P(X \leq 2) = 0.05$. This is the sum of the height of the red bars.



Error based pruning in C4.5 (J48)



Before pruning: $6(0.206) + 9(0.143) + 1(0.75) = 3.273$
After pruning: $16(0.157) = 2.521$

$U_{2.5\%}(0,1) = 0.75$ because $P(X \leq 0) = 0.25$ for a Bernoulli trial with probability of success 0.75.

Tree after pruning has lower predicted error, so we prune.

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

Distribute Instances

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

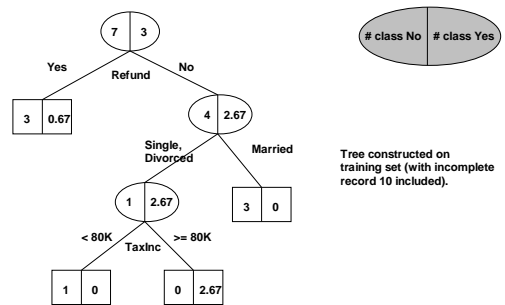
Refund	
Yes	No
Class=Yes	0
Class=No	3
Cheat=Yes	2
Cheat=No	4

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

Refund	
Yes	No
Class=Yes	0 + 3/9
Class=No	3
Class=Yes	2 + 6/9
Class=No	4

Probability that Refund=Yes is 3/9
 Probability that Refund=No is 6/9
 Assign record to the left child with weight = 3/9 and to the right child with weight = 6/9

Distribute Instances

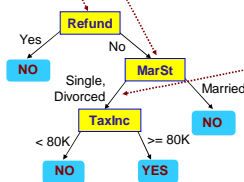


Classify Instances

New record:

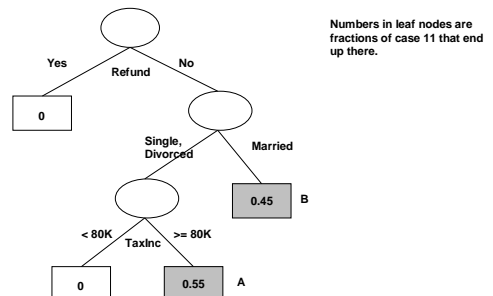
Tid	Refund	Marital Status	Taxable Income	Class
11	No	?	85K	?

	Married	Single	Divorced	Total
Class=No	3	1	0	4
Class=Yes	0	1.67	1	2.67
Total	3	2.67	1	6.67



Probability that Marital Status = Married is 3/6.67
 Probability that Marital Status = {Single, Divorced} is 3.67/6.67

Classify Instances



Classify Instances

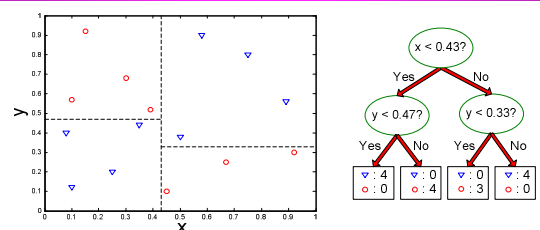
For the final prediction for case 11, take the weighted average of the predictions in node A and B.

$$P_A(\text{class=yes}) = 1, P_B(\text{class=yes}) = 0.$$

Weighted average:

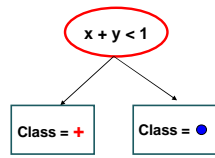
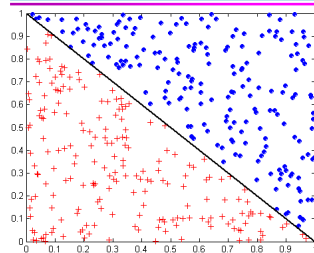
$$P(\text{class=yes}) = \frac{0.55 \cdot 2.67 \cdot 1 + 0.45 \cdot 3 \cdot 0}{0.55 \cdot 2.67 + 0.45 \cdot 3} \approx 0.52$$

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

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:

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

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

Metrics for Performance Evaluation...

	PREDICTED CLASS		
ACTUAL CLASS		Class=Yes	Class=No
	Class=Yes	a (TP)	b (FN)
	Class=No	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
ACTUAL CLASS	Class=Yes	C(Yes Yes)	C(No Yes)
	Class=No	C(Yes No)	C(No No)

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

Computing Cost of Classification

Cost Matrix	PREDICTED CLASS		
	C(i j)	+	-
ACTUAL CLASS	+	-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

True positive rate: $TP/(TP+FN)$, Fraction of positive examples that is predicted to be positive.

False positive rate: $FP/(FP+TN)$, Fraction of negative examples that is predicted to be positive.

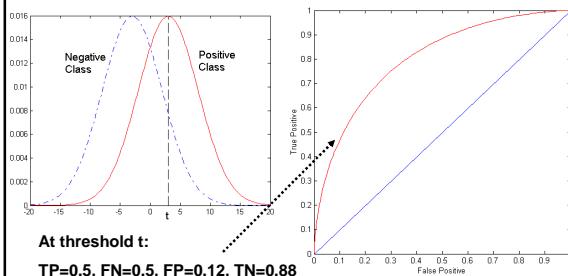
	PREDICTED CLASS		
ACTUAL CLASS		Class=Yes	Class=No
	Class=Yes	a (TP)	b (FN)
	Class=No	c (FP)	d (TN)

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

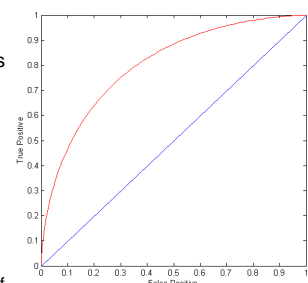
- 1-dimensional data set containing 2 classes (positive and negative)
- any points located at $x > t$ is classified as positive



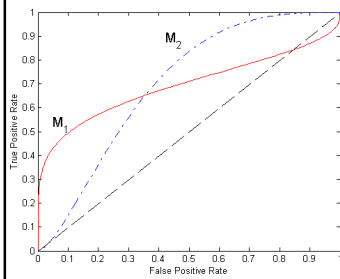
ROC Curve

(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 an ROC curve

Instance	$P(+)$	True class	FPR	TPR
1	0.95	+	0	1/5
2	0.93	+	0	2/5
3	0.87	-	1/5	2/5
4	0.85	-		
5	0.85	-		
6	0.85	+	3/5	3/5
7	0.76	-	4/5	3/5
8	0.53	+	4/5	4/5
9	0.43	-	1	4/5
10	0.25	+	1	1

How to construct an ROC curve

