

# CS256 – Midterm Exam Study Guide

By: Zayd Hammoudeh

## Chapter #04 – Classification: Basic Concepts, Decision Trees, and Model Evaluation

<b>Classification</b> Task of assigning objects to one of several predefined categories.	<b>Training Set</b> A collection of records. Each <b>record</b> contains a set of attributes one of which is the <b>class</b> .	<b>Model</b> A function from the value of record attributes to the class attribute.	<b>Test Set</b> A collection of records used to determine the accuracy of the classification model.	<b>Example Classification Techniques</b> <ol style="list-style-type: none"> <li>1. Neural Networks</li> <li>2. Decision Tree</li> <li>3. Rule Based Classifier</li> <li>4. Memory Based Reasoning</li> <li>5. Support Vector Machines</li> <li>6. Naïve Bayes and Bayesian Belief Networks</li> </ol>
---	--	--	--	---

<b>Induction</b> Using a training set to generate a model.  <b>Deduction</b> Process of applying a model to a training set.  <b>Decision Tree Induction</b> <ul style="list-style-type: none"> <li>• <b>Greedy Strategy</b></li> <li>• <b>Key Decision #1:</b> Attribute to expand next</li> <li>• <b>Key Decision #2:</b> When to stop expanding</li> </ul>	<b>Hunt's Decision Tree Induction Algorithm:</b> <ul style="list-style-type: none"> <li>• Let <math>D_t</math> be the set of training records that reach a node <math>t</math>. <ol style="list-style-type: none"> <li>1. If <math>D_t</math> contains records that <b>all belong to the same class <math>y_t</math></b>, then <math>t</math> is a leaf node with class value <math>y_t</math>.</li> <li>2. If <math>D_t</math> is an <b>empty set</b>, then <math>t</math> is a leaf node with default value <math>y_d</math>.</li> <li>3. If <math>D_t</math> contains <b>records that belong to more than one class and there are no attributes left</b>, then <math>t</math> is a leaf node with default value <math>y_d</math>.</li> <li>4. If <math>D_t</math> contains <b>records that belong to more than one class</b>, then use an attribute test to split the data into smaller subsets. Recursively apply the same procedure above.</li> </ol> </li> </ul>	<b>Attribute Types</b> <ul style="list-style-type: none"> <li>• <b>Binary</b> – Attribute with exactly two possible values.</li> <li>• <b>Nominal</b> – Two or more class values with no intrinsic Order</li> <li>• <b>Ordinal</b> – Two or more class values that can be ordered or ranked</li> <li>• <b>Continuous</b> – Quantitative attribute that can be measured along a continuum.</li> </ul>
--	--	--

<b>Splitting Nominal and Ordinal Attributes</b> <ul style="list-style-type: none"> <li>• <b>Binary</b> – Divides attribute values into two subsets. <b>This requires the additional step of finding optimal partitioning.</b></li> <li>• <b>Multi-way</b> – Use as many partitions as distinct values.</li> </ul>	<b>Splitting Based on Continuous Attributes</b> <ul style="list-style-type: none"> <li>• <b>Discretization</b> – Form an ordinal categorical attribute. <ul style="list-style-type: none"> <li>◦ <b>Static</b> – Discretize once at the beginning</li> <li>◦ <b>Dynamic</b> – Ranges can be found by equal interval bucketing, equal frequency bucketing, or clustering.</li> </ul> </li> <li>• <b>Binary Decision</b> (<math>A &lt; v</math> or <math>A &gt; v</math>) – Consider all possible splits and find the best cut. <ul style="list-style-type: none"> <li>◦ <b>Binary Decision Procedure:</b> Go between each training set record value and calculate the GINI index if the splitting point was at that value. <b>Select the splitting point with the lowest <math>GINI_{SPLIT}</math> value.</b> <ul style="list-style-type: none"> <li>▪ <b>Computationally inefficient <math>O(n)</math></b> – where <math>n</math> is the number of records.</li> </ul> </li> </ul> </li> </ul>	<b>Homogeneity/Low Impurity</b> – Extent to which nodes in the decision tree have the same class value/distribution.  <b>Nodes with high levels of homogeneity (i.e. low levels of impurity) are preferred.</b>
---	--	---

### Impurity Measures

For all of these metrics, a lower score is generally preferable.

<b>GINI Index</b> $GINI(t) = 1 - \sum_{j=1}^{n_c} (p(j t))^2$ <ul style="list-style-type: none"> <li>• <math>t</math> – Node in the decision tree</li> <li>• <math>j</math> – Class value</li> <li>• <math>n_c</math> – Number of class values</li> <li>• <math>p(j t)</math> – Probability (i.e. relative frequency) of class value <math>j</math> in node <math>t</math></li> </ul> <b>Minimum Value:</b> 0 when: $\exists j(p(j t) = 1)$ <b>Maximum Value:</b> $1 - \frac{1}{n_c}$ when: $\forall j(p(j t) = \frac{1}{n_c})$	<b><math>GINI_{SPLIT}</math></b> $GINI_{SPLIT} = \sum_{i=1}^k \frac{n_i}{n} \cdot GINI(i)$ <ul style="list-style-type: none"> <li>• <math>i</math> – Child node</li> <li>• <math>n</math> – Number of records in parent node. Note: <math display="block">n = \sum_{i=1}^k n_i</math> </li> <li>• <math>n_i</math> – Number of child nodes (i.e. attribute partitions)</li> <li>• <math>GINI(i)</math> – GINI index value of node <math>i</math>.</li> </ul> <b>Minimum Value:</b> 0 when: $\forall i(GINI(i) = 0)$ <b>Maximum Value:</b> $1 - \frac{1}{n_c}$ when: $\forall i(GINI(i) = 1 - \frac{1}{n_c})$	<b>Entropy</b> $Entropy(t) = - \sum_{j=1}^{n_c} p(j t) \cdot \log_2(p(j t))$ <ul style="list-style-type: none"> <li>• <math>t</math> – Node in the decision tree</li> <li>• <math>j</math> – Class value</li> <li>• <math>n_c</math> – Number of class values</li> <li>• <math>p(j t)</math> – Probability (i.e. relative frequency) of class value <math>j</math> in node <math>t</math></li> </ul> <b>Minimum Value:</b> 0 when: $\exists j(p(j t) = 1)$ <b>Maximum Value:</b> $\log_2(n_c)$ when: $\forall j(p(j t) = \frac{1}{n_c})$
--	---	---

Information Gain		Classification Error
$GAIN_{SPLIT}(t) = Entropy(p) - \left( \sum_{i=1}^k \frac{n_i}{n} \cdot Entropy(i) \right)$ <ul style="list-style-type: none"> <li><math>p</math> – Parent node in the decision tree</li> <li><math>i</math> – Child node in the decision tree</li> <li><math>k</math> – Number of child nodes</li> <li><math>n_i</math> – Number of records in child node <math>i</math></li> <li><math>n</math> – Number of records in parent node <math>p</math></li> </ul> $n = \sum_{i=1}^k n_i$ <p><b>Key Note:</b> A higher <math>GAIN_{SPLIT}</math> is preferable unlike with the other metrics where a lower value was better.</p> <p><b>Disadvantage of Information Gain:</b> Tends to prefer splits that result in a large number of partitions, each being small but pure (i.e. overfitting)</p>	<p><b>Normalizing for Split Size</b></p> $GainRATIO_{Split} = \frac{Gain_{SPLIT}(t)}{SplitINFO}$ $SplitINFO = - \sum_{i=1}^k \frac{n_i}{n} \cdot \log_2 \left( \frac{n_i}{n} \right)$ <p><b>SplitINFO penalizes a large split by reducing the gain.</b></p>	$Error(t) = 1 - \max_j (p(j t))$ <ul style="list-style-type: none"> <li><math>t</math> – Node in the decision tree</li> <li><math>j</math> – Class value</li> <li><math>p(j t)</math> – Probability (i.e. relative frequency) of class value <math>j</math> in node <math>t</math></li> </ul> <p><b>Minimum Value:</b> 0 when:  <math>\exists j(p(j t) = 1)</math></p> <p><b>Maximum Value:</b> <math>1 - \frac{1}{n_c}</math> when:  <math>\forall j \left( p(j t) = \frac{1}{n_c} \right)</math></p>

### Stopping Criteria for Decision Tree Induction

<p><b>Three Stopping Criteria for Decision Tree Induction</b></p> <ul style="list-style-type: none"> <li>When all records in a node have the same class value</li> <li>When all records in a node have similar attribute values.</li> <li>Early Termination</li> </ul>	<p><b>Underfitting</b> – When a model is too simple, both training and test errors are large.</p>	<p><b>Overfitting</b> – When a model becomes too complex (e.g. too large a tree), the test error begins to increase even though the training error decreases.</p> <ul style="list-style-type: none"> <li><b>Result:</b> Training error is <b>NOT</b> representative for generalization error.</li> </ul>	<p><b>Causes of Overfitting</b></p> <ul style="list-style-type: none"> <li>Noise</li> <li>Insufficient training records (i.e. lack of representative samples)</li> </ul>
--	---	--	--

Generalization Error Estimation		
<p><b>Resubstitution Error</b> Error on the <b>training</b> set.</p> <p>Single Leaf Node Error: <math>e(t)</math> Total Resubstitution Error: <math>e(T)</math></p> $e(T) = \sum e(t)$	<p><b>Generalization Error</b> Error on the <b>testing</b> data.</p> <p>Single Leaf Node Error: <math>e'(t)</math> Total Generalization Error: <math>e'(T)</math></p> $e'(T) = \sum e'(t)$	<p><b>Optimistic Estimation</b> Training error is equal to the testing error.  <math display="block">\sum e(t) = \sum e'(t)</math></p> <p><b>Pessimistic Estimation</b> Assign a penalty term to ea.  <math display="block">e'(t) = e(t) + 0.5</math></p> <p><b>Total Pessimistic Error</b>  <math display="block">e'(T) = \sum (e(t)) + N \cdot 0.5</math></p> <p><math>N</math> – Number of leaf nodes.</p> <p><b>Reduced Error Pruning</b> Use a validation set to estimate the generalization error.</p>

<p><b>Occam's Razor</b></p> <p>Given two models with similar generalization errors, one should prefer the simpler model over the more complex model.</p> <p><b>This is because more complex model has a greater chance of fitting accidentally by errors in the data.</b></p>	<p><b>Pre-pruning (Early Stopping Rule)</b></p> <ul style="list-style-type: none"> <li>Stop the induction algorithm before it becomes a full tree.</li> <li><b>Typical Stopping Rules:</b> <ul style="list-style-type: none"> <li>All remaining records have the same class value</li> <li>All attribute values are the same.</li> </ul> </li> <li><b>More restrictive conditions:</b> <ul style="list-style-type: none"> <li>Number of instances is below a user-specified threshold.</li> <li>Expanding the current node does not improve impurity measures (e.g. GINI Index, Information Gain)</li> <li>Class distribution of instances are independent of available features.</li> </ul> </li> </ul>	<p><b>Post-pruning (Early Stopping Rule)</b></p> <ul style="list-style-type: none"> <li>Grow the decision tree to its entirety.</li> <li>Trim nodes in the tree in a <b>bottom-up fashion</b>.</li> <li>Only trim nodes if by trimming the estimate of the generalization error improves.</li> <li>New leaf node's class label is determined from the majority class of instances in the merged node.</li> </ul>
---	--	--

## 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!**

Diagram illustrating a decision tree structure for post-pruning. The root node is labeled "A?". It branches into four child nodes: A1, A2, A3, and A4. Each child node has associated class counts for "Class = Yes" and "Class = No".

Node	Class = Yes	Class = No
A1	8	4
A2	3	4
A3	4	1
A4	5	1

## Examples of Post-pruning

– Optimistic error?

Don't prune for both cases

Case 1:

Diagram illustrating Case 1 for optimistic error pruning. The root node branches into two child nodes. The left child node has counts: C0: 11, C1: 3. The right child node has counts: C0: 2, C1: 4.

– Pessimistic error?

Don't prune case 1, prune case 2

– Reduced error pruning?

Depends on validation set

Case 2:

Diagram illustrating Case 2 for reduced error pruning. The root node branches into two child nodes. The left child node has counts: C0: 14, C1: 3. The right child node has counts: C0: 2, C1: 2.

## Handling Missing Attribute Values

### Issues Associated with Missing Attribute Values

- Affects how impurity measures are computed
- Affects how to distribute instances with missing value to child nodes.
- Affects how to test instance with missing value is classified.

### Computing Impurity Measure

- Calculate entropies (i.e. information gain) with element with missing value **EXCLUDED**.
- Multiply by scalar of elements included over total number of elements (in below example 9 elements included over 10 total elements hence 0.9):

### 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

- Split the missing record between the two child nodes
- Percentage of child node that goes to each child is portion to the relative frequency of that attribute value.

### 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	Cheat=Yes 2
Class=No 3	Cheat=No 4

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

  

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

### Classifying New/Unseen Records with Missing Data

- Pick the most likely of child nodes and use continue down that portion of the tree.

### 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

**Data Fragmentation** – At each level of the tree, the number of instances gets smaller. At leaf nodes, the number of instances could be too small to be statistically significant.

### Tree Induction: NP Hard

#### Alternate Strategies

- Bottom Up Tree Generation
- Bidirectional Tree Generation
  - Inside-out Bidirectional
  - Outside-in Bidirectional

**Decision Boundary** – Borderline between two neighboring regions of different classes. In non-oblique decision trees, this is parallel to access since it involves a single attribute at a time.

**Oblique Decision Tree** – Test condition in a node may involve multiple attributes.

- Advantage** – Most expressive decision tree
- Disadvantage** – Finding optimal test condition is computationally expensive.

**Expressiveness** – Decision trees do not generalize well to certain types of functions including a parity function which would require a complete tree.

**Tree Replication** – In a decision tree, a subtree may appear in multiple branches. This leads to unnecessary memory usage.

## Performance Evaluation

- Focus on the predictive capability of a model.

### Confusion Matrix

		Predicted Class	
		Class = Yes	Class=No
Actual Class	Class = Yes	a	b
	Class=No	c	d

- a – True Positive (TP)
- b – False Negative (FN)
- c – False Positive (FP)
- d – True Negative (TN)

### Accuracy

$$Accuracy = \frac{A + D}{A + B + C + D}$$

$$Accuracy = \frac{TP + TN}{TP + FN + FP + TN}$$

- Accuracy only tells part of the story.
  - Example: Two Class Problem
    - Number of Class 0 Examples: 9990
    - Number of Class 1 Examples: 10
    - If the model predicts everything as class 0, its accuracy is 99.9% but it cannot detect any class 1.

### Cost Matrix

Actual Class	Predicted Class	
	Class = Yes	Class=No
	Class = Yes	C(Y Y)
Class=No	C(Y N)	C(N N)

- $C(j|k)$  – Cost of predicting class “j” given the actual class is “k”.

$$TotalCost = a \cdot C(Y|Y) + b \cdot C(N|Y) + c \cdot C(Y|N) + d \cdot C(N|N)$$

- Cost matrix can be a better performance evaluation as it accounts for different costs of depending on the type of error.

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

- Precision** – Accuracy of positive predictions. Biased towards  $C(Y|Y)$  &  $C(Y|N)$ .
  - a – True positive.
  - c – False positive.

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

- Precision** – Accuracy of records with positive class value. Biased towards  $C(Y|Y)$  &  $C(N|Y)$ .
  - a – True positive.
  - b – False negative.

$$F\_Measure(F) = \frac{2 \cdot r \cdot p}{r + p} = \frac{a + c}{2 \cdot a + b + c}$$

- F-Measure** – Biased two all except  $C(N|N)$  (i.e. true negative)
  - r – Recall
  - p – Precision
  - c – False Positive

$WA = \frac{w_1 \cdot a + w_4 \cdot d}{w_1 \cdot a + w_2 \cdot b + w_3 \cdot c + w_4 \cdot d}$ <p> <math>n</math> – Number of instances covered by rule  <math>n_c</math> – Number of positive instances covered by rule. </p>	<p><b>Proportionality of Cost and Accuracy</b></p> <ul style="list-style-type: none"> <li>Cost and accuracy are proportional if:</li> </ul> $\mathcal{C}(Y Y) = \mathcal{C}(N N) \text{ and } \mathcal{C}(Y N) = \mathcal{C}(N Y)$	<p><b>Sample Size and Model Performance</b></p> <ul style="list-style-type: none"> <li><b>Learning Curve</b> – Shows how model accuracy changes (and varies) with sample size.</li> <li><b>Effects of Small Sample Size:</b> <ul style="list-style-type: none"> <li>Bias in the estimate</li> <li>Variance in the estimate.</li> </ul> </li> </ul>
--	--	--

## Methods for Model Comparison

<p><b>Holdout</b> – Reserve 2/3 of labelled examples for training and 1/3 for testing.</p> <p><b>Disadvantages</b></p> <ul style="list-style-type: none"> <li>Uses on a subset of the labelled examples when training the model.</li> <li>Model <b>dependent on the composition of the training and test sets</b>.</li> <li><b>Training and test sets are not independent</b> since come from same original set. If one class value is over- or under-represented in either set, it will skew the results.</li> </ul>	<p><b>Random Subsampling</b> – Repeats the whole out method multiple times with replacement.</p> <p><b>Disadvantages:</b></p> <ul style="list-style-type: none"> <li>Still <b>uses only a subset of the labelled examples</b> to build the model.</li> <li><b>No control of how many times each record appears in the training and test sets</b>. If a particular record is always in the training set, it may skew the model.</li> </ul> <p><b>Accuracy of <math>k</math> Random Subsamplings</b></p> $acc_{sub} = \frac{1}{k} \cdot \sum_{i=1}^k acc_i$ <ul style="list-style-type: none"> <li><math>k</math> – Number of iterations</li> <li><math>acc_i</math> – Accuracy of the <math>i^{th}</math> iteration.</li> </ul>	<p><b>Cross Validation</b> – Partition the labelled dataset into <math>k</math> disjoint subsets.</p> <ul style="list-style-type: none"> <li><b>k-Fold</b> – Train on <math>k-1</math> partitions and test on the remaining one.</li> <li><b>Leave-One-Out</b> – The number of partitions equals the number of training samples.</li> </ul> <p><b>Accuracy of <math>k</math>-Fold Cross Validation</b></p> $acc_{sub} = \frac{1}{k} \cdot \sum_{i=1}^k acc_i$ <ul style="list-style-type: none"> <li><math>k</math> – Number of iterations</li> <li><math>acc_i</math> – Accuracy of the <math>i^{th}</math> iteration.</li> </ul> <p><b>Disadvantages:</b></p> <ul style="list-style-type: none"> <li><b>Computationally expensive</b> as process is repeated <math>k</math> times.</li> <li>Depending on size of partition (e.g. 1 for Leave-One-Out), <b>accuracy from iteration to iteration</b> can vary significantly.</li> </ul>
<p><b>Bootstrap</b> –</p>		

Minimum Description Length	
----------------------------	--

# Chapter #05 – Additional Classification Techniques

## Rule-Based Classifiers

<p><b>Classifies records using a collection of “if...then...” rules.</b></p> <p><b>Form of Rule:</b></p> <p><math>(Condition) \rightarrow y</math></p> <ul style="list-style-type: none"> <li><b>Condition (Antecedent, LHS)</b> – Conjunction of attributes.</li> <li><b>y (Consequent, RHS)</b> – Class value.</li> </ul>	<p><b>Cover</b> – A rule <math>r</math> covers an instance <math>x</math> if the attributes of <math>x</math> satisfy the condition (LHS) of the rule.</p> <p><b>Coverage of a Rule</b> – Fraction of records that satisfy the antecedent of a rule.</p> <p><b>Accuracy of a Rule</b> – For records covered by a rule, it is the fraction of records that have the matching class value.</p>	<p><b>Mutually Exclusive Rule Set</b> – Rules in the set are independent of each other such that <b>each record is covered by at most one rule.</b></p> <p><b>Exhaustive Rule Set</b> – A set of rules that covers every possible combination of attribute values. Hence, <b>each record is covered by at least one rule.</b></p>	<p><b>Decision Tree</b> – Can be used to form a mutually exclusive, exhaustive rule set.</p> <p>Rules in a decision tree can be simplified.</p> <p><b>Effects of rule simplification:</b></p> <ul style="list-style-type: none"> <li><b>Problem #1:</b> Rules become non-mutually exclusive.</li> <li><b>Solution:</b> <ul style="list-style-type: none"> <li><b>Ordered Rule Set</b> – Rules ordered from highest to lowest priority. Records classified according to highest priority rule they satisfy.</li> <li><b>Unordered Rule Set</b> – Voting scheme</li> </ul> </li> <li><b>Problem #2:</b> Rules become non-exhaustive.</li> <li><b>Solution:</b> Use a default class.</li> </ul>
---	--	---	--

<p><b>Rule Ordering Schemes</b></p> <p><b>Rules Based Ordering</b> – Individual rules are ranked based off their quality.</p> <ul style="list-style-type: none"> <li><b>Advantage:</b> Ensures each record is classified by the “best rule” covering it.</li> <li><b>Disadvantage:</b> Interpreting lower priority rules becomes more difficult as they are negations of higher priority rules.</li> </ul> <p><b>Class-Based Ordering</b> – Rules that belong to the same class appear together.</p> <ul style="list-style-type: none"> <li><b>Advantage:</b> Simplifies rule ordering.</li> <li><b>Disadvantage:</b> May allow a lower quality rule to have higher priority than a higher quality one.</li> </ul>	<p><b>Direct Method for Rule Building</b> – Extract rules directly from the data.</p> <ul style="list-style-type: none"> <li><b>Examples:</b> RIPPER, CN2, Holte’s 1R</li> </ul> <p><b>Indirect Method for Rule Building</b> – Extract rules from other classification models (e.g. decision tree, neural network, etc.)</p> <ul style="list-style-type: none"> <li><b>Examples:</b> C4.5rules</li> </ul>	<p><b>Sequential Covering Algorithm</b></p> <ol style="list-style-type: none"> <li>Start with an empty rule set.</li> <li>Grow a rule using the “Learn-One-Rule” function.</li> <li>Remove training records covered by the rule.</li> <li>Repeat steps #2 and #3 until stopping criterion is met.</li> </ol> <p><b>Aspects of Sequential Covering</b></p> <ol style="list-style-type: none"> <li>Rule Growing</li> <li>Instance Elimination</li> <li>Rule Evaluation</li> <li>Stopping Criterion</li> <li>Rule Pruning</li> </ol>	<p><b>Rule Growing Strategies</b></p> <ol style="list-style-type: none"> <li><b>General to Specific</b> <ol style="list-style-type: none"> <li>Example: Ripper</li> </ol> </li> <li><b>Specific to General</b></li> </ol> <p><b>CN2 Algorithm</b></p> <ol style="list-style-type: none"> <li>Start from an empty rule</li> <li>Add conjuncts that minimize the entropy measure.</li> <li>Determine the rule consequent by taking majority class of covered instances.</li> </ol>
--	---	---	--

<p><b>Instance Elimination</b></p> <ul style="list-style-type: none"> <li><b>Reason for Eliminating Instances</b> – Otherwise next rule is identical to previous rule.</li> <li><b>Reason for Removing Positive Instances</b> – To ensure future rules are different.</li> <li><b>Reason for Removing Negative Instances</b> – Prevent underestimating accuracy of the rule.</li> </ul>	<p><b>Stopping Criterion</b></p> <p>Compute the information gain with the rule. If the gain is insignificant, discard the rule.</p>	<p><b>Rule Pruning</b></p> <ul style="list-style-type: none"> <li><b>Similar to post-pruning of decision trees.</b></li> <li>Uses reduced error pruning. <ul style="list-style-type: none"> <li>Remove one of the conjuncts of the rule.</li> <li>Compare error rate on validation set before and after pruning.</li> <li>If error improves, remove the conjunct.</li> </ul> </li> </ul>	<p><b>Rule Simplification</b> – Used to reduce the likelihood of overfitting.</p>
---	---	--	---

## Rule Evaluation Metrics

<p><b>Accuracy</b> = <math>\frac{n_c}{n}</math></p> <p><math>n</math> – Number of instances covered by rule  <math>n_c</math> – Number of positive instances covered by rule.</p>	<p><b>Laplace</b> = <math>\frac{n_c + 1}{n + k}</math></p> <p><math>n</math> – Number of instances covered by rule  <math>n_c</math> – Number of positive instances covered by rule.  <math>k</math> – Number of classes  <b>Used to ensure greater coverage for a rule.</b></p>	<p><b>m_estimate</b> = <math>\frac{n_c + p \cdot k}{n + k}</math></p> <p><math>n</math> – Number of instances covered by rule  <math>n_c</math> – Number of positive instances covered by rule.  <math>k</math> – Number of classes  <math>p</math> – Prior probability of positive class.</p>
---	--	--

<p><b>FOIL Information Gain</b></p> <p><math>Gain(R0, R1) = p_1 \cdot \left( \log_2 \left( \frac{p_1}{p_1 + n_1} \right) + \log_2 \left( \frac{p_0}{p_0 + n_0} \right) \right)</math></p> <p><math>R0</math> – Initial Rule  <math>R1</math> – Modified version of <math>R0</math> with added conjunct  <math>t</math> – Number of positive instances covered by both <math>R0</math> and <math>R1</math>  <math>p_0</math> – Positive instances covered by <math>R0</math>  <math>n_0</math> – Negative instances covered by <math>R0</math>  <math>p_1</math> – Positive instances covered by <math>R1</math>  <math>n_1</math> – Negative instances covered by <math>R1</math></p>	<p><b>RIPPER Algorithm</b></p> <ol style="list-style-type: none"> <li>For two classes, define one class as positive class and other as negative class. <ol style="list-style-type: none"> <li>In two class problem, <b>negative class is the default class.</b></li> </ol> </li> <li>In multi-class problem, create list of classes ordered by increasing prevalence. <ol style="list-style-type: none"> <li>Select smallest as first as positive class and rest are negative class.</li> <li>Learn rules for the smallest class first.</li> <li>Repeat with next smallest class.</li> </ol> </li> </ol>	<p><b>RIPPER Algorithm – Growing a Rule</b></p> <ol style="list-style-type: none"> <li>Start from an empty rule set.</li> <li>Add conjuncts as long as they improve <b>FOIL Information Gain</b> (i.e. <b>General-to-Specific</b>).</li> <li>Stop adding conjuncts when the rule starts covering negative examples.</li> <li>Begin pruning the rule immediately (i.e. before generating new rules) using Reduced Error Pruning.</li> <li>Delete conjuncts to maximize <math>v</math> as defined by: <math display="block">v = \frac{p - n}{p + n}</math> </li> </ol> <p><math>p</math> – Number of positive instances covered by the rule.  <math>n</math> – Number of negative instances covered by the rule.</p>
---	--	--

RIPPER Algorithm – Building the Rule Set	C4.5rules – Indirect Method	
<ol style="list-style-type: none"> <li>1. <b>Use Sequential Covering</b> <ol style="list-style-type: none"> <li>a. Find the rule that best covers the current set of positive examples.</li> <li>b. Eliminate both positive and negative examples covered by the rule.</li> <li>c. <b>Uses ordered rule set with class based ordering.</b></li> </ol> </li> <li>2. Each time a rule is added to the rule set, compute the new description length. <b>Example Stopping Conditions:</b> <ol style="list-style-type: none"> <li>a. Stopping growing the rule set if the new <b>rule increases the description length of the rule set by more than <math>d</math></b> (e.g. 64) <b>bits</b>.</li> <li>b. Stop if the error rate of the rule <b>on the validation set</b> is more than 50%.</li> </ol> </li> </ol>	<ol style="list-style-type: none"> <li>1. Start from an <b>unpruned</b> decision tree.</li> <li>2. For each rule <math>r: A \rightarrow y</math>,           <ol style="list-style-type: none"> <li>a. Consider an alternative rule <math>r': A' \rightarrow y</math> where <math>A'</math> is obtained by <b>removing one of the conjuncts</b> of <math>A</math></li> <li>b. <b>Keep the rule with the lowest pessimistic error rate</b> (assuming it is less than the original).</li> <li>c. Repeat until it is no longer possible to improve the generalization error.</li> </ol> </li> </ol>	<ol style="list-style-type: none"> <li>3. Use class-based ordering for the rule set (i.e. group by the rule consequent).</li> <li>4. Compute the description length of each class and order the rules by increasing description length.</li> </ol> <p><b>DescriptionLength = <math>L(\text{error}) + g \cdot L(\text{model})</math></b></p> <ul style="list-style-type: none"> <li>• <b><math>L(\text{error})</math></b> – Number of bits required to encode misclassified examples.</li> <li>• <b><math>L(\text{model})</math></b> – Number of bits required to encode the model.</li> <li>• <b><math>g</math></b> – Tuning parameter whose default is 0.5 and takes into account the presence of redundant attributes in the rule set.</li> </ul>

## Nearest Neighbor Classifiers

<p><b>Instance-Based Classifier</b> – Stores all training records and uses the training records directly to predict the class label of unseen records.</p> <p><b>Rote-Learner</b> – Memorizes the entire training set and performs classification only if attributes of a record match one the training examples exactly.</p> <p><b>Nearest Neighbor</b> – Uses <math>k</math> “closest” training records (i.e. nearest neighbors) for performing classification.</p>	<p><b>Nearest Neighbor Classifier Requirements</b></p> <ol style="list-style-type: none"> <li>1. <b>Set of stored labelled records.</b></li> <li>2. <b>Distance metric to compute distance between records.</b></li> <li>3. <b>Value of <math>k</math>, the number of nearest neighbors to retrieve.</b></li> </ol>	<p><b>Classifying an Unseen Record</b></p> <ol style="list-style-type: none"> <li>1. Compute the distance to all other training records.</li> <li>2. Identify the <math>k</math> nearest neighbors.</li> <li>3. Use class labels of nearest neighbors to determine the class label of unknown records</li> </ol>
---	---	--

<p><b>Voronoi Diagram</b> – Used to depict the decision boundaries for a Nearest Neighbor classifier.</p> <p><b>Euclidean Distance</b></p> $d(p, q) = \sqrt{\sum_i (p_i - q_i)^2}$ <p><b>Manhattan Distance</b></p> $d(p, q) = \sum_i  p_i - q_i $	<p><b>Determining the Class from the Nearest Neighbor List</b></p> <p><b>Option #1</b> – Take the majority vote among the <math>k</math>-Nearest Neighbors.</p> <p><b>Option #2</b> – Weight the vote according to the distance using the weight factor:</p> $\text{WeightFactor} = \frac{1}{d^2}$	<p><b>Effect of the Value of <math>k</math></b></p> <ul style="list-style-type: none"> <li>• <b><math>k</math> is too Small</b> – Underfitting and the classifier becomes sensitive to noise points.</li> <li>• <b><math>k</math> is too Large</b> – Overfitting and the neighborhood make include points from other classes.</li> </ul>	<p><b>Attribute Value Scaling Issues</b></p> <ul style="list-style-type: none"> <li>• Attributes may have to be scaled to normalize for different attribute ranges and values.</li> <li>• This is done to prevent one of the attributes dominating the distance measure.</li> </ul>
--	--	--	---

## PEBLs

<p><b>PEBLs</b></p> <ul style="list-style-type: none"> <li>• Nearest neighbor algorithm that works with both continuous and nominal features.</li> <li>• Each record is assigned a weight factor.</li> <li>• Number of nearest neighbors, <math>k = 1</math></li> </ul> <p><b>Weighted Euclidean Distance</b></p> $d(p, q) = \sqrt{\sum_i w_i \cdot (p_i - q_i)^2}$ <p><math>w_i</math> – Weight of parameter <math>i</math></p>	<p><b>Distance Between Nominal Attributes Value Difference Metric</b></p> $d(v_1, v_2) = \sum_i \left  \frac{n_{1,i}}{n_1} - \frac{n_{2,i}}{n_2} \right $ <p> <math>v_1</math> – Attribute value 1  <math>v_2</math> – Attribute value 2  <math>i</math> – The <math>i^{th}</math> class value.  <math>n_{1,i}</math> – Number of records with attribute value 1 and class value <math>i</math>  <math>n_{2,i}</math> – Number of records with attribute value 2 and class value <math>i</math>  <math>n_1</math> – Total number of records with attribute value 1  <math>n_2</math> – Total number of records with attribute value 2         </p>	<p><b>Similarity Function Used in PEBLS</b></p> $d(X, Y) = w_X \cdot w_Y \sum_{i=1}^k d(X_i, Y_i)^2$ <p> <math>X</math> &amp; <math>Y</math> – Two records  <math>w_X</math> – Distance weighting factor for record <math>X</math>  <math>w_Y</math> – Distance weighting factor for record <math>Y</math>. If <math>Y</math> is an unseen record, then <math>w_Y = 1</math>  <math>d(X_i, Y_i)</math> – Distance between records <math>X</math> and <math>Y</math> in the <math>i^{th}</math> dimension.         </p> <p><b><math>w_X = \frac{\text{Number of Times } X \text{ is Used in Prediction}}{\text{Number of Times } X \text{ Predicts Correctly}}</math></b></p> <p>If <math>w_X \cong 1</math>, then <math>X</math> makes an accurate prediction most of the time. If <math>w_X &gt; 1</math>, then <math>X</math> does not make reliable predictions.</p>
--	--	---



## Bayesian Classifiers

<p><b>Condition Probability Review</b>  <math>P(C A)</math> – Probability of <math>C</math> given <math>A</math>.</p> $P(C A) = \frac{P(C \cap A)}{P(A)}$ $P(A C) = \frac{P(C \cap A)}{P(C)}$	<p><b>Bayes Theorem</b></p> $P(A C) = \frac{P(C A) \cdot P(A)}{P(C)}$ <p><b>Bayes Classifier</b> – A probabilistic framework for solving classification problems.</p>	<p><b>Bayes Theorem Example</b></p> $P(SN M) = 0.5$ $P(M) = \frac{1}{50000}$ $P(SN) = \frac{1}{20}$ <p>Hence:</p> $P(M SN) = \frac{P(SN M) \cdot P(M)}{P(SN)} = \frac{0.5 \cdot \frac{1}{50000}}{\frac{1}{20}} = 0.0002$	<p><b>Requirement of Naïve Bayesian Classifiers</b></p> <ul style="list-style-type: none"> <li>Consider <b>each attribute</b> (<math>A_1, A_2, \dots, A_n</math>) <b>as independent random variables</b>.</li> <li>Consider the <b>class (<math>C</math>) label as a random variable</b>.</li> </ul> <p><b>Goal</b> is to find:  <math>P(C A_1, A_2, \dots, A_n)</math></p>
---	---	--	---

<p><b>Multi-Attribute Bayes Theorem</b></p> $P(C A_1, A_2, \dots, A_n) = \frac{P(C A_1, A_2, \dots, A_n) \cdot P(C)}{P(A_1, A_2, \dots, A_n)}$ <p><b>Classification Approach:</b> Select the value of <math>C</math> that maximizes the above equation.</p>	<p><b>Naïve Bayesian Simplification</b></p> $P(C A_1, A_2, \dots, A_n) = P(C A_1) \cdot P(C A_2) \cdot \dots \cdot P(C A_n)$ <p><b>This equations assumes independence among each <math>A_i</math> attribute when the class is given.</b></p>	<p><b>Estimating the Class Probability</b></p> $P(C) = \frac{N_C}{N}$ <p><math>C</math> – Class value  <math>N_C</math> – Number of training records with class value <math>C</math>  <math>N</math> – Total number of training records.</p>
---	---	--

<p><b>Estimating the Attribute-Class Probability</b></p> $P(A_i C) = \frac{ A_{i,C} }{N_C}$ <ul style="list-style-type: none"> <li><math>C</math> – Class value</li> <li><math>N_C</math> – Number of training records with class value <math>C</math></li> <li><math> A_{i,C} </math> – Number of training records with class value <math>C</math> and attribute value <math>A_i</math>.</li> </ul>	<p><b>Handling Continuous Variables</b></p> <p><b>Option #1: Discretize the continuous range into bins.</b> This makes a series of ordinal attribute values. Conditional probability is estimated by the number of records that fall in each bin.</p> <p><b>Simplest Approach:</b> Use a two-way (i.e. binary) split.</p>	
	<p><b>Standard Deviation (<math>\sigma</math>)</b></p> $\sigma = \sqrt{\frac{\sum_{i=1}^n (x_i - \mu)^2}{n - 1}}$	<p><b>Option #2:</b> Assume attribute follows a normal distribution and use that mean (<math>\mu</math>) and standard deviation (<math>\sigma</math>) to estimate the conditional probability.</p> $P(A_i = a_i   C = c_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} \cdot e^{-\frac{(a_i - \mu_{ij})^2}{2\sigma_{ij}^2}}$ <p><math>a_i</math> – Continuous value for attribute  <math>c_j</math> – Class value  <math>\mu_{ij}</math> – Mean for value for records with attribute <math>A_i</math> and class value <math>c_j</math>  <math>\sigma_{ij}</math> – Standard deviation for value for records with attribute <math>A_i</math> and class value <math>c_j</math></p>

### Handling Zero Value Conditional Probabilities – Further Conditional Probability Estimation

<p><b>Original/Standard</b></p> $P(A_i C) = \frac{ A_{i,C} }{N_C}$ <p><math> A_{i,C} </math> – Number of records with attribute value <math>A_i</math> and class value <math>C</math>  <math>N_C</math> – Number of records with class value <math>C</math></p>	<p><b>Laplace</b></p> $P(A_i C) = \frac{ A_{i,C}  + 1}{N_C + k}$ <p><math> A_{i,C} </math> – Number of records with attribute value <math>A_i</math> and class value <math>C</math>  <math>N_C</math> – Number of records with class value <math>C</math>  <math>k</math> – Number of classes</p>	$P(A_i C) = \frac{ A_{i,C}  + mp}{N_C + m}$ <p><math> A_{i,C} </math> – Number of records with attribute value <math>A_i</math> and class value <math>C</math>  <math>N_C</math> – Number of records with class value <math>C</math>  <math>p</math> – User specified “prior probability.” Most important when <math> A_{i,C}  = 0</math>. Between 0 and 1.  <math>m</math> – Equivalent sample size. Used balance between <math>p</math> and <math>\frac{ A_{i,C} }{N_C}</math></p>
---	---	--

	<b>ROC Curve</b>	
	<ul style="list-style-type: none"> <li>Used to illustrate the performance of a binary classifier.</li> <li>Two Dimensional <ul style="list-style-type: none"> <li><b>X-Axis</b> – False Positive Rate</li> <li><b>Y-Axis</b> – True Positive Rate</li> </ul> </li> </ul>	

## Comparison of Classification Algorithms

<p style="text-align: center;"><b>Decision Tree Algorithm</b></p> <p><b>Advantages</b></p> <ul style="list-style-type: none"> <li>• <b>Inexpensive to construct</b></li> <li>• <b>Extremely fast at classifying unknown records.</b></li> <li>• <b>Easy to interpret for small sized trees.</b></li> <li>• <b>Accuracy is comparable to other classification techniques for many simple datasets.</b> (Since everything comes right from the data)</li> </ul> <p><b>Disadvantages</b></p> <ul style="list-style-type: none"> <li>• <b>May not generalize well for certain types of functions</b> (e.g. Parity function requires a complete tree)</li> <li>• <b>May be insufficient for modelling continuous variables</b> that do not allow oblique nodes.</li> </ul>	<p style="text-align: center;"><b>Rule Based Classifiers</b></p> <p><b>Advantages</b></p> <ul style="list-style-type: none"> <li>• <b>As highly expressive as decision trees</b> <ul style="list-style-type: none"> <li>◦ A decision tree can be expressed via rules based classifier).</li> <li>◦ Allows for more complex models than a decision tree by allowing multiple rules to trigger on a single rule.</li> </ul> </li> <li>• <b>Easy to interpret.</b></li> <li>• <b>Easy to generate.</b></li> <li>• <b>Can classify new records quickly.</b></li> <li>• <b>Performance comparable to decision trees</b></li> <li>• <b>Well suited for handling data sets with imbalanced class distributions.</b></li> </ul>	<p style="text-align: center;"><b>Rule Based Classifiers</b></p> <p><b>Advantages</b></p> <ul style="list-style-type: none"> <li>• <b>Lazy Learner</b> – Does not require the building of a complex model.</li> <li>• <b>Can create complex decision boundaries</b> unlike rule-based and decision trees which generally create rectilinear boundaries.</li> </ul> <p><b>Disadvantages</b></p> <ul style="list-style-type: none"> <li>• <b>Each unseen records are computationally expensive</b> (since must be compared to all training records).</li> <li>• <b>Susceptible to wrong prediction without appropriate proximity measure and preprocessing is done.</b></li> <li>• Uses local data to make classification decisions so potentially <b>susceptible to noise</b>.</li> </ul>
<p style="text-align: center;"><b>Bayesian Classifier</b></p> <p><b>Advantages</b></p> <ul style="list-style-type: none"> <li>• <b>Robust to isolated noise points.</b></li> <li>• <b>Handles missing values by ignoring them in the probability estimate calculations.</b></li> <li>• <b>Robust to irrelevant attributes.</b></li> </ul> <p><b>Disadvantages</b></p> <ul style="list-style-type: none"> <li>• <b>Independence assumption may not hold for some attributes</b> (in such cases, must use a technique known as Bayesian Belief Networks).</li> </ul>		