

# 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>.</li> <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> </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: $\forall j \left( p(j t) = \frac{1}{n_c} \right)$	<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:</li> </ul> $n = \sum_{i=1}^k n_i$ <ul style="list-style-type: none"> <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: $\forall i \left( GINI(i) = 1 - \frac{1}{n_c} \right)$	<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: $\forall j \left( p(j t) = \frac{1}{n_c} \right)$
--	--	---

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 (including non-representative training set)</li> </ul>
--	---	--	--

Resubstitution Error Error on the <b>training</b> set.	Generalization Error Error on the <b>testing</b> data.	Generalization Error Estimation		
		<p><b>Optimistic Estimation</b> Training error is equal to the testing error.</p> $\sum e(t) = \sum e'(t)$	<p><b>Pessimistic Estimation</b> Assign a penalty term to ea.</p> $e'(t) = e(t) + 0.5$ <p><b>Total Pessimistic Error</b></p> $e'(T) = \sum(e(t)) + N \cdot 0.5$ <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>Single Leaf Node Error:</b> <math>e(t)</math></p> <p><b>Total Resubstitution Error:</b> <math>e(T)</math></p> $e(T) = \sum e(t)$	<p><b>Single Leaf Node Error:</b> <math>e'(t)</math></p> <p><b>Total Generalization Error:</b> <math>e'(T)</math></p> $e'(T) = \sum e'(t)$			

<p><b>Occam's Razor</b> 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>
--	--	--

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

## Miscellaneous

<p><b>Decision Tree Algorithm</b></p> <p><b>Advantages</b></p> <ul style="list-style-type: none"><li>• Inexpensive to construct</li><li>• Extremely fast at classifying unknown records.</li><li>• Easy to interpret for small sized trees.</li><li>• Accuracy is comparable to other classification techniques for many simple datasets. (Since everything comes right from the data)</li></ul>		
--	--	--