



**OPIM-410/672:  
Decision Support Systems  
Spring 2008**

**Data Mining Project Report:  
Check List**

**Student(s):** \_\_\_\_\_

**Submission Date:** \_\_\_\_\_

**Total (out of 50 marks)**

\_\_\_ / 50



The example in this tick sheet is under the assumption that you are mining a data set to determine which customers are likely to purchase caravan insurance. You will be mining a different data set, but caravan insurance is used as an example for the purposes of explanation.

<b>Business Objective:</b>	___ / 4
- <u>Objective</u> of the mining study specified.	
- <u>Justification</u> of the business value given.	
- <u>Novelty</u> - what has been done toward this objective in the past	
- <u>Novelty</u> - why my approach is novel	

<b>Data Preparation:</b>	___ / 5
- Possible <u>sources</u> of the data used specified (what internal and external databases may it have come from?)	
- Explained how the data has been <u>prepared</u> by others (e.g. was the data <i>pre-binned</i> ), and what preparation steps you took (if any). Described the implications of this preparation (e.g. in terms of biases it may have introduced).	
- Specified, in an appendix, the <u>independent (input / predictor) variables</u> and their types (e.g. numeric, categorical, binary, etc.).	
- Mentioned what <u>other variables you may have used</u> (and their sources), had they been available.	
- Specified which variable is the <u>dependent (output / predicted) variable</u> .	
- Specified how you divided the data into <u>training and test sets</u> , and why this separation is important.	
- For numeric output data: Performed an exploratory data analysis: <ul style="list-style-type: none"> <li>o produce scatter plots for every combination of input-variable with the output variable (e.g. if trying to predict income, always have income on Y axis, and then produce different scatter plots with each input variables on your X-axis: e.g. age, number of children, etc.)</li> <li>o for each input attribute: found min, max, average, standard deviation, 25<sup>th</sup> and 75<sup>th</sup> percentiles (e.g. use <code>Min()</code>, <code>Max()</code>, <code>Avg()</code>, <code>Stdev()</code>, <code>Percentile()</code> function in Excel), and found outliers.</li> </ul>	

<b>Data Mining:</b>	<b>___ / 15</b>
- Used a <u>data mining technique/algorithm</u> (such as Decision Tree induction, rule induction, or genetic algorithm) to induce a set of predictive rules or formulae. Specified which data mining technique (algorithm) was used (and what parameter values were used, if applicable).	/2
- Described the <u>relevance</u> of the data mining technique employed, and the <u>advantages and disadvantages</u> as compared to other data mining techniques.	/2
- Described what <u>other data mining techniques</u> might have been applied to this data, and what benefits could have been achieved by using them.	/1
- Showed the <u>results</u> of the data mining – specific <u>patterns discovered and interpreted</u> (e.g. showed the decision tree induced, and listed the IF ... THEN ... rules induced; or for neural nets showed the prediction formula that was suggested by the algorithm). E.g. if trying to determine who purchases caravan insurance, did you discover that hedonists or people with fire insurance policies are more likely to buy caravan insurance ?	/2
- Showed a <u>confusion matrix</u> (giving the various error rates) for each model.	/2
- Construct a <u>cost matrix</u> with hypothetical costs for your various classification decisions (true and false positives, true and false negatives), being sure to give row and column headings meaningful to your application (like "Mail / Do not mail" rather than "positive / negative").	/2
- Using your cost matrix, calculate the <u>threshold probability</u> (i.e. 'cut-off response rate') which gives Expected Cost(Predict Caravan Buyer) > Expected Cost(Predict Not Buyer).	/2
- If you used an entropy-based decision tree algorithm (e.g. C4.5 as in the free CTree tool on the internet): Use the threshold probability to <u>re-label</u> your decision tree - i.e. indicate how the most <i>profitable</i> prediction differs from the most <i>probable</i> prediction for each node in the tree. Note that decision tree built in a profit-driven way do not need to be re-labeled (re-labeling will not change the output predictions decision trees built in a profit-driven way).	/2

<b>Analysis of Data Mining Results</b>	<b>___ / 12</b>
This section is predominantly for a categorical output variable (e.g. predicting whether a person is a customer or not): adapt as appropriate for numeric output variables.	
- For each model, showed a <u>Lift Chart</u> . Explained briefly how the lift chart was created and what the meaning of the Lift Chart is.	
- <u>Estimated some costs and revenues</u> for each of your prediction types (true-positive, true-negative, false-positive, false-negative), and showed a <u>Gain Chart</u> depicting how profit varies with a change in the number of prospects contacted.	
- Explained the Gain Chart (how it was created and what it means) and specified how you would decide <u>how many prospects to contact</u> .	
- Explained what <u>other ways of determining cut-off</u> (i.e. number of prospects to mail) are possible ?	
- Showed and explained a <u>Classification Errors Graph</u> (see slides on webcafe)	
- Explain what <u>threshold values</u> are and why you needed to adjust them (improves recall).	
- Did a <u>Sensitivity Analysis</u> to show how the optimal number of prospects to contact varies with changes in costs and revenues.	
- Perform a ROC analysis	
- Plot learning curves	
- Explained why overall accuracy (overall error) is not the best metric for model quality. Described <u>what other factors influence model quality</u> and give an assessment of your model(s).	

<b>Knowledge Exploitation:</b>	____ / 6
- Described how you propose to <u>exploit</u> the knowledge gained.	
- Described what systems you would <u>integrate</u> with.	
- Drew a <u>diagram</u> showing their <u>data mining process</u> from preparation to mining to exploitation.	
- Provided a list of <u>recommendations</u> (i.e. specific actions the business should take such as particular system development projects, further data mining studies, data warehousing projects, etc.). Briefly described specific <u>goals</u> and justified each recommendation.	
<b>Post-Implementation Review: WEKA or other</b>	____ / 3
- Provide an <u>assessment</u> of the software tool used for the project	
- Mentioned <u>bugs and weaknesses</u> found in software tool.	
- Suggested <u>desirable improvements and new features</u> for software tool.	
- Suggested <u>alternatives</u> to software tool (e.g. other data mining tools like Clementine, which may provide better sampling facilities).	
<b>Report Structuring and Formatting:</b>	____ / 5
- <u>headings, bulleting, bolding, and diagrams</u> used.	
- <u>diagrams</u> labeled and numbered and appropriately referred to in text.	
- <u>title page, table of contents, and list of references</u> included.	
- <u>pages numbered and with header</u> giving report title and student name.	
- <u>professional presentation</u> .	
<b>Extra Credit:</b>	____ / 6
- created a new data set that oversampled the Response (i.e. Buyer) rows from your original data set, using the threshold probability to decide the correct proportion of Buyers-vs-NonBuyers in your new data set. (That is, constructed a new data set that is a stratified sample from the original data set, but with a greater proportion of Buyers-to-NonBuyers - the exact proportion of Buyers-to-NonBuyers can be decided using the threshold probability and the formula shown in lecture on Cost Sensitive Learning). Reconstructed the decision tree using the new data set and compare the results to the decision tree you obtained by re-labeling your original decision tree.	
- performed <u>additional analyses</u> on the data beyond what was required.	
- specified how various mining techniques might be <u>combined</u> to improve results.	
- evidence of massive effort.	
- ultra-professional presentation.	