

AI Academy: Introduction to Data Mining

Week 5 Seminar

Seminar 5 contains 2 questions.

1. Consider two models, A and B, for the classification of 100 patients. ”+” indicates positive while ”-” label indicates negative. Both models have the accuracy of 80%. Model A reports 30 positive patients correctly but incorrectly predicted 10 positive patients as negative. Model B reports 25 positive patients correctly but incorrectly predicted 5 positive patients as negative.

(a) (2 points) Fill in the Confusion Matrix for BOTH Model A and Model B.

Model A		<i>Predicted</i>		Model B		<i>Predicted</i>	
		+	-			+	-
<i>Actual</i>	+			<i>Actual</i>	+		
	-				-		

(b) (2 points) Compute **precision** and **recall** for Model A and Model B respectively.

- (c) (2 points) Produce a **cost (earning)** matrix assuming that: a true positive *earns* \$50; a true negative *earns* \$10; a false positive *loses* \$10; and a *false* negative *loses* \$100.

\$ Value		<i>Predicted</i>	
		+	-
<i>Actual</i>	+		
	-		

- (d) (2 points) Given the cost (earning) matrix above, compute **the expected earning** for models A and B respectively. **Which** model is better w.r.t. the expected earning?

2 Decision Stumps, & Cross Validation

Consider the following dataset (9 instances) with **2 binary attributes** (x_1 and x_2), and a **class attribute** y , shown in Table 1. For this question, we will consider a **Decision Stump** classifier. A Decision Stump is a decision tree with a max depth of 1 (i.e. only one split before classification).

Table 1: Data

ID	x_1	x_2	Class
1	True	False	+
2	False	False	-
3	True	True	-
4	False	False	+
5	True	True	-
6	False	True	-
7	False	False	+
8	False	True	+
9	True	False	-

1. By hand, evaluate the Decision Stump classifier, calculating the confusion matrix and testing accuracy (show your work by labeling each data object with the predicted class). You should be able to eyeball which split is best, but if you can't, use the GINI index. If the two splits end up being identical, go with x_1 . If you end up with a split where one class has an even distribution, default to a class assignment of positive.

Use the following evaluation methods:

- (a) A holdout test dataset consisting of last 4 instances
- (b) 3-fold cross-validation, using the following folds with IDs: [1,2,3], [4,5,6], [7,8,9] respectively.
- (c) Leave one out cross validation (LOOCV)