Final examination Machine Learning and Data Mining

ONLY Slides are permitted.

NB: Presentation will be taken into account when given a mark.

Marks are also subject to modification.

Two Parts that should be processed on different sheets of papers. Part I (Mephu Nguifo E.) and Part II (Antoine V.).

Part I

Exercise 1:	Classification	(6 pts)		
 Briefly outline the topic of your homework Outline the major steps of naive bayes classification. Is post-pruning preferable to prepruning in decision tree classification? Explain. Given the dataset below, explain how to build and evaluate a classifier. 				
@attribute @attribute @attribute @attribute	contact-lenses e age e spectacle-prescrip e astigmatism e tear-prod-rate e contact-lenses	<pre>{young, pre-presbyopic, presbyog {myope, hypermetrope} {no, yes} {reduced, normal} {soft, hard, none}</pre>	pic}	
<pre>@data % 24 instances young, myope, no, reduced, none young, myope, no, normal, soft</pre>				
pre-presbyopic, myope, no, reduced, none				
<pre>m presbyopic,hypermetrope,yes,reduced,none presbyopic,hypermetrope,yes,normal,none</pre>				

Exercise 2: Association rules (10 pts)

Let $\{a, b, c, d, e\}$ be a set of items; $\{t_1, t_2, t_3, t_4, t_5, t_6\}$ be a set of transactions described as:

 $t_1 = b, c, d$ $t_2 = a, b, c, d, e$ $t_3 = a, b, c, e$ $t_4 = a, b, d, e$ $t_5 = b, c, e$ $t_6 = a, b, d, e$

The minimum support threshold is to set to 3/6. Let U is a frequent itemset; V is a proper subset of U; and V is not empty.

<u>Theorem</u>: If a rule $V \rightarrow U-V$ does not satisfy the confidence threshold, then any rule $V' \rightarrow U-V'$, where V' is a subset of V, must not satisfy the confidence threshold as well.

1.	Prove the theorem above .	(1 pt)
2.	Find all frequent itemsets using Apriori algorithm. Provide details.	(2 pts)
3.	Draw the lattice of frequent itemsets.	(1 pt)
4.	Find all frequent closed itemsets. Explain.	(1,5 pt)
5.	Find one frequent closed itemset that has two minimal generators. Provide details.	(1,5 pt)
6.	Find all frequent maximal itemsets.	(1 pt)
7.	Give the negative border of frequent itemsets.	(1 pt)
8.	Let the minimal confidence equal to 66%, list three valid rules. Provide details.	(1 pt)

November 2014 - 1/2 -

Part II MANDATORY: Use another copie

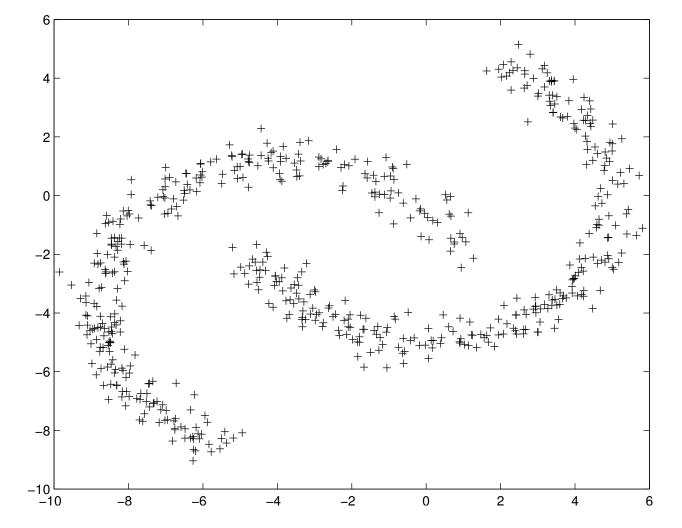
Exercise 3: Clustering

(4 pts)

1. Use the complete-link agglomerative clustering with the Euclidean distance to group the data described as follow:

Show, for each epoch, the dissimilarity matrix and the dendrogram.

- 2. Cut the final dendrogram in order to get two clusters and give the set of objects for each group.
- 3. Name two algorithms enable to find correctly the groups for the dataset below. Explain briefly your choices.



November 2014 - 2/2 -