

School of Computing and Information Systems
The University of Melbourne
COMP30027 MACHINE LEARNING (Semester 1, 2019)

Tutorial exercises: Week 9

1. Let's revisit the logic behind the **voting** method of classifier combination (used in **Bagging**, **Random Forests**, simple **Stacking**, and **Boosting** to some extent): Let's make a few assumptions (some of which we'll try to relax later):
 - (1) We have a two-class problem;
 - (2) The test (and training) instances are roughly evenly divided between the two classes;
 - (3) Our classifiers predict the test instances roughly in proportion to the distribution of the classes;
 - (4) We are building an ensemble out of two classifiers;
 - (5) The errors between the two classifiers are **uncorrelated**.
 - (a) First, let's assume our two classifiers both have an **error rate** of $e = 0.4$, calculated over 1000 instances.
 - i. Build the **confusion matrix** for these classifiers, based on the assumptions above.
 - ii. On the table overleaf, indicate the number of instances in the (count) column for the first two systems — a couple of values have been filled out; for example, there are 180 instances where the actual class is A, and both systems predicted A.
 - iii. Assuming that the voting ties are broken randomly, what the the expected error rate of the voting ensemble?
 - (b) What if we add a third classifier, also with error rate 0.4? Fill in the rest of the table, and determine the error rate of this ensemble. Why has adding a third system caused it to improve?
 - (c) Now consider two classifiers, one (1) with $e_1 = 0.1$ and the second with $e_2 = 0.2$.
 - i. Build the two confusion matrices.
 - ii. Fill out the second table overleaf. Some values are given. Determine the expected error rate of the ensemble.
 - iii. Add another system with $e_3 = 0.2$; does the error rate improve this time?
 - iv. What if the errors between the systems were very highly correlated instead? What will happen to the error rate then? What do you think would happen if we added many more highly correlated classifiers to the ensemble?
 - (d) (Extension) Find general forms for the rightmost values in the tables:
 - i. for N instances and error rates $e_{1,2,3}$;
 - ii. and, instead of the true labels being evenly divided between the two classes, a fraction α of the instances are class A, and $(1 - \alpha)$ are class B;
 - iii. and, instead of the classifier making predictions in the ratio of the true labels, it is potentially biased, predicting class A for a fraction β of the instances, and $(1 - \beta)$ for class B [Hint: the A-A cell in the confusion matrix should be $\frac{N}{2}(\alpha + \beta - e)$]
 - (e) Why can't we easily relax assumption (1) with the information given?

Predictions (all $e = 0.4$)					
	1	2	(count)	3	(count)
A	A	A	180	A	108
				B	
		B		A	
				B	
	B	A		A	
				B	
		B		A	
				B	

Predictions ($e_1 = 0.1, e_2, e_3 = 0.2$)					
	1	2	(count)	3	(count)
A	A	A	360	A	
				B	72
		B		A	
				B	
	B	A	40	A	
				B	
		B		A	
				B	
B	A	A		A	
				B	
		B		A	8
				B	
	B	A		A	
				B	
		B		A	
				B	