

School of Computing and Information  
Systems The University of Melbourne  
COMP30027 Machine Learning (Semester 1, 2021)

Tutorial: Week 12

1. [OPTIONAL] Consider the following dataset:

<i>id</i>	<i>apple</i>	<i>ibm</i>	<i>lemon</i>	<i>sun</i>	<b>label</b>
A	4	0	1	1	fruit
B	5	0	5	2	fruit
C	2	5	0	0	comp
D	1	2	1	7	comp
E	2	0	3	1	?
F	1	0	1	0	?

Treat the problem as an unsupervised machine learning problem (excluding the *id* and *label* attributes) and calculate the clusters according to (hard) **k-means** with  $k = 2$ , using the Manhattan distance, and instances A and F as the seeds.

- Repeat the previous question using “soft” k-means, and the “stiffness”  $\beta = 1$ .
- What is logic behind the EM algorithm, when used for clustering?
  - Explain the significance of the “E” step, and the “M” step.
  - Identify the “E” and “M” steps in GMM methods.
- Revise the concept of Unsupervised and Supervised evaluation for clustering evaluation
  - Explain the two main concepts that we use to measure the goodness of a clustering structure without respect to external information.
  - Explain the two main concepts that we use to measure how well do cluster labels match externally supplied class labels.
- Revise the difference between **supervised**, **semi-supervised** and **unsupervised** machine learning. When do we use semi-supervised learning?
  - What is self-training?
  - What is the logic behind active learning, and what are some methods to choose instances for the oracle?
- One of the strategies for Query sampling was query-by-committee (QBC). Using the equation below, which captures vote entropy, determine the instance that our active learner would select first.

$$x_{VE}^* = \underset{x}{\operatorname{argmax}} \left( - \sum_{y_i} \frac{V(y_i)}{C} \log_2 \frac{V(y_i)}{C} \right)$$

Respectively  $y_i$ ,  $V(y_i)$ , and  $C$  are the possible labels, the number of “votes” that a label receives from the classifiers, and the total number of classifiers.

classifier	Instance 1			Instance 2			Instance 3		
	$y_1$	$y_2$	$y_3$	$y_1$	$y_2$	$y_3$	$y_1$	$y_2$	$y_3$
$c_1$	0.2	0.7	0.1	0.2	0.7	0.1	0.6	0.1	0.3
$c_2$	0.1	0.3	0.6	0.2	0.6	0.2	0.21	0.21	0.58
$c_3$	0.8	0.1	0.1	0.05	0.9	0.05	0.75	0.01	0.24
$c_4$	0.3	0.5	0.2	0.1	0.8	0.1	0.1	0.28	0.62