# COMP3314_2C Machine Learning
## Homework:

Release date: March 12, 2024

Due date: 11:59pm, March 22, 2024

1. Consider a Perceptron with 2 inputs and 1 output. Let the weights of the Perceptron be w1=1 and w2=1 and let the bias be w0=−1.5. Calculate the output of the following inputs:(0, 0), (1, 0), (0, 1), (1, 1). (12 points)

2. Suppose that the following are a set of point in two classes:

   • Class1: (1,1), (1,2), (2,1)

   • Class2: (0,0), (1,0), (0,1)

   (1) Plot them and find the optimal separating line. (10 points)

   (2) What are the support vectors, and what is the meaning? (14 points)

3. Suppose that the probability of five events are P(first) = 0.5, P(second) = P(third) = P(fourth) = P(fifth) = 0.125. Calculate the entropy and write down in words what this means. (14 points)

4. Suppose we collect data for a group of students in a postgraduate machine learning class with features x1 = hours studies, x2 = undergraduate GPA and label y = receive an A. We fit a logistic regression and produce estimated weights as follows: w0=−6, w1=0.05, w2=1.

   (1) Estimate the probability that a student who studies for 40h and has an undergraduate GPA of 3.5 gets an A in the class. (10 points)

   (2) How many hours would the student in part 1. need to study to have a 50% chance of getting an A in the class? (10 points)

5. Given the following dataset:

| V | W | X | Y |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 |
| 1 | 0 | 0 | 1 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 0 |

Your task is to build a decision tree for classifying variable $Y$. (You can think of the dataset as replicated many times, i.e. overfitting is not an issue here).

(1) Compute the information gains $IG(Y|V)$, $IG(Y|W)$ and $IG(Y|X)$. Remember, information gain is defined as

$$IG(D_p) = I_G(D_p) - \sum_{j=1}^{m} \frac{N_j}{N_p} I_G(D_j)$$

where

$$I_G(t) = 1 - \sum_{i=1}^{c} p(i|t)^2$$

$c$ is the class number, $D_p$ and $D_j$ are the dataset of the parent and $j$-th child node. $I_G$ is gini impurity. $N_p$ is the total number of samples at the parent node and $N_j$ is the number of samples in the $j$-th child node. (10 points)

(2) Write down the entire decision tree with gini impurity. (20 points)