

**National University of Singapore
School of Computing
CS3244 Machine Learning**

Tutorial 7: Bayesian Inference

Issue: April 7, 2022

Due: April 11, 2022

Important Instructions:

- *Your solutions for this tutorial must be TYPE-WRITTEN.*
- *Make TWO copies of your solutions: one for you and one to be SUBMITTED TO THE TUTOR IN CLASS. Your submission in your respective tutorial class will be used to indicate your CLASS ATTENDANCE. Late submission will NOT be entertained.*
- *Indicate your NAME, STUDENT NUMBER, and TUTORIAL GROUP in your submitted solution.*
- *YOUR SOLUTION TO QUESTION BL 11 will be GRADED for this tutorial.*
- *You may discuss the content of the questions with your classmates. But everyone should work out and write up ALL the solutions by yourself.*

BL 11 In the solution to question TM 6.1 in Tutorial 5, I have shown you the use of the “incremental” version of Bayes’ rule. One may wonder how this ‘incremental’ version can be derived from the original Bayes’ Theorem. In this question, you are asked to derive the “incremental” version of Bayes’ rule.

Specifically, let $h \in H$. Using the Bayes’ Theorem, we know that

$$P(h|D_1, D_2) = \frac{P(D_1, D_2|h)P(h)}{P(D_1, D_2)}.$$

By assuming the conditional independence of data D_1 and data D_2 given hypothesis h , give a step-by-step derivation of the following “incremental” version of Bayes’ rule:

$$P(h|D_1, D_2) = \frac{P(D_2|h)P(h|D_1)}{\sum_{h \in H} P(D_2|h)P(h|D_1)}$$

and state any result/assumption that you have used in each step.

TM 6.5 Consider the Minimum Description Length principle applied to the hypothesis space H consisting of conjunctions of up to n Boolean attributes (e.g., *Sunny* \wedge *Warm*). Assume each

hypothesis is encoded simply by listing the attributes present in the hypothesis where the number of bits needed to encode any one of the n Boolean attributes is $1 + \log_2 n$ (i.e., $\log_2 n$ bits to indicate which of the n Boolean attributes is present in the hypothesis and 1 bit to indicate its corresponding Boolean value). Suppose that the encoding of an example given the hypothesis uses zero bits if the example is consistent with the hypothesis and uses $1 + \log_2 m$ bits otherwise (i.e., 1 bit to indicate the correct classification and $\log_2 m$ bits to indicate which of the m examples was misclassified).

- Write down the expression for the quantity to be minimized according to the Minimum Description Length principle.
- Is it possible to construct a set of training data such that a consistent hypothesis exists, but MDL chooses an inconsistent hypothesis? If so, give such a set of training data. If not, explain why not.
- Give probability distributions for $P(h)$ and $P(D|h)$ such that the above MDL algorithm outputs MAP hypotheses.

BL 12a (Final Exam AY2020/21) Fig. 1a and Fig. 1b below show perceptron units A and B . Fig. 1c below shows a network C of perceptron units with two hidden layers of two units each, while Fig. 1d below shows a network F of perceptron units with a hidden layer of two units. They are based on the following structure:

- Perceptron units A and B have one (Boolean) output unit each for producing the output o_A and the output o_B , respectively. Similarly, networks C and F of perceptron units have one (Boolean) output unit each for producing the output o_C and the output o_F , respectively.
- There should be four input units (i.e., one input unit for each of the four (Boolean) input attributes x_1, x_2, x_3, x_4).
- A Boolean is **-1 if false**, and **1** if true.
- The activation function of every (non-input) unit is a **-1 to 1 step function** (refer to page 6 of the “Neural Networks” lecture slides), including that of the output unit.
- The weights w_A (i.e., hypothesis) of perceptron unit A and the weights w_B (i.e., hypothesis) of perceptron unit B are indicated in Fig. 1a and Fig. 1b, respectively. The weights w_C (i.e., hypothesis) of network C of perceptron units and the weights w_F (i.e., hypothesis) of network F of perceptron units are indicated in Fig. 1c and Fig. 1d, respectively.
- A bias input is of value 1 and is not considered a hidden unit.

One of the four perceptron networks in Fig. 1 has been used to generate a dataset of 4 Boolean input attributes x_1, x_2, x_3, x_4 and a Boolean target output t_d with the following 3 noise-free training examples of the form $d = \langle (x_1, x_2, x_3, x_4), t_d \rangle$:

$$D = \{d_1 = \langle (-1, -1, -1, 1), 1 \rangle, d_2 = \langle (1, 1, -1, -1), 1 \rangle, d_3 = \langle (1, 1, 1, 1), -1 \rangle\}.$$

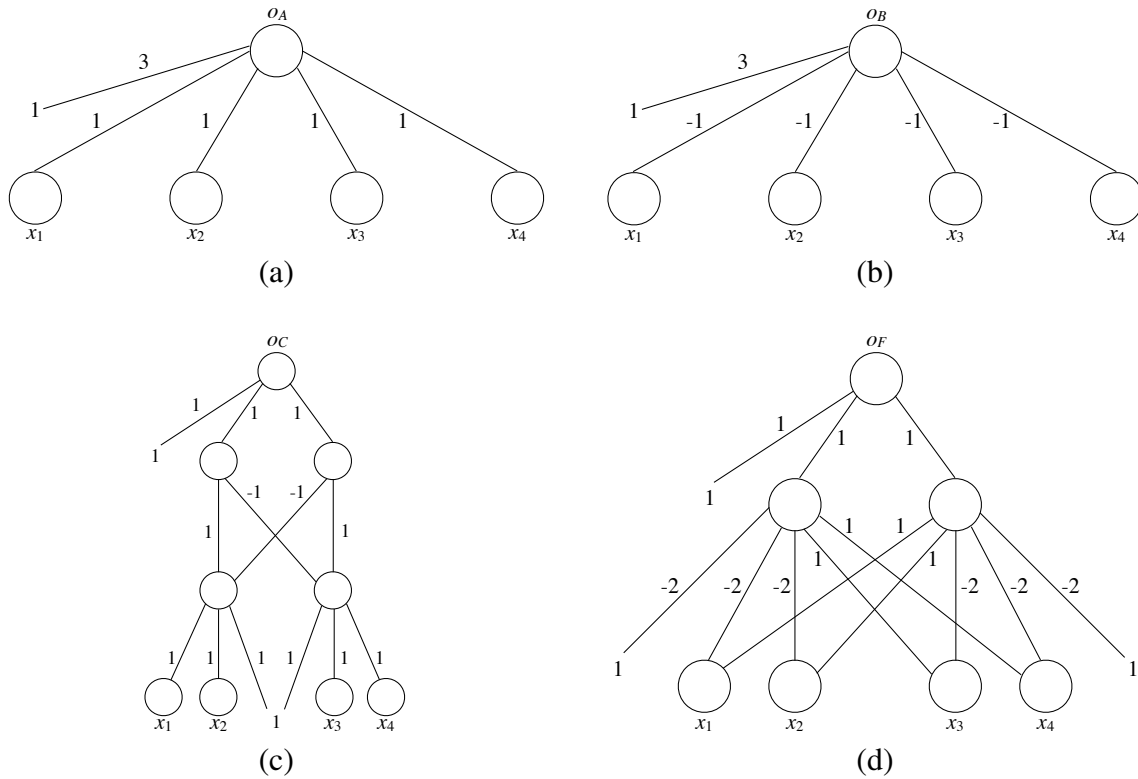


Figure 1: Perceptron networks: (a) perceptron unit A , (b) perceptron unit B , (c) network C of perceptron units, and (d) network F of perceptron units.

Suppose that the **prior beliefs** of hypotheses/weights \mathbf{w}_A , \mathbf{w}_B , \mathbf{w}_C , and \mathbf{w}_F are equal and they sum to 1.

Using Bayes' Theorem, calculate the posterior beliefs $P(\mathbf{w}_A|D)$, $P(\mathbf{w}_B|D)$, $P(\mathbf{w}_C|D)$, and $P(\mathbf{w}_F|D)$. Show the steps of your derivation. **No marks will be awarded for not doing so.**

We assume that the input instances $\mathbf{x}_d = (x_1, x_2, x_3, x_4)$ for $d \in D$ are fixed. Therefore, in deriving an expression for $P(D|\mathbf{w}_A)$, $P(D|\mathbf{w}_B)$, $P(D|\mathbf{w}_C)$, or $P(D|\mathbf{w}_F)$, we only need to consider the probability of observing the target outputs t_d for $d \in D$ for these fixed input instances \mathbf{x}_d for $d \in D$.

Furthermore, we assume that the training examples are conditionally independent given the hypothesis/weights of any perceptron network in Fig. 1.

Using the posterior beliefs calculated above, compute the **Bayes-optimal classification** for the new input instance $\mathbf{x}_{d_4} = (-1, -1, -1, -1)$. Show the steps of your derivation. **No marks will be awarded for not doing so.**