

# Exercise 2

Elements of Machine Learning, 2020, by Jens Petersen

The assignments in EML must be completed individually and written individually in English. Group discussions are allowed and encouraged, but in such cases you should list the group members in your handin. Your handin should include your solution to the exercises in a single pdf (not zipped), including code in a separate zip-file. Code used to solving the Combining multiple learners part should be attached as an appendix.

## Exercises

### From the textbooks (40 points)

Solve the following exercises

- PRML<sup>1</sup> 8.16, 8.27, ESL<sup>2</sup> 8.4, 10.1

### Graphical Models (30 points)

A standard example in graphical model literature is the student network (Figure 1).

1. What is  $p(\text{Intelligence} = 1 | \text{Letter} = 1, \text{SAT} = 1)$ ?
2. Convert the directed graph in Figure 1 to a moral graph
3. With the Letter node as root and the factor graph shown in Figure 2, Appendix A contains Sum-Product messages in the forward and backward passes and Appendix B contains Max-Sum messages in the forward pass. Give the missing messages (indicated by a question mark). Note you do not need to complete the messages relating to the Max-Sum backward pass. What are the marginal probabilities and most likely configuration of each node?

### Combining multiple learners (30 points)

The goal of this exercise is to predict whether an email is spam or not using various ways of combining multiple learners. For this you will need the spam dataset described in ESL, which can be downloaded from <https://archive.ics.uci.edu/ml/datasets/spambase>. We will use classification trees as the base learner.

---

<sup>1</sup>Pattern Recognition and Machine Learning, Christopher M. Bishop

<sup>2</sup>Elements of Statistical Learning, Trevor Hastie, Robert Tibshirani, and Jerome Friedman

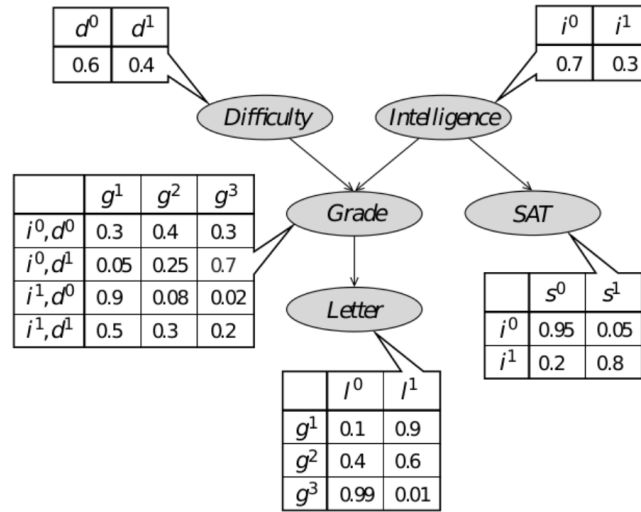


Figure 1: The student network. The model describes the relationship between class difficulty (0 = low, 1 = high), student intelligence (0 = not intelligent, 1 = intelligent), the grade the student gets in class (1 = good, 2 = average, 3 = bad), the student's score in the SAT exam (0 = low, 1 = high), and the recommendation letter the student gets from the professor (0 = not a good letter, 1 = good letter).

1. Split the dataset into non-overlapping training, validation, and test sets. Motivate your choice of splitting ratios, given the exercises below (6 points).
2. Implement, train and compare results of classifying spam (24 points)
  - (a) using a basic classification tree method. You may use the decision tree implementation in scikit-learn for this <https://scikit-learn.org/stable/modules/tree.html>.
  - (b) using bagging.
  - (c) using a boosting algorithm, such as AdaBoost.M1.

You are expected to describe, motivate, and compare possible methodological choices, including how you arrive at values for hyper parameters such as the number of iterations of AdaBoost, the number of base learners and the tree depth for the base learners. The bagging and boosting part of the exercise should be implemented by yourself, that is, it should not rely on existing implementations such as scikit-learn's ensemble methods <https://scikit-learn.org/stable/modules/ensemble.html>, however, you may use the base learner from 2a. Write up a comparison of the three approaches, what are your conclusions?

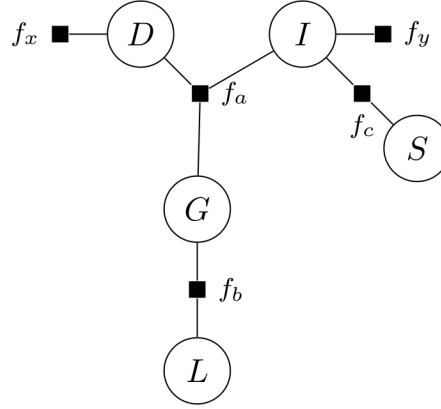


Figure 2: Factor graph corresponding to the student network with variables abbreviated by their first letters, where  $f_x(D) = p(D)$ ,  $f_y = p(I)$ ,  $f_a(D, I, G) = p(G|D, I)$ ,  $f_b(G, L) = p(L|G)$ , and  $f_c(I, S) = p(S|I)$

## Appendix A - sum product messages

Using the factor graph, we can compute the forward pass of the sum-product algorithm using the node  $L$  as root:

1.  $\mu_{f_x \rightarrow D}(D) = f_x(D)$
2.  $\mu_{f_y \rightarrow I}(I) = f_y(I)$
3.  $\mu_{S \rightarrow f_c}(S) = 1$
4.  $\mu_{f_c \rightarrow I}(I) = \sum_S f_c(S, I)$
5.  $\mu_{D \rightarrow f_a}(D) = \mu_{f_x \rightarrow D}(D)$
6.  $\mu_{I \rightarrow f_a}(I) = ?$
7.  $\mu_{f_a \rightarrow G}(G) = ?$
8.  $\mu_{G \rightarrow f_b}(G) = \mu_{f_a \rightarrow G}(G)$
9.  $\mu_{f_b \rightarrow L}(L) = \sum_G f_b(L, G) \mu_{G \rightarrow f_b}(G)$

The backward step, starts from the root and propagates back to the leaves:

1.  $\mu_{L \rightarrow f_b}(L) = 1$
2.  $\mu_{f_b \rightarrow f_G}(G) = \sum_L f_b(L, G) \mu_{L \rightarrow f_b}(L)$
3.  $\mu_{G \rightarrow f_a}(G) = \mu_{f_b \rightarrow G}(G)$
4.  $\mu_{f_a \rightarrow D}(D) = ?$

$$5. \mu_{D \rightarrow f_x}(D) = ?$$

$$6. \mu_{f_a \rightarrow I}(I) = \sum_{G,D} f_a(I, G, D) \mu_{D \rightarrow f_a}(D) \mu_{G \rightarrow f_a}(G)$$

$$7. \mu_{I \rightarrow f_y}(I) = \mu_{f_a \rightarrow I}(I)$$

$$8. \mu_{I \rightarrow f_c}(I) = \mu_{f_a \rightarrow I}(I)$$

$$9. \mu_{f_c \rightarrow S}(S) = \sum_I f_c(I, S) \mu_{I \rightarrow f_c}(I)$$

Sum product forward pass messages in detail

$$\mu_{f_x \rightarrow D}(D) = \frac{D \mid p(D)}{0 \mid 0.6 \atop 1 \mid 0.4} \quad (1)$$

$$\mu_{f_y \rightarrow I}(I) = \frac{I \mid p(I)}{0 \mid 0.7 \atop 1 \mid 0.3} \quad (2)$$

$$\mu_{S \rightarrow f_c}(S) = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad (3)$$

$$\mu_{f_c \rightarrow I}(I) = \sum_S f_c(S, I) = \sum_S \frac{I \quad S \mid f_c(S, I) = p(S|I)}{0 \quad 0 \mid 0.95 \atop 0 \quad 1 \mid 0.05 \atop 1 \quad 0 \mid 0.2 \atop 1 \quad 1 \mid 0.8} = \frac{I \mid \mu_{f_c \rightarrow I}(I)}{0 \mid 1 \atop 1 \mid 1} \quad (4)$$

$$\mu_{D \rightarrow f_a}(D) = \mu_{f_x \rightarrow D}(D) = \frac{D \mid p(D)}{0 \mid 0.6 \atop 1 \mid 0.4} \quad (5)$$

$$\mu_{I \rightarrow f_a}(I) = ? \quad (6)$$

$$\mu_{f_a \rightarrow G}(G) = ? \quad (7)$$

$$\mu_{G \rightarrow f_b}(G) = \mu_{f_a \rightarrow G}(G) = \frac{G \mid \mu_{f_a \rightarrow G}(G)}{1 \mid 0.362 \atop 2 \mid 0.2884 \atop 3 \mid 0.3496} \quad (8)$$

$$\mu_{f_b \rightarrow L}(L) = \sum_G f_b(L, G) \mu_{G \rightarrow f_b}(G) = \sum_G \frac{G \quad L \mid f_b(L, G)}{1 \quad 0 \mid 0.1 \atop 2 \quad 0 \mid 0.4 \atop 3 \quad 0 \mid 0.99 \atop 1 \quad 1 \mid 0.9 \atop 2 \quad 1 \mid 0.6 \atop 3 \quad 1 \mid 0.01} * \frac{G \mid \mu_{f_a \rightarrow G}(G)}{1 \mid 0.362 \atop 2 \mid 0.2884 \atop 3 \mid 0.3496} \quad (9)$$

$$= \frac{G \quad L \mid f_b(L, G) * \mu_{f_a \rightarrow G}(G)}{1 \quad 0 \mid 0.1 * 0.362 = 0.0362 \atop 2 \quad 0 \mid 0.4 * 0.2884 = 0.11536 \atop 3 \quad 0 \mid 0.99 * 0.3496 = 0.346104 \atop 1 \quad 1 \mid 0.9 * 0.362 = 0.3258 \atop 2 \quad 1 \mid 0.6 * 0.2884 = 0.17304 \atop 3 \quad 1 \mid 0.01 * 0.3496 = 0.003496} = \frac{L \mid \mu_{f_b \rightarrow L}(L)}{0 \mid 0.497664 \atop 1 \mid 0.502336}$$

Sum product backward pass messages in detail

$$\mu_{L \rightarrow f_b}(L) = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad (1)$$

$$\mu_{f_b \rightarrow f_G}(G) = \sum_L f_b(L, G) \mu_{L \rightarrow f_b}(L) = \sum_L \begin{array}{c|c|c} G & L & f_b(L, G) \\ \hline 1 & 0 & 0.1 \\ 2 & 0 & 0.4 \\ 3 & 0 & 0.99 \\ 1 & 1 & 0.9 \\ 2 & 1 & 0.6 \\ 3 & 1 & 0.01 \end{array} = \begin{array}{c|c} G & \mu_{f_b \rightarrow G}(G) \\ \hline 1 & 1 \\ 2 & 1 \\ 3 & 1 \end{array} \quad (2)$$

$$\mu_{G \rightarrow f_a}(G) = \mu_{f_b \rightarrow G}(G) = \begin{array}{c|c} G & \mu_{f_b \rightarrow G}(G) \\ \hline 1 & 1 \\ 2 & 1 \\ 3 & 1 \end{array} \quad (3)$$

$$\mu_{f_a \rightarrow D}(D) = ? \quad (4)$$

$$\mu_{D \rightarrow f_x}(D) = ? \quad (5)$$

$$\mu_{f_a \rightarrow I}(I) = \sum_{G,D} f_a(I, G, D) \mu_{D \rightarrow f_a}(D) \mu_{G \rightarrow f_a}(G)$$

$$= \sum_{G,D} \begin{array}{c|c} \begin{array}{ccc} I & D & G \\ \hline 0 & 0 & 1 \\ 0 & 0 & 2 \\ 0 & 0 & 3 \\ 0 & 1 & 1 \\ 0 & 1 & 2 \\ 0 & 1 & 3 \\ 1 & 0 & 1 \\ 1 & 0 & 2 \\ 1 & 0 & 3 \\ 1 & 1 & 1 \\ 1 & 1 & 2 \\ 1 & 1 & 3 \end{array} & \begin{array}{c} f_a(G, I, D) \\ \hline 0.3 \\ 0.4 \\ 0.3 \\ 0.05 \\ 0.25 \\ 0.7 \\ 0.9 \\ 0.08 \\ 0.02 \\ 0.5 \\ 0.3 \\ 0.2 \end{array} \end{array} \begin{array}{c|c} \begin{array}{cc} G & D \\ \hline 1 & 0 \\ 2 & 0 \\ 3 & 0 \\ 1 & 1 \\ 2 & 1 \\ 3 & 1 \end{array} & \begin{array}{c} \mu_{D \rightarrow f_a}(D) * \mu_{G \rightarrow f_a}(G) \\ \hline 0.6 \\ 0.6 \\ 0.6 \\ 0.4 \\ 0.4 \\ 0.4 \end{array} \end{array} \quad (6)$$

$$= \sum_{G,D} \begin{array}{c|c} \begin{array}{ccc} I & D & G \\ \hline 0 & 0 & 1 \\ 0 & 0 & 2 \\ 0 & 0 & 3 \\ 0 & 1 & 1 \\ 0 & 1 & 2 \\ 0 & 1 & 3 \\ 1 & 0 & 1 \\ 1 & 0 & 2 \\ 1 & 0 & 3 \\ 1 & 1 & 1 \\ 1 & 1 & 2 \\ 1 & 1 & 3 \end{array} & \begin{array}{c} f_a(G, I, D) * (\mu_{D \rightarrow f_a}(D) * \mu_{G \rightarrow f_a}(G)) \\ \hline 0.3 * 0.6 = 0.18 \\ 0.4 * 0.6 = 0.24 \\ 0.3 * 0.6 = 0.18 \\ 0.05 * 0.4 = 0.02 \\ 0.25 * 0.4 = 0.1 \\ 0.7 * 0.4 = 0.28 \\ 0.9 * 0.6 = 0.54 \\ 0.08 * 0.6 = 0.048 \\ 0.02 * 0.6 = 0.012 \\ 0.5 * 0.4 = 0.2 \\ 0.3 * 0.4 = 0.12 \\ 0.2 * 0.4 = 0.08 \end{array} \end{array} = \begin{array}{c|c} \begin{array}{c} I \\ \hline 0 \\ 1 \end{array} & \begin{array}{c} \mu_{f_a \rightarrow I}(I) \\ \hline 1 \\ 1 \end{array} \end{array}$$

$$\mu_{I \rightarrow f_c}(I) = \mu_{f_y \rightarrow I}(I) \mu_{f_a \rightarrow I}(I) = \begin{array}{c|c} \begin{array}{c} I \\ \hline 0 \\ 1 \end{array} & \begin{array}{c} \mu_{f_a \rightarrow I}(I) \\ \hline 1 \\ 1 \end{array} \end{array} * \begin{array}{c|c} \begin{array}{c} I \\ \hline 0 \\ 1 \end{array} & \begin{array}{c} \mu_{f_y \rightarrow I}(I) \\ \hline 0.7 \\ 0.3 \end{array} \end{array} = \begin{array}{c|c} \begin{array}{c} I \\ \hline 0 \\ 1 \end{array} & \begin{array}{c} \mu_{I \rightarrow f_c}(I) \\ \hline 0.7 \\ 0.3 \end{array} \end{array} \quad (7)$$

$$\mu_{f_c \rightarrow S}(S) = \sum_I f_c(I, S) \mu_{I \rightarrow f_c}(I) = \sum_I f_c(I, S) \begin{array}{c|c} \begin{array}{c} I \\ \hline 0 \\ 1 \end{array} & \begin{array}{c} \mu_{I \rightarrow f_c}(I) \\ \hline 0.7 \\ 0.3 \end{array} \end{array}$$

$$= \sum_I \begin{array}{c|c} \begin{array}{cc} I & S \\ \hline 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{array} & \begin{array}{c} \mu_{I \rightarrow f_c}(I) \\ \hline 0.95 \\ 0.2 \\ 0.05 \\ 0.8 \end{array} \end{array} * \begin{array}{c|c} \begin{array}{c} I \\ \hline 0 \\ 1 \end{array} & \begin{array}{c} \mu_{I \rightarrow f_c}(I) \\ \hline 0.7 \\ 0.3 \end{array} \end{array} = \begin{array}{c|c} \begin{array}{c} S \\ \hline 0 \\ 1 \end{array} & \begin{array}{c} \mu_{f_a \rightarrow S}(S) \\ \hline 0.725 \\ 0.275 \end{array} \end{array} \quad (8)$$

## Appendix B - max sum forward messages

The forward and backward messages for the max-sum algorithm are similar to the sum-product one (Appendix A), but instead of marginalizing we take the maximum, and instead of multiplying we take the sum:

1.  $\mu_{f_x \rightarrow D}(D) = \log(f_x(D))$
2.  $\mu_{f_y \rightarrow I}(I) = \log(f_y(I))$
3.  $\mu_{S \rightarrow f_c}(S) = 0$
4.  $\mu_{f_c \rightarrow I}(I) = \max_S \log(f_c(S, I))$
5.  $\mu_{D \rightarrow f_a}(D) = \mu_{f_x \rightarrow D}(D)$
6.  $\mu_{I \rightarrow f_a}(I) = ?$
7.  $\mu_{f_a \rightarrow G}(G) = ?$
8.  $\mu_{G \rightarrow f_b}(G) = \mu_{f_a \rightarrow G}(G)$
9.  $\mu_{f_b \rightarrow L}(L) = \max_G \log(f_b(L, G)) + \mu_{G \rightarrow f_b}(G)$



Max sum forward pass messages in detail (please excuse the use of equalities with rounded off values)

$$\mu_{f_x \rightarrow D}(D) = \log(f_x(D)) = \frac{D}{\begin{array}{c|c} & \mu_{f_x \rightarrow D}(D) \\ \hline 0 & \log(0.6) \\ 1 & \log(0.4) \end{array}} \quad (1)$$

$$\mu_{f_y \rightarrow I}(I) = \log(f_y(I)) = \frac{I}{\begin{array}{c|c} & \mu_{f_y \rightarrow I}(I) \\ \hline 0 & \log(0.7) \\ 1 & \log(0.3) \end{array}} \quad (2)$$

$$\mu_{S \rightarrow f_c}(S) = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad (3)$$

$$\begin{aligned} \mu_{f_c \rightarrow I}(I) &= \max_S \log(f_c(S, I)) = \max_S \begin{array}{c|c|c} I & S & f_c(S, I) = p(S|I) \\ \hline 0 & 0 & 0.95 \\ 0 & 1 & 0.05 \\ 1 & 0 & 0.2 \\ 1 & 1 & 0.8 \end{array} \\ &= \begin{array}{c|c|c} I & \mu_{f_c \rightarrow I}(I) & \arg \max_S \\ \hline 0 & \log(0.95) & 0 \\ 1 & \log(0.8) & 1 \end{array} \\ &= \begin{array}{c|c|c} I & \mu_{f_c \rightarrow I}(I) & \arg \max_S \\ \hline 0 & -0.051 & 0 \\ 1 & -0.223 & 1 \end{array} \end{aligned} \quad (4)$$

$$\mu_{D \rightarrow f_a}(D) = \mu_{f_x \rightarrow D}(D) = \frac{D}{\begin{array}{c|c} & \mu_{f_x \rightarrow D}(D) \\ \hline 0 & \log(0.6) \\ 1 & \log(0.4) \end{array}} \quad (5)$$

$$\mu_{I \rightarrow f_a}(I) = ? \quad (6)$$

$$\mu_{f_a \rightarrow G}(G) = ? \quad (7)$$

$$\mu_{G \rightarrow f_b}(G) = \mu_{f_a \rightarrow G}(G) = \begin{array}{c|c} G & \mu_{f_a \rightarrow G}(G) \\ \hline 1 & \log(0.1296) \\ 2 & \log(0.1596) \\ 3 & \log(0.1862) \end{array} = \begin{array}{c|c} G & \mu_{f_a \rightarrow G}(G) \\ \hline 1 & -2.043 \\ 2 & -1.835 \\ 3 & -1.681 \end{array} \quad (8)$$

$$\begin{aligned}
\mu_{f_b \rightarrow L}(L) &= \max_G (\log(f_b(L, G)) + \mu_{G \rightarrow f_b}(G)) \\
&= \max_G \left( \begin{array}{c|c|c} G & L & \log(f_b(L, G)) \\ \hline 1 & 0 & \log(0.1) \\ 2 & 0 & \log(0.4) \\ 3 & 0 & \log(0.99) \\ 1 & 1 & \log(0.9) \\ 2 & 1 & \log(0.6) \\ 3 & 1 & \log(0.01) \end{array} + \begin{array}{c|c} G & \mu_{G \rightarrow f_b}(G) \\ \hline 1 & \log(0.1296) \\ 2 & \log(0.1596) \\ 3 & \log(0.1862) \end{array} \right) \\
&= \max_G \begin{array}{c|c|c} G & L & \log(f_b(L, G)) + \mu_{G \rightarrow f_b}(G) \\ \hline 1 & 0 & \log(0.1 * 0.1296) = \log(0.013) \\ 2 & 0 & \log(0.4 * 0.1596) = \log(0.0638) \\ 3 & 0 & \log(0.99 * 0.1862) = \log(0.184338) \\ 1 & 1 & \log(0.9 * 0.1296) = \log(0.11664) \\ 2 & 1 & \log(0.6 * 0.1596) = \log(0.09576) \\ 3 & 1 & \log(0.01 * 0.1862) = \log(0.001862) \end{array} \tag{9} \\
&= \begin{array}{c|c|c} L & \mu_{f_b \rightarrow L}(L) & \arg \max_G \\ \hline 0 & \log(0.184338) & 3 \\ 1 & \log(0.11664) & 1 \end{array} = \begin{array}{c|c|c} L & \mu_{f_b \rightarrow L}(L) & \arg \max_G \\ \hline 0 & -1,691 & 3 \\ 1 & -2,148 & 1 \end{array}
\end{aligned}$$