



Figure 1: Bayes Net.

Labs
Machine Learning Course
 Fall 2016

EPFL
 School of Computer and Communication Sciences
Martin Jaggi & Rüdiger Urbanke
mlo.epfl.ch/page-136795.html
epfmlcourse@gmail.com

Problem Set 14, Dec 22th, 2016 (Various)

Goals. The goal of this exercise is to provide you with some further theory questions so that you can prepare better for your final exam. We might add some questions later on. So you might want to always check the latest version.

Problem 1 (Markov blanket):

In the course notes we introduced the concept of a *Markov blanket*.

Let X_i be a node in a given Bayes net and let Z be its Markov blanket. Let X_j be any other node not equal to X_i or contained in Z . Show that X_i is independent of X_j given Z by showing that X_j is *D-separated* from X_i by Z .

Solution: Consider any (undirected path) from X_i to X_j . Note that any such path must either pass through a parent or a child of X_i .

Let us first assume that the path passes through a parent. Let this parent node be X , $X \in Z$. Then this parent must be either tail-to-tail or head-to-tail with respect to this path. In either case it blocks this path.

Next assume that the path goes through a child. If the path is head-to-tail with respect to this child, again this child blocks this path and we are done. If not, then the path must be head-to-head. In this case there must be a co-parent of X_i that also lies on this path. This co-parent is also in Z by definition and it must be either head-to-tail or tail-to-tail with respect to this path. In either case it blocks this path and we are done.

Problem 2 (Bayes Net):

Consider the example shown in Figure 1. We will not work out the answers here, but it might be a good additional exercise to work out the answer yourself.

- Is X_1 independent of X_2 given X_3 ?
- Is X_1 independent of X_2 given X_5 ?

Solution:

- Is X_1 independent of X_2 given X_3 ? NO: Without conditioning, X_4 blocks this path, but X_3 is a descendant of X_4 and it is in the conditioning and hence X_4 does not block the path given X_3 .
- Is X_1 independent of X_2 given X_5 ? YES: There is only one path from X_1 to X_2 and X_5 is tail-to-tail with respect to this path and hence blocks it.

Problem 3 (Finite differences (from ETHZ course by Hofmann)):

During lab 13 you implemented the back-propagation algorithm. It is easy to get some indices wrong and so it is important to have a good way of checking the derivatives.

1. One common way to check that the computation of a derivative is correct is to use finite differences. Considering only the i -th dimension of the parameter vector w , show that finite difference yields an error $O(\epsilon)$, i.e.

$$\nabla f(w_i) = \frac{f(w_i + \epsilon) - f(w_i)}{\epsilon} + O(\epsilon), \quad (1)$$

where $\epsilon \in \mathbb{R}^+ \leq 1$.

2. The accuracy of the finite difference method can be improved significantly by using symmetrical central differences.

$$\nabla f(w_i) \approx \frac{f(w_i + \epsilon) - f(w_i - \epsilon)}{2\epsilon} \quad (2)$$

What approximation error do we get using Eq. 2?

Solution: (1) Using Taylor expansion, we get

$$f(w_i + \epsilon) = f(w_i) + \epsilon \nabla f(w_i) + O(\epsilon^2). \quad (3)$$

Re-organizing the terms and using $\frac{O(\epsilon^2)}{\epsilon} = O(\epsilon)$ yields

$$\nabla f(w_i) = \frac{f(w_i + \epsilon) - f(w_i)}{\epsilon} + O(\epsilon). \quad (4)$$

(2) For the second equation, we again use a Taylor expansion of $f(w_i + \epsilon)$ and $f(w_i - \epsilon)$ around w_i , but this time up to the third-order, i.e.

$$f(w_i + \epsilon) = f(w_i) + \epsilon \nabla f(w_i) + \frac{1}{2} \epsilon^2 \nabla^2 f(w_i) + O(\epsilon^3), \quad (5)$$

and

$$f(w_i - \epsilon) = f(w_i) - \epsilon \nabla f(w_i) + \frac{1}{2} \epsilon^2 \nabla^2 f(w_i) + O(\epsilon^3) \quad (6)$$

Subtracting Eq. 6 from Eq. 5 yields

$$\nabla f(w_i) = \frac{f(w_i + \epsilon) - f(w_i - \epsilon)}{2\epsilon} + O(\epsilon^2) \quad (7)$$

Problem 4 (Weight initialization (from ETHZ course by Hofmann)):

How shall we initialize the weights in a neural network. One point to consider is that we would like the nodes to work in a proper range. In more detail, assume that we have set the weights initially in such a way that the expected input to a node is around 100 and we use the tanh as an activation function. Then we will be operating in a regime where the activation function is very flat and it might need a lot of training to change the weights so that the expected input is closer to the “interesting” range $[-1, 1]$ where this activation function.

Suppose we have an input $X = (x_1, \dots, x_n)$ where $x_i \in \mathbb{R}^d$ and a linear neuron with random weights W that outputs a number Y as:

$$Y = W_1 x_1 + W_2 x_2 + \dots + W_n x_n$$

We assume that x_i and W_i are all independent and identically distributed.

1. Show that $\text{Var}[W_i x_i] = \mathbb{E}[x_i]^2 \text{Var}[W_i] + \mathbb{E}[W_i]^2 \text{Var}[x_i] + \text{Var}[W_i] \text{Var}[x_i]$.
2. Assuming that inputs and weights both have mean 0, show that $\text{Var}[Y] = n \text{Var}[W_i] \text{Var}[x_i]$.
3. Now assume that we want the variance of the output to be the same as the variance of the input. What is the condition for $\text{Var}[W_i]$.

Solution: (1) We first compute the variance of $W_i x_i$ is

$$\begin{aligned}
 \text{Var}[W_i x_i] &= \mathbb{E}[(W_i x_i - \mathbb{E}[W_i x_i])^2] \\
 &= \mathbb{E}[W_i^2 x_i^2] - (\mathbb{E}[W_i x_i])^2 \\
 &= \mathbb{E}[W_i^2] \mathbb{E}[x_i^2] - \mathbb{E}[W_i]^2 \mathbb{E}[x_i]^2 \\
 &= (\text{Var}[W_i] + \mathbb{E}[W_i]^2)(\text{Var}[x_i] + \mathbb{E}[x_i]^2) - \mathbb{E}[W_i]^2 \mathbb{E}[x_i]^2 \\
 &= \mathbb{E}[x_i]^2 \text{Var}[W_i] + \mathbb{E}[W_i]^2 \text{Var}[x_i] + \text{Var}[W_i] \text{Var}[x_i]
 \end{aligned} \tag{8}$$

(2) We know that $\mathbb{E}[W_i] = \mathbb{E}[x_i] = 0$ so Eq. 8 simplifies to

$$\text{Var}[W_i x_i] = \text{Var}[W_i] \text{Var}[x_i].$$

We thus get:

$$\text{Var}[Y] = \text{Var}\left[\sum_i W_i x_i\right] = n \text{Var}[W_i x_i] = n \text{Var}[W_i] \text{Var}[x_i].$$

(3) if we want the variance of the input and output to be the same, this implies $n \text{Var}[W_i] = 1$, which means the variance of the weights should be

$$\text{Var}[W_i] = \frac{1}{n}.$$

This way of initializing the weights of a neural network is known as Xavier initialization and is commonly used in existing deep network libraries such as Caffe or Tensorflow.