

## CS 475 Machine Learning: Homework 5

## Graphical Models

Due: Tuesday November 25, 2014, 11:59pm

100 Points Total

Version 1.2

Make sure to read from start to finish before beginning the assignment.

## 1 Programming (60 points)

In this assignment you will implement a loopy belief propagation (BP) algorithm for calculating marginal probabilities in a loopy MRF (more specifically, a factor graph). You will implement loopy belief propagation using provided functions for each factor (there is no learning in the homework, only inference). As background, we first summarize inference on a linear chain factor graph using message passing (sum product.) We will then explain how this algorithm can be adapted to a loopy graph structure. Since this new structure is not a DAG, we no longer have a guarantee that we will get the correct answer. Therefore, loopy BP is an approximate inference algorithm.

### 1.1 Message Passing on a Chain Factor Graph

In this section we briefly introduce the sum product algorithm on a chain factor graph for computing the marginal of each variable on the chain.

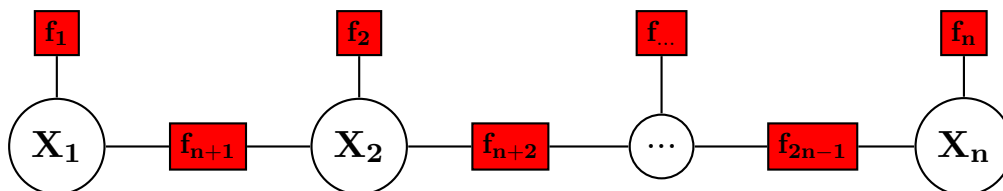


Figure 1: The chain factor graph.

Each of the variables in our chain (the  $x_i$ ) are  $k$ -ary discrete variables.

We can compute the marginal probabilities for each  $x_i$  in this chain using the sum product algorithm. We review the algorithm here. Our presentation is adapted from section 8.4.4 of Bishop. For more details see the book (<http://research.microsoft.com/en-us/um/people/cmbishop/prml/Bishop-PRML-sample.pdf>).

Our goal is to find the marginal probability of a node in the factor graph. Due to the linear structure of the factor graph, the sum product algorithm lets us write this marginal probability as:

$$p(x) = \prod_{s \in N(x)} \left[ \sum_{X_s} F_s(x, X_s) \right] \quad (1)$$

where:

$N(x)$  is the set of all factor node neighbors of  $x$

$F_s(x, X_s)$  represents the product of all the factors “downstream” of  $f_s$

Part of the above equation will be used many times, so for the sake of computation and intuition, we will define the sum term as a “message”  $\mu$ :

$$\mu_{f_s \rightarrow x}(x) := \sum_{X_s} F_s(x, X_s) \quad (2)$$

Messages can be defined using recursion:  $F_s(x, X_s) = f_s(x, x_1, \dots, x_M) \prod_{m \in N(f_s) \setminus x} \mu_{x_m \rightarrow f_s}(x_m)$ ,

where  $X_s = \{x_m | m \in N(f_s) \setminus x\}$ . Our base case is:

$$\mu_{f \rightarrow x}(x) := f(x) \text{ iff the only neighbor of } f \text{ is } x \quad (3)$$

The sum product algorithm defines the message from a variable node to a factor node as the product of the messages it receives from its “downstream” factors:

$$\mu_{x \rightarrow f_s} := \prod_{l \in N(x) \setminus s} \mu_{f_l \rightarrow x}(x) \quad (4)$$

As before, there is a base case for this equation:

$$\mu_{x \rightarrow f}(x) := 1 \text{ iff the only neighbor of } x \text{ is } f$$

We can now find the marginal probabilities using Eq. (1). While our presentation focused on chains, the Sum Product algorithm applies to any tree structured factor graph. While seemingly simple, the details may be a bit opaque. To help clarify them, we will describe the message passing procedure on our specific chain graph in Fig. (1).

Let's say, we want to compute  $p(x_2)$ . According to Eq. (1,2), we have  $p(x_2) = \mu_{f_{n+1} \rightarrow x_2}(x_2) \mu_{f_{n+2} \rightarrow x_2}(x_2) \mu_{f_2 \rightarrow x_2}(x_2)$ . So we need to compute messages  $\mu_{f_{n+1} \rightarrow x_2}$ ,  $\mu_{f_{n+2} \rightarrow x_2}$  and  $\mu_{f_2 \rightarrow x_2}$ . From Eq. (2,3), these messages can be computed as:

$$\mu_{f_{n+1} \rightarrow x_2}(x_2) = \sum_{x_1} \mu_{x_1 \rightarrow f_{n+1}}(x_1) f_{n+1}(x_1, x_2) \quad (5)$$

$$\mu_{f_{n+2} \rightarrow x_2}(x_2) = \sum_{x_3} \mu_{x_3 \rightarrow f_{n+2}}(x_3) f_{n+2}(x_2, x_3) \quad (6)$$

$$\mu_{f_2 \rightarrow x_2}(x_2) = f_2(x_2) \quad (7)$$

where Eq. (7) is already a base case while Eq. (5,6) are not. We continue to expand  $\mu_{x_1 \rightarrow f_{n+1}}(x_1)$  in Eq. (5) and  $\mu_{x_3 \rightarrow f_{n+2}}(x_3)$  in Eq. (6), according to Eq. (4) and (3):

$$\mu_{x_1 \rightarrow f_{n+1}}(x_1) = \mu_{f_1 \rightarrow x_1}(x_1) = f_1(x_1) \quad (8)$$

$$\mu_{x_3 \rightarrow f_{n+2}}(x_3) = \mu_{f_3 \rightarrow x_3}(x_3) \mu_{f_{n+3} \rightarrow x_3}(x_3) = f_3(x_3) \mu_{f_{n+3} \rightarrow x_3}(x_3) \quad (9)$$

where Eq. (8) is already a base case while Eq. (9) is not. So we continue to expand  $\mu_{f_{n+3} \rightarrow x_3}(x_3)$  in Eq. (9) and so on, until we reach the other end of the chain.

To summarize, to compute the marginal of a variable in the chain, we need to compute the messages starting from both ends of the chain one by one, until that variable has collected all its incoming messages, where each message is computed only once and then is finalized.

## 1.2 Message Passing on a Loopy Factor Graph

In this assignment you are asked to compute the marginal probability of a variable in a loopy factor graph. The graph is shown in Fig. 2. This graph is exactly the same as our linear chain, except that node  $x_N$  is connected to  $x_1$ .

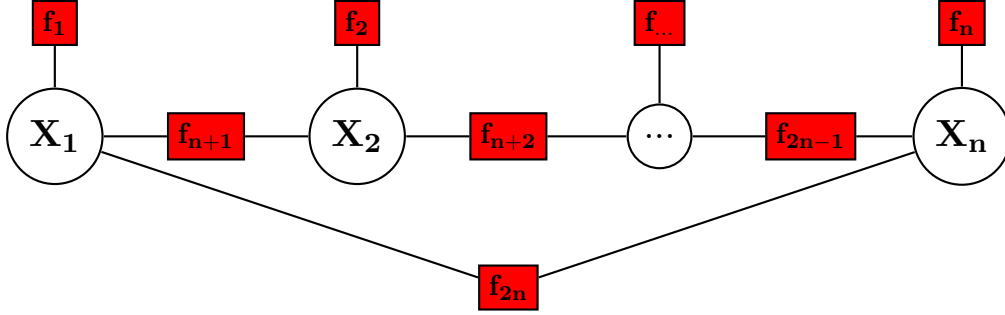


Figure 2: The loopy factor graph.

As before, each of the variables in our loop (the  $x_i$ ) are  $k$ -ary discrete variables.

Unfortunately, the procedures described in Sum Product Algorithm work for tree structured graphs only and cannot be directly applied in a loopy factor graph. If we directly applied the algorithm, we would introduce infinite recursion since there would be no base case.

Therefore, we modify the algorithm to adapt it to loopy graphs. Our change will be to start passing incomplete messages, messages that are missing information. This will obviate the problem of not having a base case to start from. We will do this by specifying an order for computing messages and then updating messages in turn. However, since messages are missing information, we cannot stop once all messages have reached their final destination. Instead, we will continue passing messages throughout the graph for  $T$  iterations, or more generally, until the values of the messages (and the marginals) have converged. Note that this new algorithm has no guarantee of correctness or convergence, but provides approximate marginals. In many cases, the results provide a good approximation to actual marginals.

The modified algorithm falls in a broader class of algorithms called *Loopy Belief Propagation*, where “Belief” here refers to how a node thinks part of the graph should be, i.e., the part that is in the opposite direction of where the message is passed. Beliefs are stored in the messages and are passed in the graph in order to update other messages. Each message is updated  $T$  times.

The modified algorithm is as follows:

1. Initialize  $\mu_{x_1 \rightarrow f_{n+1}}(x_1) = 1$  and  $\mu_{x_1 \rightarrow f_{2n}}(x_1) = 1$  for all values  $x_1$  can take.
2. For  $t$  from 1 to  $T$ 
  - (a) For  $i$  from 1 to  $n$ 
    - i. Compute  $\mu_{f_{n+i} \rightarrow x_{(i+1)\%n}}(x_{(i+1)\%n})$  from Eq. (2).
    - ii. Compute  $\mu_{x_{1+i\%n} \rightarrow f_{n+1+i\%n}}(x_{1+i\%n})$  from Eq. (4).
  - (b) For  $i$  from  $n$  to 1
    - i. Compute  $\mu_{f_{n+i} \rightarrow x_i}(x_i)$  from Eq. (2).
    - ii. Compute  $\mu_{x_i \rightarrow f_{n+(i-2)\%n+1}}(x_i)$  from Eq. (4).

After messages are passed as shown above, we will compute the marginal  $p(x_i)$  for any  $i$  according to Eq. (1,2).

### 1.3 Implementation

For the factor graph in this assignment, there are unary ( $f_1$  to  $f_n$ ) and binary ( $f_{n+1}$  to  $f_{2n}$ ) factors, and each factor  $f_i$  is associated with a potential function  $\psi_i$ . Each unary factor will have a potential function  $\psi_i(a)$  which returns a real non-negative value corresponding to the potential when variable  $x_i$  takes value  $a$ . Each factor node between two variable nodes will have a potential function  $\psi_i(a, b)$  which returns a value corresponding to the potential when variable  $x_{i-n}$  takes value  $a$  and node  $x_{i-n+1}$  takes value  $b$  (except for the factor  $f_{2n}$  where  $x_n$  takes value  $a$  and  $x_1$  takes value  $b$ ). Recall that every  $x_i$  can take values from 1 to  $k$ .

We will provide you code that gives you values of  $\psi_i(a)$  and  $\psi_i(a, b)$ . They will be in the class `cs475.loopMRF.LoopMRFPotentials` and have the signatures:

```
public double potential(int i, int a)
public double potential(int i, int a, int b)
```

There are  $n$  nodes in this loop, so the value of `i` must be between 1 and  $n$  (inclusive) in the first method and between  $n + 1$  and  $2n$  (inclusive) in the second method. Since every  $x_i$  can take values from 1 to  $k$  (inclusive), you must only call this function with values for `a` and `b` between 1 and  $k$  (inclusive). You will be able to get values for  $n$  and  $k$  by calling the following functions in `cs475.loopMRF.LoopMRFPotentials`:

```
public int loopLength() // returns n
public int numXValues() // returns k
```

These values will be read into `cs475.loopMRF.LoopMRFPotentials` from a text file that must be provided in the constructor:

```
public LoopMRFPotentials(String data_file, int iterations)
```

We are providing you with a sample of this data file, `sample_mrf_potentials.txt`. The format is "`n k`" on the first line and either "`i a potential`" or "`i a b potential`" on subsequent lines. Feel free to try out new loops to get different probability distributions, just make sure it contains all the needed potential values.

Your code will work by calculating these messages given the value of the potential functions between the variable nodes in the loop. For details on how to do this, you can refer to your notes from class, or see Bishop's examples in the book.

#### 1.3.1 Command Line Arguments

We have already added 2 command line options in `cs475.loopMRF.LoopMRFPotentials` for you.

```
registerOption("data", "String", true, "The data to use.");
registerOption("iterations", "int", true, "The number of iterations. default: 50");
```

## 1.4 What You Need to Implement

We have provided you with class `LoopyBP`, with one method left blank that you will need to implement:

```
public class LoopyBP {
    public double[] marginalProbability(int x_i) {
        // TODO
    }
}
```

The method should return a double array where the  $j$ th element is the probability that  $x_i = j$ . The length of this array should be  $k + 1$  and you should leave the 0 index as 0. These are probabilities so don't forget to normalize to sum to 1.

## 1.5 How We Will Run Your Code

We will run your code by providing you with a single command line argument which is the data file:

```
java cs475.loopMRF.LoopMRFTester -data mrf_potentials.txt -iterations T
```

Note that we will use new data files with different values of  $n$  and  $k$ , so make sure your code works for any reasonable input.

Your output should just be the results of the print statements in the code given. **Do not print anything else in the version you hand in.**

## 2 Analytical (40 points)

1. (10 points) Consider the Bayesian Network given in Figure 3. Are the sets **A** and **B** d-separated given set **C** for each of the following definitions of **A**, **B** and **C**? Justify each answer.

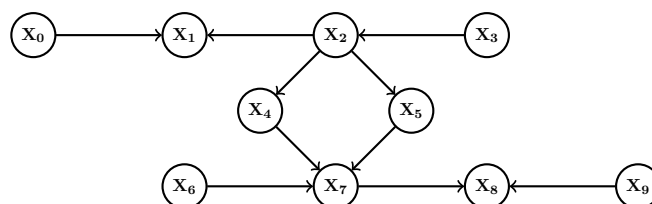


Figure 3: A directed graph

- $\mathbf{A} = \{X_3\}$ ,  $\mathbf{B} = \{X_4\}$ ,  $\mathbf{C} = \{X_1, X_2\}$
- $\mathbf{A} = \{X_5\}$ ,  $\mathbf{B} = \{X_4\}$ ,  $\mathbf{C} = \{X_2, X_7\}$
- $\mathbf{A} = \{X_4\}$ ,  $\mathbf{B} = \{X_6\}$ ,  $\mathbf{C} = \{X_8\}$
- $\mathbf{A} = \{X_5\}$ ,  $\mathbf{B} = \{X_4\}$ ,  $\mathbf{C} = \{X_2\}$

e.  $\mathbf{A} = \{X_6, X_9\}$ ,  $\mathbf{B} = \{X_8\}$ ,  $\mathbf{C} = \{X_3, X_1\}$

Now assume that Figure 3 is a Markov Random Field, where each edge is undirected (just drop the direction of each edge.) Re-answer each of the above questions with justifications for your answers.

**2. (10 points)** In this problem you will use CRFs to find the Noun Phrases (NPs) in a sentence. A sentence  $\mathbf{x}$  with label sequence  $\mathbf{y}$  is composed of  $N$  words, where each word is denoted by  $x_i$ . The label sequence  $\mathbf{y}$  is composed of  $N$  labels, where each label  $y_i \in \Gamma = \{B, I, O\}$ .  $B$ ,  $I$ ,  $O$  denote *Beginning* of an NP, *Inside* an NP, and *Outside* an NP.

Consider the following sentence:

Current economic weakness may be a result of high energy prices.

For this sentence, the labels would be:

Current [B] economic [I] weakness [I] may [O] be [O] a [B] result [I] of [I] high [I] energy [I] prices [I].

Now, consider the following CRF model for the above problem:

$$P(\mathbf{y} \mid \mathbf{x}; \mathbf{w}) = \frac{1}{Z(\mathbf{x}; \mathbf{w})} \exp \left\{ \mathbf{w}^\top \sum_{i=1}^N f(x_i, y_i, y_{i-1}) \right\} \quad (10)$$

where:

$$Z(\mathbf{x}; w) = \sum_{\mathbf{y} \in \Gamma^N} \exp \left\{ \mathbf{w}^\top \sum_{i=1}^N f(x_i, y_i, y_{i-1}) \right\} \quad (11)$$

where  $\mathbf{w} \in \mathbb{R}^m$ , and  $f : \Sigma \times \Gamma \times \Gamma \rightarrow \mathbb{R}^m$  is the feature vector defined over  $\Sigma$ , the English vocabulary.

- There are many features that may be useful for this task. Suggest one such feature.
- Describe a polynomial time algorithm which given  $\mathbf{w}$ ,  $f$  and the sentence  $\mathbf{x}$ , calculates the probability that a particular subsequence (a subset of the  $y_i$ s in  $\mathbf{y}$ ) is an NP.
- Consider the task of learning  $w$ , the parameters of our model. Write the conditional log-likelihood function of the entire data  $X$  and  $Y$ . Write the update rule for  $\mathbf{w}$  using a gradient ascent.

**3. (10 points)** Consider a HMM with a Gaussian mixture model. The probability of observation  $x_t$  given state  $z_t$  being  $i$  is given by a mixture of Gaussians:

$$p(x_t \mid z_t = i, \theta) = \sum_{c=1}^C w_{ic} \mathcal{N}(x_t \mid \mu_{ic}, \Sigma_{ic}) \quad (12)$$

where we have  $N$  states ( $i = 1, \dots, N$ ), and  $T$  observations in an instance ( $t = 1, \dots, T$ ), where each observation is real valued ( $x \in \mathbb{R}^M$ .)  $\theta$  is the set of all parameters including the transition probabilities  $\{p(z_t \mid z_{t-1})\}$ , parameters of the class-conditional densities  $\{\mu_{ic}, \Sigma_{ic}\}$ , initial state probabilities  $\{\pi_i\}$  and the mixture weights  $\{w_{ik}\}$ .

Given a set of unlabeled training examples  $\{X^{(j)}\}_{j=1}^P$  where  $X^{(j)} = [x_1^{(j)}, \dots, x_T^{(j)}]$ , the task is to learn the parameters  $\theta$  using EM algorithm.

- (a) How many parameters are there in  $\theta$ ? State how many parameters there are for each type.
  - (b) Many applications have high-dimensional data which makes this approach computationally inefficient. One possible relaxation is called *tied-mixture HMM*, where the choice of state only influences the mixture weights  $w_{ic}$ , not the Gaussian parameters. Since the Gaussians are independent of the state we need to estimate fewer parameters. Derive the  $E$  and  $M$  steps for this model.
  - (c) Consider the original Gaussian HMM model (not the tied-mixture model). Suppose the training data provided the states  $z_t$ . Describe the training algorithm for estimating  $\theta$ .
- 4. (10 points)** Consider 3 random variables  $\{X_i\}_{i=1}^3$ , drawn independently from a Bernoulli distribution with parameter  $\theta$ . Let's define  $Z_1 = X_1 \otimes X_2$  and  $Z_2 = X_2 \otimes X_3$  where  $\otimes$  is the XOR operator.
- (a) Calculate  $p(X_1, X_2 \mid Z_1)$  and  $p(X_2, X_3 \mid Z_2)$  in terms of  $\theta$ .
  - (b) Describe (or draw) the graphical model associated with this model (contains  $X_1, X_2, X_3, Z_1, Z_2$ ). Write the factorized joint distribution according to this graph structure.
  - (c) Find  $\theta$  such that  $Z_1$  is independent of  $X_1$ , and  $Z_2$  is independent of  $X_3$ . Can you see this independence in the DAG that you drew/described in the previous part? Why?

### 3 What to Submit

In each assignment you will submit two things.

1. **Code:** Your code as a zip file named `library.zip`. **You must submit source code (.java files)**. We will run your code using the exact command lines described above, so make sure it works ahead of time. Remember to submit all of the source code, including what we have provided to you.
2. **Writeup:** Your writeup as a **PDF file** (compiled from latex) containing answers to the analytical questions asked in the assignment. Make sure to include your name in the writeup PDF and use the provided latex template for your answers.

Make sure you name each of the files exactly as specified (`library.zip` and `writeup.pdf`).

To submit your assignment, visit the "Homework" section of the website (<http://www.cs475.org/>.)

### 4 Questions?

Remember to submit questions about the assignment to the appropriate group on the class discussion board: <http://bb.cs475.org>.