

# Probabilistic Graphical Models 10-708

## Homework 1: Due January 29, 2014 at 4 pm

**Directions.** This homework assignment covers the material presented in Lectures 1-3. You must complete all four problems to obtain full credit. To submit your assignment, please upload a pdf file containing your writeup and a zip file containing your code to Canvas by 4 pm on Wednesday, January 29th. We highly encourage that you type your homework using the  $\text{\LaTeX}$  template provided on the course website, but you may also write it by hand and then scan it.

### 1 Fundamentals [25 points]

This question will refer to the graphical models shown in Figures 1 and 2, which encode a set of independencies among the following variables: Season (S), Flu (F), Dehydration (D), Chills (C), Headache (H), Nausea (N), Dizziness (Z). Note that the two models have the same skeleton, but Figure 1 depicts a directed model (Bayesian network) whereas Figure 2 depicts an undirected model (Markov network).

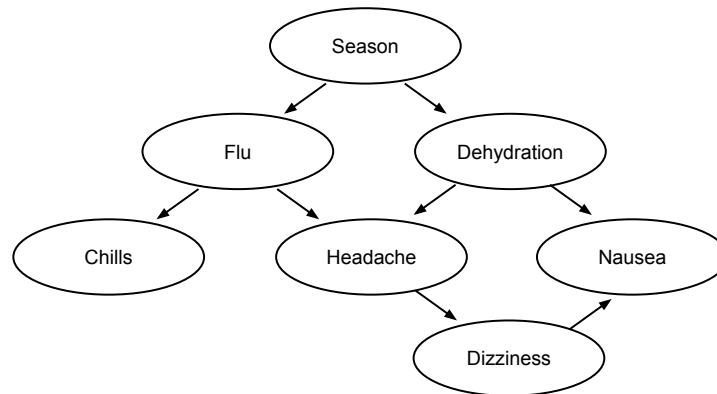


Figure 1: A Bayesian network that represents a joint distribution over the variables Season, Flu, Dehydration, Chills, Headache, Nausea, and Dizziness.

#### Part 1: Independencies in Bayesian Networks [12 points]

Consider the model shown in Figure 1. Indicate whether the following independence statements are true or false according to this model. Provide a very brief justification of your answer (no more than 1 sentence).

1. Season  $\perp$  Chills

False: influence can flow along the path Season  $\rightarrow$  Flu  $\rightarrow$  Chills, since Flu is unobserved

2. Season  $\perp$  Chills | Flu  
True: influence cannot flow through Flu, since it is observed; there are no other paths linking Season and Chills
3. Season  $\perp$  Headache | Flu  
False: influence can flow along the path Season  $\rightarrow$  Dehydration  $\rightarrow$  Headache, since Dehydration is unobserved
4. Season  $\perp$  Headache | Flu, Dehydration  
True: since both Flu and Dehydration are observed, influence cannot flow along any path that links Season and Headache
5. Season  $\perp$  Nausea | Dehydration  
False: influence can flow along the path formed by Season  $\rightarrow$  Flu  $\rightarrow$  Headache  $\rightarrow$  Dizziness  $\rightarrow$  Nausea, since Flu, Headache, and Dizziness are unobserved
6. Season  $\perp$  Nausea | Dehydration, Headache  
True: influence cannot flow along the path Season  $\rightarrow$  Dehydration  $\rightarrow$  Nausea, since Dehydration is observed; influence cannot flow along the path Season  $\rightarrow$  Flu  $\rightarrow$  Headache  $\rightarrow$  Dizziness  $\rightarrow$  Nausea, since Headache is observed; influence cannot flow along the path Season  $\rightarrow$  Flu  $\rightarrow$  Headache  $\leftarrow$  Dehydration  $\rightarrow$  Nausea, even though there is an observed v-structure centered at Headache, because Dehydration is observed
7. Flu  $\perp$  Dehydration  
False: influence can flow along the path Flu  $\leftarrow$  Season  $\rightarrow$  Dehydration, since Season is unobserved
8. Flu  $\perp$  Dehydration | Season, Headache  
False: influence can flow along the path Flu  $\rightarrow$  Headache  $\leftarrow$  Dehydration, since this is a v-structure and Headache is observed
9. Flu  $\perp$  Dehydration | Season  
True: influence cannot flow through Season, which is observed, nor through Headache or Nausea, since both form v-structures and both are unobserved
10. Flu  $\perp$  Dehydration | Season, Nausea  
False: influence can flow along the path Flu  $\rightarrow$  Headache  $\rightarrow$  Dizziness  $\rightarrow$  Nausea  $\leftarrow$  Dehydration, since Headache and Dizziness are unobserved and there is a v-structure at Nausea, which is observed
11. Chills  $\perp$  Nausea  
False: influence can flow along the path Chills  $\leftarrow$  Flu  $\leftarrow$  Season  $\rightarrow$  Dehydration  $\rightarrow$  Nausea, since Flu, Season, and Dehydration are all unobserved
12. Chills  $\perp$  Nausea | Headache  
False: influence can flow along the path Chills  $\leftarrow$  Flu  $\rightarrow$  Headache  $\leftarrow$  Dehydration  $\rightarrow$  Nausea, since there is a v-structure at Headache, which is observed

## Part 2: Factorized Joint Distributions [4 points]

1. Using the directed model shown in Figure 1, write down the factorized form of the joint distribution over all of the variables,  $P(S, F, D, C, H, N, Z)$ .

$$P(S, F, D, C, H, Z, N) = P(S) P(F|S) P(D|S) P(C|F) P(H|F, D) P(Z|H) P(N|D, Z)$$

2. Using the undirected model shown in Figure 2, write down the factorized form of the joint distribution over all of the variables, assuming the model is parameterized by one factor over each node and one over each edge in the graph.

$$\frac{1}{Z} \phi_1(S) \phi_2(F) \phi_3(D) \phi_4(C) \phi_5(H) \phi_6(N) \phi_7(Z) \cdot \phi_8(S, F) \phi_9(S, D) \phi_{10}(F, C) \phi_{11}(F, H) \phi_{12}(D, H) \phi_{13}(D, N) \phi_{14}(H, Z) \phi_{15}(N, Z)$$

### Part 3: Evaluating Probability Queries [7 points]

Assume you are given the conditional probability tables listed in Table 1 for the model shown in Figure 1. Evaluate each of the probabilities queries listed below, and show your calculations.

1. What is the probability that you have the flu, when no prior information is known?

This translates to  $P(\text{Flu} = \text{true})$

$$\begin{aligned} &P(F = \text{true}) \\ &= \sum_s P(F = \text{true}, S = s) \\ &= \sum_s P(F = \text{true} | S = s)P(S = s) \\ &= P(F = \text{true} | S = \text{wint})P(S = \text{wint}) + P(F = \text{true} | S = \text{summ})P(S = \text{summ}) \\ &= 0.4 \cdot 0.5 + 0.1 \cdot 0.5 = 0.25 \end{aligned}$$

2. What is the probability that you have the flu, given that it is winter?

	$P(S = \text{winter})$	$P(S = \text{summer})$
	0.5	0.5

	$P(F = \text{true}   S)$	$P(F = \text{false}   S)$
$S = \text{winter}$	0.4	0.6
$S = \text{summer}$	0.1	0.9

	$P(D = \text{true}   S)$	$P(D = \text{false}   S)$
$S = \text{winter}$	0.1	0.9
$S = \text{summer}$	0.3	0.7

	$P(C = \text{true}   F)$	$P(C = \text{false}   F)$
$F = \text{true}$	0.8	0.2
$F = \text{false}$	0.1	0.9

	$P(H = \text{true}   F, D)$	$P(H = \text{false}   F, D)$
$F = \text{true}, D = \text{true}$	0.9	0.1
$F = \text{true}, D = \text{false}$	0.8	0.2
$F = \text{false}, D = \text{true}$	0.8	0.2
$F = \text{false}, D = \text{false}$	0.3	0.7

	$P(Z = \text{true}   H)$	$P(Z = \text{false}   H)$
$H = \text{true}$	0.8	0.2
$H = \text{false}$	0.2	0.8

	$P(N = \text{true}   D, Z)$	$P(N = \text{false}   D, Z)$
$D = \text{true}, Z = \text{true}$	0.9	0.1
$D = \text{true}, Z = \text{false}$	0.8	0.2
$D = \text{false}, Z = \text{true}$	0.6	0.4
$D = \text{false}, Z = \text{false}$	0.2	0.8

Table 1: Conditional probability tables for the Bayesian network shown in Figure 1.

This translates to  $P(\text{Flu} = \text{true} \mid \text{Season} = \text{winter})$

$$P(F = \text{true} \mid S = \text{wint}) = 0.4$$

3. What is the probability that you have the flu, given that it is winter and that you have a headache?

This translates to  $P(\text{Flu} = \text{true} \mid \text{Season} = \text{winter}, \text{Headache} = \text{true})$

$$\begin{aligned} & P(F = \text{true} \mid S = \text{wint}, H = \text{true}) \\ &= \frac{P(F = \text{true}, S = \text{wint}, H = \text{true})}{P(S = \text{wint}, H = \text{true})} \\ &= \frac{\sum_d P(F = \text{true}, S = \text{wint}, H = \text{true}, D = d)}{\sum_{f,d} P(F = f, S = \text{wint}, H = \text{true}, D = d)} \\ &= \frac{\sum_d P(H = \text{true} \mid F = \text{true}, D = d)P(F = \text{true} \mid S = \text{wint})P(D = d \mid S = \text{wint})P(S = \text{wint})}{\sum_{f,d} P(H = \text{true} \mid F = f, D = d)P(F = f \mid S = \text{wint})P(D = d \mid S = \text{wint})P(S = \text{wint})} \\ &= \frac{0.9 \cdot 0.4 \cdot 0.1 \cdot 0.5 + 0.8 \cdot 0.4 \cdot 0.9 \cdot 0.5}{0.018 + 0.144} \\ &= \frac{0.018 + 0.144 + 0.024 + 0.081}{0.018 + 0.144 + 0.024 + 0.081} = 0.61 \end{aligned}$$

4. What is the probability that you have the flu, given that it is winter, you have a headache, and you know that you are dehydrated?

This translates to  $P(\text{Flu} = \text{true} \mid \text{Season} = \text{winter}, \text{Headache} = \text{true}, \text{Dehydration} = \text{true})$

$$\begin{aligned} & P(F = \text{true} \mid S = \text{wint}, H = \text{true}, D = \text{true}) \\ &= \frac{P(F = \text{true}, S = \text{wint}, H = \text{true}, D = \text{true})}{P(S = \text{wint}, H = \text{true}, D = \text{true})} \\ &= \frac{P(F = \text{true}, S = \text{wint}, H = \text{true}, D = \text{true})}{\sum_f P(F = f, S = \text{wint}, H = \text{true}, D = \text{true})} \\ &= \frac{P(H = \text{true} \mid F = \text{true}, D = \text{true})P(F = \text{true} \mid S = \text{wint})P(D = \text{true} \mid S = \text{wint})P(S = \text{wint})}{\sum_f P(H = \text{true} \mid F = f, D = \text{true})P(F = f \mid S = \text{wint})P(D = \text{true} \mid S = \text{wint})P(S = \text{wint})} \\ &= \frac{0.9 \cdot 0.4 \cdot 0.1 \cdot 0.5}{0.9 \cdot 0.4 \cdot 0.1 \cdot 0.5 + 0.8 \cdot 0.6 \cdot 0.1 \cdot 0.5} = \frac{0.018}{0.018 + 0.024} = 0.43 \end{aligned}$$

5. Does knowing you are dehydrated increase or decrease your likelihood of having the flu? Intuitively, does this make sense?

Knowing that you are dehydrated decreases the likelihood that you have the flu. This makes sense because the headache symptom is “explained away” by the dehydration.

## Part 4: Bayesian Networks vs. Markov Networks [2 points]

Now consider the undirected model shown in Figure 2.

1. Are there any differences between the set of marginal independencies encoded by the directed and undirected versions of this model? If not, state the full set of marginal independencies encoded by both models. If so, give one example of a difference.

There are no differences, because neither model encodes any marginal independencies at all.

2. Are there any differences between the set of conditional independencies encoded by the directed and undirected versions of this model? If so, give one example of a difference.

There are several differences. One example is that in the Markov network, we have  $\text{Flu} \perp \text{Dehydration} \mid \text{Season}, \text{Headache}$ . However, this is not the case in the Bayesian network because observing Headache creates an active v-structure at  $\text{Flu} \rightarrow \text{Headache} \leftarrow \text{Dehydration}$ .

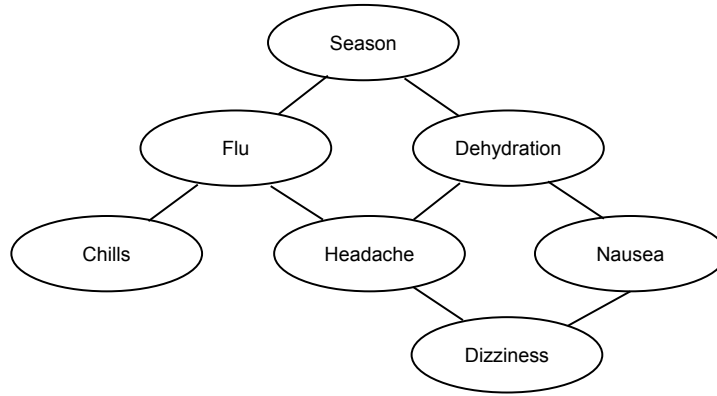


Figure 2: A Markov network that represents a joint distribution over the variables Season, Flu, Dehydration, Chills, Headache, Nausea, and Dizziness.

## 2 Bayesian Networks [25 points]

### Part 1: Constructing Bayesian Networks [8 points]

In this problem you will construct your own Bayesian network (BN) for a few different modeling scenarios described as word problems. By standard convention, we will use shaded circles to represent observed quantities, clear circles to represent random variables, and uncircled symbols to represent distribution parameters.

In order to do this problem, you will first need to understand plate notation, which is a useful tool for drawing large BNs with many variables. Plates can be used to denote repeated sets of random variables. For example, suppose we have the following generative process:

- Draw  $Y \sim \text{Normal}(\mu, \Sigma)$
- For  $m = 1, \dots, M$ :

$$\text{Draw } X_m \sim \text{Normal}(Y, \Sigma)$$

This BN contains  $M + 1$  random variables, which includes  $M$  repeated variables  $X_1, \dots, X_M$  that all have  $Y$  as a parent. In the BN, we draw the repeated variables by placing a box around a single node, with an index in the box describing the number of copies; we've drawn this in Figure 3.

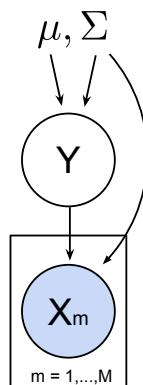


Figure 3: An example of a Bayesian network drawn with plate notation.

For each of the modeling scenarios described below, draw a corresponding BN. Make sure to label your nodes using the variable names given below, and use plate notation if necessary.

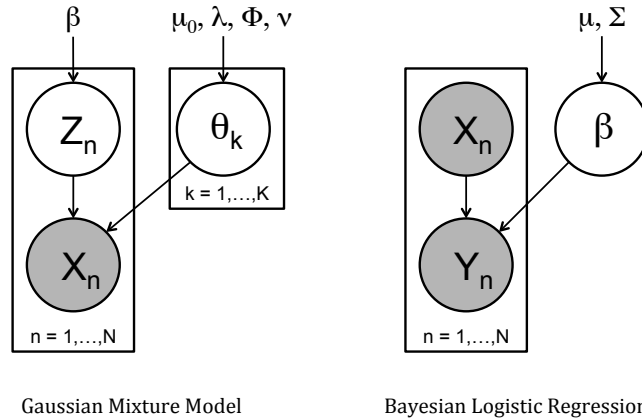
- (Gaussian Mixture Model). Suppose you want to model a set of clusters within a population of  $N$  entities,  $X_1, \dots, X_N$ . We assume there are  $K$  clusters  $\theta_1, \dots, \theta_K$ , and that each cluster represents a vector and a matrix,  $\theta_k = \{\mu_k, \Sigma_k\}$ . We also assume that each entity  $X_n$  “belongs” to one cluster, and its membership is given by an assignment variable  $Z_n \in \{1, \dots, K\}$ .

Here’s how the variables in the model relate. Each entity  $X_n$  is drawn from a so-called “mixture distribution,” which in this case is a Gaussian distribution, based on its individual cluster assignment and the entire set of clusters, written  $X_n \sim \text{Normal}(\mu_{Z_n}, \Sigma_{Z_n})$ . Each cluster assignment  $Z_n$  has a prior, given by  $Z_n \sim \text{Categorical}(\beta)$ . Finally, each cluster  $\theta_k$  also has a prior, given by  $\theta_k \sim \text{Normal-invWishart}(\mu_0, \lambda, \Phi, \nu) = \text{Normal}(\mu_0, \frac{1}{\lambda}\Sigma) \cdot \text{invWishart}(\Phi, \nu)$ .

- (Bayesian Logistic Regression). Suppose you want to model the underlying relationship between a set of  $N$  input vectors  $X_1, \dots, X_N$  and a corresponding set of  $N$  binary outcomes  $Y_1, \dots, Y_N$ . We assume there is a single vector  $\beta$  which dictates the relationship between each input vector and its associated output variable.

In this model, each output is drawn with  $Y_n \sim \text{Bernoulli}(\text{invLogit}(X_n\beta))$ . Additionally, the vector  $\beta$  has a prior, given by  $\beta \sim \text{Normal}(\mu, \Sigma)$ .

The correct graphical models are shown below. Note that for Bayesian logistic regression, it’s also correct to draw  $\{X_n\}$  as a set of fixed “parameters” since they are technically not random variables.



## Part 2: Inference in Bayesian Networks [12 points]

In this problem you will derive formulas for inference tasks in Bayesian networks. Consider the Bayesian network given in Figure 4.

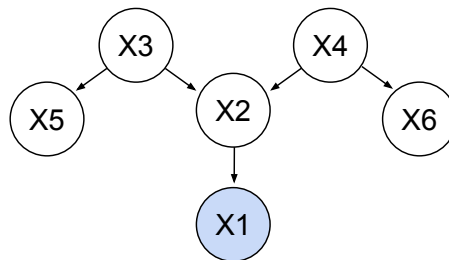


Figure 4: A Bayesian network over the variables  $X_1, \dots, X_6$ . Note that  $X_1$  is observed (which is denoted by the fact that it’s shaded in) and the remaining variables are unobserved.

For each of the following questions, write down an expression involving the variables  $X_1, \dots, X_6$  that could be computed by directly plugging in their local conditional probability distributions.

First, give expressions for the following three posterior distributions over a particular variable given the observed evidence  $X_1 = x_1$ .

1.  $P(X_2 = x_2 | X_1 = x_1)$   

$$= \frac{P(X_1 = x_1 | X_2 = x_2) \sum_{X_3} \sum_{X_4} P(X_2 = x_2 | X_3, X_4) P(X_3) P(X_4)}{\sum_{X_2} P(X_1 = x_1 | X_2) \sum_{X_3} \sum_{X_4} P(X_2 | X_3, X_4) P(X_3) P(X_4)}$$
2.  $P(X_3 = x_3 | X_1 = x_1)$   

$$= \frac{P(X_3 = x_3) \sum_{X_2} P(X_1 = x_1 | X_2) \sum_{X_4} P(X_2 | X_3 = x_3, X_4) P(X_4)}{\sum_{X_2} P(X_1 = x_1 | X_2) \sum_{X_3} \sum_{X_4} P(X_2 | X_3, X_4) P(X_3) P(X_4)}$$
3.  $P(X_5 = x_5 | X_1 = x_1)$   

$$= \frac{\sum_{X_2} P(X_1 = x_1 | X_2) \sum_{X_3} \sum_{X_4} P(X_3) P(X_4) P(X_5 = x_5 | X_3) P(X_2 | X_3, X_4)}{[\sum_{X_3} P(X_5 = x_5 | X_3) P(X_3)] [\sum_{X_2} P(X_1 = x_1 | X_2) \sum_{X_3} \sum_{X_4} P(X_2 | X_3, X_4) P(X_3) P(X_4)]}$$

Second, give expressions for the following three conditional probability queries. Note that these types of expressions are useful for the inference algorithms that we'll learn later in the class.

4.  $P(X_2 = x_2 | X_1 = x_1, X_3 = x_3, X_4 = x_4, X_5 = x_5, X_6 = x_6)$   

$$= \frac{P(X_1 = x_1 | X_2 = x_2) P(X_2 = x_2 | X_3 = x_3, X_4 = x_4) P(X_3 = x_3) P(X_4 = x_4)}{\sum_{X_2} P(X_1 = x_1 | X_2) P(X_2 | X_3 = x_3, X_4 = x_4) P(X_3 = x_3) P(X_4 = x_4)}$$
5.  $P(X_3 = x_3 | X_1 = x_1, X_2 = x_2, X_4 = x_4, X_5 = x_5, X_6 = x_6)$   

$$= \frac{P(X_2 = x_2 | X_3 = x_3, X_4 = x_4) P(X_5 = x_5 | X_3 = x_3) P(X_3 = x_3) P(X_4 = x_4)}{\sum_{X_3} P(X_2 = x_2 | X_3, X_4 = x_4) P(X_5 = x_5 | X_3) P(X_3) P(X_4 = x_4)}$$
6.  $P(X_5 = x_5 | X_1 = x_1, X_2 = x_2, X_3 = x_3, X_4 = x_4, X_6 = x_6)$   

$$= P(X_5 = x_5 | X_3 = x_3)$$

### Part 3: On Markov Blankets [5 points]

In this problem you will prove a key property of Markov blankets in Bayesian networks. Recall that the Markov blanket of a node in a BN consists of the node's children, parents, and coparents (i.e. the children's other parents). Also recall that there are four basic types of two-edge trails in a BN, which are illustrated in Figure 5: the causal trail (head-to-tail), evidential trail (tail-to-head), common cause (tail-to-tail), and common effect (head-to-head).

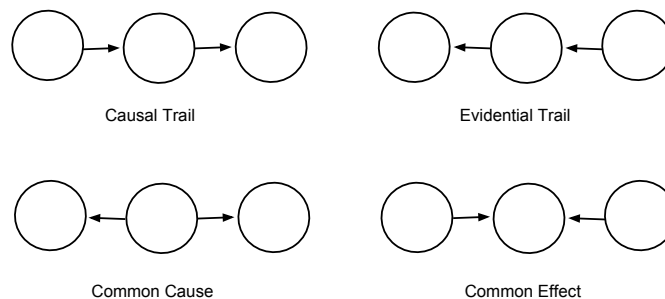


Figure 5: Illustration of the four basic types of two-edge trails in a BN.

Using the four trail types, prove the following property of BNs: *given its Markov blanket, a node in a Bayesian network is conditionally independent of every other set of nodes.*

Proof: The Markov blanket for a node  $X$  consists of its children, parents, and coparents. Assume that we condition on the Markov blanket of  $X$ . In order for  $X$  to be conditionally *dependent* on any other set of nodes  $S$ , there must exist an active trail between  $X$  and a member of  $S$ . We will show there does not exist an active trail. First, the active trail cannot include the edge between  $X$  and any of its parents. If it did, this would imply either an evidential trail or a common cause (starting at node  $X$ , going through the parent), and both of these do not yield an active trail when the parent is conditioned upon. Secondly, the active trail cannot include the edge between  $X$  and any of its children. This is because it would imply either a casual trail or a common effect. In the first case, a casual trail (starting at node  $X$ , going through the child) would not yield an active trail when the child is conditioned upon. In the second case, a common effect (starting at node  $X$ , going through the child, and ending at a coparent) would yield an active trail; however, this implies either an evidential trail or a common cause (starting at the child of  $X$ , going through the coparent), and both of these do not yield an active trail when the coparent is conditioned upon. Therefore, in all cases,  $X$  is d-separated from any set  $S$  given its Markov blanket, and is therefore conditionally independent of any set  $S$ .

### 3 Restricted Boltzmann Machines [25 points]

Restricted Boltzmann Machines (RBMs) are a class of Markov networks that have been used in several applications, including image feature extraction, collaborative filtering, and recently in deep belief networks. An RBM is a bipartite Markov network consisting of a visible (observed) layer and a hidden layer, where each node is a binary random variable. One way to look at an RBM is that it models latent factors that can be learned from input features. For example, suppose we have samples of binary user ratings (like vs. dislike) on 5 movies: Finding Nemo ( $V_1$ ), Avatar ( $V_2$ ), Star Trek ( $V_3$ ), Aladdin ( $V_4$ ), and Frozen ( $V_5$ ). We can construct the following RBM:

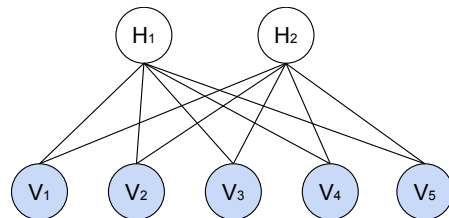


Figure 6: An example RBM with 5 visible units and 2 hidden units.

Here, the bottom layer consists of visible nodes  $V_1, \dots, V_5$  that are random variables representing the binary ratings for the 5 movies, and  $H_1, H_2$  are two hidden units that represent latent factors to be learned during training (e.g.,  $H_1$  might be associated with Disney movies, and  $H_2$  could represent the adventure genre). If we are using an RBM for image feature extraction, the visible layer could instead denote binary values associated with each pixel, and the hidden layer would represent the latent features. However, for this problem we will stick with the movie example. In the following questions, let  $V = (V_1, \dots, V_5)$  be a vector of ratings (e.g. the observation  $v = (1, 0, 0, 0, 1)$  implies that a user likes only Finding Nemo and Aladdin). Similarly, let  $H = (H_1, H_2)$  be a vector of latent factors. Note that all the random variables are binary and take on states in  $\{0, 1\}$ . The joint distribution of a configuration is given by

$$p(V = v, H = h) = \frac{1}{Z} e^{-E(v,h)} \tag{1}$$



where

$$E(v, h) = - \sum_{ij} w_{ij} v_i h_j - \sum_i a_i v_i - \sum_j b_j h_j$$

is the energy function,  $\{w_{ij}\}, \{a_i\}, \{b_j\}$  are model parameters, and

$$Z = Z(\{w_{ij}\}, \{a_i\}, \{b_j\}) = \sum_{v, h} e^{-E(v, h)}$$

is the partition function, where the summation runs over all joint assignments to  $V$  and  $H$ .

1. [7 pts] Using Equation (1), show that  $p(H|V)$ , the distribution of the hidden units conditioned on all of the visible units, can be factorized as

$$p(H|V) = \prod_j p(H_j|V) \quad (2)$$

where

$$p(H_j = 1|V = v) = \sigma \left( b_j + \sum_i w_{ij} v_i \right)$$

and  $\sigma(s) = \frac{e^s}{1+e^s}$  is the sigmoid function. Note that  $p(H_j = 0|V = v) = 1 - p(H_j = 1|V = v)$ .

$$\begin{aligned} p(H = h|V = v) &= \frac{p(v, h)}{p(v)} = \frac{p(v, h)}{\sum_h p(v, h)} \\ &= \frac{\exp(\sum_i a_i v_i) \exp(\sum_{ij} w_{ij} v_i h_j + \sum_j b_j h_j)}{\exp(\sum_i a_i v_i) \sum_h \exp(\sum_{ij} w_{ij} v_i h_j + \sum_j b_j h_j)} \\ &= \frac{\prod_j \exp(\sum_i w_{ij} v_i h_j + b_j h_j)}{\sum_h \prod_j \exp(\sum_i w_{ij} v_i h_j + b_j h_j)} \\ &= \frac{\prod_j \exp(\sum_i w_{ij} v_i h_j + b_j h_j)}{\prod_j \sum_{h_j} \exp(\sum_i w_{ij} v_i h_j + b_j h_j)} \\ &= \prod_j \frac{\exp(\sum_i w_{ij} v_i h_j + b_j h_j)}{1 + \exp(\sum_i w_{ij} v_i + b_j)} \\ &= \prod_j p(h_j|v) \end{aligned}$$

and thus  $p(H_j = 1|V = v) = \sigma(\sum_i w_{ij} v_i + b_j)$ . Note during the derivation the sum and product exchanges in the denominator because  $\sum_h \prod_j \exp(\sum_i w_{ij} v_i h_j + b_j h_j) = \sum_{h_1} \dots \sum_{h_n} f(h_1) \dots f(h_n)$  where  $f(h_j) = \exp(\sum_i w_{ij} v_i h_j + b_j h_j)$  so the sum can be pushed into products.

2. [3 pts] Give the factorized form of  $p(V|H)$ , the distribution of the visible units conditioned on all of the hidden units. This should be similar to what's given in part 1, and so you may omit the derivation.

By symmetry, we have

$$p(V|H) = \prod_i p(V_i|H)$$

and

$$p(V_i = 1|H = h) = \sigma \left( a_i + \sum_j w_{ij} h_j \right)$$

3. [2 pts] Can the marginal distribution over hidden units  $p(H)$  be factorized? If yes, give the factorization. If not, give the form of  $p(H)$  and briefly justify.

No. The form of  $p(H)$  is given by:

$$\begin{aligned}
 p(h) &= \sum_v p(v, h) \\
 &\propto \sum_v \exp(-E(v, h)) \\
 &= \sum_v \exp\left(\sum_{ij} w_{ij} v_i h_j + \sum_i a_i v_i + \sum_j b_j h_j\right) \\
 &= \exp\left(\sum_j b_j h_j\right) \sum_v \exp\left(\sum_{i,j} w_{ij} v_i h_j + \sum_i a_i v_i\right) \\
 &= \exp\left(\sum_j b_j h_j\right) \sum_v \exp\left(\sum_i \left(\sum_j w_{ij} v_i h_j + a_i v_i\right)\right) \\
 &= \left(\prod_j \exp(b_j h_j)\right) \left(\sum_v \prod_i \exp\left\{v_i \left(\sum_j w_{ij} h_j + a_i\right)\right\}\right) \\
 &= \left(\prod_j \exp(b_j h_j)\right) \left(\prod_i \sum_{v_i} \exp\left\{v_i \left(\sum_j w_{ij} h_j + a_i\right)\right\}\right) \\
 &= \left(\prod_j \exp(b_j h_j)\right) \left(\prod_i 1 + \exp\left\{\sum_j w_{ij} h_j + a_i\right\}\right)
 \end{aligned}$$

Since the second term is a product over visible units  $i$ , not hidden unit  $j$ ,  $p(h)$  does not factorize.

4. [4 pts] Based on your answers so far, does the distribution in Equation (1) respect the conditional independencies of Figure (6)? Explain why or why not. Are there any independencies in Figure 6 that are not captured in Equation (1)?

Since RBM is a full bipartite graph (all nodes are connected from the nodes on the other side), the only independencies implied by the graph are the ones shown in part 1 and 2. Thus the answer to the two parts are Yes and No.

5. [7 pts] We can use the log-likelihood of the visible units,  $\log p(V = v)$ , as the criterion to learn the model parameters  $\{w_{ij}\}, \{a_i\}, \{b_j\}$ . However, this maximization problem has no closed form solution. One popular technique for training this model is called “contrastive divergence” and uses an approximate gradient descent method. Compute the gradient of the log-likelihood objective with respect to  $w_{ij}$  by showing the following:

$$\begin{aligned}
 \frac{\partial \log p(V = v)}{\partial w_{ij}} &= \sum_h p(H = h|V = v) v_i h_j - \sum_{v,h} p(V = v, H = h) v_i h_j \\
 &= \mathbb{E}[V_i H_j | V = v] - \mathbb{E}[V_i H_j]
 \end{aligned}$$

where  $\mathbb{E}[V_i H_j | V = v]$  can be readily evaluated using Equation (2), but  $\mathbb{E}[V_i H_j]$  is tricky as the expectation is taken over not just  $H_j$  but also  $V_i$ .

Hint 1: To save some writing, do not expand  $E(v, h)$  until you have  $\frac{\partial E(v, h)}{\partial w_{ij}}$ .

Hint 2: The partition function,  $Z$ , is a function of  $w_{ij}$ .

$$\begin{aligned}
\frac{\partial \log p(v)}{\partial w_{ij}} &= \frac{\partial}{\partial w_{ij}} \log \sum_h p(v, h) \\
&= \frac{\partial}{\partial w_{ij}} \log \frac{\sum_h \exp(-E(v, h))}{Z} \\
&= \frac{Z}{\sum_h \exp(-E(v, h))} \frac{\partial}{\partial w_{ij}} \frac{\sum_h \exp(-E(v, h))}{Z} \\
&= \frac{Z}{\sum_h \exp(-E(v, h))} \left( \frac{1}{Z} \sum_h \frac{\partial}{\partial w_{ij}} \exp(-E(v, h)) - \frac{\sum_h \exp(-E(v, h))}{Z^2} \frac{\partial Z}{\partial w_{ij}} \right) \\
&= - \sum_h \frac{\exp(-E(v, h))}{\sum_h \exp(-E(v, h))} \frac{\partial E(v, h)}{\partial w_{ij}} - \frac{1}{Z} \frac{\partial}{\partial w_{ij}} \sum_{v, h} \exp(-E(v, h)) \\
&= \sum_h p(h|v) v_i h_j + \sum_{v, h} \frac{\exp(-E(v, h))}{Z} \frac{\partial E(v, h)}{\partial w_{ij}} \\
&= \sum_h p(h|v) v_i h_j - \sum_{v, h} p(v, h) v_i h_j
\end{aligned}$$

6. [2 pts] After training, suppose  $H_1 = 1$  corresponds to Disney movies, and  $H_2 = 1$  corresponds to the adventure genre. Which  $w_{ij}$  do you expect to be positive, where  $i$  indexes the visible nodes and  $j$  indexes the hidden nodes? List all of them.

$w_{11}, w_{41}, w_{51}, w_{22}, w_{32}$

## 4 Image Denoising [25 points]

This is a programming problem involving Markov networks (MNs) applied to the task of image denoising. Suppose we have an image consisting of a 2-dimensional array of pixels, where each pixel value  $Z_i$  is binary, i.e.  $Z_i \in \{+1, -1\}$ . Assume now that we make a noisy copy of the image, where each pixel in the image is flipped with 10% probability. A pixel in this noisy image is denoted by  $X_i$ . We show the original image and the image with 10% noise in Figure 7.

Given the observed array of noisy pixels, our goal is to recover the original array of pixels. To solve this problem, we model the original image and noisy image with the following MN. We have a latent variable  $Z_i$  for each noise-free pixel, and an observed variable  $X_i$  for each noisy pixel. Each variable  $Z_i$  has an edge leading to its immediate neighbors (to the  $Z_i$  associated with pixels above, below, to the left, and to the right, when they exist). Additionally, each variable  $Z_i$  has an edge leading to its associated observed pixel  $X_i$ . We illustrate this MN in Figure 8.

Denote the full array of latent (noise-free) pixels as  $\mathbf{Z}$  and the full array of observed (noisy) pixels as  $\mathbf{X}$ . We define the energy function for this model as

$$E(\mathbf{Z} = \mathbf{z}, \mathbf{X} = \mathbf{x}) = h \sum_i z_i - \beta \sum_{\{i, j\}} z_i z_j - \nu \sum_i z_i x_i \quad (3)$$

where the first and third summations are over the entire array of pixels, the second summation is over all pairs of latent variables connected by an edge, and  $h \in \mathbb{R}$ ,  $\beta \in \mathbb{R}_+$ , and  $\nu \in \mathbb{R}_+$  denote constants that must be chosen.

Using the binary image data saved in `hw1.images.mat`, your task will be to infer the true value of each pixel (+1 or -1) by optimizing the above energy function. To do this, initialize the  $Z_i$ 's to their noisy values, and then iterate through each  $Z_i$  and check whether setting it's value to +1 or



Figure 7: The original binary image is shown on the left, and a noisy version of the image in which a randomly selected 10% of the pixels have been flipped is shown on the right.

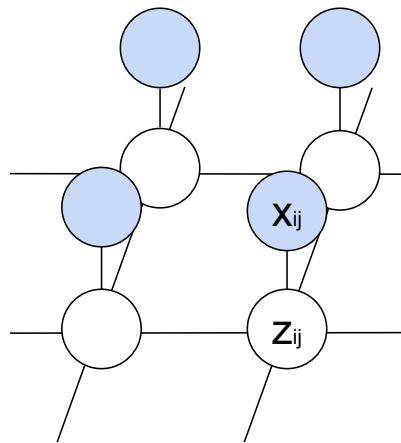


Figure 8: Illustration of the Markov network for image denoising.

$-1$  yields a lower energy (higher probability). Repeat this process, making passes through all of the pixels, until the total energy of the model has converged. You must specify values of the constants  $h \in \mathbb{R}$ ,  $\beta \in \mathbb{R}_+$ , and  $\nu \in \mathbb{R}_+$ .

Report the error rate (fraction of pixels recovered incorrectly) that you achieve by comparing your denoised image to the original image that we provide, for three different settings of the three constants. Include a figure of your best denoised image in your writeup. Also make sure to submit a zipped copy of your code. The TAs will give a “special prize” to the student who is able to achieve the lowest error on this task.

Hint 1: When evaluating whether  $+1$  or  $-1$  is a better choice for a particular pixel  $Z_i$ , you do not need to evaluate the entire energy function, as this will be computationally very expensive. Instead, just compute the contribution by the terms that are affected by the value of  $Z_i$ .

Hint 2: If you’d like to try and compete to achieve the best performance, you can work to find good parameters, or even modify the algorithm in an intelligent way (be creative!). However, if you come up with a modified algorithm, you should separately report the new error rate you achieve, and also turn in a second .m file (placed in your zipped code directory) with the modified algorithm.

[For a solution, please see the code in the zipped directory.](#)