

# 10-701 Introduction to Machine Learning

## Homework 5

*Due Dec 1, 11:59 am*

---

### Rules:

1. Homework submission is done via CMU Autolab system. Please compile your writeup and code in a single pdf file and submit to <https://autolab.cs.cmu.edu/courses/10701-f15>. Please let us know asap if you have trouble accessing Autolab.
  2. Like conference websites, repeated submission is allowed. So please feel free to refine your answers. We will only grade the latest version.
  3. Autolab may allow submission after the deadline, note however it is because of the late day policy. Please see course website for policy on late submission.
  4. Please typeset your homework using appropriate software such as L<sup>A</sup>T<sub>E</sub>X. We will not accept scanned copies of handwritten papers.
  5. You are allowed to collaborate on the homework, but you should write up your own solution and code. Please indicate your collaborators in your submission.
-

# 1 Gaussian Graphical Model (10 Points) (Hao)

In this problem, we will explore the connections between multivariate Gaussian distribution and undirected graphical models.

Given  $\mathbf{X} = \{X_1, \dots, X_d\} \in \mathbb{R}^d$  and  $P(\mathbf{X}|\mu, \Sigma) = \mathcal{N}(\mu, \Sigma)$ , and let  $\Omega = \Sigma^{-1}$  be the precision matrix (inverse covariance matrix). We can represent the distribution as an undirected graph. Specifically, let  $G = (V, E)$  be a (complete) graph with  $V = \{1, 2, \dots, d\}$ .

1. Following the notations of undirected graphical models, we can write

$$P(\mathbf{X}|\mu, \Sigma) \propto \prod_{(i,j) \in E} \psi_{ij}(x_i, x_j) \prod_{i \in V} \psi_i(x_i)$$

Derive  $\psi_{ij}(x_i, x_j)$  and  $\psi_i(x_i)$ .

2. Also, we can rewrite the above as a product of factors with respect to edges, i.e.,

$$P(\mathbf{X}|\mu, \Sigma) \propto \prod_{(i,j) \in E} \psi'_{ij}(x_i, x_j)$$

where  $\psi'_{ij}(x_i, x_j) = \psi_{ij}(x_i, x_j) \times \psi_i(x_i)^{\frac{1}{n(i)}} \times \psi_j(x_j)^{\frac{1}{n(j)}}$  and  $n(i)$  be the number of neighbors of  $i$ . Use this formulation to argue the following equivalence,

$$(i, j) \notin E \iff X_i \perp\!\!\!\perp X_j | X_{V \setminus \{i, j\}}$$

# 2 Sampling Methods (30 Points)(Hao)

## 2.1 Inverse sampling

1. One of the most widely used sampling methods is the *inverse sampling*. To sample from a distribution  $p(y)$ , we first sample a random variable  $z$  from the uniform distribution over  $(0, 1)$ , then transform  $z$  using  $y = h^{-1}(z)$  where  $h(y)$  is defined as,

$$h(y) = \int_{-\infty}^y p(\hat{y}) d\hat{y}$$

Prove that this yields a sample  $y$  which follows the distribution  $p(y)$ . Give one drawback of the inverse sampling method.

2. Given a random variable  $z$  which is uniformly distributed over  $(0, 1)$ , find a transformation  $y = g(z)$  such that  $y$  has a Cauchy distribution given by  $y \sim \frac{1}{\pi} \frac{1}{1+y^2}$ .

## 2.2 Rejection sampling

Sometimes directly drawing samples from the desired distribution  $p(z)$  is hard, but the value of  $\tilde{p}(x)$  can readily be evaluated, where  $p(x) = \frac{1}{Z_{\tilde{p}}} \tilde{p}(x)$  with  $Z_{\tilde{p}}$  as a normalization constant. In this case, we refer to another sampling method called *rejection sampling*.

The procedures of rejection sampling are described as follows: we first introduce a proposal distribution  $q(z)$ , from which we can readily draw samples. We next introduce a (smallest) constant  $k$  whose value is chosen such that  $kq(z) \geq \tilde{p}(z)$  for all values of  $z$ . To sample from  $p(z)$ , we first generate a number  $z_0$  from the distribution  $q(z)$ . Next, we generate a number  $u_0$  from the uniform distribution over  $[0, kq(z_0)]$ . Finally, if  $u_0 > \tilde{p}(z_0)$  then the sample is rejected, otherwise  $u_0$  is retained. Write down the probability that given  $u_0$  that  $u_0$  is accepted, and further prove that the above procedure yields a sample  $u_0$  which follows the distribution  $p(z)$ . Also, give one drawback of this method.

## 2.3 Markov chain Monte Carlo

Markov chain Monte Carlo (MCMC) methods are a class of algorithms for sampling from a probability distribution based on constructing a Markov chain that has the desired distribution as its equilibrium distribution. The state of the chain after a number of steps is then used as a sample of the desired distribution. The quality of the sample improves as a function of the number of steps.

1. *Metropolis-Hastings* (MH) algorithm is a MCMC method for obtaining a sequence of random samples from a probability distribution for which direct sampling is difficult. This sequence can be used to approximate the distribution. Similar to the rejection sampling, MH assumes that  $\tilde{p}(x)$ , the unnormalized target distribution, is easy to evaluate, and samples from the target distribution  $p(x)$  using the proposal distribution  $q$ . In particular, at step  $\tau$ , in which the current state is  $z^{(\tau)}$ , we draw  $z^*$  from the distribution  $q(z|z^{(\tau)})$  and accept it with probability

$$A = \min \left( 1, \frac{\tilde{p}(z^*)q(z^{(\tau)}|z^*)}{\tilde{p}(z^{(\tau)})q(z^*|z^{(\tau)})} \right)$$

Repeating the above procedures thus forms a Markov chain. Show that  $p(x)$  is the stationary distribution of the Markov chain. Once again, describe a major drawback of MH algorithm. (Hint: See the pp.29 in the slides of lecture 18 about how to prove  $p(x)$  is a stationary distribution of a Markov chain).

2. In class we showed the procedures of Gibbs sampling for target distribution  $p(\mathbf{x}) = p(x_1, \dots, x_d)$ : for each  $j \in \{1, \dots, d\}$ , draw  $t \sim p(x_j|\text{rest})$  and set  $x_j = t$ . Show that  $p(\mathbf{x})$  is the stationary distribution of the Markov chain defined by this procedure.
3. Finally, prove that Gibbs sampling is a special case of the MH algorithm. Will Gibbs sampling suffer the drawback you described for MH algorithm?

## 3 Expectation Maximization (EM) and Variational Inference (VI) (35 Points) (Zhiting)

Here we give a general treatment of the EM algorithm, and introduce VI which can be derived similarly <sup>1</sup>.

Hint: the Kullback-Leibler (KL) divergence between two probability distributions  $q(\mathbf{x})$  and  $p(\mathbf{x})$  is defined as  $KL(q||p) = -\int_{\mathbf{x}} q(\mathbf{x}) \ln \frac{p(\mathbf{x})}{q(\mathbf{x})}$ ; and we have  $KL(q||p) \geq 0$  and  $=$  holds if and only if  $q = p$ .

### 3.1 EM

Let  $\mathbf{x}$  be observed variables,  $\mathbf{z}$  latent variables, and  $\theta$  the parameters. The EM algorithm maximizes the likelihood function  $p(\mathbf{x}|\theta) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}|\theta)$  through a two-stage iterative optimization.

1. Introduce a distribution  $q(\mathbf{z})$  over the latent variables. Show that for any choice of  $q(\mathbf{z})$ , the following decomposition holds

$$\ln p(\mathbf{x}|\theta) = \mathcal{L}(q, \theta) + KL(q||p), \tag{1}$$

where

$$\begin{aligned} \mathcal{L}(q, \theta) &= \int_{\mathbf{z}} q(\mathbf{z}) \ln \left\{ \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z})} \right\}, \\ KL(q||p) &= - \int_{\mathbf{z}} q(\mathbf{z}) \ln \left\{ \frac{p(\mathbf{z}|\mathbf{x}, \theta)}{q(\mathbf{z})} \right\}, \end{aligned}$$

so that  $\mathcal{L}(q, \theta)$  is the lower bound on  $\ln p(\mathbf{x}|\theta)$ .

---

<sup>1</sup>PGM (10708) will cover more advanced topics on sampling and VI. You are welcome to take the course :)

2. In the E-step, the lower bound  $\mathcal{L}(q, \theta)$  is maximized w.r.t  $q(\mathbf{z})$  while fixing  $\theta$ . Show that  $\mathcal{L}(q, \theta)$  is maximized when  $q(\mathbf{z}) = p(\mathbf{z}|\mathbf{x}, \theta)$ .
3. In the subsequent M-step, the lower bound  $\mathcal{L}(q, \theta)$  is maximized w.r.t  $\theta$  while fixing  $q(\mathbf{z})$ . This step will necessarily increase the log likelihood  $\ln p(\mathbf{x}|\theta)$ . Explain why.

### 3.2 VI

For simplicity, here we omit the parameter  $\theta$  and only consider the latent variables  $\mathbf{z}$  and the observed variables  $\mathbf{x}$ . We repeat Eq(1) as follows

$$\ln p(\mathbf{x}) = \mathcal{L}(q) + KL(q||p).$$

From the above we know the lower bound  $\mathcal{L}(q)$  is maximized when  $q(\mathbf{z})$  equals the posterior  $p(\mathbf{z}|\mathbf{x})$ . However sometimes the true posterior is intractable so we have to resort to some approximation, one popular technique of which is the variational inference. The main idea is to restrict the family of distributions  $q(\mathbf{z})$  and find the member of this family for which the lower bound  $\mathcal{L}(q)$  is largest.

Let  $\mathbf{z} = \{z_1, \dots, z_K\}$ . Here we restrict the  $q$  family by assuming the distribution factorizes as

$$q(\mathbf{z}) = \prod_{k=1}^K q_k(z_k).$$

To find the optimal  $q$  in this family, we maximize  $\mathcal{L}(q)$  w.r.t each of the factors in turn.

1. Show that for any  $k \in \{1, \dots, K\}$ ,

$$\mathcal{L}(q) = \int q_k(z_k) \mathbb{E}_{j \neq k} [\ln p(\mathbf{x}, \mathbf{z})] dz_k - \int q_k(z_k) \ln q_k(z_k) dz_k + const,$$

where  $\mathbb{E}_{j \neq k} [\ln p(\mathbf{x}, \mathbf{z})]$  denotes the expectation of  $\ln p(\mathbf{x}, \mathbf{z})$  w.r.t the  $q$  distribution over all  $z_j$  for  $j \neq k$  (i.e.,  $\prod_{j \neq k} q_j$ ); and *const* is some constant irrelevant to  $q$ .

2. Keeping  $\{q_j\}_{j \neq k}$  fixed, show that the optimal  $q_k(z_k)$  is given by

$$\ln q_k^*(z_k) = \mathbb{E}_{j \neq k} [\ln p(\mathbf{x}, \mathbf{z})] + const.$$

## 4 HMM (20 Points) (Zhiting)

**1.1** Consider a HMM with 6 states (plus a start and end states) and an alphabet  $\{A, C, G, T\}$ . Table 1 lists the transition and emission probabilities, and Figure 1 shows the state diagram.

	0	$S_1$	$S_2$	$S_3$	$S_4$	$S_5$	$S_6$	f	A	C	G	T
0	0	1	0	0	0	0	0	0				
$S_1$	0	0	1	0	0	0	0	0	0.2	0.3	0	0.5
$S_2$	0	0	0	0.3	0	0.7	0	0	0.6	0.1	0.2	0.1
$S_3$	0	0	0	0	1	0	0	0	0.7	0	0.1	0.2
$S_4$	0	0	0	0	0	0	0	1	0.2	0.3	0.4	0.1
$S_5$	0	0	0	0	0	0	1	0	0.3	0.3	0.3	0.1
$S_6$	0	0	0	0	0	0	0	1	0.5	0.3	0	0.2

Table 1: The transition and emission probabilities.

Let  $z$  denote latent variables and  $x$  denote observed variables. Place  $<$ ,  $>$ , or  $=$  between the two components of each of the following pairs. Justify your answer.

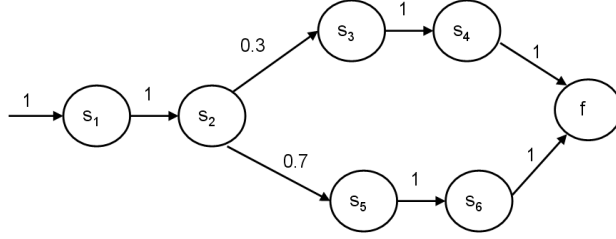


Figure 1: The state diagram of the HMM.

1.  $P(x_1 = T, x_2 = C, x_3 = A, x_4 = T, z_3 = S_3, z_4 = S_4)$   
 $P(x_1 = T, x_2 = C, x_3 = A, x_4 = T | z_3 = S_3, z_4 = S_4)$
2.  $P(x_1 = T, x_2 = C, x_3 = A, x_4 = T)$   
 $P(x_1 = T, x_2 = C, x_3 = A, x_4 = T | z_3 = S_3, z_4 = S_4)$
3.  $P(x_1 = T, x_2 = C, x_3 = A, x_4 = T, z_1 = S_1, z_2 = S_2)$   
 $P(x_1 = T, x_2 = C, x_3 = A, x_4 = T | z_1 = S_1, z_2 = S_2)$
4.  $P(x_1 = T, x_2 = C, x_3 = A, x_4 = T)$   
 $P(x_1 = T, x_2 = A, x_3 = A, x_4 = G)$

1.2 Prove that  $p(x_1, \dots, x_i, z_i) = p(x_i | z_i) \sum_{z_{i-1}} p(x_1, \dots, x_{i-1}, z_{i-1}) p(z_i | z_{i-1})$ .

## 5 Bayesian Networks (Bonus: 10 Points) (Zhiting)

Given  $n$  random variables labeled as  $1, 2, \dots, n$ , how many Bayesian networks (which is DAG) can these variables form? Give a lower bound and upper bound, respective, as tight as you can. Justify your answer.