

Problem Set 3 Solutions

Example problem 1 solution (just use algebra and the right definitions)

Problem 1. At convergence, we have

$$b_i^*(x_i) = \prod_{k \in N(i)} \psi_{ki}^*(x_i) m_{ki}^*(x_i)$$

$$b_{ij}^*(x_i, x_j) = \underbrace{\psi_{ij}^*(x_i, x_j)}_{\psi_{ij}(x_i, x_j)} \prod_{k \in N(i) \setminus j} m_{ki}^*(x_i) \prod_{k \in N(j) \setminus i} m_{kj}^*(x_j)$$

① The ~~max~~ product messages are defined by $m_{ij}(x_j) = \max_{x_i} \left\{ \phi_i(x_i) \psi_{ij}(x_i, x_j) \prod_{k \in N(i) \setminus j} m_{ki}(x_i) \right\}$

At convergence, $\max_{x_j} b_{ij}^*(x_i, x_j) \propto \max_{x_j} \left\{ \phi_i(x_i) \phi_j(x_j) \psi_{ij}(x_i, x_j) \prod_{k \in N(i) \setminus j} m_{ki}^*(x_i) \prod_{k \in N(j) \setminus i} m_{kj}^*(x_j) \right\}$

$$\propto \phi_i(x_i) \prod_{k \in N(i) \setminus j} m_{ki}^*(x_i) \max_{x_j} \left\{ \phi_j(x_j) \psi_{ij}(x_i, x_j) \prod_{k \in N(j) \setminus i} m_{kj}^*(x_j) \right\}$$

$$\propto \phi_i(x_i) \left(\prod_{k \in N(i) \setminus j} m_{ki}^*(x_i) \right) m_{ji}^*(x_i)$$

$$\propto \phi_i(x_i) \prod_{k \in N(i)} m_{ki}^*(x_i) \propto b_i^*(x_i)$$

② The sum-product case is similar to the above, except we replace \max_{x_j} with \sum_{x_j}

The sum product messages are defined as $m_{ji}(x_i) = \sum_{x_j} \phi_j(x_j) \psi_{ij}(x_i, x_j) \prod_{k \in N(j) \setminus i} m_{kj}(x_j)$

At convergence, $\sum_{x_j} b_{ij}^*(x_i, x_j) \propto \sum_{x_j} \phi_i(x_i) \phi_j(x_j) \psi_{ij}(x_i, x_j) \prod_{k \in N(i) \setminus j} m_{ki}^*(x_i) \prod_{k \in N(j) \setminus i} m_{kj}^*(x_j)$

$$\propto \phi_i(x_i) \left(\prod_{k \in N(i) \setminus j} m_{ki}^*(x_i) \right) \sum_{x_j} \left\{ \phi_j(x_j) \psi_{ij}(x_i, x_j) \prod_{k \in N(j) \setminus i} m_{kj}^*(x_j) \right\}$$

$$\propto \phi_i(x_i) \prod_{k \in N(i)} m_{ki}^*(x_i) \propto b_i^*(x_i)$$

③ Taking log of both sides, we have equivalently, $\max_x \log p(x) = \lim_{T \rightarrow 0} \log \left\{ \left(\sum_x p(x)^{\frac{1}{T}} \right)^T \right\}$

Using exponential family representation for $p(x)$, i.e. $p_\theta(x) = \exp\{\langle \theta, \phi(x) \rangle - A(\theta)\}$, the RHS simplifies to

$$\lim_{T \rightarrow 0} T \log \sum_x \exp\left\{\left\langle \frac{\theta}{T}, \phi(x) \right\rangle - \frac{A(\theta)}{T}\right\}, \text{ so equivalently } \max_x \log p(x) = \lim_{T \rightarrow 0} \frac{1}{T} \log \sum_x \exp\left\{\left\langle \frac{\theta}{T}, \phi(x) \right\rangle - \frac{A(\theta)}{T}\right\}$$

$$= \lim_{T \rightarrow 0} T \log \exp\left\{\frac{A(\theta)}{T}\right\} \sum_x \exp\left\{\left\langle \frac{\theta}{T}, \phi(x) \right\rangle\right\} = \lim_{T \rightarrow 0} T \left\{ \log \exp\left\{\frac{A(\theta)}{T}\right\} + \log \sum_x \exp\left\{\left\langle \frac{\theta}{T}, \phi(x) \right\rangle\right\} \right\} = \lim_{T \rightarrow 0} T A\left(\frac{\theta}{T}\right) - A(\theta)$$

where we used the definition of the log partition function $A(\frac{\theta}{T})$

This suggests that we can run sum-product on the distribution parameterized by θ/T (i.e., obtained by rescaling the parameters of the original distribution $p(x)$) with small T , making use of the variational representation of $A(\theta/T)$, and treat the resulting sum-product marginals as approximate MAP marginals.

Note the above solution to part 3 uses the exponential family representation of $p(x)$ (which includes log-linear models covered in class so far), and the notation $A(\theta)$ is the same as $\log Z(\theta)$; the rescaling $\theta \rightarrow \theta/T$, in the notations of log-linear model, corresponds to redefining the potential $\psi(x_c)$ (for every clique c) as $\exp((\log \psi(x_c))/T) = \psi(x_c)^{(1/T)}$.

Problem 2, Part 1.

The key equations for Gibbs sampling for MRFs are in section 12.3.3 of the textbook (Koller & Friedman):

Markov blanket

Gibbs sampling is particularly easy to implement in the many graphical models where we can compute the transition probability $P(X_i | \mathbf{x}_{-i})$ (in line 5 of the algorithm) very efficiently. In particular, as we now show, this distribution can be done based only on the *Markov blanket* of X_i . We show this analysis for a Markov network; the application to Bayesian networks is straightforward. Recalling definition 4.4, we have that:

$$\begin{aligned} P_{\Phi}(X) &= \frac{1}{Z} \prod_j \phi_j(D_j) \\ &= \frac{1}{Z} \prod_{j : X_i \in D_j} \phi_j(D_j) \prod_{j : X_i \notin D_j} \phi_j(D_j). \end{aligned}$$

Let $\mathbf{x}_{j,-i}$ denote the assignment in \mathbf{x}_{-i} to $D_j - \{X_i\}$, noting that when $X_i \notin D_j$, $\mathbf{x}_{j,-i}$ is a

full assignment to D_j . We can now derive:

$$\begin{aligned} P(x'_i | \mathbf{x}_{-i}) &= \frac{P(x'_i, \mathbf{x}_{-i})}{\sum_{x''_i} P(x''_i, \mathbf{x}_{-i})} \\ &= \frac{\frac{1}{Z} \prod_{C_j \ni X_i} \phi_j(x'_i, \mathbf{x}_{j,-i}) \prod_{C_j \not\ni X_i} \phi_j(x'_i, \mathbf{x}_{j,-i})}{\frac{1}{Z} \sum_{x''_i} \prod_{C_j \ni X_i} \phi_j(x''_i, \mathbf{x}_{j,-i}) \prod_{C_j \not\ni X_i} \phi_j(x''_i, \mathbf{x}_{j,-i})} \\ &= \frac{\prod_{C_j \ni X_i} \phi_j(x'_i, \mathbf{x}_{j,-i}) \prod_{C_j \not\ni X_i} \phi_j(\mathbf{x}_{j,-i})}{\sum_{x''_i} \prod_{C_j \ni X_i} \phi_j(x''_i, \mathbf{x}_{j,-i}) \prod_{C_j \not\ni X_i} \phi_j(\mathbf{x}_{j,-i})} \\ &= \frac{\prod_{C_j \ni X_i} \phi_j(x'_i, \mathbf{x}_{j,-i})}{\sum_{x''_i} \prod_{C_j \ni X_i} \phi_j(x''_i, \mathbf{x}_{j,-i})}. \end{aligned} \tag{12.23}$$

This last expression uses only the factors involving X_i , and depends only on the instantiation in \mathbf{x}_{-i} of X_i 's Markov blanket. In the case of Bayesian networks, this expression reduces to a formula involving only the CPDs of X_i and its children, and its value, again, depends only on the assignment in \mathbf{x}_{-i} to the Markov blanket of X_i .

First, we need to choose an initial feasible configuration for Gibbs sampling. One solution is to go through every random variable x_i (edge in this problem), collect the colors of its neighbors, and assign x_i the first unused color (e.g., if $K=5$, and neighbors of x_i are colored 5, 3 and 4, then x_i should be assigned color 1). The process stops whenever the neighbors of x_i have used up all the colors, which indicates K is too small; otherwise a feasible coloring is produced.

The implementation of Gibbs sampling is then straightforward: in each iteration,

1. Pick a random variable x_i (e.g. in round-robin fashion);
2. Fix x_i 's Markov blanket configuration to its current assignment, loop through every state x_i' of x_i and multiply together the values of containing factors to get the likelihood for $x_i=x_i'$;
3. Normalize the likelihood of x_i' to obtain the complete conditional distribution $p(x_i|x_{\setminus i})$;
4. Sample a new configuration of x_i from $p(x_i|x_{\setminus i})$, then move on to another r.v. and repeat.

Recall that in our problem, the distribution takes the form $p(x) = \{\prod_e \phi(x_e)\} \{\prod_c \psi(x_c)\}$, where e ranges over the edges of the original graph A and c ranges over the vertices of A that connects two or more edges. So in step 2 calculation of the conditional distribution $p(x_e | x_{\setminus e})$, the factors involving x_e are singleton potential $\phi(x_e)=\exp(w_{\{x_e\}})$, and clique coloring potentials $\psi(x_c)$ for all c containing e ($\psi(x_c)=0$ if the current assignment of x_e conflicts with the rest of x_c , 1 otherwise) .

At the end (after burnin), simply collect all the samples and calculate the frequency of the event $x_e=k$ to estimate the marginal probability $p(x_e=k)$.

Problem 2, Part 2.

Your table may look something like this (by Changbin):

	its_2_6	its_2_10	its_2_14	its_2_18
burnin_2_6	0	0.03125	0.054565	0.085323
burnin_2_10	0	0.10156	0.067322	0.090385
burnin_2_14	0	0.094727	0.083435	0.084759
burnin_2_18	0.17188	0.098633	0.095703	0.093063

The general conclusion is that for large number of burnin and sampling iterations, the probability in question (edge (a,d) colored with 4) is independent of the initial choice of assignment (and should approach something close to 0.09).

Example problem 3 solution (by Changbin):

Based on the assumption in this question, the probability distribution is:

$$p(x) = \frac{1}{Z} \prod_{(i,j) \in E} \exp(J x_i x_j)$$

log likelihood is:

$$\begin{aligned} l(J) &= \log \left(\prod_m p(x^m | J) \right) = \sum_m \log \left(\frac{1}{Z(J)} \prod_{(i,j) \in E} \exp(J x_i^m x_j^m) \right) \\ &= \sum_m -\log(Z(J)) + \sum_m \sum_{(i,j) \in E} \log(\exp(J x_i^m x_j^m)) \\ &= \sum_m \sum_{(i,j) \in E} J x_i^m x_j^m - M * \log(Z(J)) \\ &= \sum_{m=1}^5 J(x_1^m x_2^m + x_1^m x_3^m + x_2^m x_3^m) - 5 * \log(Z(J)) \\ &= 3J - J + 3J + 3J - J - 5 * \log(Z(J)) \\ &= 7J - 5 * \log(Z(J)) \end{aligned}$$

$$\begin{aligned} Z(J) &= \sum_{x_1, x_2, x_3} \prod_{(i,j) \in E} \exp(J x_i x_j) \\ &= \sum_{x_1, x_2, x_3} \exp \left(\sum_{(i,j) \in E} J x_i x_j \right) \\ &= (\exp(3J))_{1,1,1} + (\exp(-J))_{1,1,-1} + (\exp(-J))_{1,-1,1} + (\exp(-J))_{-1,1,1} \\ &\quad + (\exp(-J))_{-1,1,-1} + (\exp(-J))_{-1,-1,1} + (\exp(3J))_{-1,-1,-1} \\ &= 2e^{3J} + 6e^{-J} \end{aligned}$$

$$\text{Thus: } l(J) = 7J - 5 * \log(2e^{3J} + 6e^{-J})$$

To get the maximum, we have

$$\frac{\partial l}{\partial J} = 7 - 5 * \frac{6e^{3J} - 6e^{-J}}{2e^{3J} + 6e^{-J}} = 0$$

then

$$\begin{aligned} 14e^{3J} + 42e^{-J} - 30e^{3J} + 30e^{-J} &= 0 \\ 9e^{-J} &= 2e^{3J} \\ J &= \frac{\log(9/2)}{4} \end{aligned}$$