

# Homework 4 of Statistical Machine Learning

Wang Yikai, 2017310740

June 29, 2018

## 1 Probabilistic Graphical Models

### 1.1 Conditional Queries in a Bayesian Network

1. The probabilistic graph of the model is shown in Figure 1. There are:

$$\begin{aligned} P(G_1) &= (0.5, 0.5), \quad P(G_i|G_1) = \begin{pmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{pmatrix}, \quad (i = 2, 3) \\ p(X_i|G_i = 1) &= \mathcal{N}(X_i|\mu = 55, \sigma^2 = 10), \quad (i = 1, 2, 3) \\ p(X_i|G_i = 2) &= \mathcal{N}(X_i|\mu = 65, \sigma^2 = 10), \quad (i = 1, 2, 3) \end{aligned}$$

The **Markov blanket** of  $G_3$  contains  $G_1$  and  $X_3$ .

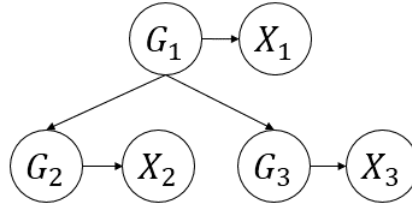


Figure 1: Probabilistic graph

2. Suppose  $X_2$  is observed to be 50. In Python, uses `scipy.stats.norm.pdf([50], 55, np.sqrt(10))`, there is  $P(X_2 = 50|G_2 = 1) = 0.03614$ . Similarly,  $P(X_2 = 50|G_2 = 2) = 1.6 \times 10^{-6}$ . The posterior belief  $P(G_1 = 1|X_2 = 50)$  is:

$$\begin{aligned} P(G_1 = 1|X_2 = 50) &= \frac{P(G_1 = 1, X_2 = 50)}{P(X_2 = 50)} \\ &\propto P(G_1 = 1, X_2 = 50) \\ &= \sum_{i=1}^2 P(G_1 = 1, G_2 = i, X_2 = 50) \\ &= \sum_{i=1}^2 P(G_1 = 1)P(G_2 = i|G_1 = 1)P(X_2 = 50|G_2 = i) \\ &= 0.5 \times 0.9 \times 0.03614 + 0.5 \times 0.1 \times 1.6 \times 10^{-6} = 0.0163 \end{aligned}$$

$$\begin{aligned}
P(G_1 = 2|X_2 = 50) &\propto P(G_1 = 2, X_2 = 50) \\
&= \sum_{i=1}^2 P(G_1 = 2)P(G_2 = i|G_1 = 2)P(X_2 = 50|G_2 = i) \\
&= 0.5 \times 0.1 \times 0.03614 + 0.5 \times 0.9 \times 1.6 \times 10^{-6} = 0.0018
\end{aligned}$$

As  $P(G_1 = 1|X_2 = 50) + P(G_1 = 2|X_2 = 50) = 1$ , there are:

$$P(G_1 = 1|X_2 = 50) = \frac{0.0163}{0.0163 + 0.0018} = 0.901$$

$$P(G_1 = 2|X_2 = 50) = \frac{0.0018}{0.0163 + 0.0018} = 0.099$$

3. Suppose we observe both  $X_2 = 50$  and  $X_3 = 50$ . Then  $P(G_1|X_2, X_3)$  is:

$$\begin{aligned}
P(G_1 = 1|X_2 = 50, X_3 = 50) &= \frac{P(G_1 = 1, X_2 = 50, X_3 = 50)}{P(X_2 = 50, X_3 = 50)} \propto P(G_1 = 1, X_2 = 50, X_3 = 50) \\
&= \sum_{i=1}^2 \sum_{j=1}^2 P(G_1 = 1)P(G_2 = i|G_1 = 1)P(X_2 = 50|G_2 = i)P(G_3 = j|G_1 = 1)P(X_3 = 50|G_3 = j) \\
&= P(G_1 = 1) \sum_{i=1}^2 [P(G_2 = i|G_1 = 1)P(X_2 = 50|G_2 = i)] \sum_{j=1}^2 [P(G_3 = j|G_1 = 1)P(X_3 = 50|G_3 = j)] \\
&= 0.5 \times (0.9 \times 0.03614 + 0.1 \times 1.6 \times 10^{-6})^2 = 5.29 \times 10^{-4} \\
P(G_1 = 2|X_2 = 50, X_3 = 50) &= 0.5 \times (0.1 \times 0.03614 + 0.9 \times 1.6 \times 10^{-6})^2 = 6.53 \times 10^{-6}
\end{aligned}$$

As  $P(G_1 = 1|X_2 = 50, X_3 = 50) + P(G_1 = 2|X_2 = 50, X_3 = 50) = 1$ , there are:

$$P(G_1 = 1|X_2 = 50, X_3 = 50) = \frac{5.29 \times 10^{-4}}{5.29 \times 10^{-4} + 6.53 \times 10^{-6}} = 0.988$$

$$P(G_1 = 2|X_2 = 50, X_3 = 50) = \frac{6.53 \times 10^{-6}}{5.29 \times 10^{-4} + 6.53 \times 10^{-6}} = 0.012$$

## 1.2 Conditional Random Fields

1. The undirected graph and the factor graph of the CRF are shown as Figure 2, 3:

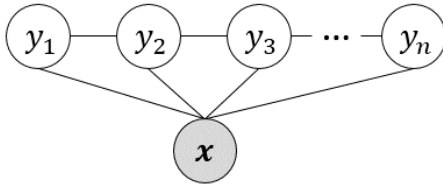


Figure 2: Undirected Graph

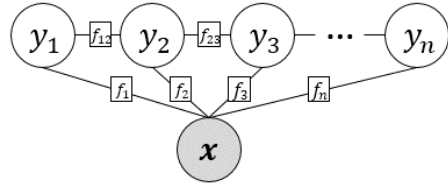


Figure 3: Factor Graph

2. The cliques of CRF is illustrated as Figure 4.

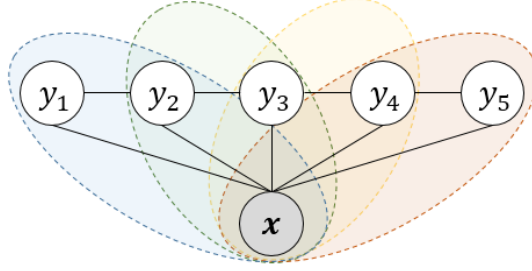


Figure 4: Cliques of CRF

The junction tree of CRF is illustrated as Figure 5.

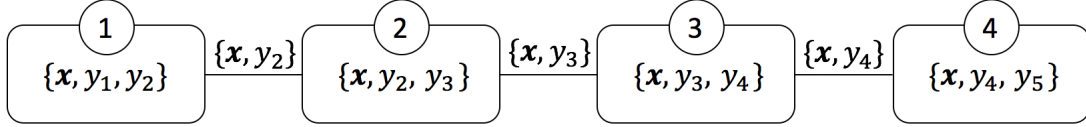


Figure 5: Junction Tree of CRF

Denote each clique in the junction tree as  $C_i, i = 1, 2, 3, 4$ . First, set initial factors at each cluster as products. For example,  $\psi_1(\mathbf{x}, y_1, y_2) = \phi(\mathbf{x}, y_1)\phi(y_1, y_2)$ .

By running the sum-product algorithm on the junction tree, we can get  $p(y_3, \mathbf{x}; \mathbf{w})$ :

- In  $C_1$ , eliminate  $y_1$  by sending  $\delta_{1 \rightarrow 2}(\mathbf{x}, y_2) = \sum_{y_1} \psi_1(\mathbf{x}, y_1, y_2)$  to  $C_2$ ;
- In  $C_2$ , eliminate  $y_2$  by sending  $\delta_{2 \rightarrow 3}(\mathbf{x}, y_3) = \sum_{y_2} \delta_{1 \rightarrow 2}(\mathbf{x}, y_2) \psi_2(\mathbf{x}, y_2, y_3)$  to  $C_3$ ;
- In  $C_4$ , eliminate  $y_5$  by sending  $\delta_{4 \rightarrow 3}(\mathbf{x}, y_4) = \sum_{y_5} \psi_4(\mathbf{x}, y_4, y_5)$  to  $C_3$ ;
- In  $C_3$ , eliminate  $y_4$  to get  $p(y_3, \mathbf{x}; \mathbf{w}) = \sum_{y_4} \delta_{2 \rightarrow 3}(\mathbf{x}, y_3) \delta_{4 \rightarrow 3}(\mathbf{x}, y_4) \psi_3(\mathbf{x}, y_3, y_4)$ .

After getting  $p(y_3, \mathbf{x}; \mathbf{w})$ , there is  $p(\mathbf{x}; \mathbf{w}) = \sum_{y_3} p(y_3, \mathbf{x}; \mathbf{w})$ . Then by  $p(y_3 | \mathbf{x}; \mathbf{w}) = \frac{p(y_3, \mathbf{x}; \mathbf{w})}{p(\mathbf{x}; \mathbf{w})}$ , we can get  $p(y_3 | \mathbf{x}; \mathbf{w})$ .

3. Given the following parametric form of CRF:

$$p(\mathbf{y} | \mathbf{x}; \mathbf{w}) = \frac{1}{Z(\mathbf{x}; \mathbf{w})} \exp \left\{ \mathbf{w}^T \sum_{i=2}^n \mathbf{f}(\mathbf{x}, y_i, y_{i-1}) \right\}$$

where

$$Z(\mathbf{x}; \mathbf{w}) = \sum_{\mathbf{y}' \in \mathcal{L}^n} \exp \left\{ \mathbf{w}^T \sum_{i=2}^n \mathbf{f}(\mathbf{x}, y'_i, y'_{i-1}) \right\}$$

To learn the parameters of the CRF, we maximize the conditional log likelihood of training data  $\mathcal{D} = \{(\mathbf{x}^i, \mathbf{y}^i)\}_{i=1}^N$  over parameters  $\mathbf{w}$  using gradient descent:

$$\mathcal{L}(\mathbf{w}) = \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \log p(\mathbf{y}|\mathbf{x}; \mathbf{w}) = \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \left\{ \mathbf{w}^T \sum_{i=2}^n \mathbf{f}(\mathbf{x}, y_i, y_{i-1}) - \log Z(\mathbf{x}; \mathbf{w}) \right\}$$

Compute its gradient:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \left\{ \sum_{i=2}^n \mathbf{f}(\mathbf{x}, y_i, y_{i-1}) - \frac{1}{Z(\mathbf{x}; \mathbf{w})} \frac{\partial Z(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}} \right\}$$

where  $\frac{1}{Z(\mathbf{x}; \mathbf{w})} \frac{\partial Z(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}}$  is:

$$\begin{aligned} \frac{1}{Z(\mathbf{x}; \mathbf{w})} \frac{\partial Z(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}} &= \frac{1}{Z(\mathbf{x}; \mathbf{w})} \frac{\partial \sum_{\mathbf{y}' \in \mathcal{L}^n} \exp \left\{ \mathbf{w}^T \sum_{i=2}^n \mathbf{f}(\mathbf{x}, y'_i, y'_{i-1}) \right\}}{\partial \mathbf{w}} \\ &= \frac{1}{Z(\mathbf{x}; \mathbf{w})} \sum_{\mathbf{y}' \in \mathcal{L}^n} \exp \left\{ \mathbf{w}^T \sum_{i=2}^n \mathbf{f}(\mathbf{x}, y'_i, y'_{i-1}) \right\} \sum_{i=2}^n \mathbf{f}(\mathbf{x}, y'_i, y'_{i-1}) \\ &= \sum_{\mathbf{y}' \in \mathcal{L}^n} \frac{1}{Z(\mathbf{x}; \mathbf{w})} \exp \left\{ \mathbf{w}^T \sum_{i=2}^n \mathbf{f}(\mathbf{x}, y'_i, y'_{i-1}) \right\} \sum_{i=2}^n \mathbf{f}(\mathbf{x}, y'_i, y'_{i-1}) \\ &= \sum_{\mathbf{y}' \in \mathcal{L}^n} p(\mathbf{y}'|\mathbf{x}; \mathbf{w}) \sum_{i=2}^n \mathbf{f}(\mathbf{x}, y'_i, y'_{i-1}) \\ &= \mathbb{E}_{p(\mathbf{y}'|\mathbf{x}; \mathbf{w})} \left[ \sum_{i=2}^n \mathbf{f}(\mathbf{x}, y'_i, y'_{i-1}) \right] \end{aligned}$$

Therefore  $\frac{\partial \mathcal{L}}{\partial \mathbf{w}}$  is:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{w}} &= \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \left\{ \sum_{i=2}^n \mathbf{f}(\mathbf{x}, y_i, y_{i-1}) - \mathbb{E}_{p(\mathbf{y}'|\mathbf{x}; \mathbf{w})} \left[ \sum_{i=2}^n \mathbf{f}(\mathbf{x}, y'_i, y'_{i-1}) \right] \right\} \\ &= \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \sum_{i=2}^n \left\{ \mathbf{f}(\mathbf{x}, y_i, y_{i-1}) - \mathbb{E}_{p(\mathbf{y}'|\mathbf{x}; \mathbf{w})} [\mathbf{f}(\mathbf{x}, y'_i, y'_{i-1})] \right\} \end{aligned}$$

4. Choose  $y_n$  as the root node, and apply forward propagation( $O(n)$ ) and backward propagation( $O(n)$ ). Therefore the whole process of belief propagation is  $O(n)$ .

Forward propagation:

$$\begin{aligned} \alpha_0(y_0|\mathbf{x}) &= 1(y_0 = \text{start}) \\ \alpha_i(y_i|\mathbf{x}) &= \delta_i^T(\mathbf{X}, y_i, y_{i-1}) \alpha_{i-1}(y_{i-1}|\mathbf{x}), \quad i = 1, \dots, n+1 \end{aligned}$$

where

$$\delta_i(\mathbf{X}, y_i, y_{i-1}) = \exp \left\{ \mathbf{w}^T \mathbf{f}(\mathbf{x}, y_i, y_{i-1}) \right\}$$

Backward propagation:

$$\begin{aligned}\beta_{n+1}(y_{n+1}|\mathbf{x}) &= 1(y_{n+1} = \text{stop}) \\ \beta_i(y_i|\mathbf{x}) &= \beta_{i+1}(y_{i+1}|\mathbf{x})\delta_{i+1}^T(\mathbf{x}, y_{i+1}, y_i), \quad i = 1, \dots, n+1\end{aligned}$$

Here  $y_i$  has three probable values, which means  $\alpha_i$  and  $\beta_i$  are both vectors with 3 dimensions and  $\delta_i$  is a  $3 \times 3$  matrix.

According to the definition:

$$\begin{aligned}Z(\mathbf{x}; \mathbf{w}) &= \alpha_n^T(\mathbf{x}) \cdot \mathbf{1} = \mathbf{1}^T \cdot \beta_1(\mathbf{x}) \\ p(y_i|\mathbf{x}) &= \frac{\alpha_i^T(y_i|\mathbf{x})\beta_i(y_i|\mathbf{x})}{Z(\mathbf{x}; \mathbf{w})} \\ p(y_i, y_{i-1}|\mathbf{x}) &= \frac{\alpha_{i-1}^T(y_i|\mathbf{x})\delta_i(\mathbf{x}, y_i, y_{i-1})\beta_i(y_i|\mathbf{x})}{Z(\mathbf{x}; \mathbf{w})}\end{aligned}$$

Therefore the expectation is:

$$\sum_{i=2} \mathbb{E}_{p(\mathbf{y}'|\mathbf{x}, \mathbf{w})}[f(\mathbf{x}, y'_i, y'_{i-1})] = \sum_{i=2}^n \sum_{y'_i, y'_{i-1}} p(y'_i, y'_{i-1}|\mathbf{x}) f(\mathbf{x}, y'_i, y'_{i-1})$$

## 2 Deep Generative Models: Class-conditioned VAE

1. The variational lower bound of Class-conditioned VAE is:

$$\begin{aligned}\log p_\theta(x|y) &\geq \log p_\theta(x|y) - D_{KL}[q_\phi(z|x, y) \| p_\theta(z|x, y)] \\ &= \log p_\theta(x|y) - \mathbb{E}_{z \sim q_\phi(z|x, y)} \left[ \log \frac{q_\phi(z|x, y)}{p_\theta(z|x, y)} \right] \\ &= \log p_\theta(x|y) - \mathbb{E}_{z \sim q_\phi(z|x, y)} [\log q_\phi(z|x, y) - \log p_\theta(x, z|y) + \log p_\theta(x|y)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x, y)} [\log p_\theta(x, z|y) - \log q_\phi(z|x, y)] \\ &= \mathcal{L}(\theta, \phi)\end{aligned}$$

Design the variational posterior as:

$$q_\phi(z|x, y) = N(z | \mu_z(x, y; \phi), \sigma_z^2(x, y; \phi))$$

2. The code is submitted, simply run main.py is ok.
3. Results(in the next page)

The generated result is illustrated in Figure 6.

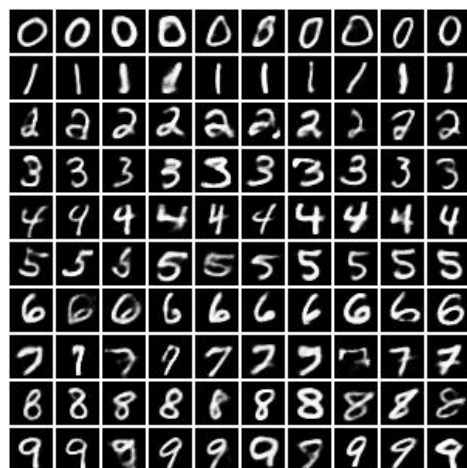


Figure 6: Generated Result

The variational lower bound of the training and testing process is illustrated in Figure 7.

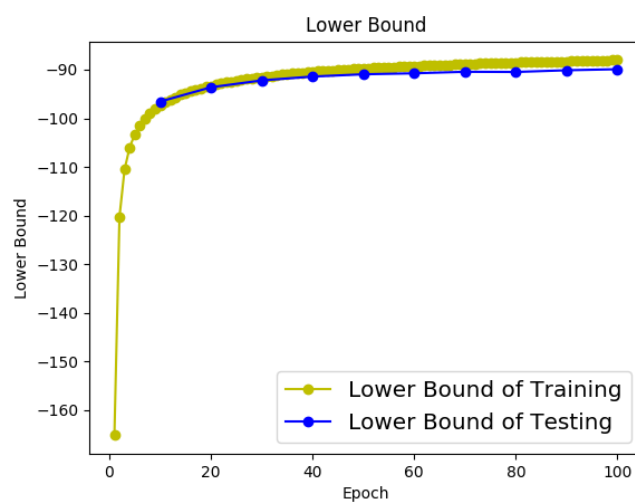


Figure 7: Lower Bound