# Assignment 4: Graphical Model and Unsupervised Learning

AIGS/CSED515 Machine Learning

Instructor: Jungseul Ok

jungseul@postech.ac.kr

Due: 20:00pm Dec 16, 2020

**Remarks**

- Group study and open discussion via LMS board are encouraged, however, assignment that your hand-in must be **of your own work**, and **hand-written** unless you're asked a coding task.

- Submit a scanned copy of your answer on LMS online in **a single PDF file**, to which you also print and add your code if apply, i.e., no zip file, just a single PDF containing everything.

- Delayed submission may get some penalty in score: 5% off for delay of $0 \sim 4$ hours; 20% off for delay of $4 \sim 24$ hours; and delay longer than 24 hours will not be accepted.

1. [6 pt] (An application of belief propagation) Consider an integer programming (IP) of $x_1, ..., x_5$ with linear objective and constraints in the followings:

$$\underset{x_1,...,x_5 \in \{0,1\}}{\text{maximize}} \quad x_1 + 2x_2 + 3x_3 + 2x_4 + 2x_5$$

$$\text{subject to} \quad x_1 + x_2 + x_3 \leq 1$$
$$x_3 + x_4 \leq 1$$
$$x_4 + x_5 \leq 1$$

In order to solve the IP, we can formulate a maximum a posterior (MAP) problem of the joint probability of $x_1, ..., x_5$ in the following factorized form:

$$p(x_1, ..., x_5) = \frac{1}{Z} f_a(x_1) f_b(x_2) f_c(x_3) f_d(x_4) f_e(x_5) f_A(x_1, x_2, x_3) f_B(x_3, x_4) f_C(x_4, x_5) ,$$

where $Z$ is the normalization constant, $f_a(x_1) = \exp(x_1)$, $f_b(x_2) = \exp(2x_2)$, $f_c(x_3) = \exp(3x_3)$, $f_d(x_4) = \exp(2x_4)$, $f_e(x_5) = \exp(2x_5)$,

$$f_A(x_1, x_2, x_3) = \begin{cases} 1 & \text{if } x_1 + x_2 + x_3 \leq 1 \\ 0 & \text{otherwise} \end{cases} ,$$

$$f_B(x_3, x_4) = \begin{cases} 1 & \text{if } x_3 + x_4 \leq 1 \\ 0 & \text{otherwise} \end{cases} ,$$

$$f_C(x_4, x_5) = \begin{cases} 1 & \text{if } x_4 + x_5 \leq 1 \\ 0 & \text{otherwise} \end{cases} .$$

Note that only configuration of $x_1, ..., x_5$ verifies all the constraints in the IP has non-zero probability, which is proportional to the exponential of the objective value of the IP. Hence, the MAP configuration is a solution to the integer programming.

(a) [2pt] <u>Draw</u> the factor graph corresponding to the joint probability $p(x_1, ..., x_5)$.

(b) [4pt] <u>Solve</u> the IP using the max-product belief propagation algorithm. <u>What</u> is the optimal value?

2. [11pt] (Graph learning) Consider 4 binary random variables $X_1, X_2, X_3, X_4 \in \{0, 1\}$. Assume we have 100 observations and the following table shows the number of counts of observations. We want to learn the structure of random variables $X_i$'s from our observations using the Chow-Liu tree algorithm.

| $X_1$ | $X_2$ | $X_3$ | $X_4$ | Count | $X_1$ | $X_2$ | $X_3$ | $X_4$ | Count |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 2 | 1 | 0 | 0 | 0 | 7 |
| 0 | 0 | 0 | 1 | 5 | 1 | 0 | 0 | 1 | 5 |
| 0 | 0 | 1 | 0 | 2 | 1 | 0 | 1 | 0 | 7 |
| 0 | 0 | 1 | 1 | 5 | 1 | 0 | 1 | 1 | 12 |
| 0 | 1 | 0 | 0 | 2 | 1 | 1 | 0 | 0 | 10 |
| 0 | 1 | 0 | 1 | 4 | 1 | 1 | 0 | 1 | 10 |
| 0 | 1 | 1 | 0 | 6 | 1 | 1 | 1 | 0 | 10 |
| 0 | 1 | 1 | 1 | 8 | 1 | 1 | 1 | 1 | 5 |

(a) [4pt] Compute the marginal probability $p(X_i)$ for each $i \in \{1, 2, 3, 4\}$ and $p(X_i, X_j)$ for all $i \neq j \in \{1, 2, 3, 4\}$.

| $x$ | 0 | 1 |
|---|---|---|
| $p(X_1 = x)$ | | |
| $p(X_2 = x)$ | | |
| $p(X_3 = x)$ | | |
| $p(X_4 = x)$ | | |

| $(x, y)$ | $(0, 0)$ | $(0, 1)$ | $(1, 0)$ | $(1, 1)$ |
|---|---|---|---|---|
| $p(X_1 = x, X_2 = y)$ | | | | |
| $p(X_1 = x, X_3 = y)$ | | | | |
| $p(X_1 = x, X_4 = y)$ | | | | |
| $p(X_2 = x, X_3 = y)$ | | | | |
| $p(X_2 = x, X_4 = y)$ | | | | |
| $p(X_3 = x, X_4 = y)$ | | | | |

(b) [2pt] <u>Compute</u> the mutual information $I(X_i, X_j)$ for all $i \neq j \in \{1, 2, 3, 4\}$.

| | |
|---|---|
| $I(X_1, X_2)$ | |
| $I(X_1, X_3)$ | |
| $I(X_1, X_4)$ | |
| $I(X_2, X_3)$ | |
| $I(X_2, X_4)$ | |
| $I(X_3, X_4)$ | |

(c) [2pt] Performing Chow-Liu algorithm, <u>draw</u> a Bayesian network $T_1$ rooted from@$X_1$.

(d) [2pt] Performing Chow-Liu algorithm, <u>draw</u> a Bayesian network $T_3$ rooted from $X_3$.

(e) [1pt] Let $p_T$ denote the probability corresponding to Bayesian network $T$. <u>Compute</u> the difference $\mathrm{KL}(p\|p_{T_1}) - \mathrm{KL}(p\|p_{T_3})$ where Bayesian networks $T_1$ and $T_3$ are obtained in Problems 2c and 2d, resp.

3. [12pt] (K-means) Given a dataset $\mathcal{D} = \{x^{(i)}\}_{i \in [N]}$ of $N$ data points in $\mathbb{R}^2$, we want to partitioning them into $K$ clusters using $K$-means algorithm. Let $\mu_k \in \mathbb{R}^2$ denote the center of cluster $k \in [K]$. Then, the $K$-means algorithm aims at optimizing:

$$\min_{\{r_{ik}\},\{\mu_k\}} \sum_{i \in [N]} \sum_{k \in [K]} \frac{1}{2} r_{ik} \|x^{(i)} - \mu_k\|_2^2 \qquad (1a)$$

$$\text{s.t.} \quad r_{ik} \in \{0, 1\} \quad \forall i \in [N], \forall k \in [K] \quad \text{and;} \qquad (1b)$$

$$\sum_{k \in [K]} r_{ik} = 1 \quad \forall i \in [N]. \qquad (1c)$$

(a) [2pt] Given fixed cluster centers $\{\mu_k\}_{k \in [K]}$, obtain the optimal $r_{ik}$ for (1). Justify your solution.

(b) [2pt] Given fixed $\{r_{ik}\}_{i \in [N], k \in [K]}$, verifying (1b) and (1c), obtain the optimal cluster center $\mu_k$ for (1). Justify your solution.

(c) [4pt] We want to check the convergence of $K$-means algorithm which alternates Problems (3a) and (3b). Describe the algorithm. Let $L_t$ be the loss (1a) after $t$-th iteration of $K$-means algorithm. Check if $L_t$ is monotonically increasing in $t$. Using the following theorem (a part of monotone convergence theorem), check the convergence of $K$-means algorithm in terms of loss function. Can we guarantee that $K$-means algorithm converges to the global optimality?

**Theorem 1.** If $(a_t)_{t \in \mathbb{N}}$ is a _monotone_ sequence of real numbers, i.e., if $a_t \leq a_{t+1}$ for every $t \geq 1$, or $a_t \geq a_{t+1}$ for every $t \geq 1$, then this sequence has a finite limit if and only if the sequence is _bounded_.

(d) [4pt] Complete Kmeans.py which performing $K$-means algorithm aforementioned. For the given dataset, after how many updates does the algorithm converge? What cost function value does it converge to? What are the obtained centers?

5

4. [20pt] (Generative Adversarial Newtorks) Consider the following max-min problem for a dataset $\mathcal{D}$ consisting of $x$'s:

$$\max_{\theta} \min_{w} -\sum_{x \in \mathcal{D}} \log p_w(y = 1 \mid x) - \sum_{z \in \mathcal{Z}} \log(1 - p_w(y = 1 \mid G_\theta(z))) + \frac{C}{2}\|w\|_2^2 . \quad (2)$$

Here the generator $G_\theta(z)$ parameterized by $\theta$ transforms noise $z \in \mathcal{Z}$ into artificial data. The discriminator $p_w(y \mid x)$ parameterized by $w$ checks if $x$ is artificial or not, where $y = 1$ indicates that $x$ is real, and $y = 0$ indicates that $x$ is artificial. The hyper-parameter $C \geq 0$ controls impact of regularization. Note that solving (2) is challenging mainly due to the objective is neither convex in $w$ nor concave in $\theta$ in general. We will check if the cost function is is convex in $w$ for specific choice of the discriminator model. To do so, we use several facts:

**Fact1.** A function $f(w)$ is convex in $w$ if Hessian[1] of $f(w)$ is positive semi-definite[2].

**Fact2.** A sum of convex functions is also convex.

(a) [2pt] Suppose that we model the discriminator as follows:

$$p_w(y = 1 \mid x) = \frac{1}{1 + \exp(w^\top x)} .$$

Using this, <u>write down</u> the resulting cost function for (2).

(b) [2pt] <u>Obtain</u> Hessian of (A) $= \frac{C}{2}\|w\|_2^2 - w^\top b$ in $w$. <u>Check</u> if (A) is convex, and justify your answer.

(c) [2pt] <u>Obtain</u> Hessian of (B) $= \log(1 + \exp(w^\top b))$ in $w$. <u>Check</u> if (B) is convex, and justify your answer.

(d) [2pt] <u>Check</u> if the cost function obtained in Problem 4a is convex, and justify your answer.

(e) [2pt] Introducing auxiliary variables $\xi_x = w^\top x$ and $\xi_z = w^\top G_\theta(z)$, consider the following optimization (for the discriminator):

$$\min_{w} \quad \sum_{x \in \mathcal{D}} \log(1 + \exp \xi_x) + \sum_{z \in \mathcal{Z}} \log(1 + \exp(\xi_z)) - \sum_{z \in \mathcal{Z}} w^\top G_\theta(z) + \frac{C}{2}\|w\|_2^2 \quad (3a)$$

$$\text{s.t.} \quad \xi_x = w^\top x \quad \forall x \in \mathcal{D} \quad (3b)$$

$$\xi_z = w^\top G_\theta(z) \quad \forall z \in \mathcal{D} \quad (3c)$$

<u>Write</u> the Lagrangian for this optimization, where $\lambda_x$ and $\lambda_z$ are Lagrange multipliers corresponding to (3b) and (3c), resp.

---

[1]https://en.wikipedia.org/wiki/Hessian_matrix
[2]https://en.wikipedia.org/wiki/Definite_symmetric_matrix

(f) [2pt] <u>Obtain</u> the value of

$$\min_{w} \frac{C}{2}\|w\|_2^2 - w^\top b$$

in terms of $b$ and $C \geq 0$.

(g) [2pt] <u>Obtain</u> the value of

$$\min_{\xi} \quad \lambda\xi + \log(1 + \exp \xi)$$

in terms of $\lambda$ assuming $-1 \leq \lambda \leq 0$.

(h) [4pt] Combining Problems 4e, 4f, and 4g and using $H(a) = a\log(-a) - (1 + a)\log(1 + a)$, <u>obtain</u> dual function $g(\lambda)$ for (3). For training the discriminator, we can replace the original minimization over $w$ described in (2) with the dual maximization over valid values of $\lambda$. Using this, <u>write down</u> an alternative of GAN training in (2), in which we have a max-max problem instead of the max-min problem. Note that such an alternative training in max-max form can help to bypass challenges from fining a saddle-point, i.e., solving the max-min problem.

(i) [2pt] <u>Complete</u> `GAN.py`, which is an implementation of the alternative training of GAN obtained in Problem 4h with the $\log D$ trick in the lecture. (Hint: use `target1` and `target2`)