

Problem 4 (20pt): [Bayesian Network] Suppose we are given 5 random variables, A_1, A_2, B_1, B_2, B_3 . These random variables have dependence relations as follows: B_1 depends on A_1 , A_2 depends on B_1 . B_2 depends on A_2 and A_1 . B_3 depends on A_2 . All 5 random variables are binary, i.e., $A_i, B_j \in \{0, 1\}, i = 1, 2; j = 1, 2, 3$.

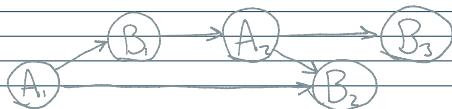
(1) [5pt] Draw the corresponding bayesian network.

(2) [5pt] Based on the bayesian network in (1), write down the joint distribution $p(A_1, A_2, B_1, B_2, B_3)$.

(3) [5pt] How many independent parameters are needed to fully specify the joint distribution in (2).

(4) [5pt] Suppose we do not have any independence assumption, write down one possible factorization of $p(A_1, A_2, B_1, B_2, B_3)$, and state how many independent parameters are required to describe joint distribution in this case.

1) Draw the bayesian network.



2) Write Joint Distribution $p(A_1, A_2, B_1, B_2, B_3)$

$$\begin{aligned} p(A_1, A_2, B_1, B_2, B_3) &= p(B_1 | A_1, A_2, B_2, B_3) \cdot p(A_1, A_2, B_2, B_3) \\ &= p(B_1 | A_1, A_2, B_2) \cdot p(B_1 | A_1, A_2, B_3) \cdot p(A_1, A_2, B_2) \\ &= p(B_1 | A_1, A_2, B_2) \cdot p(B_1 | A_1, A_2, B_3) \cdot p(B_2 | A_1, A_2) \cdot p(A_1, A_2) \\ &\quad + p(B_1 | A_1, A_2, B_2) \cdot p(B_1 | A_1, A_2, B_3) \cdot p(B_3 | A_1, A_2) \cdot p(A_1, A_2) \end{aligned}$$

3) # of independent parameters needed to specify the joint distribution?

$$\begin{aligned} p(B_1 | A_1, A_2, B_2, B_3) &\rightarrow 2^4 = 16 \text{ parameters} \\ p(B_1 | A_1, A_2, B_2) &\rightarrow 2^3 = 8 \text{ parameters} \\ p(B_1 | A_1, A_2) &\rightarrow 2^2 = 4 \text{ parameters} \\ p(A_2 | A_1) &\rightarrow 2^1 = 2 \text{ parameters} \\ p(A_1) &\rightarrow 2^0 = 1 \text{ parameter} \end{aligned}$$

32 total parameters needed

4) One possible factorization + state # of independent parameters needed.

$$p(A_1, A_2, B_1, B_2, B_3) = p(B_1 | A_2) \cdot p(B_2 | A_1, A_2) \cdot p(A_2 | B_1) \cdot p(B_1 | A_1) \cdot p(A_1)$$

$$\begin{aligned} p(B_1 | A_2) &\rightarrow 2^1 = 2 \text{ parameters} \\ p(B_2 | A_1, A_2) &\rightarrow 2^2 = 4 \text{ parameters} \\ p(A_2 | B_1) &\rightarrow 2^1 = 2 \text{ parameters} \\ p(B_1 | A_1) &\rightarrow 2^1 = 2 \text{ parameters} \\ p(A_1) &\rightarrow 2^0 = 1 \text{ parameter} \end{aligned}$$

11 total parameters needed

Problem 2 (15pt): [K-means and gradient descent] Recall the loss function for k-means clustering with k clusters, sample points x_1, x_2, \dots, x_n , and centers $\mu_1, \mu_2, \dots, \mu_k$:

$$L = \sum_{j=1}^k \sum_{x_i \in S_j} \|x_i - \mu_j\|^2$$

, where S_j refers to the set of data points that are closer to μ_j than to any other cluster mean.

(1) [5pt] Instead of updating μ_j by computing the mean, let's minimize L with batch gradient descent while holding the sets S_j fixed. Derive the update formula for μ_1 with learning rate ϵ .

(2) [2pt] Derive the update formula for μ_1 with stochastic gradient descent on a single sample point x_i . Use learning rate ϵ .

(3) [8pt] In this part, we will connect the batch gradient descent update equation with the standard k-means algorithm. Recall that in the update step of the standard algorithm, we assign each cluster center to be the mean of the data points closest to that center. It turns out that a particular choice of the learning rate ϵ (which may be different for each cluster) makes the two algorithms (batch gradient descent and the standard k-means algorithm) have identical update steps. Let's focus on the update for the first cluster, with center μ_1 . Calculate the value of ϵ so that both algorithms perform the same update for μ_1 .

1) Minimize L w/ batch gradient descent, holding sets S_j fixed.
Derive update formula for μ_1 w/ learning rate ϵ .

$$\begin{aligned} L &= \sum_{i=1}^n \sum_{x_i \in S_1} \|x_i - \mu_1\|^2 \\ &\rightarrow \text{for } \mu_1 : L = \sum_{x_i \in S_1} \|x_i - \mu_1\|^2 \end{aligned}$$

We know L is a quadratic function of μ_1 where $x_i \in S_1$ are constants.
For gradient descent, we need to calculate the gradient of L .

$$\begin{aligned} \frac{\partial L}{\partial \mu_1} &= 2 \sum_{i=1}^n (x_i - \mu_1) \cdot \frac{\partial}{\partial \mu_1} (x_i - \mu_1) = -2 \sum_{i=1}^n (x_i - \mu_1) \\ &\quad \star \text{ because current is negative, gradient is negative} \\ &\Rightarrow \text{Update Formula: } \mu_1 = \mu_1 + \epsilon / 2 \sum_{i=1}^n (x_i - \mu_1) \\ &\quad \text{where } \epsilon \text{ is a chosen learning rate (e.g. 0.001) + N} \\ &\quad \text{is the total number of points assigned to cluster 1.} \end{aligned}$$

2) Derive update formula for μ_1 w/ stochastic gradient descent.

Similar to the above, we need to calculate the gradient of L . However, since we are now using SGD, only need to use one sample point x_i per update. Thus, we can remove the summation.

$$\begin{aligned} \frac{\partial L}{\partial \mu_1} &= -2(x_i - \mu_1) \quad x_i \text{ is a single point in } S_1 \text{ where } i=1 \\ &\Rightarrow \text{Update Formula: } \mu_1 = \mu_1 + \epsilon / 2(x_i - \mu_1) \end{aligned}$$

3) Calculate the value of ϵ so that both algorithms perform the same update for μ_1 .

Recall that for K-means clustering, the formula to calculate the mean (aka centroid) for each cluster of points is:

$$\mu_k = \frac{\sum_{x_i \in S_k} x_i}{\sum_{x_i \in S_k}}, \text{ where } x_i \text{ is a binary indicator that } = 1 \text{ if the point } x_i \text{ is assigned to cluster } k + 0 \text{ otherwise.} \\ \Rightarrow \sum_{x_i \in S_k} x_i \text{ is just the number of points in cluster } k, N$$

Since we are considering 1 cluster $k=1$:

$$\mu_1 = \frac{\sum_{x_i \in S_1} x_i}{N}$$

We can then just set this equal to the update formula we got in 1) + solve for ϵ .

Problem 3 (10pt): [Latent variable model and GMM] Consider the discrete latent variable model where the latent variable z use 1-of-K representation. The distribution for latent variable z is defined as:

$$p(z_k = 1) = \pi_k$$

where $\{\pi_k\}$ satisfy: $0 \leq \pi_k \leq 1$ and $\sum_{k=1}^K \pi_k = 1$. Suppose the conditional distribution of observation x given particular value for z is Gaussian:

$$p(x|z_k = 1) = \mathcal{N}(x|\mu_k, \Sigma_k)$$

(1) [2pt] Write down the compact form of $p(z)$ and $p(x|z)$.

(2) [3pt] Show that the marginal distribution $p(x)$ has the following form:

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k)$$

(3) [5pt] If we want to find the MLE solution for parameters π_k, μ_k, Σ_k in such model, what algorithm should we use? Briefly describe its major difference compared to K-means algorithm.

1) Write compact form of $p(z) + p(x|z)$.

The prior, $p(z)$, where z is a K -dimensional latent variable having 1-of-K representation (one element $z_u = 1$, rest = 0) and $\sum z_u = 1$, can be written as $p(z_u = 1) = \pi_u$
 $\star \sum_{u=1}^K \pi_u = 1$

We get this from $p(z) = \prod_{u=1}^K \pi_u^{z_u}$

Since every $z_u = 0$ except for one, the above will simplify:
e.g. $u=3 \rightarrow z = (0, 0, 1) \Rightarrow z_u = 3$

$$\Rightarrow \prod_{u=1}^3 \pi_u^{z_u} = \pi_1^0 \cdot \pi_2^0 \cdot \pi_3^1 = \pi_3$$

$$\Rightarrow p(z) = \prod_{u=1}^K \pi_u^{z_u} \star \text{The } \pi_u = \text{mixing coefficient}$$

The conditional likelihood, $p(x|z)$, where x is the input data, can be written as $p(x|z_u = 1) = \mathcal{N}(x|\mu_u, \Sigma_u)$, a Gaussian distribution with mean μ_u + covariance Σ_u .

Using the same logic used to define $p(z)$ above since there can only be one $z_u = 1$:

$$p(x|z) = \prod_{u=1}^K N(x|\mu_u, \Sigma_u)^{z_u}$$

2) Show that marginal distribution has the form: $p(x) = \sum_{k=1}^K \pi_k N(x|\mu_k, \Sigma_k)$

The marginal distribution of x , $p(x)$, can be obtained by summing the joint distribution over all possible values of z .

We know from the product (aka chain) rule of probabilities:

$$p(x, z) = p(z) \cdot p(x|z)$$

Both $p(z)$ + $p(x|z)$ were obtained in 1). We can use these with marginalization (sum rule) to get the marginal distribution, $p(x)$, by summing up the terms on z .

$$p(x) = \sum_{k=1}^K p(z_k) p(x|z_k) = \sum_{k=1}^K \pi_k N(x|\mu_k, \Sigma_k)$$

3) IF MLE solution for π_k, μ_k, Σ_k wanted, which algorithm?
Describe major difference compared to K-means algorithm.

$$\mu_1 = \frac{\sum x_n}{N}$$

We can then just set this equal to the update formula we get in 1) + solve for ϵ .

$$\mu_1 + \epsilon [2\sum_{n=1}^N (x_n - \mu_1)] = \frac{\sum_{n=1}^N x_n}{N}$$

$$\rightarrow \epsilon [2\sum_{n=1}^N (x_n - \mu_1)] = \frac{\sum_{n=1}^N x_n}{N} - \mu_1$$

$$\rightarrow \epsilon = \left[\frac{\sum_{n=1}^N x_n}{N} - \mu_1 \right] / 2\sum_{n=1}^N (x_n - \mu_1)$$

$$\rightarrow \epsilon = \frac{1}{N} [\sum_{n=1}^N x_n - N\mu_1] / 2\sum_{n=1}^N (x_n - \mu_1)$$

$$\rightarrow \epsilon = \frac{\sum_{n=1}^N x_n - N\mu_1}{2N\sum_{n=1}^N (x_n - \mu_1)} \Rightarrow \boxed{\epsilon = \frac{1}{2N}}$$

- 3) If MLE solution for $T_{ik}, \mu_k, + \Sigma_k$ wanted, which algorithm? Describe major difference compared to K-means algorithm.

Since z is a latent (hidden) variable, Expectation Maximization, EM, is needed to find $z +$ the parameters $T_{ik}, \mu_k, + \Sigma_k$

The EM algorithm consists of 2 steps, the E step + M step. First, $T_{ik}, \mu_k, + \Sigma_k$ can be randomly initialized + the initial value of the log likelihood, $\ln p(x | T, \mu, \Sigma) = \sum_{n=1}^N \ln (\sum_{k=1}^K T_{nk} N(x_n | \mu_k, \Sigma_k))$, is computed.

E (Expectation) Step

Algorithm tries to infer value of z using responsibility $r(z_{nk})$, aka the posterior $p(z|x)$, which is gotten from Bayes's rule.

To simplify things, the expectation $E(z_{nk}) = \frac{\pi_k N(x | \mu_k, \Sigma_k)}{\sum_{i=1}^K \pi_i N(x | \mu_i, \Sigma_i)}$, the expected value of the indicator variable Z_{nk} under posterior $p(z|x, T, \mu, \Sigma)$, can be calculated.

M (Maximization) Step

Algorithm updates parameter $T_{ik}, \mu_k, + \Sigma_k$ based on the guess of z .

This is done by maximizing $E_z[\ln p(x, z | T, \mu, \Sigma)]$:

$$T_{ik} = \frac{n_{ik}}{N}$$

$$\mu_k = \frac{1}{n_k} \sum_{n=1}^N r(z_{nk}) x_n$$

$$\Sigma_k = \frac{1}{n_k} \sum_{n=1}^N r(z_{nk}) (x_n - \mu_k)(x_n - \mu_k)^T$$

The E + M steps repeat until convergence, such as when the parameters stop changing or the log likelihood stops changing

The difference between K-means + this algorithm is that K-means uses hard assignment, where each data point is only assigned to 1 cluster. Alternatively, this algorithm uses soft assignment, meaning posterior probabilities are used to find how much a data point is associated with a specific cluster. This adds a measure of uncertainty. K-means can be thought of as a special case of this algorithm that uses the same + infinitely small variance for each Gaussian. K-means may not be as accurate but it converges faster.