

HW4 Handwritten Assignment

Lecturer: Pei-Yuan Wu

TAs: Chun-Lin Huang(Problem 1, 2), Yuan-Chia Chang(Problem 3, 4, 5)

November 2023, Sixth Edition

Problem 1 (EM algorithm for mixture of Bernoulli model)(1.5%)

Consider the generative model parameterized by $\theta = (\pi_k, \mu_k)_{k=1}^K$, where $\pi_1, \dots, \pi_K \in [0, 1]$ satisfies $\sum_{k=1}^K \pi_k = 1$, and that $\mu_1, \dots, \mu_K \in [0, 1]^D$, so that the probability of generating a D -dimensional binary vector $\mathbf{x} = (x^{(1)}, \dots, x^{(D)}) \in \{0, 1\}^D$ is

$$p(\mathbf{x}; \theta) = \sum_{k=1}^K \pi_k \prod_{j=1}^D \mu_{kj}^{x^{(j)}} (1 - \mu_{kj})^{1-x^{(j)}}$$

In other words, with given μ_k , the elements $x^{(1)}, \dots, x^{(D)}$ are independent, where $x^{(j)}$ follows Bernoulli distribution of mean μ_{kj} . Suppose we observe training data of N binary vectors $\mathbf{x}_1, \dots, \mathbf{x}_N \in \{0, 1\}^D$, derive the E-step and M-step equations of the EM algorithm for optimizing the mixing coefficients π_k and the Bernoulli means μ_{kj} by maximum likelihood.

Problem 2 (EM algorithm for mixture of exponential model)(1.5%)

Given N samples $x_1, \dots, x_N \in [0, \infty)$, we would like to cluster them into K clusters. Assume the samples are generated according to Exponential mixture models

$$X \sim \sum_{j=1}^K \pi_j \text{Exp}(\tau_j)$$

where $\pi_1 + \dots + \pi_K = 1$, and $\text{Exp}(\tau)$ denotes the exponential distribution with probability density function

$$f_\tau(x) = \begin{cases} (1/\tau)e^{-x/\tau} & , x \geq 0 \\ 0 & , x < 0 \end{cases}$$

We would like to apply Expectation Maximization algorithm to find the maximum likelihood estimation of parameters $\theta = \{(\pi_k, \tau_k)\}_{k=1}^K$.

- (a) Please write down the E-step and M-step and show that the parameters are updated from $\theta^{(t)} = \{(\pi_k^{(t)}, \tau_k^{(t)})\}_{k=1}^K$ to $\theta^{(t+1)} = \{(\pi_k^{(t+1)}, \tau_k^{(t+1)})\}_{k=1}^K$ in the following form:

$$\tau_k^{(t+1)} = \frac{\sum_{i=1}^N \delta_{ik}^{(t)} x_i}{\sum_{i=1}^N \delta_{ik}^{(t)}}, \quad \pi_k^{(t+1)} = \frac{1}{N} \sum_{i=1}^N \delta_{ik}^{(t)}$$

- (b) What is the closed form expression of $\delta_{ik}^{(t)}$?

Problem 3 (Boosting)(0.5%)

1. Consider training a boosting classifier using decision stumps on the data set illustrated in Figure 1:



Figure 1: AdaBoost Data set

- (a) Which examples will have their weights increased at the end of the first iteration? Circle them.
 - (b) How many iterations will it take to achieve zero training error? Justify your answers.
2. Suppose AdaBoost is run on N training examples, and suppose on each round that the weighted training error ϵ_t of the t 'th weak hypothesis is at most $1/2 - \gamma$, for some number $0 < \gamma < 1/2$. After how many iterations, T , will the combined hypothesis be consistent with the N training examples, i.e., achieves zero training error? Your answer should only be expressed in terms of N and γ . (Hint: Recall that exponential loss is an upper bound for 0-1 loss. What is the training error when 1 example is misclassified?)

Problem 4 (Expectation Maximization Interpretation behind Semi-Supervised Learning)(1%)

Given N samples $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^m$ as well as their labels $y_1, \dots, y_N \in \{0, 1, \dots, K\}$. Consider the generative model where each sample \mathbf{x}_i is generated independently according to Gaussian mixture model that depends on the label y_i , as represented by random variable

$$X_i \sim \begin{cases} \sum_{j=1}^K \pi_j \mathcal{N}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) & , \text{ if } y_i = 0 \\ \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) & , \text{ if } y_i = k \neq 0 \end{cases}$$

where $\pi_1 + \dots + \pi_K = 1$, and $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes the Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$, with probability density function

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^m |\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

We would like to apply Expectation Maximization algorithm to find the maximum likelihood estimation of parameters $\theta = \{(\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\}_{k=1}^K$.

1. Please write down the E-step and M-step and show that the parameters are updated from $\theta^{(t)} = \left\{(\pi_k^{(t)}, \boldsymbol{\mu}_k^{(t)}, \boldsymbol{\Sigma}_k^{(t)})\right\}_{k=1}^K$ to $\theta^{(t+1)} = \left\{(\pi_k^{(t+1)}, \boldsymbol{\mu}_k^{(t+1)}, \boldsymbol{\Sigma}_k^{(t+1)})\right\}_{k=1}^K$ in the following form:

$$\pi_k^{(t+1)} = \frac{\sum_{i:y_i=0} \delta_{ik}^{(t)}}{\sum_{i:y_i=0} 1}$$

$$\boldsymbol{\mu}_k^{(t+1)} = \frac{\sum_{i:y_i=k} \mathbf{x}_i + \sum_{i:y_i=0} \delta_{ik}^{(t)} \mathbf{x}_i}{N_k + \sum_{i:y_i=0} \delta_{ik}^{(t)}}$$

$$\boldsymbol{\Sigma}_k^{(t+1)} = \frac{\sum_{i:y_i=k} (\mathbf{x}_i - \boldsymbol{\mu}_k^{(t+1)}) (\mathbf{x}_i - \boldsymbol{\mu}_k^{(t+1)})^T + \sum_{i:y_i=0} \delta_{ik}^{(t)} (\mathbf{x}_i - \boldsymbol{\mu}_k^{(t+1)}) (\mathbf{x}_i - \boldsymbol{\mu}_k^{(t+1)})^T}{N_k + \sum_{i:y_i=0} \delta_{ik}^{(t)}}$$

where $N_k = \sum_{i:y_i=k} 1$ is the number of samples in class k . Please show your derivations.

2. What is the closed form expression of $\delta_{ik}^{(t)}$? Please show your derivations.

Problem 5 (Label Propagation Algorithm)(1.5%)

In this problem, we will investigate label propagation algorithm by executing on a toy example. Next, we will show that the algorithm will convergence, which can be expressed analytically.

Let's consider the graph that we have seen in HW3 Problem 2.

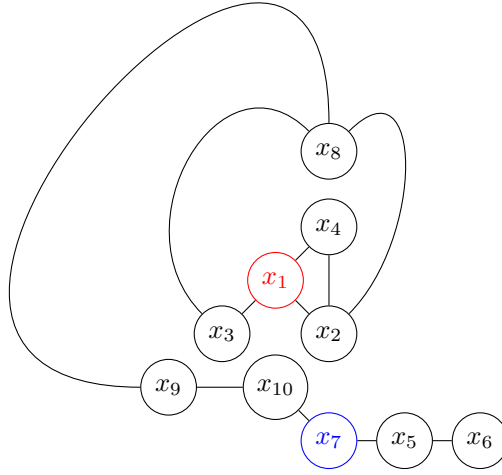


Figure 2: undirected connected graph G with labeled node

We have previously known that x_1 node is in Class 1 and x_7 node is in Class 2. Now, we want to separate these 10 nodes into Class 1 and Class 2.

Consider the transition matrix \mathbf{T} ,

$$\mathbf{T}_{i,j} = \frac{\tilde{\mathbf{W}}_{i,j}}{\sum_{k=1}^{10} \tilde{\mathbf{W}}_{k,j}}$$

, where $\tilde{\mathbf{W}}$ is the adjusted adjacency matrix of the graph G , which is defined as

$$\tilde{\mathbf{W}} = \mathbf{W} + \delta \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

. $\delta > 0$ is a small number. By adjusting the adjacency matrix, the weight of edge being connected in the original graph G is $1 + \delta$ and the weight of other edges is δ . Through the adjustment, we can prevent that the algorithm runs unsupervised due to the isolated labeled nodes. In the toy example, we set $\delta = 0.01$. $\mathbf{T}_{i,j}$ represents the probability that node j will propagate its own state to node i . For example, $\mathbf{T}_{2,1} = \mathbf{T}_{3,1} = \mathbf{T}_{4,1} \approx 0.326 \approx \frac{1}{3}$, $\mathbf{T}_{1,1} = \mathbf{T}_{5,1} = \mathbf{T}_{6,1} = \mathbf{T}_{7,1} = \mathbf{T}_{8,1} = \mathbf{T}_{9,1} = \mathbf{T}_{10,1} \approx 0.003$, which shows that node 1 will transfer its label to three neighbors with probability around 1 over 3. Also, it will transfer its label to other nodes(including itself) with probability slightly greater than zero.

1. Please write down the transition matrix \mathbf{T} .

Next, we define the label matrix sequence

$$\mathbf{Y}^t \in \mathbb{R}^{10 \times 2} \quad t = 0, 1, \dots$$

where the i_{th} row of \mathbf{Y}^t means the probability distribution of node x_i at time t . In this example, $\mathbf{Y}_{i,1}^t$ means the probability that the node x_i lies in Class 1 at time t , and $\mathbf{Y}_{i,2}^t$ means the probability that the node x lies in Class 2 at time t . We initialize $\mathbf{Y}_{1,1}^0 = 1, \mathbf{Y}_{1,2}^0 = 0$ because x_1 is labeled as Class 1. Also, $\mathbf{Y}_{7,1}^0 = 0, \mathbf{Y}_{7,2}^0 = 1$ because x_7 is labeled as Class 2. For other nodes i , we initialize $\mathbf{Y}_{i,1}^0 = \mathbf{Y}_{i,2}^0 = 0.5$.

After defining label matrix \mathbf{Y} and transition matrix \mathbf{T} , we introduce the algorithm below:

Algorithm 1 Label Propagation Algorithm in Toy Example

Input label matrix \mathbf{Y}^0 , transition matrix \mathbf{T} , tolerance level ϵ
Output node $x_i \in \{\text{Class 1, Class 2}\} \quad i = 1, \dots, 10$

- 1: **procedure** LABEL PROPAGATION($\mathbf{Y}^0, \mathbf{T}, \epsilon=10^{-8}$)
- 2: $t = 0$
- 3: **repeat**
- 4: $t = t + 1$
- 5: $\mathbf{Y}^t = \mathbf{T}\mathbf{Y}^{t-1}$ ▷ Random walk to its neighbor
- 6: $\mathbf{Y}_{i,j}^t = \mathbf{Y}_{i,j}^t / (\mathbf{Y}_{i,1}^t + \mathbf{Y}_{i,2}^t) \quad i = 1, \dots, 10, j = 1, 2$ ▷ Normalize the probability distribution
- 7: $\mathbf{Y}_{1,1}^t = 1, \mathbf{Y}_{1,2}^t = 0, \mathbf{Y}_{7,1}^t = 0, \mathbf{Y}_{7,2}^t = 1$ ▷ Clamp the labeled data
- 8: **until** $\|\mathbf{Y}^t - \mathbf{Y}^{t-1}\|_F < \epsilon$
- 9: If $\mathbf{Y}_{i,1}^t > 0.5$ then node x_i lies in Class 1, otherwise node x_i lies in Class 2, $i = 1, \dots, 10$
- 10: **end procedure**

There are three main procedures in the loop. First, every node propagates its own state to its neighbors with the transition probability. Next, we normalize the probability distribution for every node. In the last step, we clamp the probability distribution of the labeled data, which prevents the distribution of labeled data being influenced by unlabeled data and accelerates the convergence speed of the algorithm.

2. Please execute the algorithm. Write down the iteration number t , \mathbf{Y}^t , which nodes lies in Class 1 and which nodes lies in Class 2. Does the result correspond with the graph?

To show the convergence of label propagation algorithm, we consider more general case as the following statement.

Let $(x_1, y_1), \dots, (x_l, y_l)$ be labeled data, where y_i takes value in $\{1, \dots, C\}$, which indicates that x_i lies in Class y_i . Also, we have u unlabeled data $(x_{l+1}, y_{l+1}), \dots, (x_{l+u}, y_{l+u})$, where y_j is an unknown value, which lies in $\{1, \dots, C\}$.

We can construct the transition matrix \mathbf{T} ,

$$\mathbf{T}_{i,j} = \frac{\tilde{\mathbf{W}}_{i,j}}{\sum_{k=1}^{l+u} \tilde{\mathbf{W}}_{k,j}}$$

, where $\tilde{\mathbf{W}}_{i,j} > 0$, $i, j = 1, \dots, l+u$.

Also, we define the label matrix sequence

$$\mathbf{Y}^t \in \mathbb{R}^{(l+u) \times C} \quad t = 0, 1, \dots$$

where the i_{th} row of \mathbf{Y}^t means the probability distribution of node x_i at time t . $\mathbf{Y}_{i,j}^t$ means the probability that the node x_i lies in Class j at time t . For \mathbf{Y}^0 , we clamp the first l rows as following,

$$\mathbf{Y}_{i,j}^t = \mathbb{1}\{y_i = j\} \quad i = 1, \dots, l, j = 1, \dots, C$$

, which indicates that x_i must lies in Class y_i . For the other rows, we initialize

$$\mathbf{Y}_{i,j}^t = \frac{1}{C} \quad i = l+1, \dots, l+u, j = 1, \dots, C$$

We execute by the algorithm below.

Algorithm 2 Generalized Label Propagation Algorithm

Input label matrix \mathbf{Y}^0 and transition matrix \mathbf{T}

Output \mathbf{Y}^*

- 1: **procedure** GENERALIZE LABEL PROPAGATION($\mathbf{Y}^0, \mathbf{T}, t = 0$)
 - 2: **repeat**
 - 3: $t = t + 1$
 - 4: $\mathbf{Y}^t = \mathbf{T}\mathbf{Y}^{t-1}$ ▷ Random walk to its neighbor
 - 5: $\mathbf{Y}_{i,j}^t = \mathbf{Y}_{i,j}^t / \sum_{k=1}^C \mathbf{Y}_{i,k}^t \quad i = 1, \dots, l+u, j = 1, \dots, C$ ▷ Normalize the probability distribution
 - 6: $\mathbf{Y}_{i,j}^t = \mathbb{1}\{y_i = j\} \quad i = 1, \dots, l, j = 1, \dots, C$ ▷ Clamp the labeled data
 - 7: **until** \mathbf{Y}^t converges to \mathbf{Y}^*
 - 8: **end procedure**
-

Once we get \mathbf{Y}^* , we can conclude that the unlabeled data (x_i, y_i) belong to Class y_i , where

$$y_i = \arg \max_j \mathbf{Y}_{i,j}^*$$

Next, we want to calculate \mathbf{Y}^* analytically.

3. Please show that the line 4 and 5 in Algorithm 2 can be combined as

$$\mathbf{Y}^t = \bar{\mathbf{T}} \mathbf{Y}^{t-1}$$

, where $\bar{\mathbf{T}}_{i,j} = \mathbf{T}_{i,j} / \sum_{k=1}^{l+u} \mathbf{T}_{i,k}$, $i, j = 1, \dots, l+u$.

We split $\bar{\mathbf{T}}$ to $\begin{bmatrix} \bar{\mathbf{T}}_{ll} & \bar{\mathbf{T}}_{lu} \\ \bar{\mathbf{T}}_{ul} & \bar{\mathbf{T}}_{uu} \end{bmatrix}$, where $\bar{\mathbf{T}}_{ll} \in \mathbb{R}^{l \times l}$, $\bar{\mathbf{T}}_{uu} \in \mathbb{R}^{u \times u}$. Also, we split \mathbf{Y}^t to

$$\begin{bmatrix} \mathbf{Y}_L^t \\ \mathbf{Y}_U^t \end{bmatrix}, \text{ where } \mathbf{Y}_L^t \in \mathbb{R}^l, \mathbf{Y}_U^t \in \mathbb{R}^u$$

4. Please show that after an iteration,

$$\mathbf{Y}_U^t = \bar{\mathbf{T}}_{uu} \mathbf{Y}_U^{t-1} + \bar{\mathbf{T}}_{ul} \mathbf{Y}_L^{t-1}, t = 1, 2, \dots,$$

$$\mathbf{Y}_L^t = \mathbf{Y}_L^{t-1}, t = 1, 2, \dots$$

5. From the above result, we let $\mathbf{Y}_L = \mathbf{Y}_L^t$ for any t . Show that for $t \geq 1$

$$\mathbf{Y}_U^t = \bar{\mathbf{T}}_{uu}^t \mathbf{Y}_U^0 + \sum_{i=1}^t \bar{\mathbf{T}}_{uu}^{i-1} \bar{\mathbf{T}}_{ul} \mathbf{Y}_L$$

6. Please show that $\sum_{j=1}^u \bar{\mathbf{T}}_{uu,i,j} = \gamma_i$, where $0 < \gamma_i < 1$, for $i = 1, \dots, u$. Use the fact to derive $\sum_{j=1}^u \bar{\mathbf{T}}_{uu,i,j}^n \leq \gamma^n$, for $i = 1, \dots, u$, where $\gamma = \max_{i=1, \dots, u} \gamma_i$. Last, derive $\lim_{n \rightarrow \infty} \bar{\mathbf{T}}_{uu}^n = \mathbf{O}$.

7. Define $\mathbf{S}_n = \mathbf{I} + \bar{\mathbf{T}}_{uu} + \bar{\mathbf{T}}_{uu}^2 + \dots + \bar{\mathbf{T}}_{uu}^{n-1}$, $\mathbf{S}_n(\mathbf{I} - \bar{\mathbf{T}}_{uu}) = \mathbf{I} - \bar{\mathbf{T}}_{uu}^n$. Use the fact to show that $\lim_{n \rightarrow \infty} \mathbf{S}_n = (\mathbf{I} - \bar{\mathbf{T}}_{uu})^{-1}$. Combined all the result above, please show that $\mathbf{Y}_U^* = (\mathbf{I} - \bar{\mathbf{T}}_{uu})^{-1} \bar{\mathbf{T}}_{ul} \mathbf{Y}_L$. Hence, $\mathbf{Y}^* = \begin{bmatrix} \mathbf{Y}_L^* \\ \mathbf{Y}_U^* \end{bmatrix}$, which can be obtained analytically regardless of the initial value \mathbf{Y}^0 .

8. Please calculate the analytical solution \mathbf{Y}^* on the toy example above. Does the solution correspond to the iteration solution \mathbf{Y}^t ?

Version Description

1. First Edition: Finish Problem 1 to 5
2. Second Edition: Revise the data description in Problem 5 Generalized Label Propagation.
3. Third Edition: Revise Problem 5 to make the definition more robust.
4. Fourth Edition: Fix small typo at Problem 5 (5) $\mathbf{Y}^0 \rightarrow \mathbf{Y}_U^0$
5. Fifth Edition: Fix typo at Problem 5 (6)
6. Sixth Edition: Problem 4: edit π_k^{t+1}