

This assignment contains 6 questions. Please read and follow the following instructions.

- **DUE DATE: Oct 4th, 11:45 PM**
 - **TOTAL NUMBER OF POINTS: 135**
 - **NO PARTIAL CREDIT** will be given so provide concise answers.
 - Clearly list your team ID, each team member's **names and Unity IDs** at the top of your submission.
 - Submit only a **single PDF file** per group, containing your answers.
-

1. (45 points) [Song Ju] [**Expectation Maximization**]

You are running a Naive Bayes classifier using 3 binary feature variables X_1, X_2, X_3 to predict the status of a nuclear power plant Y (0 for "Normal" and 1 for "Malfunction". Note that the actual status of the nuclear power plant Y cannot be directly observed but rather it can only be estimated through sensors such as the core temperatures. In this simplified scenario, let's assume there is only one sensor, Z with binary values: 0 for "Normal" vs. 1 for "Abnormal". One day you realize some of the sensor values Z are missing. Based on the nuclear power plant's manual book, the probability of missing values is much higher when $Y = 1$ than when $Y = 0$. More specifically, the exact values from the sensor specifications are:

$$P(Z \text{ missing} | Y = 1) = .15; P(Z = 1 | Y = 1) = .85$$

$$P(Z \text{ missing} | Y = 0) = .03; P(Z = 0 | Y = 0) = .97$$

- (5 points) Draw a Bayes net that represents this problem with a node Y that is the unobserved label, a sensor node Z that is either a copy of Y or has the value "missing", and the three features X_1, X_2, X_3 .
- (5 points) What is the probability of $Y = 1$ given that our sensor data Z is missing, i.e., $P(Y = 1 | Z = \text{"missing"})$? Derive your answer in terms of $\theta_{Y=1}$, denoting $P(Y = 1)$.
- (5 points) Based on the sensor data Z , we approximate the count and probabilities

of each X_i variables under condition $Y = 1$, which are shown below.

$$\text{Count}(X_1 = 1|Y = 1) = 20,$$

$$\text{Count}(X_2 = 1|Y = 1) = 24,$$

$$\text{Count}(X_3 = 1|Y = 1) = 4,$$

$$P(X_1 = 1|Y = 1) = \frac{1}{2},$$

$$P(X_2 = 1|Y = 1) = \lambda,$$

$$P(X_3 = 1|Y = 1) = 1 - 3\lambda,$$

Assume $X_1|Y, X_2|Y, X_3|Y$ are all Bernoulli variables. Please apply log likelihood to calculate λ .

- (d) (15 points) Despite the approximation in part (c), we would like to theoretically learn the best choice of parameters for $P(Y), P(X_1|Y), P(X_2|Y)$, and $P(X_3|Y)$. Assume $Y, X_1|Y, X_2|Y, X_3|Y$ are all Bernoulli variables and let us denote θ to include the following parameters:

$$\theta_{Y=y} = P(Y = y),$$

$$\theta_{X_1=x_1|Y=y} = P(X_1 = x_1|Y = y),$$

$$\theta_{X_2=x_2|Y=y} = P(X_2 = x_2|Y = y)$$

$$\theta_{X_3=x_3|Y=y} = P(X_3 = x_3|Y = y).$$

Write the log-likelihood of X, Y and Z given θ , in terms of θ and $P(Z|Y)$, first for a single example $(X_1 = x_1, X_2 = x_2, X_3 = x_3, Z = z, Y = y)$, then for n i.i.d. examples $(X_1^i = x_1^i, X_2^i = x_2^i, X_3^i = x_3^i, Z^i = z^i, Y^i = y^i)$ for $i = 1, \dots, n$.

- (e) (15 points) Provide the E-step and M-step for performing expectation maximization of θ for this problem. For the $(t + 1)$ th iteration, in the E-step compute the distribution $Q_{t+1}(Y|Z, X)$ using

$$Q_{t+1}(Y = 1|Z, X) = E[Y|Z, X_1, X_2, X_3, \theta_t]$$

using your Bayes net from part (a) and conditional probability from part (b) for the unobserved class label Y of a single example.

In the M-step, compute:

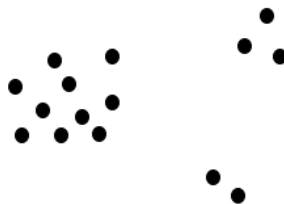
$$\theta_{t+1} = \operatorname{argmax}_{\theta} \sum_{i=1}^n \sum_y Q(Y^i = y|Z^i, X^i) \log P(X_1^i, X_2^i, X_3^i, Y^i, Z^i|\theta_t)$$

using all of the examples $(X_1^1, X_2^1, X_3^1, Y^1, Z^1), \dots, (X_1^n, X_2^n, X_3^n, Y^n, Z^n)$. Note: it is OK to leave your answers in terms of $Q(Y|Z, X)$.

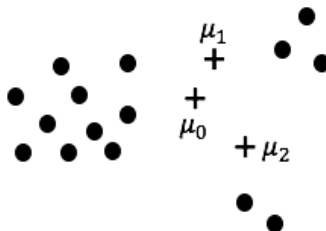
2. (20 points) [Song Ju] **Gaussian Mixture Model (GMM)**

Consider the set of training data in the graph below, let's assume it contains three clusters. For GMM, the means and variances of three Gaussians are μ_0 and σ_0, μ_1

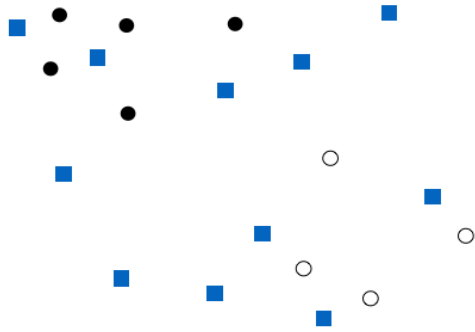
and σ_1 , and μ_2 and σ_2 , respectively. Additionally, we have π_0, π_1, π_2 to denote the mixture proportions of the three Gaussians (i.e., $p(x) = \pi_0 N(\mu_0, \sigma_0 I) + \pi_1 N(\mu_1, \sigma_1 I) + \pi_2 N(\mu_2, \sigma_2 I)$), where I is the identity matrix and $\pi_0 + \pi_1 + \pi_2 = 1$. We will also use θ to refer to the entire collection of parameters $(\mu_0, \mu_1, \mu_2, \sigma_0, \sigma_1, \sigma_2, \pi_0, \pi_1, \pi_2)$ defining the mixture model $p(x)$.



- (a) (10 points) Would K-Means ($K = 3$) and our 3-cluster GMM trained using EM produce the same cluster centers (means) for this data set above? Justify your answer. (Answer without any justification will get zero point.)
- (b) (10 points) In the following, we apply EM to train our 3-cluster GMM on the data below. The ‘+’ points indicate the current means μ_0 , μ_1 , and μ_2 of the three Gaussians after the k th iteration of EM.



- (b.1) (3 points) On the figure, draw the directions in which μ_0 , μ_1 and μ_2 will move in the next EM iteration.
- (b.2) (3 points) Will the marginal likelihood of the data, $\prod_j P(x^j | \theta)$ increase or decrease on the next EM iteration? Explain your reasoning.
- (b.3) (4 points) Will the estimate of π_0 increase or decrease on the next EM iteration? Explain your reasoning.
3. (26 points) [Farzaneh Khoshnevisan] **Semi-supervised Learning**
 Consider the following figure which contains labeled (L) (class 1: black circles, class 2: hollow circles) and unlabeled (U) (blue squares) data. In this question, you will use two semi-supervised methods: S³VM and co-training, to utilize the unlabeled data for further improvement of a SVM classifier.



- (a) (10 points) Explain how would the semi-supervised SVM (S^3VM) perform on this data as compared to the supervised SVM, by plotting the separating hyper-planes produced by both algorithms. For SVM, please draw marginal boundaries and separating hyper-plane by solid lines. Using a different color, draw the marginal boundaries and the separating hyper-plane of S^3VM using dash lines.
- (b) (10 points) In applying co-training, g_1 and g_2 are SVM classifiers, and p_i/n_i represents the number of positive or negative points to label at iteration i , respectively. In this example, assume that $p_1, n_1 = 2$ and $p_2, n_2 = 1$.
- (b.1) (5 points) Explain What is the main underlying assumption in applying co-training? How can we apply it to the above data (what are the two classifiers)?
- (b.2) (5 points) Identify the label of each point $x_i \in U$ and the final SVM separating hyper-plane after applying 2 iterations of semi-supervised co-training. Assume that at iteration i , p_i and n_i points are labeled based on the farthest distance from the separation hyper-plane.
- (c) (6 points) Compare the separating hyper-planes produced by S^3VM from part (a) and co-training from part (b.2). Explain why you think these two algorithms perform similarly/differently?
4. (15 points) [Farzaneh Khoshnevisan] **Generalized Sequential Pattern**
Consider a data sequence:

$$S = (\{A, B\}\{B, C\}\{D, E\}\{A, D\}\{B, E, F\})$$

and the following time constraints:

- $\text{min_gap}=0$ (interval between last event in e_i and first event in e_{i+1} is > 0)
- $\text{max_gap}=2$ (interval between first event in e_i and last event in e_{i+1} is ≤ 2)
- $\text{max_span}=5$ (interval between first event in e_1 and last event in e_{last} is ≤ 6)
- $w_s = 1$ (time between first and last events in e_i is ≤ 1)

For each of the sequences $w = (e_1, \dots, e_{last})$ below, determine whether they are subsequences of S , and if not, which constraint is excluding them.

(a) (3 points) $w = (\{A\}\{B\}\{C\}\{D\}\{E\})$

(b) (3 points) $w = (\{B\}\{D\}\{E\})$

(c) (3 points) $w = (\{D, E\}\{D, E\})$

(d) (3 points) $w = (\{A\}\{C, D, E\}\{A, F\})$

(e) (3 points) $w = (\{A, B, C, D\}\{E, F\})$

5. (15 points) [Farzaneh Khoshnevisan] **Generalized Sequential Pattern**

Consider the following frequent 3-sequences:

$\langle \{1, 2, 3\} \rangle, \langle \{1, 3\}\{4\} \rangle, \langle \{1\}\{2, 4\} \rangle, \langle \{1\}\{4\}\{5\} \rangle, \langle \{2, 3\}\{4\} \rangle, \langle \{2, 3\}\{5\} \rangle, \langle \{2, 4\}\{4\} \rangle, \langle \{2\}\{4, 5\} \rangle, \langle \{3\}\{4, 5\} \rangle, \langle \{3\}\{4\}\{5\} \rangle, \langle \{4\}\{4, 5\} \rangle.$

(a) (5 points) List all the candidate 4-sequences produced by the candidate generation step of the GSP algorithm.

(b) (5 points) List all the candidate 4-sequences pruned during the candidate pruning step of the GSP algorithm (assuming no timing constraints).

(c) (5 points) List all the candidate 4-sequences pruned during the candidate pruning step of the GSP algorithm (assuming $\text{max_gap} = 1$).

6. (14 points) [Farzaneh Khoshnevisan] **Decision Theory**

Imagine you want to purchase a start-up company. With 5% chance the company's value will be up and you will make \$10,000K, and 95% chance that the company will fail and you will lose \$600K. Additionally, it would cost you \$200K to hire a group of experts for consultation who will help you determine how promising the company will be. This group of the experts are known to be accurate 75% of time.

Your goal is to maximize the expected value of your decision. What, if any, the best action(s) should you take, and what is your expected value? Assume that you are risk-neutral.

Draw a decision tree that supports your conclusion and show all the probabilities and utility values for every node.