

Internet Analytics (COM-308): Final Exam

June 29, 2024

Duration: **2h45**.

Total points: **100**.

Number of pages: **19**.

Allowed documents: **class notes, lab handouts, homeworks, your own code**.

There should in general be enough room below every question for intermediate calculations and your answer. However, you are allowed to use additional sheets of paper; please **write your name on every sheet**, number them, and staple them to this document before handing in.

The use of **mobile phones, tablets, laptop computers**, and other communication devices is **prohibited**.

Last name:
First name:
SCIPER number:
Signature:

Please leave blank.

1	2	3	4	5	6	7	Total
30	12	14	14	10	10	10	100

Question 1: Multiple Choice Questions (30 points)

(30 pts) All questions have a single answer. Check the correct one. Grading:

- Correct answer: +2 points;
- Wrong answer: −1 point;
- No answer or "I don't know": 0 point.

1. Here is a basic implementation of the power method to estimate the PageRank scores using python and numpy:

```
def power_method(probas, G_matrix, S):  
    for iteration in range(S):  
        probas = probas @ G_matrix  
    return probas
```

G_matrix is the google matrix, which has been created using a teleportation probability $0 < 1 - \theta < 1$.

We note N the number of nodes of the graph, M the number of edges of the graph and S the number of iterations of the algorithm.

What is the smallest upper bound we can give on the number of operations of this function (additions and multiplications)?

- ☐ $O(SN^2)$
- ☐ $O(SNM)$
- ☐ $O(SM)$
- ☐ I don't know

2. We propose the following algorithm to reduce the number of operations without changing the result:

```
def power_method2(probas, graph_transitions, theta, S):  
    nb_nodes = len(probas[0])  
    for iteration in range(S):  
        probas = probas @ graph_transitions * theta  
        [.....]  
    return probas
```

Which line of code is missing (marked by [...])?

- ☐ `probas += (1 - theta) * np.ones((1, nb_nodes)) / nb_nodes`
- ☐ `probas = (1 - theta) * probas @ graph_transitions`
- ☐ `probas = (1 - theta) * probas + np.ones((1, nb_nodes)) @ graph_transitions`
- ☐ I don't know

3. We use the same notation as 1.

What is the smallest upper bound we can give on the number of operations of the new function `power_method2()`? (we assume that adding 0 or multiplying by 0 does not count as an operation)

- ☐ $O(SN \log(M))$
- ☐ $O(SM)$
- ☐ $O(SN)$
- ☐ I don't know

4. We consider a random graph G whose edges are all independent. G has N nodes split in two communities A and B. 50% of the nodes are in A and 50% are in B. For every pair of nodes in A, the probability of an edge is $P_{AA} = 0.1$, for every pair in B, $P_{BB} = 0.2$ and for every pair of nodes in different communities, the edge probability is $P_{AB} = 0.001$.

When N goes to infinity, which of the following is more likely to happen?

- ☐ G will have a giant component and some other connected components.
- ☐ G will be connected.
- ☐ G will not have any giant component.
- ☐ I don't know

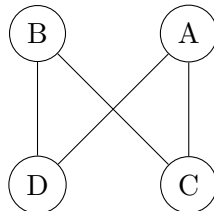
5. On the same graph G as the previous question, someone is trying to estimate the size of A and B by performing a random walk on the graph. The random walk starts from a random node. For each node visited, she records the cluster of the node (A or B).

After a very large number of steps, she tries to estimate the size of each cluster by naively comparing the number of nodes of each cluster she visited.

Which cluster will she estimate to be bigger?

- ☐ A
- ☐ B
- ☐ It depends on the starting point of the random walk.
- ☐ I don't know

6. Consider the undirected graph $G = (V, E)$ with vertices $V = \{A, B, C, D\}$ below. What is the number of Bayesian networks compatible with G (same edges as G but all oriented) where the following two conditional independences hold: $A \perp B \mid D$ and $C \perp D \mid A, B$?



- ☐ 0
- ☐ 3
- ☐ 4
- ☐ I don't know

7. Here is a simple implementation of the Latent Dirichlet Allocation (LDA) model using RDD operations in spark.

- **RDD_corpus**: An RDD of documents, where each document is represented as a tuple of (docID, list of terms). Example:

```

RDD_corpus = sc.parallelize([
    (1, ['word1', 'word2', 'word3']),
    (2, ['word2', 'word4', 'word5'])
])

```

- **num_topics**: Number of topics to be extracted from the documents.
- **max_iterations**: Number of iterations for the LDA algorithm to converge.

```

# Step 1: Initialize random topic assignments
def initialize_topics(doc):
    docID, terms = doc
    return [(docID, term, randint(0, num_topics - 1)) for term in terms]

topic_assignments = RDD_corpus.flatMap(initialize_topics)

# Step 2: Iterate to optimize the topic assignments
for iteration in range(max_iterations):
    # Calculate topic probabilities and assign new topics
    def sample_topics((docID, term, topic)):
        # Placeholder for actual topic sampling logic
        return (docID, term, randint(0, num_topics - 1))

    topic_assignments = topic_assignments.map(sample_topics)

# Step 3: Aggregate topic assignments into final model structures
A = topic_assignments.map(lambda x: (x[0], x[2])).countByValue()
B = topic_assignments.map(lambda x: (x[1], x[2])).countByValue()

```

What are A and B :

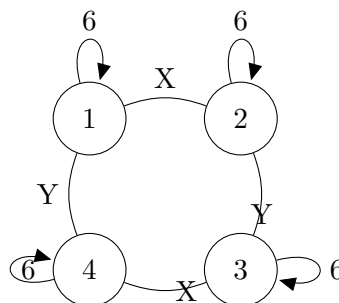
- ☐ doc topic distribution, term topic distribution
- ☐ number of docs, number of topics
- ☐ number of topics, number of docs
- ☐ I don't know

8. As seen in class, the Louvain method generates a sequence of increasingly smaller, weighted graphs containing self-loops. The general expression for modularity in this case is

$$Q = \frac{1}{2m} \sum_{c_i \in C} \sum_{u,v \in c_i} \left(w_{uv} - \frac{d_u d_v}{2m} \right), \quad (1)$$

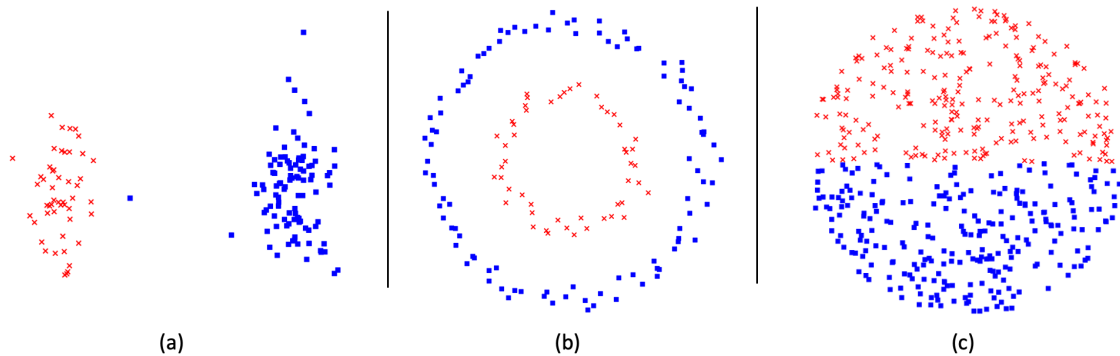
where w_{uv} is the weight of the edge between node u and v (w_{uu} for a self-loop), m is the sum of the weights of all links in the network, d_u and d_v are the sum of the weights of the links incident to node u and v respectively.

Which combination of X and Y will cause the graph to have only two communities in the next iteration?



- ☐ $X = 1, Y = 1$
- ☐ $X = 0.5, Y = 1$
- ☐ $X = 1, Y = 2$
- ☐ I don't know

9. We run the clustering methods K-means and Gaussian Mixture Models (GMM) on two-dimensional data points, with both methods using $K = 2$ clusters. In each figure below, points assigned to cluster 1 are represented by red crosses, while points assigned to cluster 2 are represented by blue squares. For each of the following figures, determine which clustering methods could have led to the cluster assignments. Note that more than one method could be correct for each figure, or none of them. For GMM, a point is assigned to the cluster with maximum posterior.



Which of the following statements is true?

- ☐ Only GMM can produce (c)
- ☐ Only K-means can produce (c)
- ☐ Both GMM and K-means can produce (c)
- ☐ I don't know

10. As in the previous question, which of the following statements is true?

- ☐ K-means can produce (a)
- ☐ GMM can produce (a) and (b)
- ☐ None of the above statements are true
- ☐ I don't know

11. Which of the following statements is true?

- ☐ The generative model of probabilistic Latent Semantic Analysis (pLSA) is:
For each document d in the corpus: (a) sample a topic z from a prior $p(z)$; (b) sample a document d according to $p(d|z)$; (c) sample each word w in d according to $p(w|z)$
- ☐ If PCA is performed on an uncorrelated data set (meaning that the covariance matrix is diagonal), the variances along the different principal directions are all equal
- ☐ If A is a real symmetric matrix with non-negative eigenvalues, then the eigenvalues and singular values of A coincide
- ☐ I don't know

12. We classify a message $M = (w_1, w_2, w_3)$ consisting of three words as ham (G) or spam (B), using a Naive Bayes classifier. Suppose the prior for both classes is positive: $P(B) > 0$, $P(G) > 0$, and the class posterior for message M is positive as well: $P(G|M) > 0$, $P(B|M) > 0$.

Now we create a new, longer message M' , by repeating M multiple times: $M' = (w_1, w_2, w_3, w_1, w_2, \dots, w_3)$. Which of the following inequalities guarantees that for M' sufficiently long, we classify M' as ham (i.e., $P(G|M') > P(B|M')$)?

- ☐ $P(G) > P(B)$
- ☐ $P(M|G) > P(M|B)$
- ☐ $P(G)P(M|G) > P(B)P(M|B)$
- ☐ I don't know

13. Which of the following assertions about the Latent Dirichlet Allocation (LDA) text model is false?

- ☐ Conditional on the topic, the word is independent of the document.
- ☐ Increasing the α parameter of the Dirichlet prior tends to lead to sparser topic distributions for documents.
- ☐ It is possible to generate words for a new document not seen during training.
- ☐ I don't know

14. Consider a SIR epidemic model with one percent of the population initially infected on a population of $N = 10'000$ people. We assume that each person meets on average 120 people every year, and patients stay infected for an average of one week.

Among the following options, what is the smallest percentage of the population that needs to be vaccinated in order to prevent the epidemic to already be a large-scale epidemic? (we assume the vaccine fully prevents transmission)

- ☐ 1/2
- ☐ 2/3
- ☐ 3/4
- ☐ I don't know

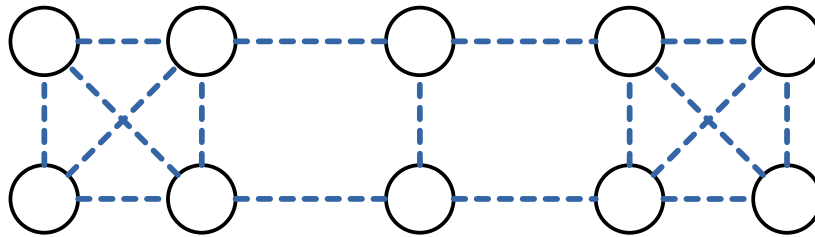
15. In relation to Skip-gram, for a specific word-context pair (w, c) , which of the following is the partial derivative of the objective function $L_{wc} = \log \sigma_s(c, w)$ with respect to the context vector v_c ? ($u_w \in \mathbb{R}^d$ and $v_c \in \mathbb{R}^d$ are the vectors of the word and context and

$$\sigma_s(c, w) = \frac{e^{u_w^T v_c}}{\sum_{c' \in V} e^{u_w^T v_{c'}}})$$

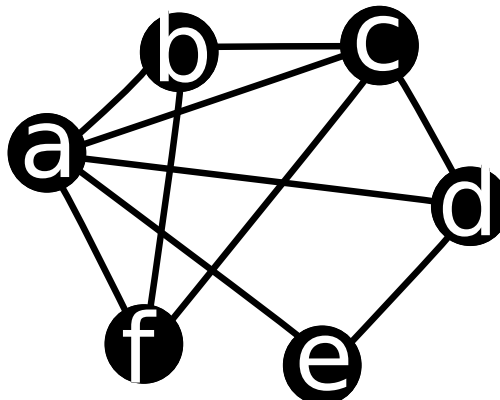
- ☐ u_w
- ☐ $u_w(1 - \sigma_s(c, w))$
- ☐ $u_w(1 - \frac{1}{\sum_{c' \in V} e^{u_w^T v_{c'}}})$
- ☐ I don't know

Question 2: Network Structure (12 points)

1. (4 pts) Find the largest possible set of strong links such that the strong triadic closure (STC) property holds for all nodes of the following graph. Mark the strong edges in the figure.



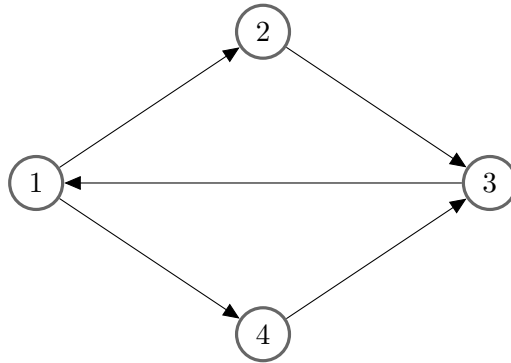
Consider the following graph. We are interested in its transitivity.



2. (2 pts) What is the clustering coefficient c_d of node d ?
3. (4 pts) What is the weighted average clustering coefficient (transitivity) of the entire graph of the previous question?
4. (2 pts) Is this graph strongly clustered relative to an equivalent-density $G(n, p)$?

Question 3: PageRank (14 points)

We perform PageRank on the following graph.



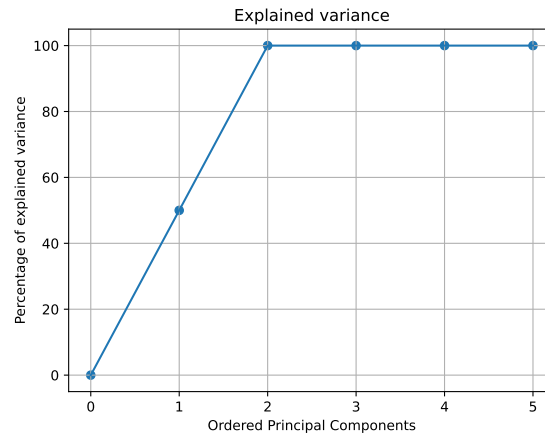
1. (1 pt) For $\theta = 1$, which nodes are dangling nodes and what are the connected components?
2. (2 pts) For which values of θ is this Markov chain ergodic? When it is not, justify why not.
3. (3 pts) Approximate the PageRank scores for $\theta = 0.999$.

Now we add nodes 5, 6, 7 to the graph as shown on the image below.

6. (3 pts) Rank the nodes by increasing PageRank scores if $\theta = 0.001$ (specify if there are equalities).

Question 4: PCA (14 points)

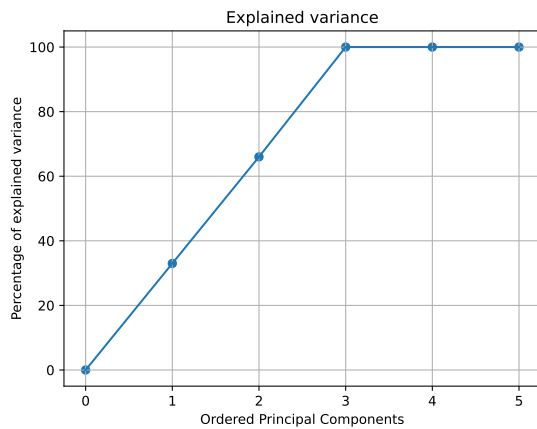
You have a dataset of N data points $x_1, \dots, x_N \in \mathbb{R}^5$. You rigorously follow the steps of PCA and you plot the cumulative variance explained by the principal components. This is the plot you obtain:



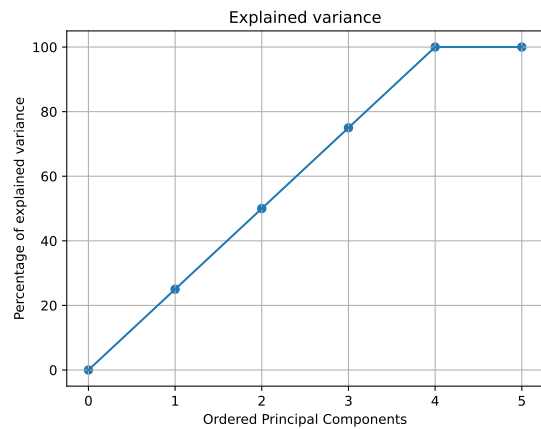
- (1 pt) From this plot, what can you deduce about the dimensionality of the data?
- (4 pts) Is it possible to observe this plot if $N = 1$, $N = 2$, $N = 3$, $N = 4$? Justify your answers.
- (3 pts) Let N_0 be the smallest N you found for the previous question. Give an example of a set of point(s) $x_1, \dots, x_{N_0} \in \mathbb{R}^5$ that explain the plot (give their coordinates).

4. (3 pts) Compute the covariance matrix for these points. What are the eigenvalues of the covariance matrix?

5. (3 pts) You don't change your dataset, but you add a new point to it. It now contains $N_0 + 1$ points $x_1, \dots, x_{N_0+1} \in \mathbb{R}^5$. Can the new plot be plot (a)? Can it be plot (b)? Justify your answers.



(a)



(b)

Question 5: Collaborative Filtering (10 points)

Here are two different datasets (a) and (b) of the ratings given by 3 users (1, 2, 3) about 4 items (A, B, C, D). The X represents a missing score.

	A	B	C	D
1	1	2	3	-1
2	0	1	2	-2
3	3	4	5	X

(a)

	A	B	C	D
1	1	2	3	-1
2	0	1	2	-2
3	0	4	5	X

(b)

- (1 pt) Call r_{ui} the score given by user u to item i . Call $(b_u), u \in \{1, 2, 3\}$ the user biases and $(b_i), i \in \{A, B, C, D\}$ the item biases. One of the two datasets follows the model $r_{ui} = b_u + b_i$, the other one follows the model $r_{ui} = b_u + b_i + q_u p_i$, where q_u and p_i are respectively a user parameter and an item parameter.
Identify which dataset follows which model. Justify your answer.

For the rest of the exercise, we focus on the model $r_{ui} = b_u + b_i$.

- (1 pt) Predict the score that user 3 would give to item D in the dataset you identified to follow the model $r_{ui} = b_u + b_i$.

3. (2 pts) Suppose we have a single user u and a single item i , and the corresponding rating r_{ui} . We are trying to find bias parameters in the usual way, which amounts to the optimization problem $(b_u^*, b_i^*) = \operatorname{argmin}_{(b_u, b_i)} (r_{ui} - b_u - b_i)^2$.
Give (all) the solution(s) to this minimization problem.

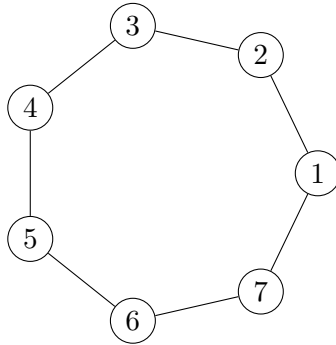
4. (3 pts) Modify the minimization problem above by adding a quadratic regularizer to the parameters.
Give (all) the solution(s) to the modified minimization problem.

5. (3 pts) We now return to the full problem with multiple u and i :

$$\min_{b_u, b_i} \sum_{u,i} (r_{ui} - b_u - b_i)^2$$
Under which condition does this problem have a unique solution?

Question 6: Random Walk and Conductance (10 points)

A particle performs a random walk on a cycle graph with nodes numbered from 1 to n , denoted as C_n . For example C_7 is as follows:



1. (2 pts) For which n is this random walk ergodic? Justify your answer.

For the rest of the exercise, consider all n such that this ergodicity condition is satisfied.

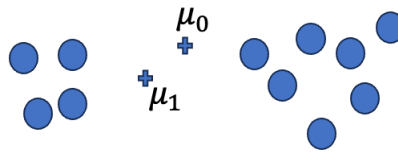
2. (3 pts) Let $\lambda_1(n)$ and $\lambda_2(n)$ be the first and second largest eigenvalues of the transition matrix of the random walk C_n . Determine $\lambda_1(n)$ and $\lambda_2(n)$ as n goes to infinity?

3. (3 pts) For a fixed n , add an edge to C_n that increases the conductance the most. Where would you add this edge? Call this new graph C'_n and compute its conductance $\Phi_{C'_n}$. Compute the ratio $\frac{\Phi_{C'_n}}{\Phi_{C_n}}$ asymptotically when n goes to infinity.
4. (2 pts) For the graph C_n , compute the expected number of steps for a particle starting at node 1 to visit node $\frac{n+1}{2}$ for the first time. Show your detailed calculations (Hint: consider this expected number of steps for every node and write a system of equations based on them).

Question 7: Clustering (10 points)

This question concerns training Gaussian Mixture Models (GMMs). Throughout, we will assume there are two Gaussian components in the GMM. We will use μ_0, μ_1, σ_0^2 , and σ_1^2 to define the means and variances of these two Gaussians, and will use π_0 and $(1 - \pi_0)$ to denote the mixture proportions of the two Gaussians (i.e., $p(x) = \pi_0 N(\mu_0, \sigma_0^2 I) + (1 - \pi_0) N(\mu_1, \sigma_1^2 I)$). We will also use θ to refer to the entire collection of parameters $\langle \mu_0, \mu_1, \sigma_0^2, \sigma_1^2, \pi_0 \rangle$ defining the mixture model $p(x)$.

1. (2 pts) Consider applying EM to train a Gaussian Mixture Model (GMM) to cluster the data below into two clusters. The “+” points indicate the current means μ_0 and μ_1 of the two Gaussian mixture components after the k^{th} iteration of EM. Draw on the figure the directions in which μ_0 and μ_1 will move during the next M-step. Justify your answer.



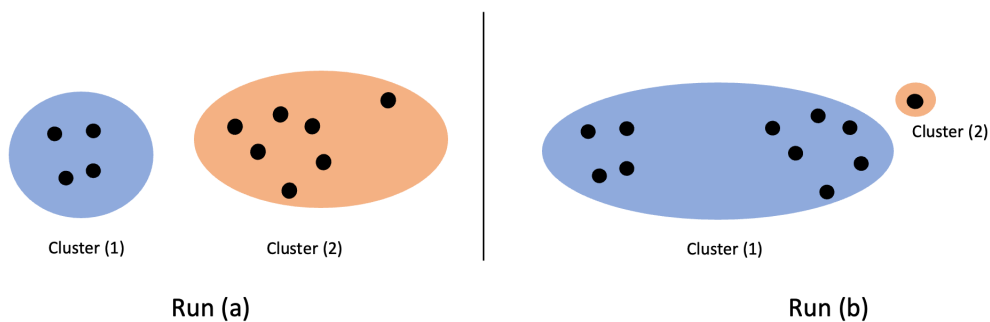
2. (1 pt) Is it possible that the marginal likelihood of the data (objective function), $\prod_j P(x_j|\theta)$, decreases on the next EM iteration? Justify your answer.

3. (2 pts) Consider the set of training data below. After running the Gaussian Mixture Model (GMM) twice on this data, with different initializations, we obtain the following results for the determinants of the covariance matrices of the probability density functions (pdf) for each cluster (we have two clusters, $K = 2$):

- For run **a**: $\det(\Sigma_1) = 5$ and $\det(\Sigma_2) = 10$
- For run **b**: $\det(\Sigma_1) = 15$ and $\det(\Sigma_2) = 0.00001$

Here, Σ_1 and Σ_2 represent the covariance matrices for the Gaussian distributions of clusters 1 and 2, respectively.

Which run, **a** or **b**, has the higher likelihood value? Provide a detailed justification for your answer.



4. (2 pts) Consider a training dataset consisting of partially labeled samples. For the first m samples x_1, x_2, \dots, x_m , we have corresponding z labels (which can be 0 or 1). The remaining samples, $x_{m+1} \dots x_{m+n}$, are unlabeled with unknown z values. Given this setup, which of the likelihood functions should be maximized to exploit both labeled and unlabeled data in GMM clustering? Justify your answer.

- (a) $(\prod_{i=1}^m P(x_i|z_i, \theta))(\prod_{i=m+1}^{m+n} P(x_i|\theta))$
 (b) $(\prod_{i=1}^m P(x_i, z_i|\theta))(\prod_{i=m+1}^{m+n} P(x_i|\theta))$
 (c) $(\prod_{i=1}^m P(z_i|\theta)P(x_i|\theta))(\prod_{i=m+1}^{m+n} P(x_i|\theta))$

5. (3 pts) How do you need to change the E-step and M-steps of the EM algorithm to fit the above objective function?

6. (0 pts, but a deep sense of satisfaction if you are right) Who wins the Euro 2024 championship?