

Mining Massive Datasets — Final Exam

1	2	3	4	5	6	7	8	9	10	11	Σ
/8	/10	/11	/10	/8	/4	/8	/6	/4	/4	/12	/48

Do not write anything above this line

Name:

Student ID:

Signature:

- Only write on the sheets given to you by supervisors. If you need more paper, ask the supervisors.
- Pages 14-16 can be used as scratch paper.
- All sheets (including scratch paper) have to be returned at the end.
- **Do not unstaple the sheets!**
- Wherever answer boxes are provided, please write your answers in them.
- Please write your student ID (*Matrikelnummer*) on every sheet you hand in.
- **Only use a black or a blue pen (no pencils, red or green pens!).**
- You are allowed to use your A4 sheet of handwritten notes (two sides). **No other materials (e.g. books, cell phones, calculators) are allowed!**
- Exam duration - 90 minutes.
- This exam consists of 11 pages, 7 problems. You can earn 48 points.

 Student ID:

1 Hidden Markov Models

Problem 1 [10 points] In this question, we will discuss hidden Markov models with **continuous** observations. We will use the notation from the lecture, where $Z_t \in \{1, \dots, K\}$ denotes the state at time t and $X_t \in \mathbb{R}$ denotes the observation at time t . The conditional probability of an observation at a state k is $\Pr(X_t | Z_t = k, \theta) = \mathcal{N}(X_t | \mu_k, \sigma_k^2)$, i.e. a Gaussian distribution parametrized by mean μ_k and variance σ_k^2 . θ is the set of parameters of the HMM, which includes the initial probabilities $\pi \in \mathbb{R}^K$, transition probability matrix $A \in \mathbb{R}^{K \times K}$ and the means and variances $\{\mu_1, \dots, \mu_K, \sigma_1^2, \dots, \sigma_K^2\}$.

- a) Write down the log-likelihood $\log \Pr(Z_1, \dots, Z_T, X_1, \dots, X_T | A, \pi, \mu_1, \dots, \mu_K, \sigma_1^2, \dots, \sigma_K^2)$ as a function of A, π, μ_k 's and σ_k^2 's.

$$\begin{aligned} \log \Pr(X_1, \dots, X_T, Z_1, \dots, Z_T) &= \log \left[\Pr(Z_1) \prod_{t=1}^{T-1} \Pr(Z_{t+1} | Z_t) \prod_{t=1}^T \Pr(X_t | Z_t) \right] \\ &= \log(\pi_{Z_1}) \left[\sum_{t=1}^{T-1} \log(A_{Z_t Z_{t+1}}) \right] \left[\sum_{t=1}^T \log(\mathcal{N}(X_t | \mu_{Z_t}, \sigma_{Z_t}^2)) \right] \end{aligned}$$

- b) Write the forwards and backwards update equations for this HMM. Explain in a single line how they are different from the updates we studied in class (discrete observations).

$$\begin{aligned} \alpha_t(k) &= \mathcal{N}(X_t | \mu_k, \sigma_k^2) \sum_{i=1}^K \alpha_{t-1}(i) A_{ik} \\ \beta_t(k) &= \sum_{i=1}^K A_{ki} \beta_{t+1}(i) \mathcal{N}(X_{t+1} | \mu_i, \sigma_i^2) \end{aligned}$$

The equations are similar in form. But in this case, the output probabilities are Gaussian rather than multinomial. The outputs are also continuous rather than discrete.

- c) You are given a sequence of observations $\{X_1, \dots, X_T\}$ and the corresponding states $\{Z_1, \dots, Z_T\}$. Are the maximum likelihood estimates of A_{ij} and π_i for this model different from the ones for HMM with discrete observations (that we studied in class)? Explain why or why not.

The update equations for A_{ij} and π_i are the same. They involve only the state transition counts and so are independent of the form chosen for emission probabilities.

- d) You are given a sequence of observations $\{X_1, \dots, X_T\}$ and the corresponding states $\{Z_1, \dots, Z_T\}$. Write down the closed-form maximum likelihood estimates for the parameters μ_k and σ_k^2 . You don't have to derive the closed-form MLE, stating it with an intuitive explanation is enough.

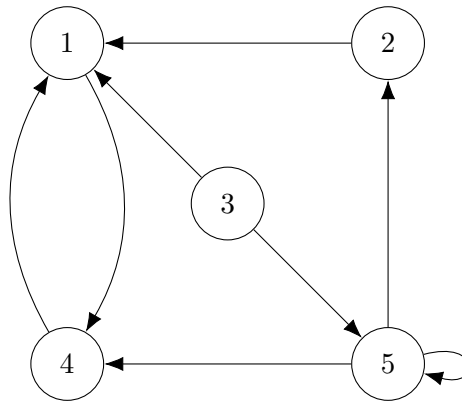
Given the states, we have K Gaussian distributions and some samples from each. Thus we can estimate the parameters of each Gaussian via sample mean and sample variance.

More precisely (optional):

$$\mu_k = \frac{\sum_{t=1}^T \mathbb{I}(Z_t = k) X_t}{\sum_{t=1}^T \mathbb{I}(Z_t = k)}$$

$$\sigma_k^2 = \frac{\sum_{t=1}^T \mathbb{I}(Z_t = k) (X_t - \mu_k)^2}{\sum_{t=1}^T \mathbb{I}(Z_t = k)}$$

2 PageRank



Problem 2 [4 points] Given the graph above and the topic-specific PageRank vector $r = [0.3662, 0.0442, 0.0800, 0.3519, 0.1578]$ using the teleport set $S = \{3, 5\}$. What is the value of β (corresponding to the teleport probability $1 - \beta$)? Justify your answer.

Node 3 is in the teleport set. Since it has no parent nodes, it can only be reached via teleportation. The value of node 3 is 0.08. Thus, the teleport probability is $1 - \beta = 2 \cdot 0.08 = 0.16$. $\beta = 0.84$.

3 Graph Clustering

Problem 3 [8 points] You are given a connected undirected unweighted graph with a set of nodes V . Let S_n^* be the partitioning minimizing the **normalized cut**, and let S_r^* be the partitioning minimizing the **ratio cut**. That is:

$$S_n^* = \arg \min_S \text{n-cut}(S) \quad S_r^* = \arg \min_S \text{r-cut}(S)$$

where $S = \{C_1, V \setminus C_1\}$ is a two-way partitioning of the nodes in the graph.

Prove or disprove that $\text{n-cut}(S_n^*) \leq \text{r-cut}(S_r^*)$.

Hint: Think about the values of the cuts when we plug in the same partitioning S .

Let $S = \{C_1, V \setminus C_1\}$ be any two way partitioning, and let $c = \text{cut}(C_1, V \setminus C_1)$. Then we have:

$$\text{n-cut}(S) = \frac{c}{\text{vol}(C_1)} + \frac{c}{\text{vol}(V \setminus C_1)} \leq \frac{c}{|C_1|} + \frac{c}{|V \setminus C_1|} = \text{r-cut}(S)$$

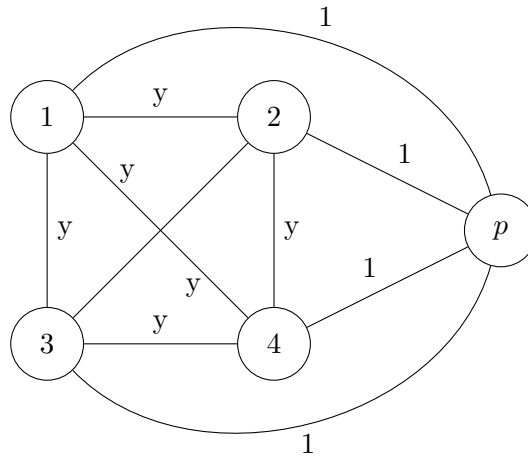
since $\text{vol}(C_1) \geq |C_1|$. Furthermore, by definition we have that $\text{n-cut}(S_n^*) \leq \text{n-cut}(S)$. Combining these facts we obtain for any S :

$$\text{n-cut}(S_n^*) \leq \text{n-cut}(S) \leq \text{r-cut}(S)$$

Finally, replacing S with S_r^* above we obtain $\text{n-cut}(S_n^*) \leq \text{r-cut}(S_r^*)$.

You are given an undirected, weighted graph with a set of nodes $V = \{1, 2, \dots, n, p\}$. It contains one clique of n nodes that is fully interconnected, each edge having weight y , and one additional node p that is connected to all nodes in the clique with weight 1.

See the example figure below for $n = 4$ as an illustration. We want to partition the nodes into **two** clusters.



- a) What is the ratio cut when one of the clusters contains only the node p ? Provide the solution as a function of y and n .

$$c_r = \frac{n}{1} + \frac{n}{n} = n + 1$$

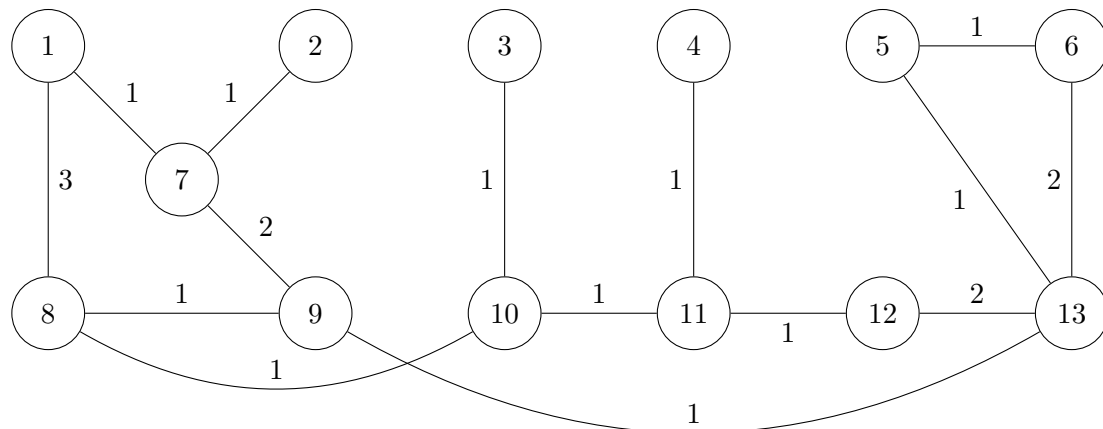
- b) What is the ratio cut if one of the clusters corresponds to $S \subseteq \{1, 2, \dots, n\}$, with $|S| = k > 0$? Provide the solution as a function of y , n , and k .

$$\begin{aligned} c_r(k, y, n) &= \frac{k(n-k)y + k}{k} + \frac{ky(n-k) + k}{n-k+1} \\ &= \frac{(n-k+1)(y(n-k) + 1) + k(y(n-k) + 1)}{n-k+1} \\ &= \frac{y(n-k+1)(n-k+1+k)}{n-k+1} \\ &= y(n+1) \end{aligned}$$

- c) When is the partitioning $\{p\}$ vs $\{1, 2, \dots, n\}$ the optimal solution for the minimum ratio cut? Specify all values of y and n , such that this holds.

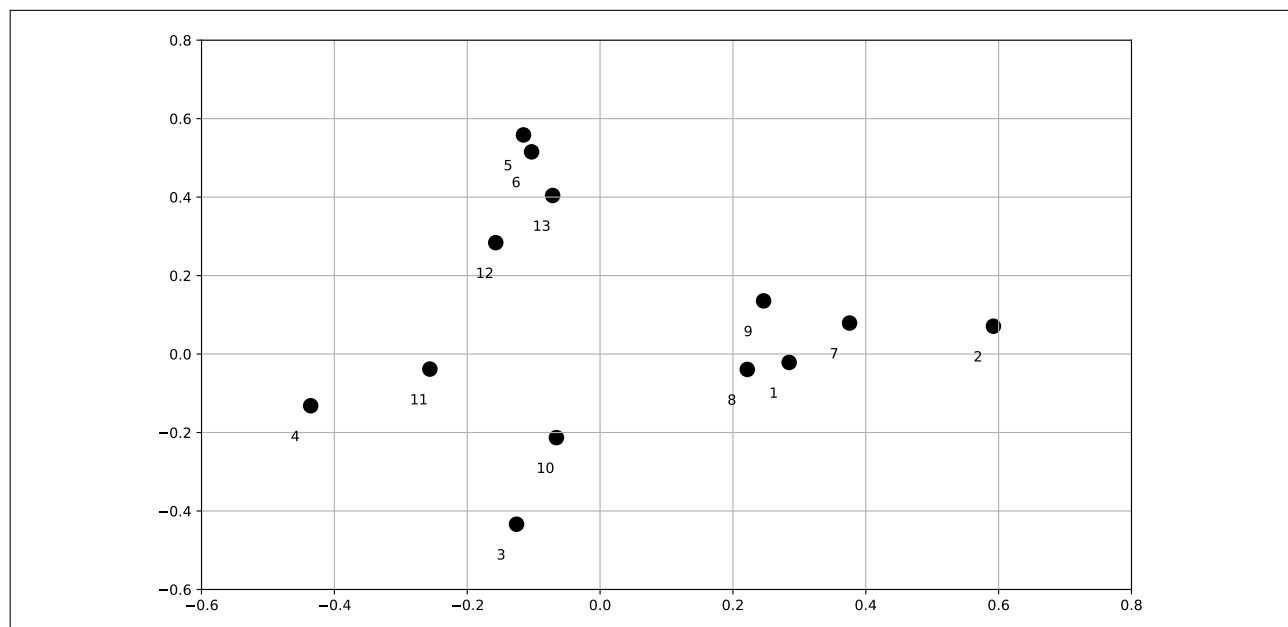
$$n + 1 > y(n + 1) \Leftrightarrow y < 1$$

Problem 4 [6 points] Spectral embedding has been applied to the following undirected weighted graph

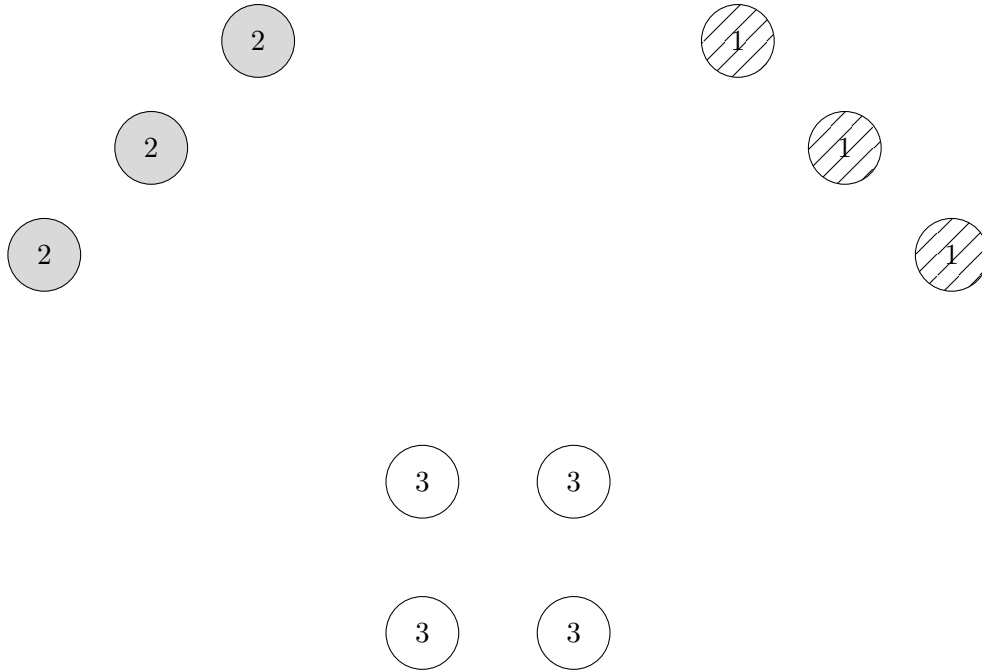


The plot below demonstrates the coordinates of the nodes in the embedding space (as defined by the second and the third eigenvectors of the unnormalized graph Laplacian).

Your task is to annotate the nodes in the plot below with their corresponding node IDs from the figure above.



Problem 5 [4 points] The graph below has been generated using the Stochastic block model with $K = 3$ communities (edges are hidden in the drawing). The community assignments of the 10 nodes in the graph are known (number inside each circle indicates the community ID z_i).



- (a) Draw the edges in the graph such that performing maximum likelihood estimation of the parameter $\boldsymbol{\eta}$ will yield

$$\boldsymbol{\eta} = \begin{bmatrix} 1 & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{4} & \frac{1}{3} & \frac{2}{3} \end{bmatrix}$$

In case you draw some edges incorrectly and need to correct your solution, please redraw the entire graph from scratch on pages 14-16. Note that there are multiple correct solutions and specifying any of them is enough.

Provide a brief explanation below.

The entries of \mathbf{B} relate to the edge counts as:

$$\mathbf{B}_{ii} = \frac{\text{\#edges within group } i}{\binom{\text{\#nodes in group } i}{2}}$$

$$\mathbf{B}_{ij} = \frac{\text{\#edges between groups } i, j}{\text{\#nodes in group } i \cdot \text{\#nodes in group } j} \quad \text{if } i \neq j$$

Therefore the number of edges between each pair of communities is

$$\mathbf{E} = \begin{bmatrix} 3 & 3 & 3 \\ 3 & 1 & 4 \\ 3 & 4 & 4 \end{bmatrix}$$

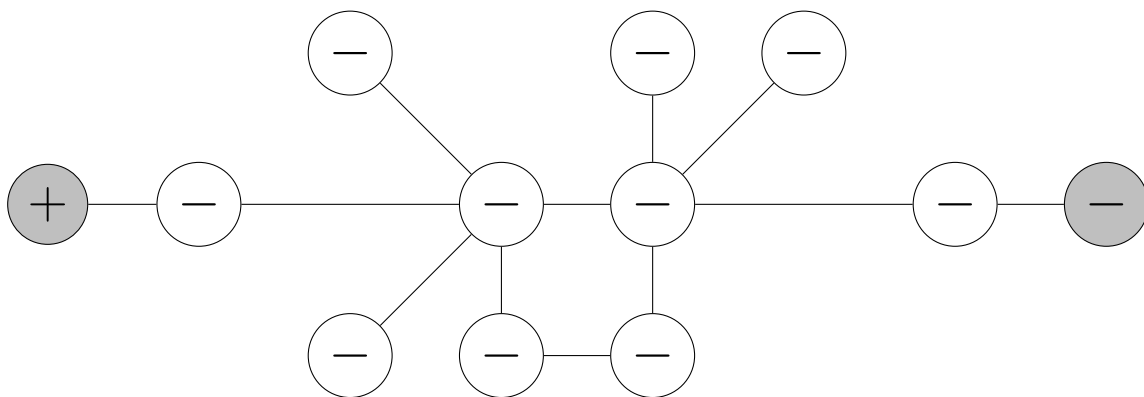
(b) What is the maximum likelihood estimate of π for the given graph and communities?

$$\pi = [3/10, 3/10, 4/10]$$

4 Label Propagation

Problem 6 [4 points] You are given the graph below with two labeled seed nodes (shaded), belonging to two classes: + and −. Specify the labels of the remaining nodes (by writing + or − inside each node) that would be obtained with an **exact** solution of the standard binary label propagation problem (i.e. assuming label smoothness).

Justify your answer. Multiple correct solutions are possible, it is enough to provide one.



Minimizing the energy \Leftrightarrow minimizing the cut. This coloring minimizes the cut.

5 Deep Learning on Graphs

Problem 7 [12 points] You are given an undirected unweighted graph G with N nodes. Your task is to instantiate the differentiable message passing framework to compute PageRank. More specifically you have to specify:

- The input features $\mathbf{x}_v \in \mathbb{R}^2$ of each node v .
- The function M that computes the message from node u to its neighbor v .
- The function U that updates the hidden representation of a node v .

such that the hidden representations in the last layer approximately recover the PageRank score of the nodes in the graph G .

Hint 1: Increasing the number of layers in this formulation should improve the approximation of the PageRank scores. Hint 2: You do not need any trainable parameters.

Let the $\mathbf{x}_v \in \mathbb{R}^2$. Now let the first coordinate of \mathbf{x}_v encode the initial value of the page rank score r_v and the second coordinate encode the degree of node d_v .

Specifically $\mathbf{x}_v = [\frac{1}{N}, d_v]$.

Since the graph is undirected the recursive page rank formula is:

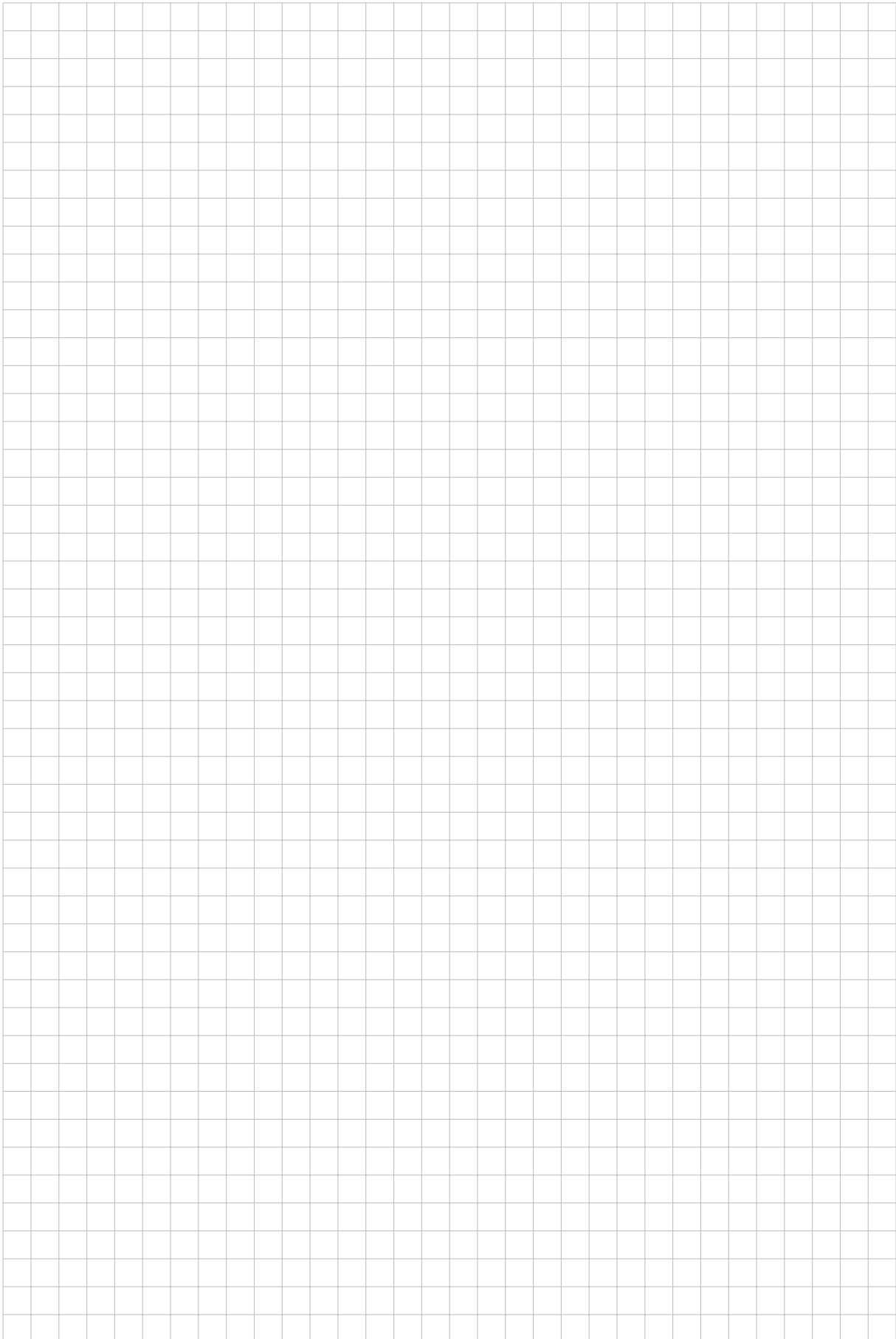
$$r_v = \sum_{u \in N(v)} \frac{r_u}{d_u}$$

Setting the function $M(h_v^{(k-1)}, h_u^{(k-1)}, E_{vu}) = \left[\frac{h_{u0}^{(k-1)}}{h_{u1}^{(k-1)}}, 0 \right]$ where we divide the first coordinate by the second coordinate and summing over the neighbors of v get the message:

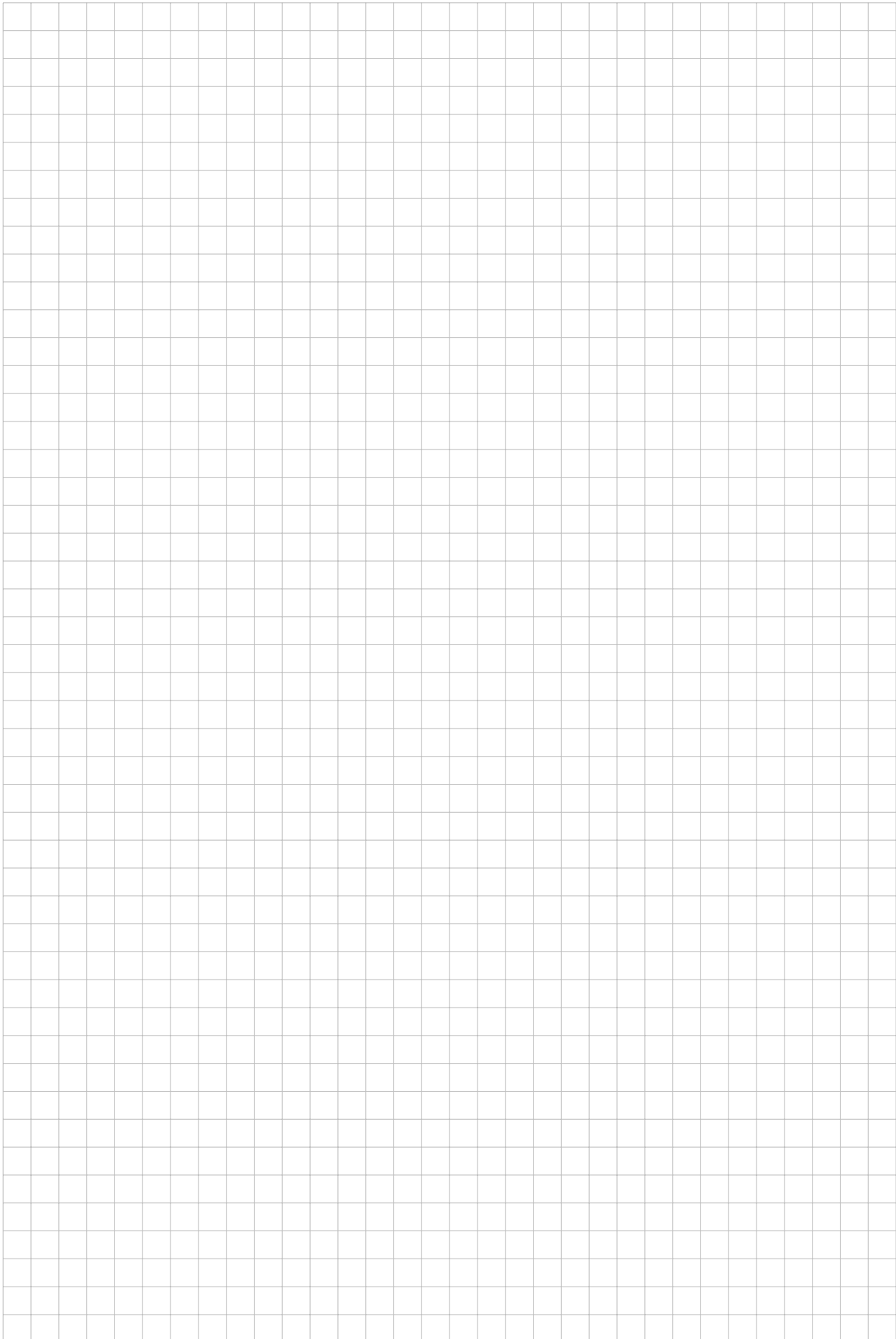
$$m_v^k = \sum_{u \in N(v)} M(h_v^{(k-1)}, h_u^{(k-1)}, E_{vu}) = \left[\sum_{u \in N(v)} \frac{h_{u0}^{(k-1)}}{h_{u1}^{(k-1)}}, 0 \right]$$

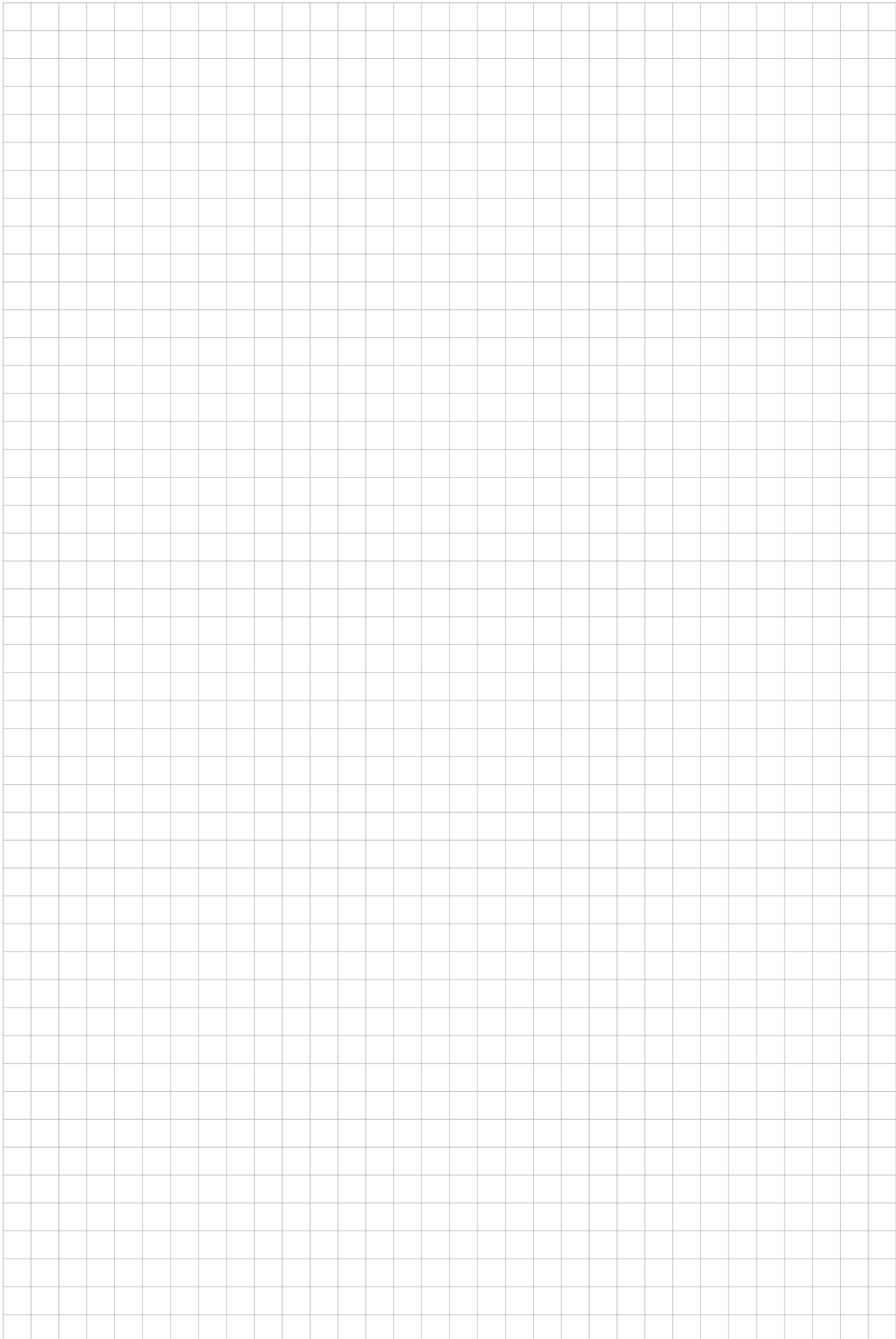
Setting the function $U(h_v^{(k-1)}, m_v^k) = \left[m_{v0}^k, h_{v1}^{(k-1)} \right]$ where the first coordinate is the first coordinate of the message and the second coordinate is the second coordinate of the previous hidden representation (to preserve the degree) gives us the complete specification.

Stacking multiple such layers is equivalent to performing power iteration.



Student ID:





Student ID: