

Q1	Q2	Q3	Q4	Total

Student Name: _____

Number: _____

Yıldız Technical University

BLM4800–Introduction to Data Mining

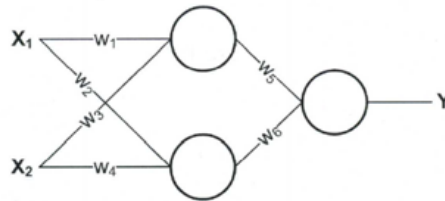
Final Exam - Spring 2022-2023

- **Duration:** 90 minutes
- **Exam information:**
 - Attempts to cheat in the exam will not be tolerated. If an attempt to cheat is discovered, it will be severely punished.
 - Read all the questions carefully before you start answering them.
 - The point value of each question is indicated next to the question.

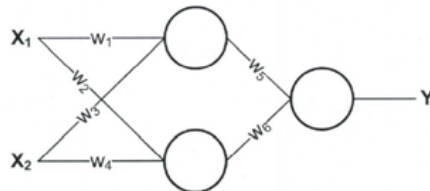
1. [20 points]. The following two-layer Neural Network estimates the target variable Y using the weights of w_2, \dots, w_6 and activation functions on bivariate data such as $X=(X_1, X_2)$. There are two activation function options;

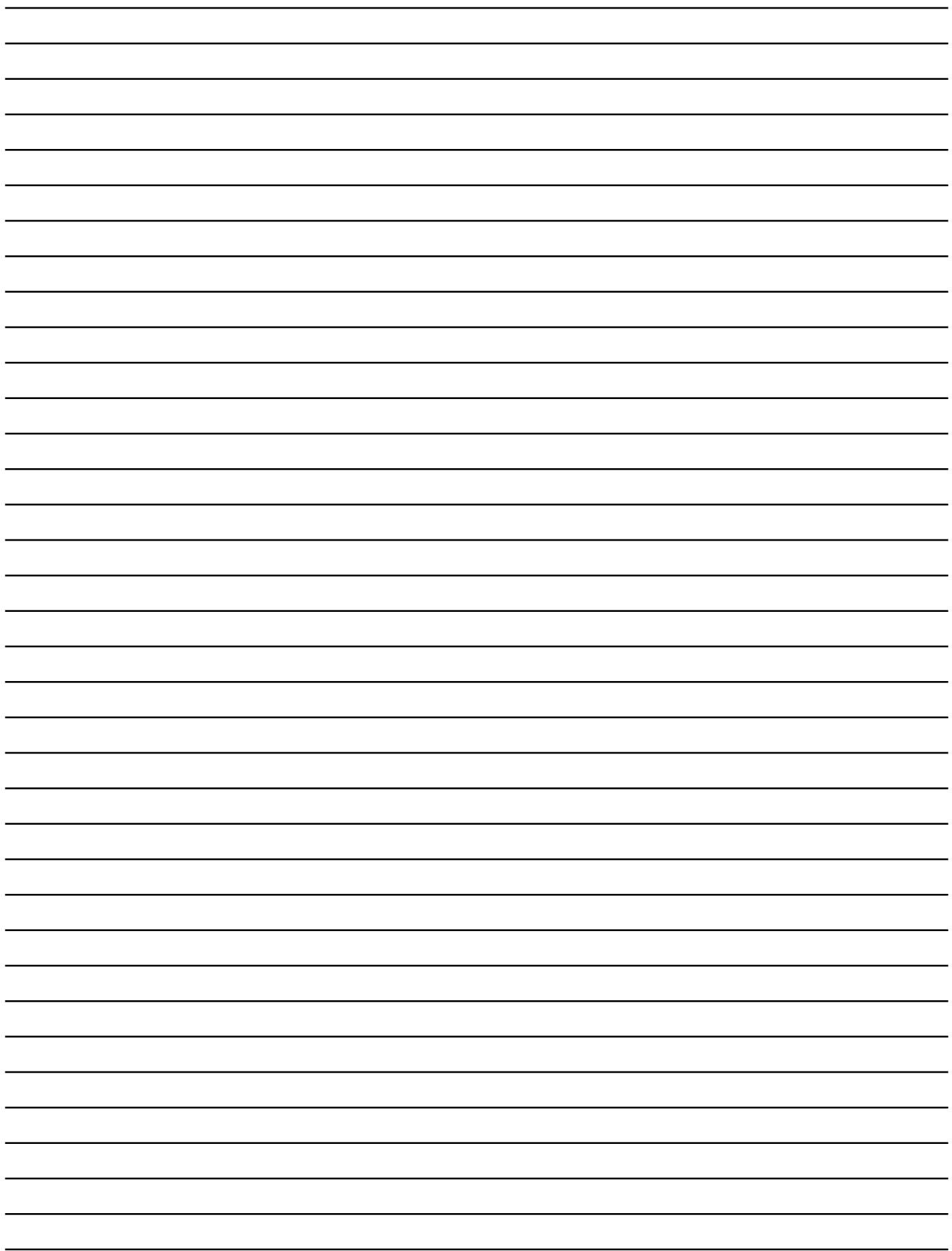
- **S:** sigmoid function $S(\alpha) = \text{sign}[\sigma(\alpha) - 0.5] = \text{sign}\left[\frac{1}{1+e^{(-\alpha)}} - 0.5\right]$
 - **L:** Linear function $L(\alpha) = c \cdot \alpha$
- In both cases $\alpha = \sum_i w_i X_i$

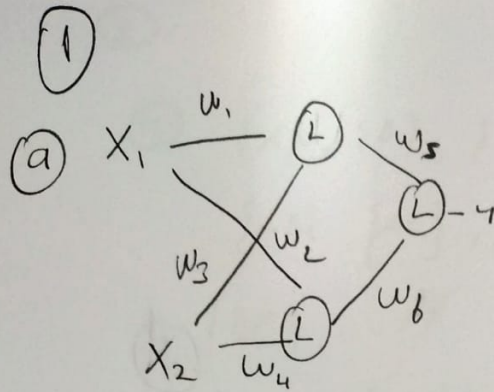
a. [10 points]. To model a linear regression of this Artificial Neural Network as $Y = \beta_1 X_1 + \beta_2 X_2$, write the appropriate activation functions in the blanks in the figure below and explain your reasoning (Write S or L inside the neurons).



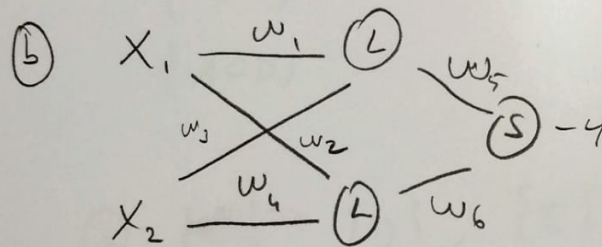
b. [10 points]. To model this Artificial Neural Network $Y = \arg \max_y P(Y=y|X)$ as a binary logistic regression $P(Y=1|X) = \left(\frac{1}{1+e^{-(\beta_1 X_1 + \beta_2 X_2)}}\right)$, Write the appropriate activation functions in the gaps in the figure and explain your reasoning. (just write S or L inside the neurons).





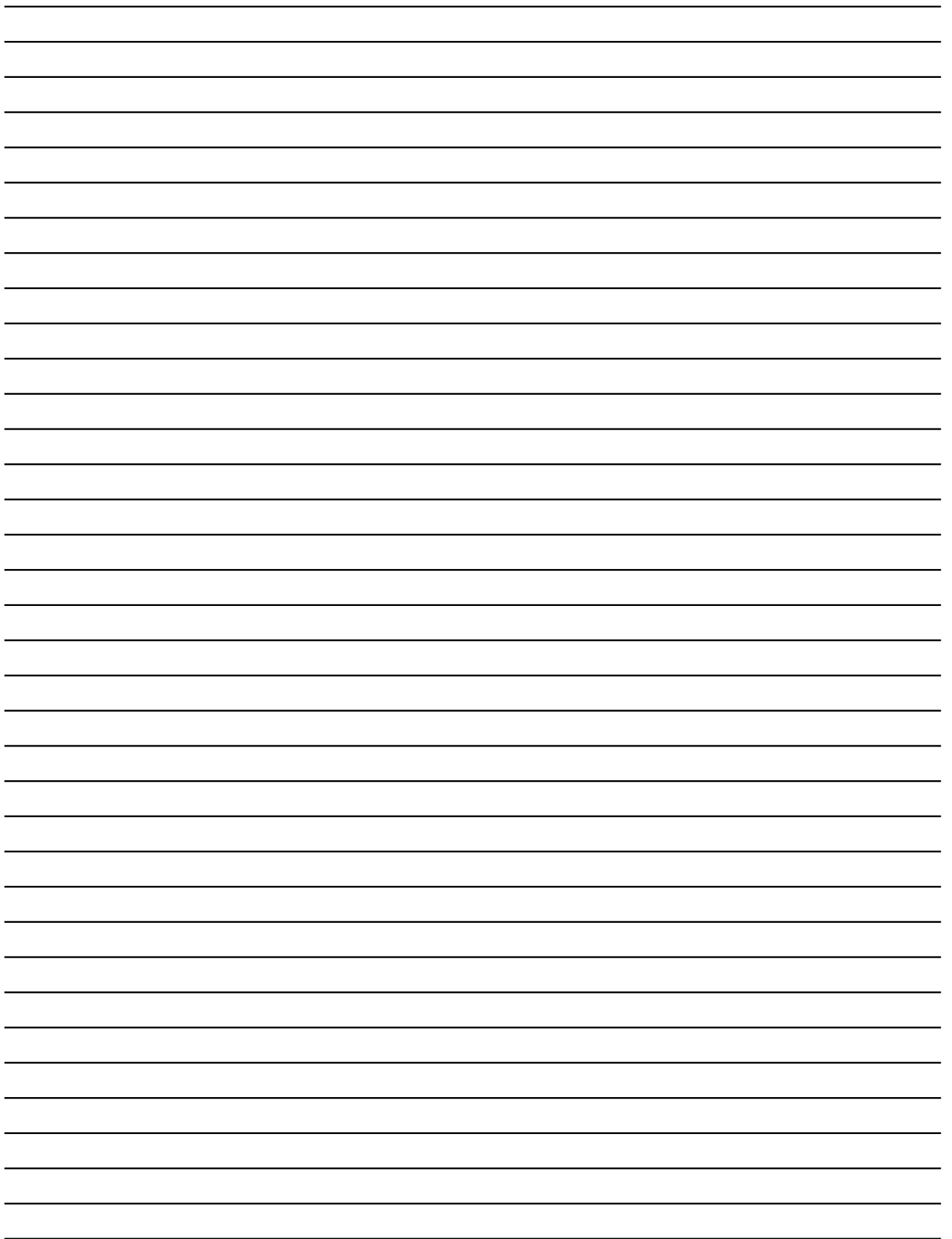


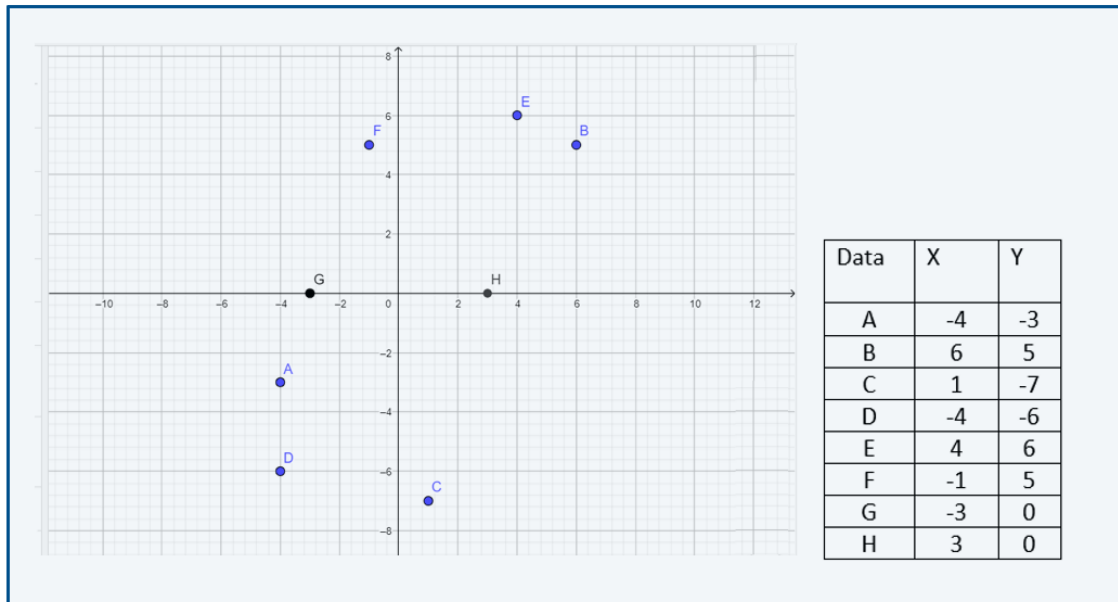
All of activation func in the model must be linear in order for the output of a NN to be linear



The provided information indicates that the desired outcome of the neural network is a binary output, achieved through the use of the sigmoid activation func. The input of the last layer follows a linear func, implying that inner layers should also consist of entirely linear activation

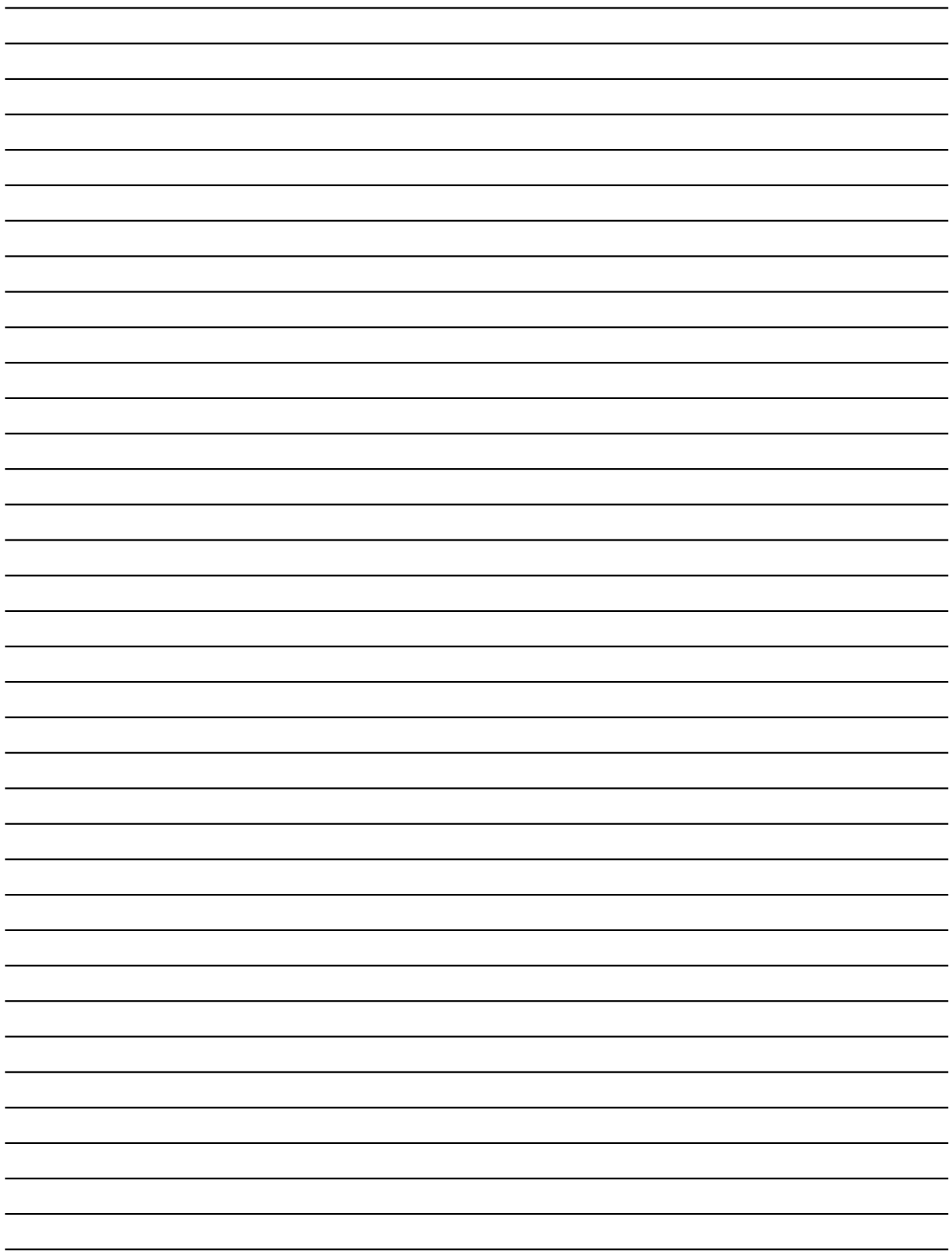
This suggests that network aims to linearly combine the input features to generate the desired binary output





2. [30 points]. Using the diagram above, answer the following questions in the context of **cluster analysis**. Suppose you are given the G data point as the initial centroid of first cluster {cluster1: G point} and the H data point as the second cluster {cluster2: H point}. Simulate the K-Means algorithm for (#k=2) assuming it uses the Euclidean distance.
- What happens to cluster assignments after ONE iteration?
 - What happens to cluster assignments after convergence? (Fill in the table below and show your calculations on the side)

Data	Cluster Assignment After First Iteration	Cluster Assignment After Second Iteration
1		
2		
3		
4		
5		
6		
7		
8		
9		
10		



(2)

(a.)

After first Loop

1. Cluster = A D F G

2. Cluster = B C E H

(b)

Second loop

$$\begin{aligned} 1. \text{Centroid} &= \left(\frac{(A_x + D_x + F_x + G_x)}{4}, \frac{(A_y + D_y + F_y + G_y)}{4} \right) \\ &= \frac{-4 - 4 - 1 - 3}{4}, \frac{-3 - 6 + 5 + 0}{4} = [-3, -1] \end{aligned}$$

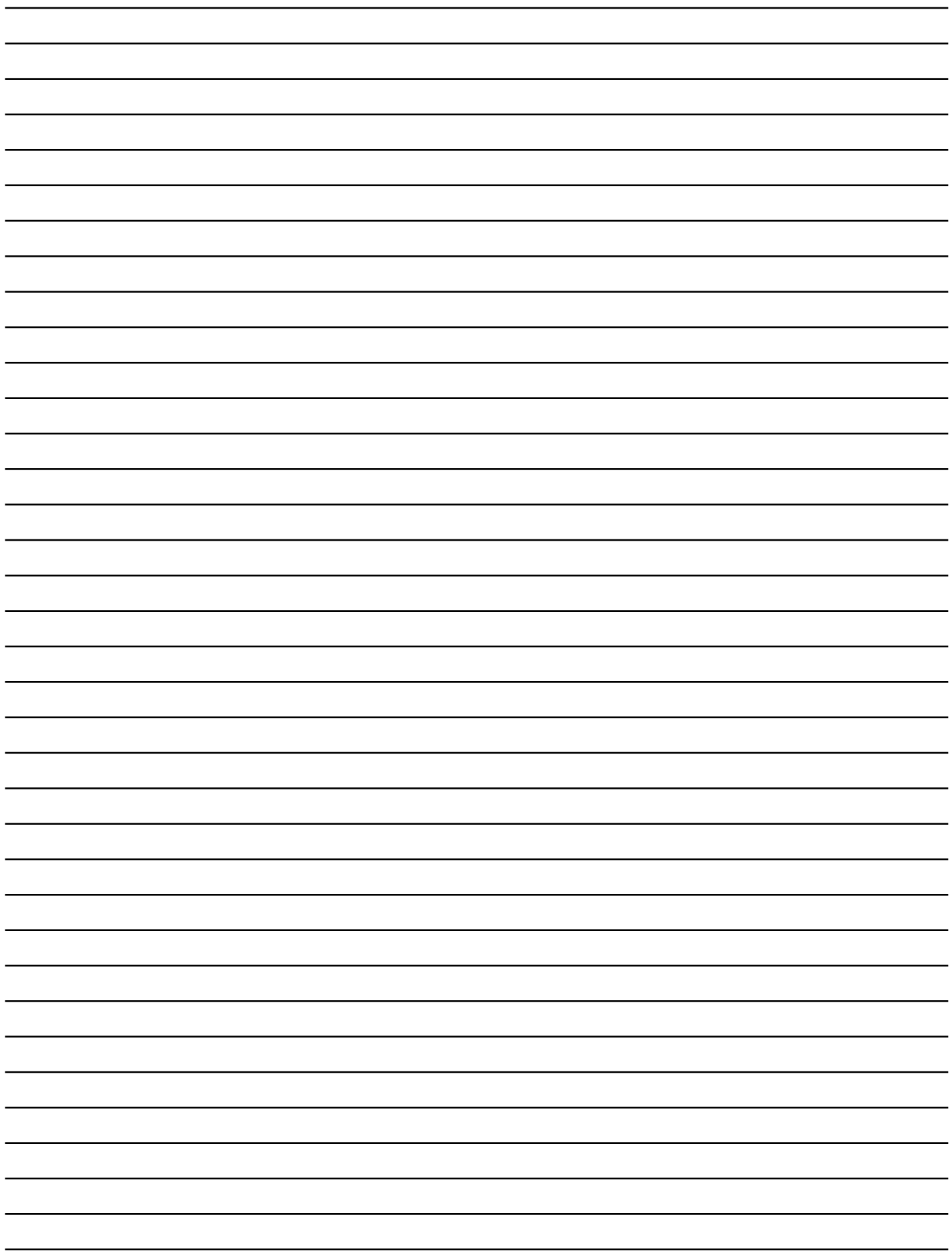
$$\begin{aligned} 2. \text{Centroid} &= \left(\frac{(B_x + C_x + E_x + H_x)}{4}, \frac{(B_y + C_y + E_y + H_y)}{4} \right) \\ &= \left(\frac{6 + 1 + 3 + 4}{4}, \frac{(5 + 6 + 7)}{4} \right) = [3.5, 7] \end{aligned}$$

So After second Loop

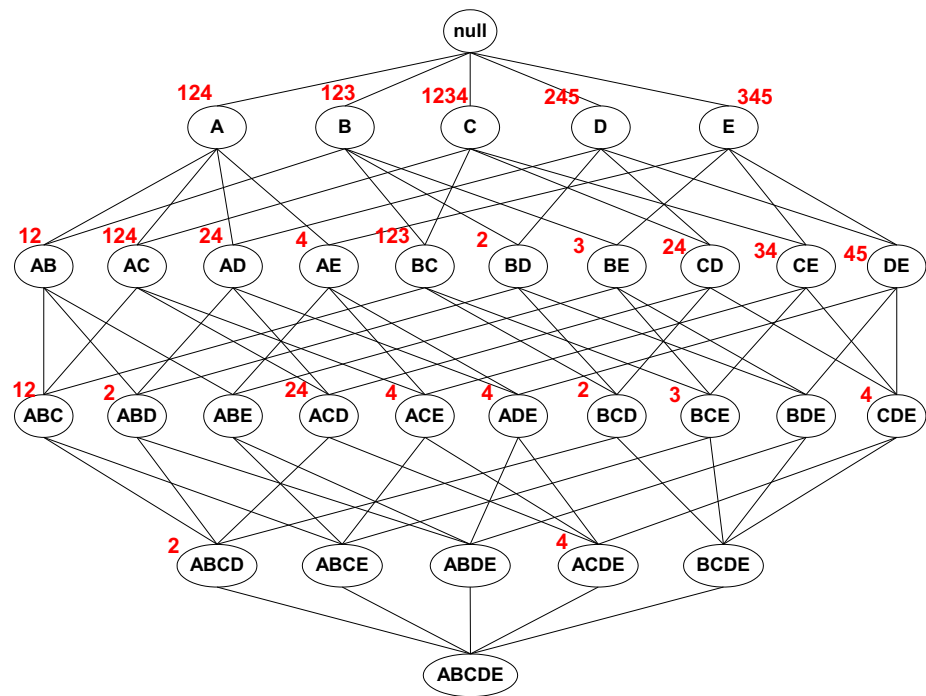
1. Cluster = A C G D

2. Cluster = B F E H

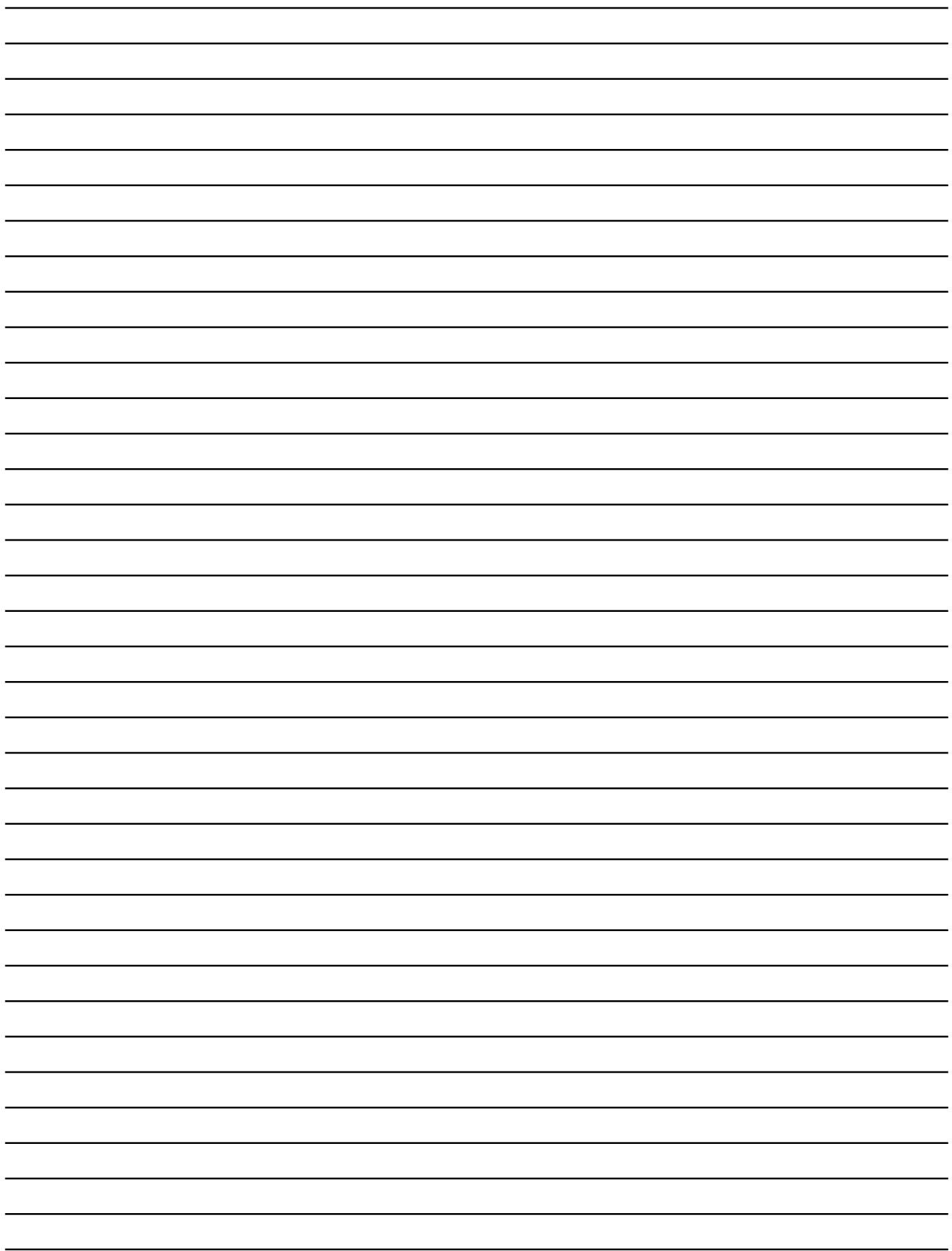
<u>Date</u>	<u>1.</u>	<u>2.</u>
1	6	6
2	4	4
3	4	6
4	6	6
5	4	4
6	6	4
7	6	6
8	4	4



ID	Products
1	ABC
2	ABCD
3	BCE
4	ACDE
5	DE



3. [26 points]. The diagram above shows the transaction number and the list of products sold within the transactions. Using this diagram, answer the following questions. Use the minimum support threshold value of 2 for these questions.
- [6 points]. Using the diagram, write the list of all Frequent itemsets below.
 - [10 points]. “An item set is a Closed Frequent itemset if none of its nearest supersets has the same support as the item set.” Using the diagram, write the Closed Frequent itemset list of all items below.
 - [10 points]. “An item set is called a Maximal Frequent itemset if none of its immediate supersets are frequent.” Using the diagram, write the Maximal Frequent itemset list of all items below.



③

①. $\{A\}, \{B\}, \{C\}, \{D\}, \{E\}$

$\{AB\}, \{AC\}, \{AD\}$

$\{BC\},$

$\{CD\}, \{CE\}$

$\{DE\}$

$\{ABC\}$

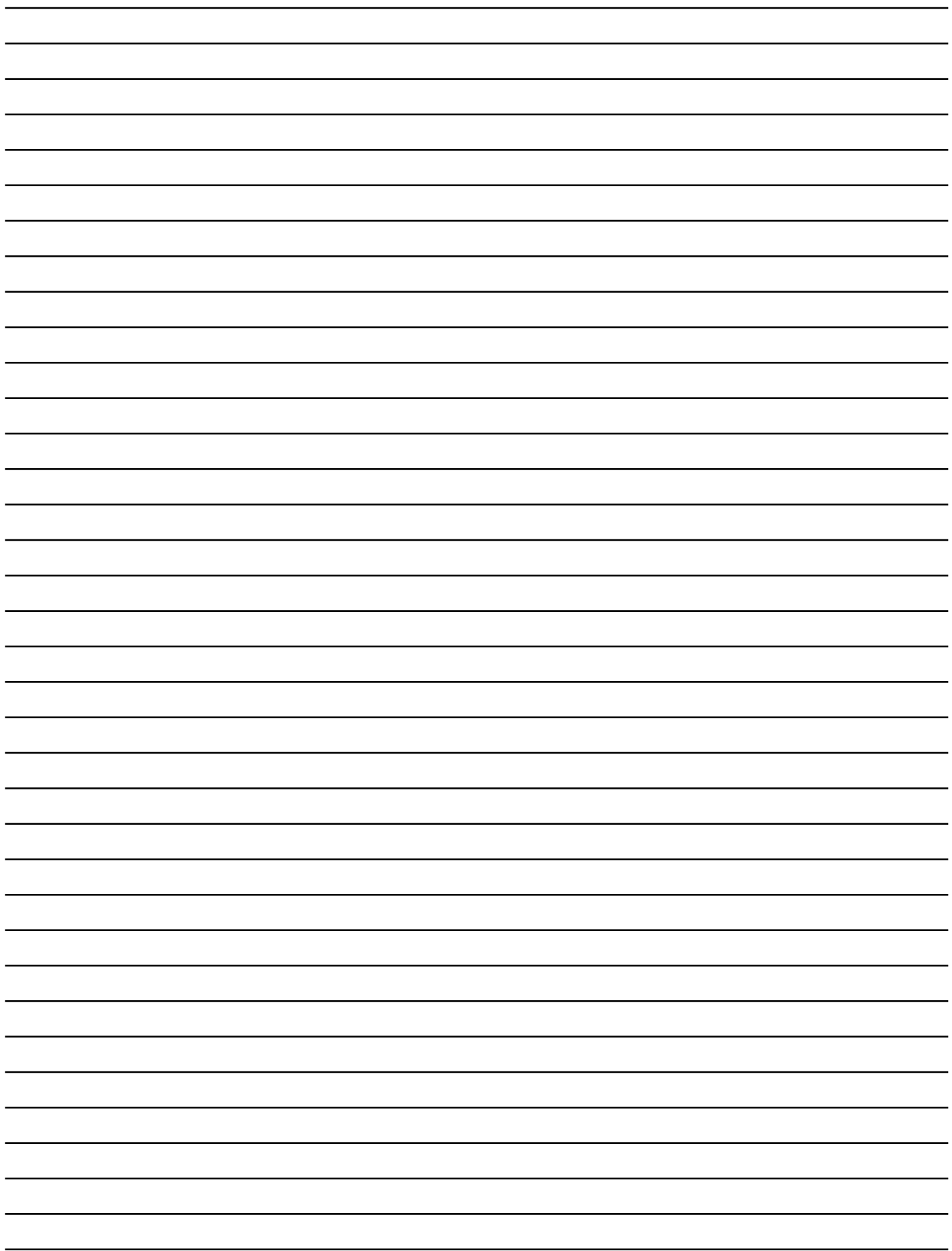
$\{ACD\}$

② $\{C\}, \{D\}, \{E\}$

$\{AC\}, \{BC\}, \{CE\}, \{DE\},$

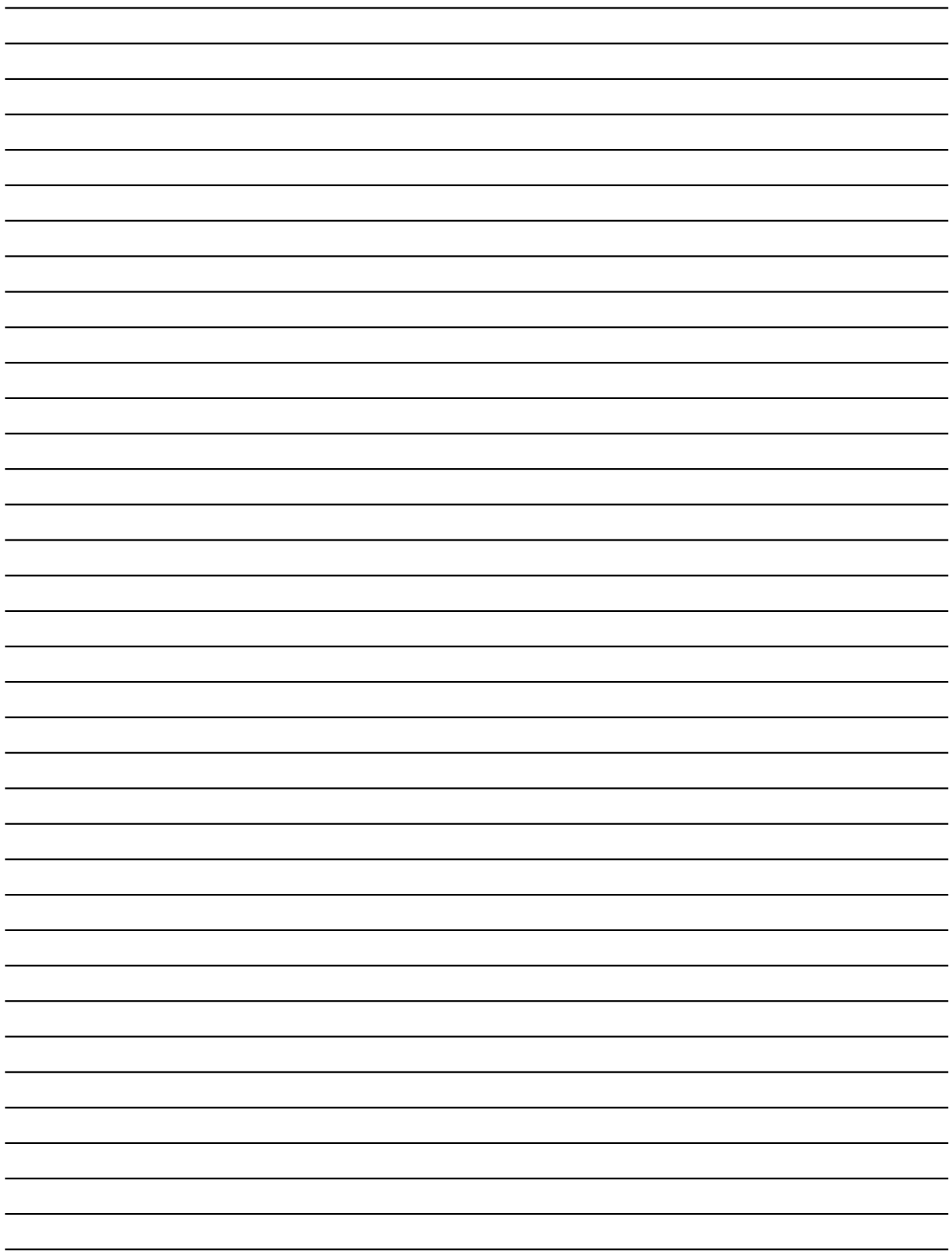
$\{ABC\}, \{ACD\}$

③ $\{CE\}, \{DE\}, \{ABC\}, \{ACD\}$



Data	X	Y
1	126	78
2	128	80
3	128	82
4	130	82
5	130	84
6	132	86

4. [24 points]. Using the data above, answer the following questions in the context of principal component analysis for data reduction. You want to reduce the data to a single size. Assume that the first Principal component is given as (0.59, 0.81). According to this;
- a. [12 points]. What is the corresponding projection point on the First Principal component line for data point #3 ($x=128$, $y=82$)?
- b. [12 points]. Assume that the second Principal component is given as (-0.81, 0.59). What is the corresponding projection point on the second Principal component line for data point #5 ($x=130$, $y=84$)?



- ④ We have 2D (X, Y) and want to reduce to 1D
by projecting it onto the 1. principal Component
with given $(0.59, 0.81)$

(X, Y) — vector (u, v)
use formula dot product

1. PCA = dividing the PC vector by magnitude

$$\text{Magnitude} = \sqrt{(0.59)^2 + (0.81)^2} \approx 1$$

$$\begin{aligned}\text{Product Dot Product} &= (128 \times 0.59) + (82 \times 0.81) \\ &= 75.52 + 66.42 \\ &= 141.94\end{aligned}$$

⑤ Dot product = $(130 \times -0.81) + (94 \times 0.59)$

$$\begin{aligned}&= -105.3 + 55.46 \\ &= -49.84\end{aligned}$$

