

# CS 181 Artificial Intelligence (Fall 2024), Final Exam

## Instructions

- Time: 10:30 – 12:30 (120 minutes)
- This exam is closed-book, but you may bring one A4-size cheat sheet. Put all the study materials and electronic devices (except a calculator) into your bag and put your bag in the front, back, or sides of the classroom.
- Two blank pieces of paper are attached, which you can use as scratch paper. Raise your hand if you need more paper.
- For multiple choice questions in Question 2–5:
  - ☐ means you should mark ALL choices that apply;
  - ☐ means you should mark exactly ONE choice;
  - When marking a choice, please fill in the bubble or square COMPLETELY (e.g., ☒ and ☒). Ambiguous answers will receive no points.
  - For each question with ☐ choices, you get half of the points for selecting a non-empty proper subset of the correct answers.

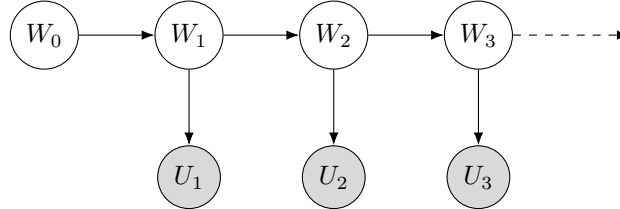
## 1 Multiple choices (10 pt)

Each question has one or more correct answers. Select all the correct answers. For each question, you get 1 point if you select all the correct answers and nothing else, 0 point if you select one or more wrong answers, and 0.5 point if you select a non-empty proper subset of the correct answers.

1	2	3	4	5
6	7	8	9	10

1. Markov Models are based on the Markov assumption. Consider the first-order Markov models with states  $X_1, X_2, X_3, \dots$ , which of the following statement(s) is/are correct?
  - A. The Markov assumption states that  $P(X_{t+1}|X_{1:t}) = P(X_{t+1}|X_t)$ .
  - B. The Markov assumption states that  $P(X_{t+1}|X_{1:t}) = P(X_{t+1}|X_{t-1})$ .
  - C. Based on the Markov assumption, the past states  $X_{1:t-1}$  are independent of the future states  $X_{t+1:T}$  given the current state  $X_t$ .
  - D. Based on the Markov assumption, the future state  $X_{t+1}$  is independent of the current state  $X_t$  given the past states  $X_{1:t-1}$ .

2. Consider the weather Hidden Markov Model (HMM) described in class. The weather state  $W_t$  can be either sunny or rainy, and the observation  $U_t$  can be either umbrella or no umbrella. Given the observations of  $U_1 = \text{umbrella}$ ,  $U_2 = \text{umbrella}$ ,  $U_3 = \text{no umbrella}$ , we want to find the most likely weather state sequence  $W_{0:3}$ . Which ONE of the following algorithms is the most suitable for this task?



- A. Forward algorithm
  - B. Viterbi algorithm
  - C. Gradient descent
  - D. Gibbs sampling
3. Which of the following is NOT a required element in the formal definition of a Markov Decision Process (MDP)?
- A. A set of states  $\mathcal{S}$  representing all possible configurations of the environment.
  - B. A set of actions  $\mathcal{A}$  representing all possible actions an agent can take.
  - C. A reward function  $R(s, a, s')$  that assigns a reward to each transition from state  $s$  to state  $s'$  via action  $a$ .
  - D. An exploration strategy  $\epsilon$ -greedy that balances exploration and exploitation.
4. Which of the following statement(s) is/are correct?
- A. In policy iteration, we do policy improvement by doing policy extraction from the values computed in policy evaluation.
  - B. The goal of an MDP is to maximize the single-step reward.
  - C. In MDP, the successor state depends on not only the current state, but also the previous states.
  - D. Negative rewards do not exist in MDP.
5. Which of the following statements about the exploration-exploitation tradeoff in Q-learning is *false*?
- A. The  $\epsilon$ -greedy strategy ensures that the agent continues to explore new actions before  $\epsilon$  becomes 0, even when it has found a good policy.
  - B. As  $\epsilon$  decreases over time, the agent will prioritize exploiting the learned Q-values over exploring new actions.
  - C.  $\epsilon$ -greedy always chooses the optimal action.
  - D. Reducing  $\epsilon$  too quickly can lead to suboptimal learning because the agent may stop exploring before finding the optimal policy.

6. Which of the following is necessary for ensuring that Q-learning converges to the optimal Q-values?
  - A. The learning rate  $\alpha$  is expected to decrease over time and eventually become sufficiently small.
  - B. The exploration policy should always select the action with the highest Q-value in every state (purely greedy policy).
  - C. The discount factor  $\gamma$  should always be set to 1 to ensure the algorithm prioritizes future rewards equally with immediate rewards.
  - D. The samples to update the Q-function require to be collected by the optimal policy.
7. You are developing two machine learning models to predict house prices based on various features such as size, location, number of bedrooms, and age of the property. After training, you observe that the model *A* performs exceptionally well on the training data but poorly on the validation set. Additionally, a simpler model *B* shows consistently poor performance on both training and validation data. Based on provided information, which of the following statement(s) is/are true?
  - A. The model *A* is overfitting on the training data and it cannot generalize well.
  - B. The model *B* is overfitting on the validation data.
  - C. For model *A*, it is possible to improve the performance on validation data by adding more training data that is representative of the true data distribution.
  - D. For model *B*, it is possible to improve the performance on the same training data by redesigning the model to be more expressive.
8. You are training a logistic regression model using gradient ascent to maximize the log-likelihood function  $l(\theta)$  of your data, where  $\theta$  is the set of model parameters to optimize. Which of the following statement(s) is/are true?
  - A. If the learning rate is set too high, gradient ascent may fail to converge to the maximum.
  - B. Gradient ascent is guaranteed to find the global maximum of  $l(\theta)$ .
  - C. The output of a logistic regression model is not always in  $(0, 1)$  (between 0 and 1).
  - D. Your logistic regression model can also be used for classification.
9. Which of the following statement(s) about self-attention mechanisms in language models is/are correct?
  - A. Self-attention mechanisms try to help the tokens pay more attention to relevant tokens.
  - B. The attention weights are determined by the similarities between the query vectors and the value vectors.
  - C. Causal attention restricts each token to compute attention solely with the tokens before it (include itself) in the sequence.
  - D. Self-attention mechanism only works for sequences of fixed length.
10. Which of the following statement(s) about unsupervised learning is/are correct?
  - A. GMM assumes that data points are generated from a mixture of a finite number of Gaussian distributions.
  - B. Unsupervised learning does not require manual labeling of data.
  - C. The EM algorithm is suitable for solving the parameter estimation problem of Gaussian mixture model (GMM).
  - D. The final result of the K-means algorithm will not be affected by the choice of the initial cluster centers.

**Solution:**

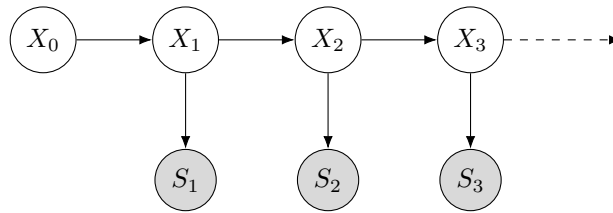
1. AC
2. B
3. D
4. A
5. C
6. A
7. ACD
8. AD
9. AC
10. ABC

## 2 Sleepy Students (10 pt)

A professor wants to know if students are getting enough sleep. Each day, the professor observes whether the students sleep in class. The professor has the following domain theory:

- The prior probability of getting enough sleep, with no observations, is 0.7.
- The probability of getting enough sleep on night  $t$  is 0.8 given that the student got enough sleep the previous night, and 0.3 if not.
- The probability of sleeping in class is 0.1 if the student got enough sleep, and 0.3 if not.

Let  $X_t$  denotes the state of getting enough sleep on night  $t$ , and  $S_t$  denotes the state of sleeping in class on day  $t$ . We have the following graphical model:



$X_0$	$P(X_0)$
0	0.3
1	0.7

$X_t$	$X_{t+1}$	$P(X_{t+1} X_t)$
0	0	0.7
0	1	0.3
1	0	0.2
1	1	0.8

$X_t$	$S_t$	$P(S_t X_t)$
0	0	0.7
0	1	0.3
1	0	0.9
1	1	0.1

### 2.1 Filtering (3 pt)

Given the observations  $s_1 = 1$ ,  $s_2 = 0$ , what is the probability that the student got enough sleep on night 2? That is, what is  $P(X_2 = 1|s_{1:2})$ ?

- ☐ 0.11
 ☐ 0.33
 ☐ 0.55
 ☐ 0.77
 ☐ 0.99

**Solution:**

C

### 2.2 Particle Filtering (5 pt)

Consider 5 samples in  $t = 0$  in the particle filtering algorithm:  $X_0^{(1)} = 0$ ,  $X_0^{(2)} = 1$ ,  $X_0^{(3)} = 0$ ,  $X_0^{(4)} = 1$ ,  $X_0^{(5)} = 1$ .

#### 2.2.1 Propagate Forward (2 pt)

What are the values of  $X_1^{(1)}, X_1^{(2)}, X_1^{(3)}, X_1^{(4)}, X_1^{(5)}$  after the propagate forward step?

To generate random samples, use as many values as needed from the table below, which we generated independently and uniformly at random from 0 to 1. Use numbers from left to right. To sample a binary variable  $W$  with probability  $P(W = 0) = p$ , select a value  $\alpha$  from the table, and choose  $W = 1$  if  $\alpha \geq p$  and  $W = 0$  otherwise.

We have done the fifth sample for you as an example.

0.123	0.722	0.830	0.476	0.261	0.593
-------	-------	-------	-------	-------	-------

- |                |                                    |                         |
|----------------|------------------------------------|-------------------------|
| 1. $X_1^{(1)}$ | <input type="radio"/> 1            | <input type="radio"/> 0 |
| 2. $X_1^{(2)}$ | <input type="radio"/> 1            | <input type="radio"/> 0 |
| 3. $X_1^{(3)}$ | <input type="radio"/> 1            | <input type="radio"/> 0 |
| 4. $X_1^{(4)}$ | <input type="radio"/> 1            | <input type="radio"/> 0 |
| 5. $X_1^{(5)}$ | <input checked="" type="radio"/> 1 | <input type="radio"/> 0 |

**Solution:**

1.  $X_1^{(1)} = 0$
2.  $X_1^{(2)} = 1$
3.  $X_1^{(3)} = 1$
4.  $X_1^{(4)} = 1$

### 2.2.2 Observe (2 pt)

Suppose that the observation is  $s_1 = 1$ . What is the weight of the sample  $X_1^{(5)}$  after the observe step?

- ☐ 0      ☐ 0.1      ☐ 0.286      ☐ 0.3      ☐ 0.714      ☐ 1

**Solution:**

B

### 2.2.3 Resample (1 pt)

Suppose after observing step of particle filtering at  $t = 2$ , the particles and its weight are as follow:

Particle	$X_2^{(1)}$	$X_2^{(2)}$	$X_2^{(3)}$	$X_2^{(4)}$	$X_2^{(5)}$
State	0	1	0	1	1
Weight	0.7	0.9	0.7	0.9	0.9

Given the particles and weights after the observe step, what is the weighted sample distribution  $P'(X_2)$  used in the resampling step? Select the correct value for  $P'(X_2 = 0)$ .

- ☐ 0      ☐ 0.134      ☐ 0.341      ☐ 0.413      ☐ 0.431      ☐ 1

**Solution:**

C

### 2.3 The Argument (2 pt)

The professor prefers the filtering algorithm to the particle filtering algorithm. He also points out that the particle filtering algorithm does not suit the problem of sleepy students. Do you agree with the professor? Please choose the most reasonable explanation from the following options.

- ☐ Yes. Since the state space is very small in this problem, it defeats the motivation of using the particle filtering algorithm.
- ☐ Yes. Since the state space is discrete in this problem, the particle filtering algorithm may produce incorrect results.
- ☐ No. Particle filtering is more efficient than filtering in most cases, including this one.
- ☐ No. Particle filtering could provide more accurate results than filtering in this problem.

**Solution:**

A

### 3 MDP and RL(10 pt)

#### 3.1 MDP (4 pt)

Consider the MDP with transition model and reward function as given in the table below. Assume the discount factor  $\gamma = 1$ , i.e., no discounting.

$s$	$a$	$s'$	$T(s, a, s')$	$R(s, a, s')$
$A$	1	$A$	0	0
$A$	1	$B$	1	0
$A$	2	$A$	1	1
$A$	2	$B$	0	0
$A$	3	$A$	0.5	0
$A$	3	$B$	0.5	0

$s$	$a$	$s'$	$T(s, a, s')$	$R(s, a, s')$
$B$	1	$A$	0.5	10
$B$	1	$B$	0.5	0
$B$	2	$A$	1	0
$B$	2	$B$	0	0
$B$	3	$A$	0.5	2
$B$	3	$B$	0.5	4

(a)  $V_k(s)$  is the expected value of starting in state  $s$  if the game ends in  $k$  more time steps, i.e.,

$$V_{k+1}(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')] .$$

Let initial  $V_0(A) = 0$ ,  $V_0(B) = 0$ . Fill in the values for  $V_1(A)$  and  $V_1(B)$  after one value iteration:

$$V_1(A) = \underline{\hspace{2cm}} \quad , \quad V_1(B) = \underline{\hspace{2cm}}$$

**Solution:**

1 5

(b) Let  $\pi_k^*(s)$  be the optimal action in state  $s$  if the game ends in  $k$  more time steps, i.e.,

$$\pi_{k+1}^*(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')] .$$

Fill in the following tables (fill in 1, 2 or 3):

$s$	$\pi_1^*(s)$
$A$	
$B$	

**Solution:**

2 1

#### 3.2 Reinforcement learning (6 pt)

Consider an unknown MDP with three states ( $A, B$  and  $C$ ) and two actions ( $\leftarrow$  and  $\rightarrow$ ). Suppose the agent chooses actions according to some policy  $\pi$  in the unknown MDP, collecting a dataset consisting of samples  $(s, a, s', r)$  representing taking action  $a$  in state  $s$  resulting in a transition to state  $s'$  and a reward of  $r$ .

$s$	$a$	$s'$	$r$
$A$	$\rightarrow$	$B$	2
$C$	$\leftarrow$	$B$	2
$B$	$\rightarrow$	$C$	-2
$A$	$\rightarrow$	$B$	4

You may assume a discount factor of  $\gamma = 1$ .

Assume that all  $Q$ -values are initialized to 0, and use a learning rate of  $\alpha = \frac{1}{2}$ .



(a) Run Q-learning on the above experience table and fill in the following  $Q$ -values:

$$Q(A, \rightarrow) = \underline{\hspace{2cm}}, \quad Q(B, \rightarrow) = \underline{\hspace{2cm}}$$

**Solution:**

2.5 -0.5

(b) Use the empirical frequency count model-based reinforcement learning method described in lectures to estimate the transition function  $\hat{T}(s, a, s')$  and reward function  $\hat{R}(s, a, s')$ . (Do not use pseudocounts; if a transition is not observed, it has a count of 0.) Write down the following quantities. You may write N/A for undefined quantities.

$$\hat{T}(A, \rightarrow, B) = \underline{\hspace{2cm}} \quad \hat{R}(A, \rightarrow, B) = \underline{\hspace{2cm}}$$

$$\hat{T}(B, \rightarrow, A) = \underline{\hspace{2cm}} \quad \hat{R}(B, \rightarrow, A) = \underline{\hspace{2cm}}$$

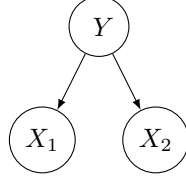
**Solution:**

1 3  
0 N/A

## 4 Supervised Machine Learning (10 pt)

### 4.1 Naïve Bayes (5 pt)

Answer questions about the following Naïve Bayes model with label  $Y \in \{0, 1\}$  and features  $X_1, X_2 \in \{0, 1\}$ .



#### 4.1.1 Play with Naïve Bayes (2 pt)

Consider the training data below:

data	1	2	3	4	5	6	7	8	9	10
$X_1$	1	1	1	0	1	1	1	0	1	0
$X_2$	1	0	1	0	1	1	1	0	0	1
$Y$	0	1	1	1	0	0	1	1	1	1

Maximum Likelihood Estimation (MLE) is a fundamental approach for parameter estimation within statistical models. In real world, Laplace Smoothing is commonly employed. With a smoothing strength of  $k = 2$ , determine the probabilities for the following scenarios using both MLE ( $P_{\text{MLE}}$ ) and Laplace Smoothing ( $P_{\text{LAP}}$ ). (If necessary, your answers can be in fraction)

(1)  $P_{\text{MLE}}(X_2 = 1|Y = 0) =$

(2)  $P_{\text{LAP}}(X_2 = 1|Y = 0) =$

**Solution:**

1,  $\frac{5}{7}$

#### 4.1.2 Naïve Bayes Prediction (2 pt)

$X_1$	$Y$	$P(X_1 Y)$
0	0	$\theta_1$
1	0	$1 - \theta_1$
0	1	0.5
1	1	0.5

$Y$	$P(Y)$
0	0.5
1	0.5

$X_2$	$Y$	$P(X_2 Y)$
0	0	$1 - \theta_2$
1	0	$\theta_2$
0	1	0.5
1	1	0.5

Table 1: CPTs for the Naïve Bayes

Now given the above CPTs, we would like to make predictions on a new sample with  $X_1 = 0$  and  $X_2 = 1$ . What is the probability of  $Y = 0$  given  $X_1 = 0$  and  $X_2 = 1$  ( $P(Y = 0|X_1 = 0, X_2 = 1)$ )?

- ☐  $\frac{\theta_1(1 - \theta_2)}{\theta_1(1 - \theta_2) + \frac{1}{4}}$
- ☐  $\frac{\theta_2(1 - \theta_1)}{\theta_2(1 - \theta_1) + \frac{1}{4}}$
- ☐  $\frac{\theta_1(1 - \theta_2)}{\theta_1(1 - \theta_2) + \frac{1}{8}}$
- ☐  $\frac{\theta_1\theta_2}{\theta_1\theta_2 + \frac{1}{4}}$

**Solution:**

D.  $\frac{\theta_1\theta_2}{\theta_1\theta_2 + \frac{1}{4}}$

#### 4.1.3 Naive Bayes Classifier (1 pt)

Which of the following statement(s) about the Naive Bayes Classifier is/are **True**?

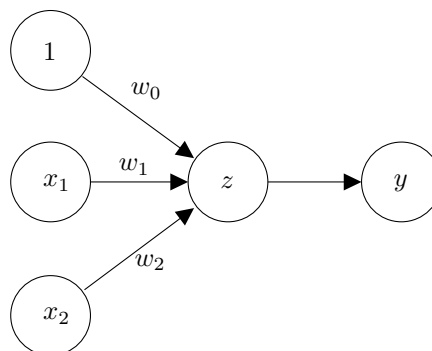
- ☐ Naive Bayes assumes that all features are dependent given the class.
- ☐ Laplace Smoothing can be used to handle zero probability issues in Naive Bayes.
- ☐ Naive Bayes can be used for multiclass classification problems.
- ☐ The Naive Bayes Classifier always outperforms other classifiers in terms of accuracy.
- ☐ None of the above.

**Solution:**

BC

#### 4.2 Perceptron and Logistic Regression (5 pt)

You are working on a binary classification problem using the Perceptron Learning Algorithm. Consider the following computation graph, for which the inputs are two features  $x_1, x_2 \in \mathbb{R}$ . The output is a binary label  $y \in \{-1, +1\}$ .



where  $z = w_0 + w_1x_1 + w_2x_2$ . Suppose we have 2 models: a perceptron and a logistic regression model. The perceptron calculates  $y$  as follows:

$$y = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

The logistic regression model first calculates the probability of  $y = +1$  as follows:

$$P(y = +1|x; w) = \frac{1}{1 + e^{-z}}$$

and then do classification by selecting  $\hat{y} = \arg \max_y P(y|x; w)$ . For binary classification, we can have the following prediction rule

$$\hat{y} = \begin{cases} 1 & \text{if } P(y = 1 | x; w) \geq 0.5, \\ -1 & \text{if } P(y = 1 | x; w) < 0.5. \end{cases}$$

The dataset consists of four training examples, each with two features  $x_1$  and  $x_2$  as detailed in Figure 1. We label the positive data (+1) as “+” symbol and the negative data (-1) as “-” symbol. The order of training is from 1 to 4 data point (in red circles), where the data points are (-2, 0), (0, 1), (2, 0), (0, -1).

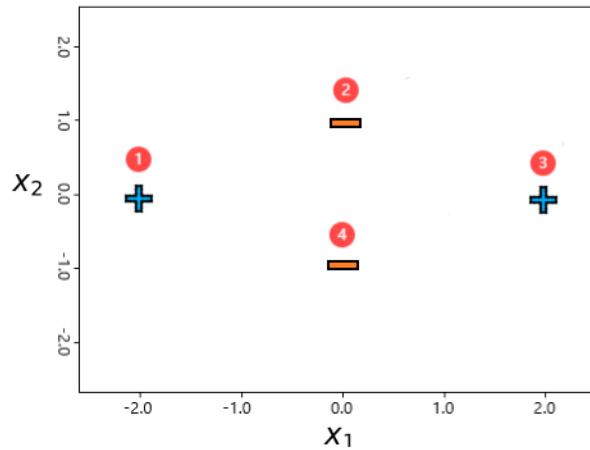


Figure 1: Dataset for perceptron.

#### 4.2.1 Classification (2 pt)

Now we perform binary perceptron update rule (for prediction, please refer to the above) to update the weight vector which is initialized as  $[w_0 \ w_1 \ w_2] = [0 \ 0 \ 0]$ , where  $w_0$  is the bias term. What is the final weight if we update in the order of  $1 \rightarrow 2 \rightarrow 3 \rightarrow 4$  of data points?

**Solution:**

**$[-1 \ 2 \ 0]$**

#### 4.2.2 Training and prediction (1 pt)

Suppose we train the perceptron on the training data points in Figure 1 in another order ( $4 \rightarrow 3 \rightarrow 2 \rightarrow 1$ ) using features  $x_1$  and  $x_2$  with the same initialized weight. Compared with the resulting weight under the order ( $1 \rightarrow 2 \rightarrow 3 \rightarrow 4$ ) in the above, which of the following statement is correct?

- ☐ They are the same, because the order of training of perceptron will never affect the final weight.
- ☐ They are different, because the order of training of perceptron will possibly affect the final weight.
- ☐ They are the same, though the order of training of perceptron will possibly affect the weight.

**Solution:**

B

#### 4.2.3 Improving Perceptrons (Logistic Regression) (1 pt)

Suppose we train the a logistic regression model on the training data points in Figure 1 using features  $x_1$  and  $x_2$ . Are we able to achieve a consistent 100% accuracy of training (correctly classify all the points)?

- ☐ Yes, because the logistic regression model uses probabilistic decision which is better than perceptron.
- ☐ No, because the dataset is not linear separable.
- ☐ Unknown, because it depends on how well the model is trained.

**Solution:**

B

#### 4.2.4 Feature Exploration (1 pt)

Based on the training data points in Figure 1 using features  $x_1$  and  $x_2$ , you are considering adding one additional feature to the dataset. For each of the following feature transformations, select it if adding it as a third feature (alongside and) would make the two classes linearly separable. Which of the following features should be selected? (Select all that apply.)

- ☐  $\|x\|_2^2$
- ☐  $\|x\|_2^3$
- ☐  $x_1 + x_2$
- ☐  $2x_1$

For the above,  $\|x\|_2$  means the  $L_2$  norm of a vector. For example, if  $x = (x_1 \ x_2)$ ,  $\|x\|_2 = \sqrt{x_1^2 + x_2^2}$ .

**Solution:**

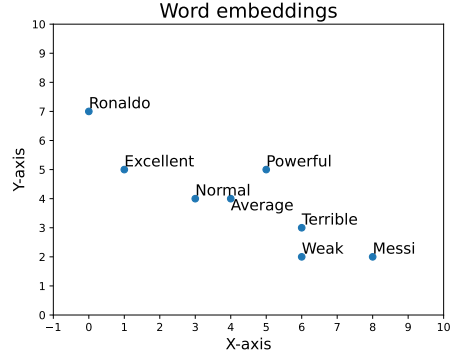
A,B

## 5 Large Language Model (10 pt)

### 5.1 Word Embeddings (6 pt)

In a large language model, a word embedding is a feature vector that represents each word in a high-dimensional space. Words with similar semantics typically have similar embeddings. As illustrated in the table below, given  $N$  word embeddings, our goal is to classify these words using the k-means algorithm.

Word	Embedding $(x_i, y_i)$
Excellent	(1,5)
Weak	(6,2)
Average	(4,4)
Messi	(8,2)
Terrible	(6,3)
Normal	(3,4)
Ronaldo	(0,7)
Powerful	(5,5)



For simplicity, all word embeddings are represented as two-dimensional vectors  $E_i = (x_i, y_i)$  where  $x_i, y_i \in \mathbb{R}$ .

We use the Euclidean distance for distance calculation between the embedding vectors, i.e.,

$$\text{dist}(E_i, E_j) = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2}$$

#### 5.1.1 Data Point Assignment (2.5 pt)

The cluster center  $c_j^{(t)}$  represents the cluster center of category  $j$  after the  $t$  iteration. In this problem, we define the number of categories  $k = 3$ , and the initial three cluster centers are  $c_0^{(0)}, c_1^{(0)}, c_2^{(0)} = (1, 5), (3, 4), (5, 5)$ .

Assign data points to the closest cluster center given the initial centers.

	(1, 5)	(6, 2)	(4, 4)	(8, 2)	(6, 3)	(3, 4)	(0, 7)	(5, 5)
$c_0^{(0)} = (1, 5)$	■	□	□	□	□	□	□	□
$c_1^{(0)} = (3, 4)$	□	□	□	□	□	■	□	□
$c_2^{(0)} = (5, 5)$	□	□	□	□	□	□	□	■

**Solution:**

	(1, 5)	(6, 2)	(4, 4)	(8, 2)	(6, 3)	(3, 4)	(0, 7)	(5, 5)
$c_0^{(0)} = (1, 5)$	■	□	□	□	□	□	■	□
$c_1^{(0)} = (3, 4)$	□	□	■	□	□	■	□	□
$c_2^{(0)} = (5, 5)$	□	■	□	■	■	□	□	■

### 5.1.2 Update Cluster Center (1.5 pt)

Assign each center to the average of its assigned points. Calculate the new cluster centers and fill in the blanks below.

(1)  $c_0^{(1)} =$

(2)  $c_1^{(1)} =$

(3)  $c_2^{(1)} =$

**Solution:**

(1)  $c_0^{(1)} = (0.5, 6)$

(2)  $c_1^{(1)} = (3.5, 4)$

(3)  $c_2^{(1)} = (6.25, 3)$

### 5.1.3 Convergence (2 pt)

In the K-means algorithm, the choice of initial centers affects convergence speed. Which of the following initial centers can make the K-means algorithm converge faster (convergence time is defined as the number of iterations required)

☐  $(1, 5), (6, 2), (4, 4)$

☐  $(6, 2), (8, 2), (0, 7)$

**Solution:**

A

## 5.2 Attention is all you need (4 pt)

For this section, we are going to explore the attention mechanism in large language models. The simplified attention weights of a sentence can be calculated by the following formula:

$$A = \text{Softmax}(QK^\top, \text{dim} = 1)$$

$A \in \mathbb{R}^{L \times L}$  is the attention weights,  $Q \in \mathbb{R}^{L \times D}$  is the query vectors,  $K \in \mathbb{R}^{L \times D}$  is the key vectors, where  $L$  represents the length of the sentence and  $D$  represents the feature dimension. The Softmax function is performed on the last dimension.

### 5.2.1 Self-Attention (2 pt)

Now we have a sentence  $S = \text{"Ronaldo is the best soccer player"}$  which consists of six words. At the same time, we calculated the  $Q$  and  $K$  of this sentence and get (In this question  $D = 2$ ):

$$S = \begin{bmatrix} \text{Ronaldo} \\ \text{is} \\ \text{the} \\ \text{best} \\ \text{soccer} \\ \text{player} \end{bmatrix}, Q = \begin{bmatrix} 1 & 2 \\ -2 & 1 \\ -3 & 1 \\ 1 & 3 \\ 2 & 2 \\ 2 & 1 \end{bmatrix}, K = \begin{bmatrix} 1 & 1 \\ -3 & 2 \\ -2 & 2 \\ 1 & 2 \\ 2 & 1 \\ 2 & 0 \end{bmatrix}$$

By calculating the attention weights, please answer which word does "Ronaldo" put the most attention weight on?

- ☐ Ronaldo
 ☐ is
 ☐ the
 ☐ best
 ☐ soccer
 ☐ player

**Solution:**

- ☐ Ronaldo
 ☐ is
 ☐ the
 ☒ best
 ☐ soccer
 ☐ player

### 5.2.2 Causal Attention (2 pt)

In large language models, causal attention is proposed to make each word only pay attention to the words before it (include itself). The specific implementation method is usually to add a mask when calculating the attention:

$$A = \text{Softmax}(QK^T + M)$$

For the above sentence of length 6, if causal attention is applied, which of the following should the mask  $M$  be?

☐

$$\begin{bmatrix} 0, & -\infty, & -\infty, & -\infty, & -\infty, & -\infty \\ -\infty, & 0, & -\infty, & -\infty, & -\infty, & -\infty \\ -\infty, & -\infty, & 0, & -\infty, & -\infty, & -\infty \\ -\infty, & -\infty, & -\infty, & 0, & -\infty, & -\infty \\ -\infty, & -\infty, & -\infty, & -\infty, & 0, & -\infty \\ -\infty, & -\infty, & -\infty, & -\infty, & -\infty, & 0 \end{bmatrix}$$



$$\bigcirc \begin{bmatrix} 0, & -\infty, & -\infty, & -\infty, & -\infty, & -\infty \\ 0, & 0, & -\infty, & -\infty, & -\infty, & -\infty \\ 0, & 0, & 0, & -\infty, & -\infty, & -\infty \\ 0, & 0, & 0, & 0, & -\infty, & -\infty \\ 0, & 0, & 0, & 0, & 0, & -\infty \\ 0, & 0, & 0, & 0, & 0, & 0 \end{bmatrix}$$

$$\bigcirc \begin{bmatrix} 0, & 0, & 0, & 0, & 0, & 0 \\ -\infty, & 0, & 0, & 0, & 0, & 0 \\ -\infty, & -\infty, & 0, & 0, & 0, & 0 \\ -\infty, & -\infty, & -\infty, & 0, & 0, & 0 \\ -\infty, & -\infty, & -\infty, & -\infty, & 0, & 0 \\ -\infty, & -\infty, & -\infty, & -\infty, & -\infty, & 0 \end{bmatrix}$$

**Solution:**

**B**