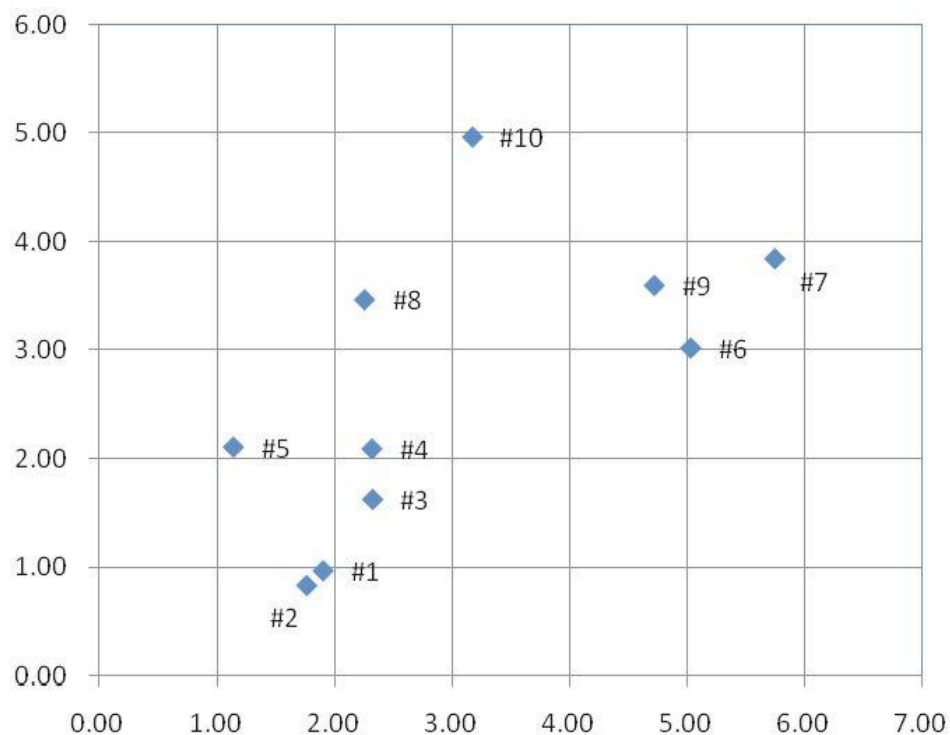


**Q11**

70 Points

Suppose you are given the following  $\langle x, y \rangle$  pairs. You will simulate the k-means algorithm to identify TWO clusters ( $k=2$ ) in the data.

Data #	x	y
1	1.90	0.97
2	1.76	0.84
3	2.32	1.63
4	2.31	2.09
5	1.14	2.11
6	5.02	3.02
7	5.74	3.84
8	2.25	3.47
9	4.71	3.60
10	3.17	4.96



Suppose you are given the initial assignment cluster center as:  
{cluster1: #1}, {cluster2: #10} – the first data point is used as the first

cluster center and the 10-th as the second cluster center.  
Assume k-means uses Euclidean distance.

### Q1.1 1

15 Points

Please simulate the k-means ( $k=2$ ) algorithm for the first iteration.  
(Upload a table that shows the Euclidean distance results and clusters for each sample in the dataset and show calculation)

▼ IMG\_20210528\_172341\_N.jpg

Download

Q1.1 I calculated distance with python codes.

	1	2	3	4	5	6	7	8	9	10
1	0	0.191	0.782	1.192	1.370	3.733	4.794	2.524	3.848	4.187
2		0	0.968	1.365	1.413	3.921	4.984	2.675	4.039	4.354
3			0	0.460	1.273	3.036	4.071	1.841	3.097	3.436
4				0	1.170	2.865	3.850	1.381	2.835	2.926
5					0	3.985	4.914	1.755	3.868	3.499
6						0	1.091	2.806	0.657	2.680
7							0	3.509	1.05	2.808
8								0	2.463	1.751
9									0	2.054
10										0

Data	Cluster Assignment
1	cluster 1
2	cluster 1
3	cluster 1
4	cluster 1
5	cluster 2
6	cluster 2
7	cluster 2
8	cluster 2
9	cluster 2

Data	cluster assign
10	cluster 2

### Q1.2 2

15 Points

Please simulate the k-means ( $k=2$ ) algorithm for the second iteration. (Upload a table that shows the Euclidean distance results and clusters for each sample in the dataset and show calculation)

▼ IMG\_20210528\_172350\_N.jpg

Download

Q1.2

$$m_1 = \left( \frac{1}{4} (1.90 + 1.76 + 2.32 + 2.31), \frac{1}{4} (0.97 + 0.84 + 1.63 + 2.09) \right)$$

$$m_2 = \left( \frac{1}{6} (1.14 + 5.02 + 5.92 + 2.25 + 4.71 + 3.17), \frac{1}{6} (2.11 + 3.02 + 3.84 + 3.47 + 3.60 + 4.96) \right)$$

$$m_1 = (2.0725, 1.3825)$$

$$m_2 = (4.402, 4.2)$$

if must be point

Data	$m_1$	$m_2$
1	0.447	4.085 $\rightarrow m_1$
2	0.626	4.274 $\rightarrow m_1$
3	0.350	3.307 $\rightarrow m_1$
4	0.746	2.974 $\rightarrow m_1$
5	1.182	3.874 $\rightarrow m_1$
6	3.371	1.332 $\rightarrow m_2$
7	4.414	1.385 $\rightarrow m_2$
8	2.095	2.272 $\rightarrow m_1$
9	3.445	0.674 $\rightarrow m_2$
10	3.742	1.447 $\rightarrow m_2$

$m_1$  is a assign cluster 1  
 $m_2$  is a assign cluster 2  
~~new  $m_1$~~

new  $m_1 = (2.108, 1.8)$   
 new  $m_2 = (2.948, 3.506)$

### Q1.3 3

15 Points

Please simulate the k-means ( $k=2$ ) algorithm for the third iteration. (Upload a table that shows the Euclidean distance results and clusters for each sample in the dataset and show calculation)

No files uploaded

### Q1.4 4

## 25 Points

What are the cluster assignments until convergence? (Upload a table that shows the Euclidean distance results and clusters for each sample in the dataset)

 No files uploaded

**Q2.2**

30 Points

Check the boxes for ALL CORRECT CHOICES.

Every question should have at least one box checked.

**Q2.1.1**

3 Points

In terms of the bias-variance decomposition, a 1-nearest neighbor classifier has \_\_\_\_\_ than a 3-nearest neighbor classifier.

☒ higher variance☐ higher bias☐ lower variance☒ lower bias**Q2.2.2**

3 Points

Which of the following are true about bagging?

☒ In bagging, we choose random subsamples of the input points with replacement

☐ The main purpose of bagging is to decrease the bias of learning algorithms.

☐ Bagging is ineffective with logistic regression because all of the learners learn exactly the same decision boundary

☒ If we use decision trees that have one sample point per leaf, bagging never gives lower training error than one ordinary decision tree

### Q2.3 3

3 Points

You've just finished training a random forest for spam classification, and it is getting abnormally bad performance on your validation set, but good performance on your training set. Your implementation has no bugs. What could be causing the problem?

☒ Your decision trees are too deep

☒ You have too few trees in your ensemble

☒ You are randomly sampling too many features when you choose a split

☒ Your bagging implementation is randomly sampling sample points without replacement

### Q2.4 4

3 Points

What strategies can help reduce overfitting in decision trees?

☒ Pruning☒ Enforce a minimum number of samples in leaf nodes☐ Make sure each leaf node is one pure class☒ Enforce a maximum depth for the tree**Q2.5 5**

3 Points

Which of the following are true of convolutional neural networks (CNNs) for image analysis?

☒ Filters in earlier layers tend to include edge detectors☒ Pooling layers reduce the spatial resolution of the image☐ They have more parameters than fully-connected networks with the same number of layers and the same numbers of neurons in each layer☐ A CNN can be trained for unsupervised learning tasks, whereas an ordinary neural net cannot**Q2.6 6**

3 Points

Neural networks

☐ optimize a convex cost function☐ always output values between 0 and 1☒ can be used for regression as well as classification☒ can be used in an ensemble**Q2.7 7**

3 Points

As the number of training examples goes to infinity, your model trained on that data will have:

☐ Lower bias☐ Higher bias☒ Same bias**Q2.8 8**

3 Points

As the number of training examples goes to infinity, your model trained on that data will have:

☒ Lower variance☐ Higher variance☐ Same variance**Q2.9 9**

3 Points

Which of the following statements about ensemble methods is true?

- ☐ Combining weak learners using bagging is good since it can reduce the variance.
- ☒ Combining strong learners using boosting is good since it can reduce the bias.
- ☒ Combining weak learners using boosting is good since it can reduce the variance.
- ☐ Combining strong learners using bagging is good since it can reduce the variance

## Q2.10 10

3 Points

Which of the following can not be used in unsupervised learning:

- ☒ Classification
- ☒ Regression
- ☐ Dimension reduction
- ☐ Clustering

## Quiz-6

● GRADED

STUDENT

MEHMET TAHA USTA



TOTAL POINTS

49 / 100 pts

QUESTION 1

1		25 / 70 pts
1.1	1	15 / 15 pts
1.2	2	10 / 15 pts
1.3	3	0 / 15 pts
1.4	4	0 / 25 pts

QUESTION 2

2		24 / 30 pts
2.1	1	3 / 3 pts
2.2	2	3 / 3 pts
2.3	3	3 / 3 pts
2.4	4	3 / 3 pts
2.5	5	3 / 3 pts
2.6	6	3 / 3 pts
2.7	7	3 / 3 pts
2.8	8	3 / 3 pts
2.9	9	0 / 3 pts
2.10	10	0 / 3 pts