The "Data Mining" Specialization        Learn More                          ✕

# Feedback — Week 1 Quiz                                    Help Center

Thank you. Your submission for this quiz was received.

You submitted this quiz on **Sat 2 May 2015 10:50 AM PDT**. You got a score of **13.00** out of **13.00**.

## Question 1

Which of the following statements are true?

| Your Answer | | Score | Explanation |
|---|---|---|---|
| ☑ Clustering analysis in unsupervised learning since it does not require labeled training data. | ✔ | 0.25 | |
| ☐ When clustering, we want to put two dissimilar data objects into the same cluster. | ✔ | 0.25 | This is false because, as discussed in the lecture, the objective of clustering is having similar objects in the same cluster. |
| ☑ Clustering analysis has a wide range of applications in tasks such as data summarization, dynamic trend detection, multimedia analysis, and biological network analysis. | ✔ | 0.25 | |
| ☐ It is impossible to cluster objects in a data stream. We must have all the data objects that we need to cluster ready before clustering can be performed. | ✔ | 0.25 | This is false because clustering algorithms can be adapted to perform clustering in a streaming fashion. |
| Total | | 1.00 / 1.00 | |

**Question Explanation**

The correct answers are: "Clustering analysis in unsupervised learning since it does not require

labeled training data." and "Clustering analysis has a wide range of applications in tasks such as data summarization, dynamic trend detection, multimedia analysis, and biological network analysis."

# Question 2

What are some common considerations and requirements for cluster analysis?

| Your Answer | | Score | Explanation |
|---|---|---|---|
| ☑ In order to perform cluster analysis, we need to have a similarity measure between data objects. | ✔ | 0.25 | |
| ☑ We need to consider the desired shape and size of clusters. | ✔ | 0.25 | |
| ☐ Cluster analysis requires a large amount of labeled training data. | ✔ | 0.25 | This is false since cluster analysis is unsupervised learning, which by definition does not require training data. |
| ☑ We need to consider the space on which the clustering is performed. In other words, we need to decide what subset of the available features we are going to consider in the similarity measure. | ✔ | 0.25 | |
| Total | | 1.00 / 1.00 | |

**Question Explanation**

The correct answers are: "In order to perform cluster analysis, we need to have a similarity measure between data objects." , "We need to consider the desired shape and size of clusters." and "We need to consider the space on which the clustering is performed. In other words, we need to decide what subset of the available features we are going to consider in the similarity measure."

# Question 3

Which of the following statements are true?

| Your Answer | Score | Explanation |
|---|---|---|

| ☐ We can only visualize the clustering results when the data is 2-dimensional. | ✔ 0.25 | This is false because we can also easily visualize clusters in 3D. HD-eye is an example of software used for visualizing higher dimensional clusters. |
|---|---|---|
| ☑ When dealing with high-dimensional data, we sometimes consider only a subset of the dimensions when performing cluster analysis. | ✔ 0.25 | |
| ☑ Agglomerative clustering is an example of a distance-based clustering method. | ✔ 0.25 | |
| ☑ Graphs, time-series data, text, and multimedia data are all examples of data types on which cluster analysis can be performed. | ✔ 0.25 | |
| Total | 1.00 / 1.00 | |

**Question Explanation**

The correct answers are: "Agglomerative clustering is an example of a distance-based clustering method.", "When dealing with high-dimensional data, we sometimes consider only a subset of the dimensions when performing cluster analysis.", and "Graphs, time-series data, text, and multimedia data are all examples of data types on which cluster analysis can be performed."

# Question 4

The following real dataset contains information about Iris setosa and versicolor.

What is the supremum distance between these two objects?

| Species | Sepal length | Sepal width | Petal length | Petal width |
|---|---|---|---|---|
| Iris setosa | 4.9 | 3.0 | 1.4 | 0.2 |
| Iris versicolor | 5.6 | 2.5 | 3.9 | 1.1 |

| Your Answer | Score | Explanation |
|---|---|---|
| ○ 2.8 | | |
| ○ 4.6 | | |

○ 7.8

◉ 2.5                              ✔            1.00

Total                                           1.00 / 1.00

**Question Explanation**

The supremum distance, corresponding to p = 2 Minkowski distance, between two objects is defined as follows:

$d_\infty(i,j) = {}^l\max_{k=1} | x_{i,k} - x_{j,k} |$

# Question 5

The following real dataset contains information about Iris setosa and versicolor.

| Cases | F1 | F2 | F3 | F4 | F5 | F6 | F7 | F8 | F9 | F10 |
|-------|----|----|----|----|----|----|----|----|----|-----|
| 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| 2 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 |

Assume all the activities are symmetric binary variables. What is the distance between Case 1 and Case 2?

| Your Answer | Score | Explanation |
|-------------|-------|-------------|
| ○ 4/7 | | |
| ○ 3/7 | | |
| ○ 3/10 | | |
| ◉ 4/10 | ✔ 1.00 | |
| Total | 1.00 / 1.00 | |

**Question Explanation**

If the binary variables are symmetric, we have:

d(i,j) = r + s / q + r + s + t = 4/10

# Question 6

The following real world dataset contains two samples from Car Evaluation Database, which was derived from a simple hierarchical decision model originally developed for the demonstration of DEX (M. Bohanec, V. Rajkovic: Expert System for Decision Making. Sistemica 1(1), pp. 145-157, 1990.). The model evaluates cars according to the following concept structure:

CAR            car acceptability

. PRICE            overall price

. . buying            buying price

. . maint            price of the maintenance

. TECH            technical characteristics

. . COMFORT        comfort

. . . doors            number of doors

. . . persons        capacity in terms of persons to carry

. . . lug_boot        the size of luggage boot

. . safety            estimated safety of the car

The attribute values are as follows:

| Attrubute | Values (categorical) |
|---|---|
| buying | v-high, high, med, low |
| maint | v-high, high, med, low |
| doors | 2, 3, 4, 5-more |
| persons | 2, 4, more |
| lug_boot | small, med, big |
| safety | low, med, high |

| Case | buying | maint | doors | persons | lug_boot | safety |
|---|---|---|---|---|---|---|
| Car 1 | med | v-high | 3 | more | small | med |
| Car 2 | high | v-high | 4 | 4 | big | med |

To calculate the distance between objects with categorical attributes, we use a set of binary attributes to represent each categorical attribute. Assume all the binary attributes are symmetric. What is the distance between Car 1 and Car 2?

| Your Answer | Score | Explanation |
|---|---|---|
| ○ 2/3 | | |
| ◉ 8/21 | ✔ | 1.00 |

○ 8/10

○ 1/3

○ 8/17

Total                                                  1.00 / 1.00

---

**Question Explanation**

Considering converting the categorical random variables into binary random variables, there will be (4 + 4 + 4 + 3 + 3 + 3 = 21) binary random variables in total. Moreover, we have the following contingency table:

|  |  | Car 1 |  |  |
|---|---|---|---|---|
|  |  | 1 | 0 | sum |
| Car 2 | 1 | 2 | 4 | 6 |
|  | 0 | 4 | 11 | 15 |
|  | sum | 6 | 15 | 21 |

If all the binary attributes are asymmetric, we have the distance between Car 1 and Car 2 as:

d(i,j) = s+r/q+s+r+t = 8/21

---

# Question 7

Given the following two short texts with punctuation removed, calculate the cosine similarity between them based on the bag of words model.

Text1: language is the source of misunderstandings

Text2: language is the soul of a nation

| Your Answer | Score | Explanation |
|---|---|---|
| ○ 0 | | |
| ○ 0.44 | | |
| ○ 0.095 | | |
| ◉ 0.617 | ✔ 1.00 | |
| Total | 1.00 / 1.00 | |

---

**Question Explanation**

We can get the vector representations for the two short texts,

T1 = (1, 1, 1, 1, 1, 1, 0, 0, 0)

T2 = (1, 1, 1, 0, 1, 0, 1, 1, 1),

where the dimensions correspond to: language, is, the, source, of, misunderstanding, soul, a, nation.

The cosine similarity can be obtained via

cos(T1,T2) = T1*T2/ ||T1|| ||T2|| = 4/ $\sqrt{6} \sqrt{7}$

# Question 8

With regard to the species of Iris setosa, we have sampled data on the features of sepal length and sepal width, as follows:

| Feature | Sepal length | Sepal width |
|---|---|---|
| Case 1 | 5.1 | 3.5 |
| Case 2 | 4.9 | 3.0 |
| Case 3 | 4.7 | 3.2 |
| Case 4 | 4.6 | 3.1 |
| Case 5 | 5.0 | 5.4 |

What is the sample correlation coefficient between sepal length and sepal width?

| Your Answer | | Score | Explanation |
|---|---|---|---|
| ○ 2.396 | | | |
| ○ 0.185 | | | |
| ◉ 0.479 | ✔ | 1.00 | |
| ○ 0.895 | | | |
| ○ 0.398 | | | |
| Total | | 1.00 / 1.00 | |

**Question Explanation**

The sample correlation coefficient can be calculated as:

$\rho_{12} = \sigma_{12}/\sigma_1 \sigma_2 = {}^n\Sigma_{i=l} (x_{il} - \mu_l)(x_{i2} - \mu) / \sqrt{{}^n\Sigma_{i=l} (x - \mu_2)^2}$

# Question 9

Considering the K-means algorithm, after the current iteration, we have 3 centroids (0, 1) (2, 1), (-1, 2). Will points (0.5, 0.5) and (-0.5, 0) be assigned to the same cluster in the next iteration?

| Your Answer | | Score | Explanation |
|---|---|---|---|
| ⦿ Yes | ✔ | 1.00 | |
| ◯ No | | | |
| Total | | 1.00 / 1.00 | |

**Question Explanation**

They will be both assigned to (0, 1).

# Question 10

Considering the K-means algorithm, if points (1, -3), (1, 1), and (-2, 2) are the only points which are assigned to the first cluster now, what is the new centroid for this cluster?

| Your Answer | | Score | Explanation |
|---|---|---|---|
| ⦿ (0, 0) | ✔ | 1.00 | |
| ◯ (0, 3) | | | |
| ◯ (0, 2) | | | |
| ◯ (-2, 1) | | | |
| Total | | 1.00 / 1.00 | |

**Question Explanation**

Calculate the average value for x and y separately. You will then find the answers are 0 and 0, and thus, the new centroid should be (0, 0).

# Question 11

K-means++ algorithm is designed for better initialization for K-means, which will take the

farthest point from the currently selected centroids. Suppose k = 2 and we have selected the first centroid is (0, 0). Among the following points (these are all the remaining points), which one should we take for the second centroid? (Here, the distance is measured by Euclidean Distance).

| Your Answer | Score | Explanation |
|---|---|---|
| ⦿ (3, 0) | ✔ 1.00 | |
| ○ (0, 1) | | |
| ○ (2, -2) | | |
| ○ (-2, 1) | | |
| Total | 1.00 / 1.00 | |

**Question Explanation**

K-means++ will take the farthest point from the currently selected centroids.

# Question 12

Considering the K-median algorithm, if points (-1, 3), (-3, 1), and (-2, -1) are the only points which are assigned to the first cluster now, what is the new centroid for this cluster?

| Your Answer | Score | Explanation |
|---|---|---|
| ○ (0, 2) | | |
| ⦿ (-2, 1) | ✔ 1.00 | |
| ○ (0, 0) | | |
| ○ (0, 3) | | |
| Total | 1.00 / 1.00 | |

**Question Explanation**

Calculate the median value for x and y separately. You will then find the answers are -2 and 1, and thus the new centroid should be (-2, 1)

# Question 13

Which of the following statements about K-medoids, K-median, and K-modes algorithms are correct?

| Your Answer | | Score | Explanation |
|---|---|---|---|
| ☑ The centroids in the K-means algorithm may not be any observed data points. | ✔ | 0.25 | |
| ☑ The K-modes algorithm is designed for categorical data. | ✔ | 0.25 | |
| ☑ The K-medoids and K-median algorithms are less sensitive to outliers than K-means. | ✔ | 0.25 | |
| ☐ The centroids in the K-medoids algorithm may not be any observed data points. | ✔ | 0.25 | |
| Total | | 1.00 / 1.00 | |

**Question Explanation**

In the K-medoids algorithm, the centroids are selected from the given data points.