

Course: Data Mining (INFOMDM) Test name: [20251106] INFOMDM - Data mining - 1-GS - USP Status: Closed Type: Exam Time left: 116 min. and 30 sec.

No filters selected; all answers are displayed

Course	Data Mining (INFOMDM)
Blueprint name	[20251106] INFOMDM - Data mining - 1-GS - USP
Blueprint type	Exam
Status	Closed
Number of questions	10
Start/End time	November 6, 2025 from 1:30 PM to 3:39 PM
Score	78 (maximum score: 100)
Number of questions answered correctly:	6
Questions mostly answered correctly:	3
Questions mostly answered incorrectly:	1

[Hide question text](#)

Question 1 Answered on: November 6, 2025 - 3:16 PM Duration: 22 min. and 20 sec. Score: 14 of 20 pts.

14 pts.

TRUE OR FALSE?

Indicate whether the following statements are true or false.

(Frequent Sequence Mining) "ai" occurs 2 times as a subsequence of "taai taai pop" (there are 2 different mappings)

- TRUE
 FALSE 2 pt. ✓

(Bayesian Networks) Two directed independence graphs are equivalent if they have the same skeleton and the same v-structures (imoralities).

- TRUE
 FALSE 2 pt. ✓

(Random Forests) Each split in a tree in a random forest is allowed to use only a random subset of the features.

- TRUE
 FALSE 2 pt. ✓

(Bias-Variance decomposition) As the training set size increases, the bias component of expected prediction error decreases.

- TRUE
 FALSE 2 pt. ✓

(Link-based classification) In link-based classification, we model the dependency between node labels and the labels of their neighboring nodes.

- TRUE
 FALSE 2 pt. ✓

(Gradient Boosting) In gradient boosting with trees, each additional tree tries to reduce the errors of the model constructed so far.

- TRUE
 FALSE 2 pt. ✓

(Cost-complexity pruning) In cost-complexity pruning, we prune in the nodes of the current pruned subtree of T_{max} that yield the smallest increase in resubstitution error per leaf reduction (i.e., per unit decrease in the number of leaves).

- TRUE
 FALSE 2 pt. ✓

(Bayesian networks) In the directed graphical model $A \rightarrow B \leftarrow C$, A and C are independent given B.

- TRUE
 FALSE 2 pt. ✓

(Link prediction) In link prediction, we only choose pairs of nodes from the training interval that are not connected by an edge. The test interval is used to provide the class labels for those pairs.

- TRUE
 FALSE 2 pt. ✓

(Frequent pattern mining) Consider an alternative tree mining scenario with just a single data tree. In this scenario, the support of a pattern tree is equal to the number of distinct occurrences of the pattern tree in the data tree. Two occurrences are considered distinct if they correspond to mapping functions φ_1 and φ_2 , where $\varphi_1(v) \neq \varphi_2(v)$ for at least one node v in the pattern tree.

Claim: the anti-monotonicity property between support and the induced subtree relationship holds in this scenario.

- TRUE
 FALSE 2 pt. ✓

Question 2 Answered on: November 6, 2025 - 1:44 PM Duration: 4 min. and 24 sec. Score: 5 of 5 pts.

5 pts.

Classification Trees: Computing Splits

As we are growing a classification tree, we encounter a node that contains the following data on numerical attribute x and binary class label y :

x	4	4	8	10	16	16	20	26
y	0	0	0	1	0	1	1	1

We use the gini-index as impurity measure.

If we use the segment borders algorithm to determine the best split on x , we need to compute the impurity reduction of the following splits (1 or more answers may be correct):

- $x \leq 6$
 $x \leq 9$
 $x \leq 13$
 $x \leq 18$
 $x \leq 23$

Question 3 Answered on: November 6, 2025 - 1:47 PM Duration: 4 min. and 3 sec. Score: 5 of 5 pts.

5 pts.

Classification Trees: Cost-Complexity Pruning

We are pruning a tree T that has been grown on $n=100$ training examples. The class variable has 3 possible values, denoted by A, B and C.Consider a node t in this tree, with the following class distribution:

class	A	B	C
number of examples	64	6	0

The branch T_t of tree T has 3 leaf nodes that together make 2 errors.What is the critical alpha value $g(t)$ for node t ? (round your answer to two decimal places)The critical value $g(t)$ for node t is: 0.04 0.04 5 pt.

Question 4 Answered on: November 6, 2025 - 1:58 PM Duration: 13 min. and 54 sec. Score: 10 of 10 pts.

10 pts.

Frequent Itemset Mining: Closed Frequent Itemsets

Consider the following transactions on items (A,B,C,D):

ti	Items
1	ABC
2	ABC
3	AB
4	BCD
5	BCD
6	CD
7	C
8	BC

We use the Apriori-Close (A-Close) algorithm to find all closed frequent itemsets with minimum support of 2.

Which of the following itemsets are level-2 generators? (1 or more answers may be correct)

- AB
 AC
 AD
 BC
 BD
 CD

Which of the following are closed frequent itemsets? (1 or more answers may be correct)

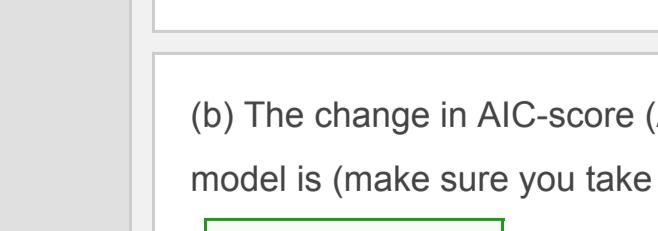
- ABC
 BD
 BC
 CD
 AC
 BCD
 D

Question 5 Answered on: November 6, 2025 - 2:04 PM Duration: 5 min. and 31 sec. Score: 5 of 10 pts.

5 pts.

Undirected Graphical Models

Consider a graphical log-linear model on variables A,B,C,D, and E with the following independence graph:



A, B, C and D are binary variables, and E has 5 possible values.

Answer the following questions:

(a) The maximum likelihood fitted counts for this model are given by:

- $\frac{n(A, B)n(A, C)n(A, D)n(B, E)}{n(A)^2 n(B)}$ 5 pt. ✓

$\frac{n(A, B, E)n(A, C)n(A, D)}{n(A)^2}$

$\frac{n(A, B)n(A, C)n(A, D)n(B, E)}{2 n(A)n(B)}$

This model does not have a closed form solution for the maximum likelihood fitted counts.

(b) The number of u-terms of this model is:

10 16 5 pt. (correct) 15 3 pt.

Question 6 Answered on: November 6, 2025 - 3:27 PM Duration: 16 min. and 29 sec. Score: 15 of 15 pts.

15 pts.

Undirected Graphical Models

The following data concerns an outbreak of food poisoning after the traditional Christmas Lunch of the personnel of the Department of Information and Computing Sciences of our University. This time the theme was Dutch cuisine. Of the food eaten, interest focused on the "Berenhap" and "Frikandel".

The variables are:

1. Berenhap (B) eaten (yes) or not eaten (no)
2. Frikandel (F) eaten (yes) or not eaten (no)
3. Sick (S) (yes) or not (no)

Questionnaires were completed by 100 of the 114 persons attending. The table of observed counts is given below.

Observed counts		Sick
Berenhap	Frikandel	no yes
no	no	22 9
yes	yes	3 12
yes	no	8 1
yes	yes	12 38

We fit the model $B \perp S \mid F$ to this data and obtain fitted counts:

Fitted counts	Sick
Berenhap	Frikandel
no	22.29
yes	3.46
no	11.54
yes	1.29

In your calculations, always use the natural logarithm. Round your final answers to two decimal places. Don't round intermediate results.

(a) The fitted count for $B = \text{yes}$, $F = \text{yes}$, $S = \text{yes}$ is equal to: 38.46 38.46 (+/- 0.1) 5 pts. ✓(b) The contribution of the cell $B = \text{no}$, $F = \text{no}$, $S = \text{no}$ to the deviance of the model $B \perp S \mid F$ is equal to: -0.58 -0.58 (+/- 0.01) 5 pts. -0.29 (+/- 0.01) 5 pts.It is given that the deviance of the model $B \perp S \mid F$ is equal to 0.21, and the critical value of χ^2 distribution with 2 degrees of freedom is equal to 9.21 for $\alpha = 0.01$.

(c) Based on this data, the best supported conclusion is that people got sick because of eating: Frikandel Frikandel 5 pt.

Question 7 Answered on: November 6, 2025 - 2:25 PM Duration: 13 min. and 50 sec. Score: 4 of 10 pts.

4 pts.

Text Classification: Multinomial Naive Bayes

You are given the following collection of song lyrics and corresponding music genre:

Words in lyrics	Genre
gone baby gone	Blues
woke up this morning	Blues
shake baby shake	Funk
shake ya funky funky ya ya	Funk

Answer the following questions:

(a) The estimate of $P(\text{baby} \mid \text{Blues})$ according to the multinomial naive Bayes model with Laplace smoothing is:

- $\frac{n(A, B)n(A, C)n(A, D)n(B, E)}{n(A)^2 n(B)}$ 5 pt. ✓

$\frac{n(A, B, E)n(A, C)n(A, D)}{n(A)^$