

COURSE CODE: 16B11BI612

MAX. MARKS: 35

COURSE NAME: Datawarehousing and mining for bioinformatics

COURSE CREDITS: 4

MAX. TIME: 2HRS

---

*Note: All questions are compulsory. Carrying mobile phone during examinations will be treated as a case of unfair means.*

---

1. Explain the difference between Inmon and Kimball models of datawarehousing. (5)
2. Explain the limitations of using single identification tree. How are these overcome by using a committee of decision tree approach? (2+6)
3. How do we use the partitioning approaches for clustering data in bioinformatics. Explain the k-means and k-medoids methods in this context. (3+3)
4. Explain briefly with steps how we perform PCA and describe its application briefly. (5)
5. For the following square matrix (5)

$$\begin{bmatrix} 3 & 0 & 1 \\ -4 & 1 & 2 \\ -6 & 0 & -2 \end{bmatrix}$$

Decide which, if any of the following vectors are eigen vectors of that matrix and give the corresponding eigen value.

$$\begin{matrix} 2 & -1 & -1 & 0 & 3 \\ (a) 2 & (b) 0 & (c) 1 & (d) 1 & (e) 2 \\ -1 & 2 & 3 & 0 & 1 \end{matrix}$$

6. Explain why gain ratio is a better measure for identification trees than entropy. Highlight the limitations of identification tree. What is meant by overfitting in ID trees and how do we circumvent it? (6)