

BIO543: Big Data Mining Healthcare

(15th April 2024, Quiz3)

Maximum Marks: 20

Duration: 20 Minutes

Name: _____

Admission No: _____

Instructions: Attempt all questions, each question carry one mark. Please tick only best possible option.

- Which of the following software is best software for big data mining?
 - Scikit learn
 - Mahout**
 - Weka
 - MySQL
- Which component of the Hadoop ecosystem is used for scripting?
 - Mahout
 - Hive
 - Pig**
 - Ambari
- Which of the following is not part of search engine (Lucene)?
 - Crawling
 - Ranking
 - Indexing
 - Embedding**
- What type of files are not extracted or analyze by Lucene search engine?
 - PDF
 - XML
 - TIFF**
 - MS Word
- What is purpose of shingling?
 - Convert nucleotides to amino acid
 - Convert documents to sets**
 - Convert SQL to NOSQL
 - Computing Jaccard similarity
- In which technique similar documents are grouped in buckets for similarity search.
 - Min-Hashing
 - LSH**
 - Curve clustering
 - BLAST
- What will be Jaccard similarity between set (1,1,0,0,1,1) and (1,0,1,1,1,1)?
 - 0.5**
 - 0.75
 - 0.75
 - 1.0
- In which technique, large sets are converted to short signature.
 - Shingling
 - RNA-seq
 - Word embedding
 - Min-Hashing**
- In LSH, if document C1 and C2 have similarity 0.80 or more, $b=20$, $r=5$, what is percent probability that C1 and C2 will be in same band.
 - 100 %
 - < 95.00%
 - > 99.96%**
 - < 90.25%

10. Which of the following is not part of cancer hallmarks?
- a) Evading growth suppressors
 - b) **Maintaining apoptosis**
 - c) Resisting cell death
 - d) Inducing angiogenesis
11. Study of all mRNA molecules in a cell is called.
- a) Genomics
 - b) Epigenomics
 - c) Proteomics
 - d) **Transcriptomics**
12. Which technique all to measure relative gene expression (e.g. disease vs healthy)?
- a) **cDNA microarray**
 - b) RNA-seq
 - c) Affymetrix
 - d) DNA-seq
13. Which of the following is used for prioritization of anticancer drugs in CancerDP?
- a) **Transcriptomics**
 - b) Glycomics
 - c) Epigenomics
 - d) Metabolomics
14. Which of the following database do not contain gene expression information?
- a) GEO
 - b) **Uniprot**
 - c) ArrayExpress
 - d) TCGA
15. In which stage you will classify a cancer patient who have single tumor that has not grown into any blood vessels?
- a) Late-stage
 - b) Stage-IV
 - c) **Stage-I**
 - d) Stage-II
16. Which of the following is simple feature not ML/DL trained features?
- a) FastText
 - b) GloVe
 - c) Word2Vec
 - d) **TF-IDF**
17. Which of the following software used large language model for predicting protein structures?
- a) **AlphaFold**
 - b) Modeller
 - c) HHpred
 - d) Robetta
18. Which of the following is not part of large language models
- a) **Transition probabilities**
 - b) Attention
 - c) Pre-trained models
 - d) Fine tuning of models
19. Which of the following approach is not used in large language model?
- a) Sentence correction
 - b) Text translation
 - c) Text completion
 - d) **Text similarity**
20. Which company developed large language model “Go AI” ?
- a) OpenAI
 - b) Meta
 - c) **DeepMind**
 - d) Google