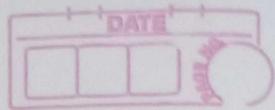


# Bio Unit-I

## Assignment - I



1) what is Bioinformatics? Explain the history of Bioinformatics.

⇒ Bioinformatics is new discipline which addresses the need, manages the need and a data heading past was massively generated by genomic research. This discipline represent the conversion of genetics biotechnology, IT and encompasses analysis as implementation of data modeling of biological phenomenon and their development of algorithms and statistics.

History →

BI history started with an  'Gregor Mendel' is known as the father of genetics. he cross fertilizes different colours of same species of flowers. he create record of the colours of flowers that produce inheritance of traits could be more easily extend it was controls by factors passed by generation to a generation.

In 1972 all work metofirst DNA molecules using ligas in same year Stanley Cohen, Annchan and Harifwost Groval producing first predominant & DNA organism.

In 1988, the human genome organisation was founded this is an international organization of scientist, involved in human genome project.

In 1993 human research lab produce physical a human map

3) Why Bioinformatics needed? Benefits of bioinformatics.

→ Bioinformatics is crucial because it combines biology, computer science, and mathematics to analyze and interpret biological data. It helps researchers make sense of vast amounts of genetic, molecular, and clinical information. This is vital for understanding diseases, drug discovery, genetic research, and personalized medicine.

Benefits:

i) Data Analysis: It enables the efficient analysis of large-scale biological data, such as DNA sequences, protein structures.

ii) Genomic Research: Bioinformatics plays a pivotal role in deciphering genomes, identifying genes and understanding genetic variations.

iii) Drug Discovery: It aids in identifying potential drug targets, predicting drug interactions.

iv) Disease Understanding: By analyzing biological data, bioinformatics help unravel the molecular mechanisms underlying diseases, leading to better diagnosis and treatment strategies.

v) Evolutionary Studies: It assists in tracing the evolutionary history of species by comparing genetic sequences, shedding light on their relationships.

4) Explain the term molecular biology.

→ Molecular biology is a specialized branch of biology and biochemistry which is specifically concerned study of various biological activities at molecular level.

- The history of molecular biology:

The term of molecular biology was coined by an american scientist warren in year 1938.

- The discovery of molecular biology began in year 1938 early 1940's and its fundamental development took place in year

1953.

- Living things are made of chemicals just as non-living things are, So a molecular biologist studies how molecules interact with one another in living organisms to perform the functions of life.

- Molecular biologists conduct experiments to investigate the structure, function, processing, regulation and evolution of biological molecules and their interactions with one another - providing micro-level insights into how life works.

~~Imp.~~  
Q6)

what is Biological Database? explain the characteristics in detail.

⇒

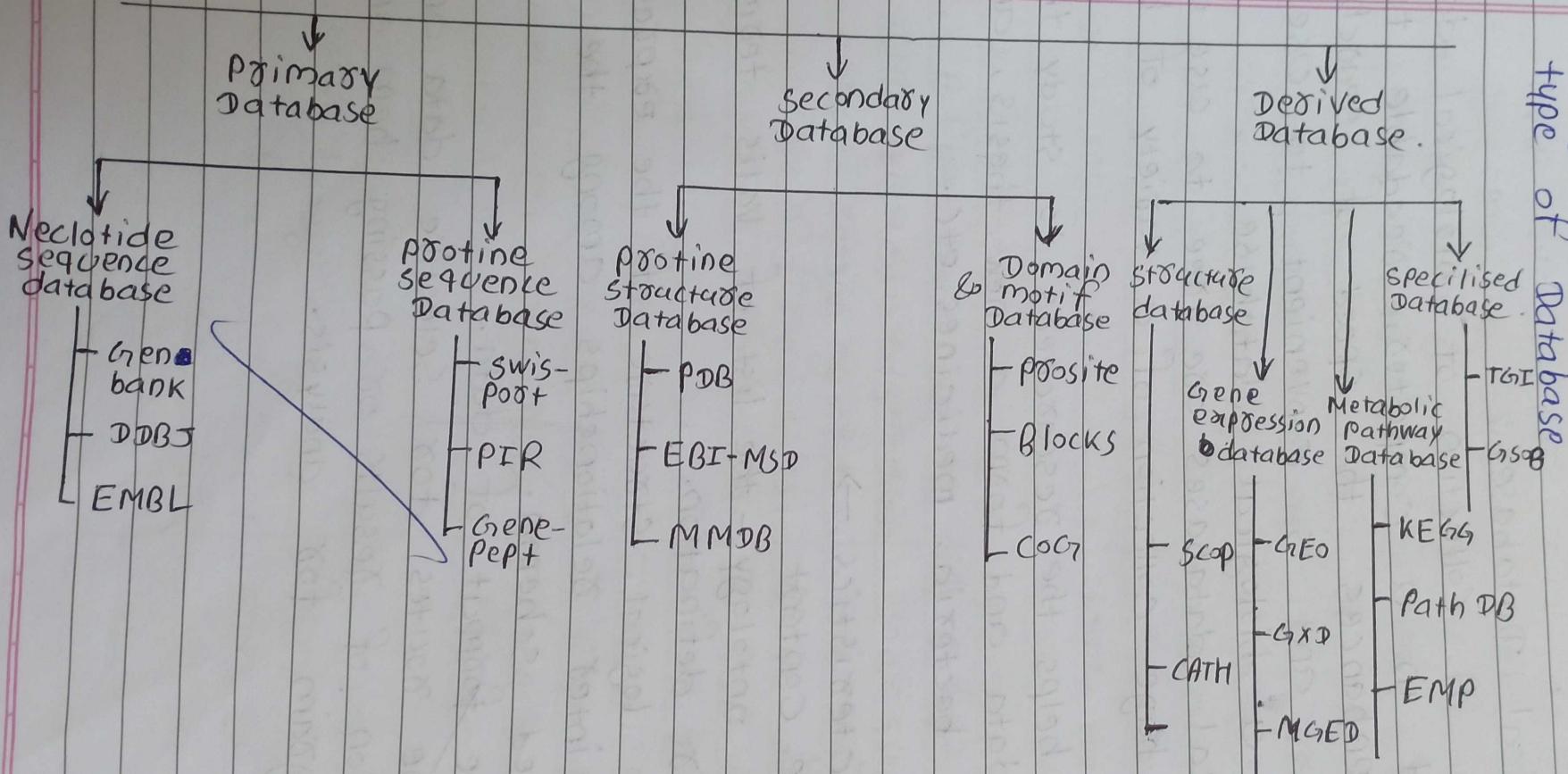
Biological Database :

- A collection of biological data arranged in a computer readable form that enhances the speed of search and retrieval and is convenient to use is called a biological database.
- Biological databases are free to use and contain a huge collection of a variety of biological data.
- It helps the researchers to study the available data and form a new thesis, antivirous, helpful bacteria, medicines, etc.

Characteristics →

- The Content
- The ontology - the list of valid term and their definition.
- The logical structure of the expressions of the inter relationships among the data called Schema.
- The format of data.
- The routes for Selective data and presentation of results or passing them on to a program for analysis.

## Biological Database



## 10) Difference between primary and secondary database

Primary Database	Secondary Database
i) It is also known as archival database.	i) It is also known as curated database or derived database.
ii) Database consisting of data derived experimentally such as nucleotide sequences and 3 dimensional structures are known as primary database.	ii) Database consisting of data derived from the analysis of primary data such as nucleotide sequences, protein structures.
iii) It contains the original experimental results are directly submitted into database by researchers	iii) It contains results of analysis of primary database and significant data in form of conserved sequences, signature sequences.
iv) Primary database has high levels of redundancy or duplication of data.	iv) This database has low levels of redundancy of data as it is curated.
v) e.g. Gen bank, DDBJ, PDB	v) e.g. swiss-port.

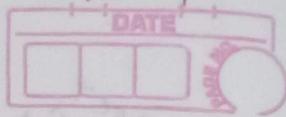
~~JMP correct the ans.~~

"> What is mean by sequence database and structure database?

→ i) Sequence Database :

A sequence database contains a collection of biological sequences, such as DNA, RNA, and protein sequences, from various organisms. These sequences are typically stored in a standardized format and can be accessed based on their identifiers, accessed based on their identifiers, accessed based on their identifiers accession numbers, or other search criteria. Sequence databases are used for performing sequence similarity searches, where a query sequence is compared against sequences in the database to find similar or homologous sequences.

e.g. GEOBank, EMBL, DDBJ.



12) Define Bioinformatics tool and explain the Sequence analysis tool in detail.

→

Bioinformatics tool are userfriendly software programs that allow researchers to analyse biological data. along with construction duration of biological database, bioinformatic also consists of development of Computational tool for Sequence, Structure and function analysis.

#### Sequence Analysis tool :

This tools are used for analysing nucleotides or protein Sequence then also use for identifying homologous sequences and understanding the evolutionary relationship between different organism. They include tools used for Sequence alignment, Sequence database searching, Motif discovering and genome assembly. Some popular Sequence analysis tool are BLAST, T-coffee, MEGA, MEME.

- 14) What is meant by function analysis tool. Explain GEO, CobRA tool box and Pathway tool.
- These tools are essential for understanding the function and relationship between different proteins and genes and identifying key pathways involved in diseases. It includes tools that are used for profiling gene expression, predicting protein-protein interaction, predicting protein sub-cellular localization and modeling metabolic pathways.

#### i) GEO :

GEO stands for Gene Expression omnibus. It's a public database maintained by the National Center for Biotechnology Information (NCBI). GEO is used for storing and sharing high-throughput gene expression and molecular abundance data. Researchers use this resource to deposit and access data from various experimental techniques like microarrays and next-generation sequencing.

#### ii) CobRA Toolbox :

The CobRA Toolbox (constraint-Based Reconstruction and Analysis Toolbox) is a computational toolkit used for function analysis of biological systems, particularly metabolic networks. It employs constraint-based modeling techniques to study how metabolic pathways function within an organism.

The CobRA Toolbox enables researchers to:

- Reconstruct Metabolic Models
- Constraint-Based Analysis
- Phenotype Prediction
- Metabolic Engineering
- Biomarker Discovery
- Drug Target Identification.

#### iii) Pathway tool :

Pathway tools are software tools used in bioinformatics for pathway analysis. They help researchers understand the complex network of biochemical reactions and interactions that occur within a biological system.

In functional analysis, pathway tools enable the following:

- Pathway Visualization
- Data Integration
- Enrichment Analysis
- Prediction of Functional Outcomes
- Comparative Analysis
- Hypothesis Generation.

15)

list out the application of Bioinformatics and explain any four.

→

Applications :

i) Genome Sequencing and Analysis :

Bioinformatics is crucial in analysing large-scale genome sequences, identifying genes, regulatory elements, and functional regions. It helps in understanding genetic variations, mutations, and their impact on diseases.

ii) Protein Structure Prediction :

Bioinformatics tools predict the 3D structure of proteins based on their amino acid sequences. This information aids in understanding protein functions, interactions, and designing drugs that target specific proteins.

iii) Phylogenetics :

Bioinformatics assists in reconstructing evolutionary relationships among species by analyzing molecular data such as DNA and protein sequences. It helps to understand the evolutionary history and relatedness between different organisms.

~~Ans~~

## Q.1. DNA ,

- - DNA is a group of molecules that is responsible for carrying and transmitting the hereditary material or the genetic instruction from parents to offsprings.
- DNA also plays crucial roles in the production of proteins.

These are 3 types of DNA.

## i) A-DNA :

It is right handed double helix similar to B-DNA form. Dehydrated DNA takes an A form protein~~s~~ts the DNA in extreme condition such as complete dryness.

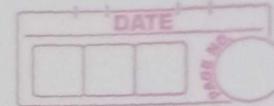
## ii) B-DNA :

It is the most common DNA conformation and right handed helix. The majority of DNA has B-type conformation under normal physiological condition.

## iii) Z-DNA :

It is left handed DNA where Double helix bits to the left in a zig zag pattern. This found in a zig zag pattern.

This found the head of start side of gene hence they need to some role in gene.



## DNA

Structure 1) DNA is double stranded molecule consisting of long chain of nucleotides.

Function 2) It transmits genetic information to make other cells and new organisms.

3) DNA is self replicating

4) Located in Nucleus of cell and in mitochondria

5) Basic Pairing : AT, CG

molecular weight

6) 2 to 6 million

7) DNA is more stable

molecule than RNA.

DNA is stable under alkaline conditions

8) DNA is vulnerable to damage by ultraviolet light

## RNA

Structure 1) RNA is single stranded molecule consisting of short chain of nucleotides

Function 3) It transfers the genetic code from nucleus to ribosomes to make proteins.

3) RNA synthesized from DNA.

4) Located in cytoplasm, nucleus and in ribosome

5) Basic Pairing : AU, GC

6) 25,000 to 2 million

7) Much more reactive

than DNA and is not stable in alkaline conditions

8) RNA is much more sensitive to damage from UV light than DNA.

Q.2 Explain regulation of gene expression.

→ Gene expression is a process by which the instruction present in DNA are converted into a functional product such as protein. During gene expression genetic code from DNA are converted into protein with help of transcription & translation.

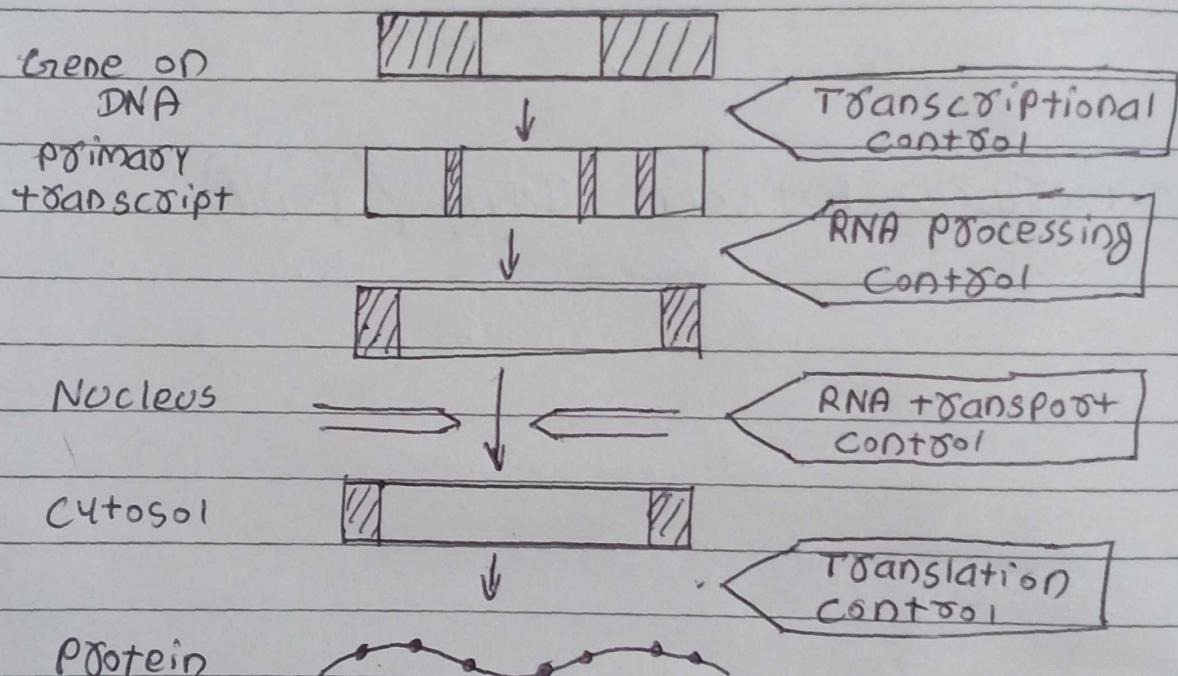


Fig. How ~~the~~ DNA is synthesis to protein

First transcription factors recognize specific sequences on DNA, acting like keys unlocking access for RNA Polymerase. RNA polymerase binds to promoter, starting the transcription process. RNA polymerase reads the DNA code, building a complementary RNA molecule (mRNA). The mRNA undergoes modifications before leaving the nucleus. Ribosomes in the cytoplasm read the mRNA code and assemble the corresponding protein.

3) what is DNA Sequencing?

→ DNA Sequencing is the method that determines the order of four nucleotide bases that make up DNA molecules and convey important genetic information.

The human genome contains 3 billion base pairs that provides the instructions for creating and maintaining human being.

Base pair structure of the DNA is best suited for the storage of vast amount of genetic information.

Data Sequencing methods:

i) Sanger Sequencing method:-

This method was discovered by English biochemist Frederick Sanger in 1970.

$N = A, T, G, C$  to prevent addition of another nucleotides.

ii) Next Generation Sequencing method:-

It is also known as massively parallel sequencing. This technology is largely separated Sanger sequences with advantages such as cost efficiency, rapidness. The NGs can determine order of millions of fragments simultaneously.

NGS is short read sequencing that requires construction of small fragment sequencing followed by deep sequencing, raw data processing, DNA sequence alignment, assembly, annotation and downstream analysis.

iii) Emerging third Generation Sequencing:  
It is known as long read sequencing including PacBio SMART sequencing Oxford Nanopore sequencing can examine billions of template of DNA and RNA and simultaneously detect variable methylation or without bias. The long read method detects more variations some of which cannot be observed by short read sequencing.

4) What is meant by Genome mapping.

→ Gene mapping refers to the process of locating genes within a genome. Gene mapping refers to the techniques used to identify genes location & distance between genes. The distance between various sites inside a gene can also be described through gene mapping.

Types of gene mapping:

Genetic linkage maps & physical maps are the two main categories of "Maps" used in genome mapping. Both maps consist of genetic markers & gene.

- Genetic linkage mapping — It shows the location of each gene on a chromosome & their relative distances from one another.

- Physical mapping — physical map always provide the actual DNA base pair distance between landmarks. It is one genome mapping approach that has a high accuracy in determining the sequences of DNA base pairs.

5) List out and explain any four implications of genomics for medical science.

→ i) Genomic medicine for personalized treatment: Genomic allows for the tailoring of medical treatment based on an individual's unique genetic makeup. e.g. knowledge of a patient's genomic profile can guide the selection of targeted therapies for cancer, improving treatment effectiveness.

ii) Disease Risk assessment & early detection:

- Genomic analysis enables the identification of genetic markers associated with various diseases including genetic predisposition.
- For instance, identifying specific genetic variations can help predict the risk of developing conditions like Alzheimer's or cardiovascular diseases, enabling proactive health care planning.

iii) Genomic Research for disease understanding and discovery :

- Genomic aids in researching the genetic basis of diseases, uncovering causative genetic mutations or alterations.
- This understanding is crucial for developing targeted treatments.

iv) Rare disease, Diagnosis and Research :

7

6) Explain the term proteomics.

→ Proteomics focuses on the study of proteins, including their identification, quantification, and the characterization and structures. This field helps in understanding proteins functions, interactions and structures.

Proteomics is the large-scale study of proteins. Proteins are complex molecules that perform essential functions in living organisms, including structure, metabolism, transport, communication and defense.

Proteomics is used to study the structure, function and interactions of proteins. Proteomics research use a variety of techniques to identify, characterize and quantify proteins. These technique includes:

- i) Mass Spectrometry.
- ii) Gel electrophoresis
- iii) Immunoassays.

Proteomics is a rapidly growing field with a wide range of applications. Proteomics research is being used to develop new diagnostic tests and treatments for diseases to understand basic of disease and to develop new drugs and therapies.

Applications:-

- i) Cancer Research
- ii) Infectious disease research
- iii) Neuroscience Research
- iv) Personalized medicine.

8'

7) What are the applications of proteomics to the medicines.

→ Here are some important applications of the proteomics in the field of medicine:-

- i) Disease Biomarker discovery.
- ii) Personalized medicine
- iii) Drug Targeted Identification.
- iv) Pharmacoproteomics.
- v) Protein - protein Interaction Analysis.
- vi) Analysis of Post-Translational Modifications.
- vii) Vaccine Development.
- viii) Drug Resistance studies.
- ix) Toxicology and adverse event Analysis.

8) Differentiate between genomics & proteomics.

	Genomics	Proteomics
Scope and focus	Study of an organism's entire genome, including genes & non-coding regions	Study of an organism's entire proteome, encompassing all expressed proteins.

	DNA Sequences	Proteins
Information Content	Genetic code, nucleotide sequences (A,G,T,C)	functional elements, protein presence, abundance, interactions, modifications.

	Understanding genetic inheritance, gene structure genotype - phenotype relationships	Understanding protein functions, cellular processes, disease biomarkers.
Purpose	Understanding genetic inheritance, gene structure genotype - phenotype relationships	Understanding protein functions, cellular processes, disease biomarkers.

	DNA Sequencing, genome mapping, PCR	Mass spectrometry, 2D gel electrophoresis, protein-protein interaction studies.
Methodology	DNA Sequencing, genome mapping, PCR	Mass spectrometry, 2D gel electrophoresis, protein-protein interaction studies.

9) Explain the term protein modeling.

→ Protein modeling, also known as protein structure modeling or protein structure prediction, is a computational approach used in bioinformatics and structural biology to predict the three-dimensional structure of a protein based on its amino acid sequence.

The primary aim or goal of protein modeling is to generate a 3D model of a protein's atomic coordinates, which can provide insight into its structure, function, and interaction with other molecules.

10) Differentiate between :-

→ Pairwise Sequence Alignment

- i) An alignment procedure comparing two biological sequences either protein, DNA or RNA.
- ii) Generally categorized as a global or local alignment method.
- iii) comparatively simple algorithm is used.
- iv) Applications:- primarily to find out conserved regions between the two sequences.
- v) e.g.: BLAST, LALIGN, EMBOSS, Needle.

Multiple sequence Alignment.

- i) An alignment comparing three or more biological sequences either proteins, DNA or RNA.
- ii) Generally categorized as global multiple sequence method.
- iii) complex sophisticated algorithm is used.
- iv) To detect regions of variability or consecutive in family of protein.
- e.g.: MUSCLE, T-coffee, MARRT.

11) Explain dynamic programming method of pairwise sequence alignment with example.

- 
- Dynamic programming is used to find the optimal alignment between two proteins or nucleic acid sequences by comparing all possible pairs of characters in the sequences.
  - Dynamic programming can be used to produce both global & local alignment. The global pairwise alignment algorithm using dynamic programming is based on Needleman-Wunch algorithm, while the dynamic programming local alignment is based on the Smith Waterman algorithm.

The method works in following three steps:-

1) Initialization:

The first step in global alignment dynamic programming approach is to create a matrix with  $M+1$  columns &  $N+1$  rows where  $M \leq N$  corresponds to the size of the sequences to be aligned.

2) Matrix fillings:

Matrix is filled with, if the characters are match then the value be 1, if characters in sequences are mismatch the value be 0 and for gap the value be -1 and for no gap value be zero.

In this way complete matrix is filled with numbers.

3) 3rd point after question no. 16.

12) Explain BLAST Method with example.

- BLAST stands for Basic Local Alignment Search tool. It is a local alignment algorithm based tool used for aligning multiple sequences and finding similarities or dissimilarities among various species.
- BLAST is a heuristic method which means that is the dynamic programming algorithm that is fasted, efficient but relatively less sensitive.

*Alg with Ed. Series*

BLAST first breaks the query sequence into smaller fragments, then it compares each fragment to the sequences in the database. BLAST extends the alignment to see if there is a longer match.

USE :- The scientists can use BLAST to search for similar genes in other organisms, the scientists can hypothesize that these genes are involved in the same disease pathway.

13) Explain K-tuple method with example.

- The K-tuple method is simple but effective way to identify sequences. It works by dividing the sequences into K-tuples and identical, then the corresponding sequences are considered to be similar.
- EX:-

$$S_1 = "ACG T A C G T"$$

$$S_2 = "A C G A C G T"$$

for,  $K=3$

$$S_1 = ACG, CGT, GTA, TAC, ACG, CGT$$

$$S_2 = ACG, CGA, GAC, ACG, CGT$$

We can compare the K-tuple for  $S_1$  &  $S_2$ . We find that the K-tuple "ACG" is present in both  $S_1$  &  $S_2$ .

This indicates sequence  $S_1$  &  $S_2$  are similar.

14) Differentiate between BLAST & K-tuple method.

→

BLAST	K-tuple
i) Find similar sequences in a database	i) A unit used in various sequence analysis task.
ii) Can be used for DNA, RNA and protein sequences.	ii) A general concept application to various sequence types.
iii) Searches for local and global sequence similarities.	iii) Represents a fixed-length sequence segment.
iv) Returns a list of similar sequences from a database along with statistical measures.	iv) Results in lists of matching K-tuples or sequence segment in analysed sequence.

15) Explain the term dynamic programming method of multiple sequence alignment.

→ Dynamic programming is a fundamental technique used in bioinformatics, including MSA in the context of MSA. Dynamic programming is often employed to efficiently compute the optimal alignment of three or more biological sequences.

Here is how dp works in MSA:

i) Problem formulation :- Given multiple sequences, goal is to align them in a way that maximizes the similarities.

ii) Scoring System :- Define a scoring system that assigns score to different possibilities such as match, mismatch, gaps, etc.

iii) Dynamic programming Table:- Create a multidimensional table to store intermediate results & ultimately compute the optimal alignment score.

iv) Filling the table :- Start filling table using a dp often based on the Needleman-Wunsch or Smith Waterman algorithm.

v) Backtracking :- After table filled, backtrack through table from last cell to reconstruct optimal alignments.

vi) Optimal Alignments :- They are set of sequences aligned in a way that maximizes the overall alignment score based on defined scoring system.

16) Explain the term cluster-W.

→ Cluster W is a tool for aligning multiple protein or nucleotide sequences. This alignment is achieved via three steps:-

- Pairwise alignment.
- Grid tree generation.
- Progressive alignment.

- Cluster W-MPI is a distributed & parallelism implication of cluster W.

- Cluster W uses progressive alignment methods align most similar sequences first and work their way down to the last similar sequences until a global alignment is created.

- Cluster W is a matrix based algorithm tools like T-coffee and Dialign are consistent base cluster W is fairly efficient algorithm competes - against other software.

*Done by:  
9/6/2023*

question

## ① Define phylogenetic analysis?

→ A phylogenetic analysis is the study of the evolutionary development of a species or a group of organisms or a particular characteristic of organisms.

question

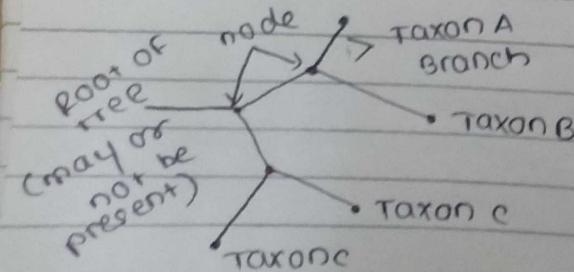
## ② List out? and explain fundamental concept of phylogenetic tree.

→ Phylogenetic tree, also known as an evolutionary tree, is a diagram that depicts the evolutionary relationship between different species or genes.

It's like a family tree for the entire living world, tracing back lineages through shared ancestry.

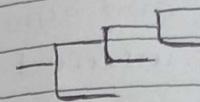
properties.

- Root, representing the oldest common ancestor of all the groups included
- Branches, representing a lineage that diverged at a specific point in time
- Leaves, represent the present day species or genes being studied.



Types.

① Rooted tree → have a designated root representing the oldest common ancestor. most commonly used



② unrooted tree → do not specify a root, focusing only on the relative relationship between groups



question

## ③ process of phylogenetic analysis.

→

## (i) choosing the right markers

- Homologous sequence → you need sequences from a common ancestor, like genes or protein segment
- informative sequence → The chosen markers should have enough variation to distinguish between lineages.

## (ii) Gathering data →

- sequencing DNA or proteins
- collecting morphological data

## (iii) Alignment → multiple sequence alignment aligns the chosen sequences accounting for insertions, deletions and mutations.

(iv) Tree building → Algorithms analyze the aligned sequences to infer evolutionary relationships.

(v) Evaluating the tree → Bootstrapping and other statistical test assess the robustness of inferred tree.

(vi) Interpreting the tree → may reflect branches

(vii) visualization and analysis

question

④ Explain tree building methods.

→

(i) distance method

(ii) character method

maximum parsimony  
maximum likelihood

(iii) bayesian inference.

(i) distance method → calculate distance between taxa and build a tree minimizing total distance. Distance are then used to construct a tree that minimizes the total distance between all pairs.

(ii) character Method → This method consider the presence or absence of specific characters in different taxa. goal is to find best tree

(iii) Bayesian Interface → considers uncertainty in data and models to find the most probable tree.

question

⑤ Define paralogs & orthologs

→

(i) paralogs → paralogs are the genes that diverge within the same species. They are multiple genes (tree leafs) per species.

genes

(ii) orthologs → are the species the diverge with species.

question

⑥ Define microarray

→ microarray is essentially a miniaturized laboratory on a chip, used to investigate thousands of biological interactions simultaneously.

Applications

- gene expression studies
- cancer research
- microbial identification

question

## ① concept of tree evaluation



### (i) Phylogenetic tree Evaluation

(ii) Decision tree Evaluation → access the predictive performance of decision tree used for tasks like classifying genes. (cross validation, accuracy etc)

(iii) Expression tree Evaluation → Evaluate tree representing gene expression patterns or regulatory networks

(iv) hierarchical clustering tree evaluation → It assess the quality of tree generated by hierarchical clustering algorithms used to group similar biological entities. (genes, proteins, samples)

(v) other tree based models

random forest  
gradient boosted tree.

question

## ② Microarray Techniques

→ microarray are like tiny lab-on-chips with thousands of mini detectors. They scan samples for specific target, like gene or proteins, using fluorescent labels. By analyzing the

Date:

P. No:

Date:

P. No:

glowing patterns, scientists gain insights into disease, how genes work, and evolution.

- detectors (probes) wait on the chip
- sample with labeled cues (RNA/DNA) is added
- detectors bind their matching cues, glowing brightly
- the chip is scanned, revealing the detector cue matches.

question

## ③ applications of microarray

- (i) Gene expression profiling
- (ii) mutation & SNP detection
- (iii) comparative genomics
- (iv) drug development
- (v) clinical diagnostics
- (vi) Agriculture & plant Research
- (vii) forensic science.