UNIT 1 BIOINFORMATICS AND BIOSTATISTICS

UNIT-I: Scope of Computers in Current Biological Research

Scope of Computers in Biological Research

- Data Analysis: Processing and analyzing large datasets, such as genomic sequences, protein structures, and biological pathways.
- Modeling and Simulation: Creating computational models of biological systems to understand complex interactions and predict behaviors.
- Bioinformatics: Using software tools for gene sequencing, molecular modeling, and database management.
- Image Analysis: Enhancing and analyzing images from microscopy and other imaging techniques.
- Machine Learning: Applying algorithms to classify biological data, predict outcomes, and discover new patterns.

Basic Operations and Architecture of Computers

Basic Operations

- Input: Receiving data from input devices like keyboards, mice, and scanners.
- Processing: Performing operations on the data using the central processing unit (CPU).
- Storage: Saving data in memory units like RAM (temporary) and hard drives (permanent).
- Output: Sending processed data to output devices like monitors, printers, and speakers.

Architecture of Computers

- CPU (Central Processing Unit): The brain of the computer, executing instructions and performing calculations.
 - ALU (Arithmetic Logic Unit): Handles arithmetic and logical operations.
 - Control Unit: Directs the operation of the processor.
 - Registers: Small, fast storage locations for immediate data processing.
- Memory:
 - Primary Memory (RAM): Volatile memory used for temporary data storage.
 - Secondary Memory (Hard Drives, SSDs): Non-volatile storage for long-term data.
- Input/Output Devices: Interface for data entry and retrieval (e.g., keyboard, mouse, display).
- Bus Systems: Pathways that transmit data between components.

Introduction to Digital Computers

Organization of Digital Computers

- Motherboard: The main circuit board connecting all components.
- Power Supply: Provides power to the computer.
- Peripherals: External devices like printers, external drives, and input devices.

Name of the Faculty: **Sri E Bharat Raju** Head &Lecturer in Biotechnology

Study material for MSc

Programming Languages

Low-Level and High-Level Languages

- Low-Level Languages:
 - Machine Language: Binary code directly executed by the CPU.
 - Assembly Language: Symbolic code translated into machine language using an assembler.
- High-Level Languages:
 - C, C++, Java, Python: More abstract and closer to human language, translated to machine code using compilers or interpreters.

Binary Number System

- Basics:
 - Binary Digits (Bits): 0 and 1.
 - Binary to Decimal Conversion: Summing powers of 2.
 - Decimal to Binary Conversion: Dividing by 2 and recording remainders.

The Soft Side of the Computer

Different Operating Systems

- Windows:
 - Features: User-friendly interface, widespread use, compatibility with numerous applications.
 - \circ Applications: Office work, gaming, general computing.
- Linux:
 - Features: Open-source, customizable, secure, various distributions (Ubuntu, Fedora, CentOS).
 - Applications: Servers, development environments, scientific computing.

Introduction to Programming in C

Basics of C Programming

- Structure of a C Program:
 - Header Files:#include <stdio.h>
 - Main Function:int main () { }
- Data Types:int, float, char, double.
- Variables and Constants: Declaring and initializing variables.
- Operators: Arithmetic, relational, logical.
- Control Structures: if, else, while, for, switch.
- Functions: Defining and calling functions, passing arguments.
- Input/Output:printf(), scanf().

Introduction to Internet and Its Applications

Basics of the Internet

- Definition: Global network connecting millions of private, public, academic, business, and government networks.
- Components: Routers, servers, protocols (TCP/IP).

Name of the Faculty: **Sri E Bharat Raju** Head &Lecturer in Biotechnology

Study material for MSc

Internet Applications

- Email: Communication through electronic messages.
- World Wide Web: Accessing and sharing information through websites and browsers.
- File Transfer Protocol (FTP): Transferring files between computers.
- Search Engines: Tools for finding information on the web (Google, Bing).
- Online Databases: Accessing scientific databases (PubMed, GenBank).
- Social media: Platforms for social interaction and information sharing.

Name of the Faculty: **Sri E Bharat Raju** Head &Lecturer in Biotechnology

Study material for MSc

UNIT-II: INTRODUCTION TO BIOINFORMATICS, GENOMICS, AND PROTEOMICS

Introduction to Bioinformatics

Definition: Bioinformatics is an interdisciplinary field that develops and applies computational methods to analyze biological data.

Key Areas:

- Data Analysis: Interpreting biological data (DNA, RNA, protein sequences).
- Data Management: Storing and retrieving large datasets.
- Algorithm Development: Creating tools for data analysis and interpretation.
- Visualization: Presenting data in understandable formats (graphs, models).

Applications:

- Genomic sequencing and annotation.
- Protein structure prediction.
- Drug discovery and development.
- Comparative genomics.

Genomics

Definition: Genomics is the study of the complete set of DNA (including all of its genes) in an organism.

Key Concepts:

- Genome Sequencing: Determining the order of nucleotides in an organism's DNA.
- Genomic Annotation: Identifying genes and other functional elements in the genome.
- Comparative Genomics: Comparing genomes between different species.

Applications:

- Identifying disease genes.
- Evolutionary biology studies.
- Personalized medicine.

Proteomics

Definition: Proteomics is the large-scale study of proteins, particularly their structures and functions.

Key Concepts:

Name of the Faculty: **Sri E Bharat Raju** Head &Lecturer in Biotechnology

Study material for MSc

- Protein Expression: Analysis of protein levels under different conditions.
- Protein-Protein Interactions: Studying how proteins interact within a cell.
- Post-translational Modifications: Investigating changes to proteins after synthesis.

Applications:

- Biomarker discovery.
- Understanding cellular processes.
- Drug target identification.

UNIT-II: Bioinformatics Tools, Databases, and Projects

Bioinformatics Tools

Online Tools

- **BLAST (Basic Local Alignment Search Tool):** Compares nucleotide or protein sequences to sequence databases.
- FASTA: Similar to BLAST, but uses different algorithms for sequence alignment.

Offline Tools

- Bioconductor: Open-source software for bioinformatics.
- **EMBOSS (European Molecular Biology Open Software Suite):** Provides tools for sequence analysis.

Biological Databases

Overview: Biological databases store and organize biological data. They are essential for retrieving, managing, and analyzing biological information.

Types of Biological Databases:

- **NCBI** (National Center for Biotechnology Information): Houses a collection of databases relevant to biotechnology and biomedicine.
- **EMBL (European Molecular Biology Laboratory):** Provides access to various biological data, including sequence data.
- GenBank: A nucleotide sequence database maintained by NCBI.
- Swiss-Prot: A manually curated protein sequence database.
- **PDB** (**Protein Data Bank**): Archive of 3D structural data of biological molecules.

Database Searching Using BLAST and FASTA

BLAST:

- Function: Finds regions of similarity between biological sequences.
- Applications: Identifying species, gene functions, and evolutionary relationships.

FASTA:

Name of the Faculty: **Sri E Bharat Raju** Head &Lecturer in Biotechnology

Study material for MSc

- Function: Aligns sequences and finds regions of similarity.
- Applications: Similar to BLAST, used for sequence alignment and database searching.

Human Genome Project

Overview:

- An international research initiative aimed at mapping and understanding all the genes of the human genome.
- Completed in 2003, it provided a complete and accurate sequence of the 3 billion DNA base pairs in the human genome.

Goals:

- Identify all the approximately 20,000-25,000 genes in human DNA.
- Determine the sequences of the 3 billion chemical base pairs that make up human DNA.
- Store this information in databases.
- Improve tools for data analysis.
- Transfer related technologies to the private sector.
- Address ethical, legal, and social issues that may arise from the project.

Name of the Faculty: **Sri E Bharat Raju** Head &Lecturer in Biotechnology

Study material for MSc

UNIT IIISEQUENCE ALIGNMENT

Introduction and Significance of Sequence Alignments

Sequence Alignment:

- **Definition:** A method of arranging sequences of DNA, RNA, or protein to identify regions of similarity.
- Purpose: To infer functional, structural, or evolutionary relationships between the sequences.

Significance:

- Identifying homologous sequences.
- Inferring functional and structural annotations.
- Studying evolutionary relationships.

Pairwise Sequence Alignment

Types:

- **Global Alignment:** Aligns entire sequences end-to-end (e.g., Needleman-Wunsch algorithm).
- Local Alignment: Aligns subsequences within larger sequences (e.g., Smith-Waterman algorithm).

Applications:

- Comparing gene or protein sequences.
- Identifying conserved domains.

Multiple Sequence Alignment (MSA)

Principles:

- Aligning three or more sequences simultaneously.
- Identifying conserved regions across multiple sequences.

Tools:

- Clustal Omega
- MUSCLE
- MAFFT

Name of the Faculty: **Sri E Bharat Raju** Head &Lecturer in Biotechnology

Study material for MSc

Applications:

- Constructing phylogenetic trees.
- Identifying conserved motifs and domains.

Gene and Genome Annotation

Introduction:

• The process of identifying elements within a genome and attaching biological information to them.

Tools Used:

- Gene Prediction: GENSCAN, AUGUSTUS.
- Genome Annotation: Ensembl, NCBI Genome Annotation Pipeline.
- Functional Annotation: BLAST, InterProScan.

Physical Map of Genomes

Definition:

• A physical map represents the physical distances between landmarks on the genome, such as genes and markers.

Techniques:

- Restriction mapping.
- Fluorescence in situ hybridization (FISH).
- Radiation hybrid mapping.

Protein Secondary Structure Prediction

Principles:

• Predicting the local structures (α -helices, β -sheets) within a protein sequence.

Tools:

- PSIPRED
- JPred
- SOPMA

Applications:

- Understanding protein function.
- Guiding experimental structure determination.

Name of the Faculty: **Sri E Bharat Raju** Head &Lecturer in Biotechnology

Study material for MSc

Protein 3D Structure Prediction

Principles:

• Predicting the three-dimensional structure of a protein from its amino acid sequence.

Methods:

- Homology modeling
- Threading (fold recognition)
- Ab initio prediction

Tools:

- SWISS-MODEL
- I-TASSER
- Rosetta

Applications:

- Drug design
- Understanding protein function and interactions

Protein Docking

Principles:

• Predicting the preferred orientation of one molecule to a second when bound to each other to form a stable complex.

Tools:

- AutoDock
- ClusPro
- HADDOCK

Applications:

- Drug discovery
- Studying protein-protein and protein-ligand interactions

Introduction to Homology Modeling

Principles:

• Predicting a protein's 3D structure based on the known structure of a homologous protein.

Steps:

Name of the Faculty: **Sri E Bharat Raju** Head &Lecturer in Biotechnology

Study material for MSc

- Identify templates
- Align target sequence with template
- Build model
- Refine and validate model

Tools:

- MODELLER
- SWISS-MODEL
- Phyre2

Applications:

- Functional annotation
- Drug design

Computer-Aided Drug Design (CADD) in Drug Discovery

Principles:

• Using computational methods to discover, design, and optimize drugs.

Types:

- Structure-Based Drug Design (SBDD): Uses the 3D structure of a target protein.
- Ligand-Based Drug Design (LBDD): Uses knowledge of other molecules that bind to the target.

Tools:

- Schrodinger Suite
- AutoDock
- MOE (Molecular Operating Environment)

Applications:

- Identifying lead compounds
- Optimizing drug candidates

Molecular Phylogeny

Concept:

• The study of evolutionary relationships among biological entities, often using genetic data.

Methods of Tree Construction:

1. Distance-Based Methods:

- UPGMA (Unweighted Pair Group Method with Arithmetic Mean)
- Neighbor-Joining

Name of the Faculty: **Sri E Bharat Raju** Head &Lecturer in Biotechnology Study material for MSc

2. Character-Based Methods:

- Maximum Parsimony
- Maximum Likelihood
- Bayesian Inference

Tools:

- MEGA (Molecular Evolutionary Genetics Analysis)
- PHYLIP (Phylogeny Inference Package)
- RAxML (Randomized Axelerated Maximum Likelihood)

Applications:

- Studying evolutionary relationships
- Understanding genetic diversity and species evolution

UNIT-IV: STATISTICAL METHODS IN DATA ANALYSIS

Brief Description and Tabulation of Data

Description of Data:

- Qualitative Data: Descriptive data (e.g., colors, labels).
- Quantitative Data: Numerical data (e.g., height, weight).
 - **Discrete Data:** Countable numbers (e.g., number of students).
 - **Continuous Data:** Measurable quantities (e.g., temperature).

Tabulation of Data:

- Organizing data into tables for better understanding.
 - Frequency Distribution Table: Shows how often each value occurs.
 - Class Interval: Range of values.
 - **Frequency:** Number of occurrences within each interval.

Graphical Representation of Data:

- Bar Chart: For categorical data.
- Histogram: For continuous data, showing frequency distribution.
- **Pie Chart:** For proportional data.
- Line Graph: For showing trends over time.
- Scatter Plot: For showing relationships between two variables.

Measures of Central Tendency

Mean (Arithmetic Average): Mean $(\bar{x}) = \sum_{i=1}^{n} x_i$

• Sensitive to extreme values (outliers).

Median:

Name of the Faculty: **Sri E Bharat Raju** Head &Lecturer in Biotechnology

Study material for MSc

- Middle value when data is ordered.
- If even number of observations, median is the average of the two middle values.
- Not affected by outliers.

Mode:

- Most frequently occurring value in a data set.
- Can have more than one mode (bimodal, multimodal).

Measures of Dispersion

Range: Range = Maximum value – Minimum value

Simplest measure of dispersion.

Variance (σ^2 for population, s^2 for sample):Sample Variances² = $\frac{\sum_{i=1}^{n} (x_i - \overline{x})^2}{n-1}$

• Measures the average squared deviation from the mean.

Standard Deviation (σ for population, s for sample): Sample Standard Deviation(s)= $\sqrt{s^2}$

• Provides a measure of the spread of data around the mean.

Simple Linear Regression and Correlation

Simple Linear Regression:

- Predicts the value of a dependent variable (Y) based on the value of an independent variable (X).
- Equation: Y = a + bXY
 - *a*: Intercept
 - b: Slope

Correlation:

- Measures the strength and direction of the relationship between two variables.
 - **Pearson Correlation Coefficient (r)**
 - \circ *r* ranges from -1 to 1.
 - \circ r=1 Perfect positive correlation.
 - \circ r = -1 Perfect negative correlation.
 - \circ r = 0 No correlation.

Name of the Faculty: **Sri E Bharat Raju** Head &Lecturer in Biotechnology

Study material for MSc

Types of Errors and Level of Significance

Types of Errors:

- **Type I Error** (*α*): Rejecting a true null hypothesis (false positive).
- **Type II Error** (β): Failing to reject a false null hypothesis (false negative).

Level of Significance (a):

- The probability of making a Type I error.
- Common levels: 0.05, 0.01.

Tests of Significance

t-test:

- Compares means between two groups.
 - Independent t-test: For two independent groups.
 - **Paired t-test:** For related groups (e.g., pre-test and post-test).

Chi-square test:

- Tests the association between categorical variables.
- **Chi-square statistic:** $\chi 2 = \frac{\sum (Oi Ei)^2}{E_i}$
 - *0i*: Observed frequency.
 - *Ei*: Expected frequency.

ANOVA (Analysis of Variance):

- Compares means among three or more groups.
- **F-statistic:** $F = \frac{between-group variance}{within-group variance}$

Name of the Faculty: **Sri E Bharat Raju** Head &Lecturer in Biotechnology

Study material for MSc