

①



Date
Page

Introduction to Bioinformatics.

Bioinformatics is an interdisciplinary field mainly involving molecular biology and genetics, computer science, mathematics, and statistics. Large scale biological problems are addressed from a computational point of view. The most common problems are modeling biological processes at the molecular level and making inferences from collected data. A bioinformatics solution usually involves the following steps.

- Collect statistics from biological data.
- Build a computational model.
- Solve a computational modeling problem.
- Test and evaluate a computational algorithm.

Sequence analysis is the analysis of DNA and protein sequence clues regarding function and includes subproblems such as identification of homologs, multiple sequence alignment, searching sequence patterns and evolutionary analysis.

Protein structure are three-dimensional data and the associated problems are structure prediction (secondary and tertiary), analysis of protein structure for clues regarding function and structural alignment.



Gene expression data is usually represented as matrices and analysis of microarray data mostly involves statistics analysis, classification and clustering approaches. Biological networks such as gene regulatory networks, metabolic pathways and protein-protein interaction networks are usually modeled as graph and graph theoretic approaches are used to solve associated problems such as construction and analysis of large scale - networks.

History

The term "Bioinformatics" was invented by Paulien Hogeweg and Ben Hesper in 1970.

- The collection of amino acid sequence was compiled in the Atlas of protein sequence and structure by the National biomedical foundation.
- This collection was edited by Margaret D. Dayhoff from 1965 to 1978.
- Dayhoff and coworkers contribute to the consolidation of amino acid sequences by developing computer software for detecting distantly related sequence.



- The EMBL established their data library in 1980. to collect, organize and distribute nucleotide sequence data and related information.
- NCBI was established in U.S.A. NCBI serves as primary information databank and provider of information.
- The National Biomedical Research Foundation established in PIR in 1984.

DNA Sequence

- The symbol used to represent DNA sequence data.
- The four bases denoted by a single letter A (Adenine), C, (Cytosine), T (Thymine) and ~~uracil~~ G (Guanine).
- But often data sequence contain ambiguities in that it is not clear as to which of the four base present at several positions.
- For example sequence data may indicate that the base present at a specific position may be either G or A, it is purine.
- Similarly, if a position may have either C or T it is pyrimidine.
- The base sequence of the two complementary strands of a DNA molecule are represented by this system of a symbol.



Amino Acid sequence of protein: -

- The amino acid conventionally represented by three letter symbols, e.g. - Ala for Alanine, Val for Valine, Valine etc.
- But in Bioinformatics they are denoted by single letter e.g. A for Alanine, C for Cysteine, D as Aspartic Acid etc.
- But some positions in protein sequences have ambiguities this sequence is comparable to that for DNA sequences.
- For e.g. it may not be clear that a position has glutamine or glutamic acid the symbol is given the symbol Z.
- The protein synthesis begin at the N-terminus and proceeds to the C-terminus.
- The amino acid sequences in databases are listed from the N-terminus to the C-terminus of the polypeptide.



Types of Sequences in Nucleotide Sequence Databases.

- The databases on DNA sequences contain a different types cDNA sequences:-
- A cDNA molecule is obtained by reverse transcription of an RNA molecule.
- The cDNA sequences, therefore represent that part of the genome that is transcribed into RNA.
- If the cDNA is obtained from mRNA, it will represent only the exon sequences of the gene sequence expressed in the concerned cell / tissue / organisms.

Genomic DNA Sequence :-

- These sequence represent the complete genome of the organisms.
- When the genome sequences is completed, it will contain the sequences of the entire genome of the organisms.
- In case of prokaryotes genome consists of usually, a single chromosome while in case of eukaryotes it relates to the nuclear DNA.

#



Expressed Sequence Tag (EST) sequences:-

- The sequences are obtained by sequencing only a part of the cDNA molecules produced using mRNA.
- These sequences are dubbed as 'Tags' because they can be used as probe for the isolation of the concerned genes from the genomic DNA.
- This approach was used by J. Craig Ventol and his group for obtaining the sequence of expressed portion of human genome.
- The EST technique generated enormous sequence data that permitted the construction of a preliminary transcript map of the human genome.

Genome Sequence Tag (GST)

GST were developed for the identifying the genes of *Plasmodium falciparum*.

- It was observed that the enzyme *Mung bean nuclease* (Mnase) cuts *P. falciparum* genomic DNA b/w genes.
- GSTs are developed by sequencing the DNA fragments on either side of the points of cuts generated by Mnase.



Organelles DNA Sequence :-

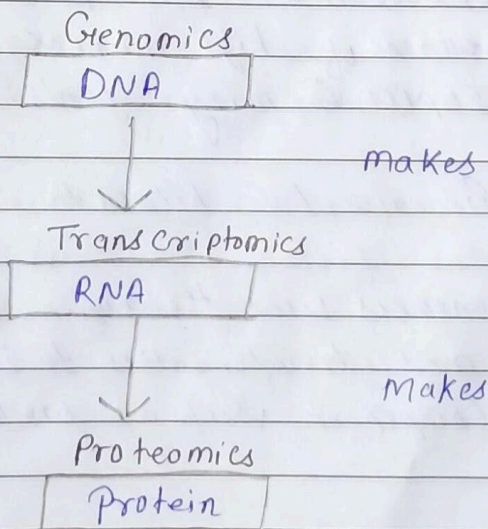
- Organelles DNA is the DNA found in mitochondria (mDNA) and chloroplast (cpDNA).
- The sequence of the data are compiled in database.

Branches of Bioinformatics

A living cell is a system whose cellular components such as genome, the gene transcript, and the proteins interact with each other, and these interactions determine the fate of the cell, e.g. - whether a stem cell is going to become a liver cell or a cancer cell.

The Three Branches of Bioinformatics

- (i) Genomics
- (ii) Transcriptomics
- (iii) Proteomics





Genomics

Genomics plays a significant role in modern biological research in which the nucleotide sequences of all the chromosomes of an organism are mapped and the location of different genes and their sequence is determined.

- This involves extensive analysis of the nucleic acids through molecular biology techniques before the data are ready for processing by computers.
- It is a science that attempts to describe a living organism in terms of the sequence of its genome.
- It was not reliable to estimate the no. of genes in an organism based on the no. of nucleotide base pair. because of the presence of high numbers of redundant copies of many genes.
- Genomics has helped to rectify this problem.
- Genomics uses technique of molecular biology and bioinformatics to identify cellular component such as proteins, rRNA, tRNA etc.



and analyse the sequence and non-coding sequence.

- The first automatic DNA Sequences was developed in 1986 by Leroy Hood.
- Haemophilus influenza was the first bacterium to be sequenced in 1995.
- Even if one can identify all the genes on a genome, the genes only indicate that at some point in time, it might be transcribed to produce Cellular Component.
- sig - A human genome contains about 30,000 to 60,000 protein coding genes, but only a subset of them is expressed in a particular cell type at a particular time.

#(i) Transcriptomics

Transcriptomics is the study of the transcriptome which include the whole set of mRNA molecules in one or a population of biological cells.

- This study help us to depict the expression level of genes, often using techniques such as DNA micro arrays, that is capable of sampling ten thousands of different mRNAs at a time.



- This kind of new technique has helped biologist to routinely monitor the gene expression between the control cells and treatment cells.
- Transcriptomics has a few limitations.
- The relative abundance of transcripts as characterized by the sequential analysis of gene expression (SAGE) is misusually experiments.

Proteomics / Proteomnce

- proteomics represents the earliest to identify a major sub-class of a cellular component, the proteins and their interactions.
- Proteomics involves the sequencing of amino acid in a protein determining its 3D structure and relating it to the function of protein.
- Before computer processing comes into the picture, extensive data, particularly through crystallography and nuclear magnetic resonance (NMR)



- With such data known as proteins, the structure and its relationship to the function of newly discovered protein.
- In such areas, bioinformatics has enormous analytical and predictive potential.
- Metabolic proteins such as haemoglobin and insulin have been subjected to intensive proteomic investigation.
- The term 'proteomics' was coined to make an analogy with genomics.
- Scientists feel that the bioinformatics of protein is crucial, to understand the cellular components and the interactions completely.

Aims of Bioinformatics

- The various important ways in which bioinformatics can be used, such as include data acquisition tool and database development, data analysis and data integration, and data management.

Data Acquisition

- Data Acquisition is primarily concerned with accessing and storing data generated directly from the biological experiments
- The data generated by various sequencing projects have to be retrieved in the appropriate format, and capable of being linked to all information related to the DNA samples.
- The data are organized in different databases so that the researchers can access existing information.

Tool and Database development

- Many laboratories generate large volume of data such as DNA sequence, gene expression information, 3D molecular structure and highly-throughput screening.
- Consequently, they must develop effective databases for storing and quickly accessing data. The other aim is to develop tools and resources that aid in the analysis of data.

Data Analysis

- The third aim is to use these tool to analyse the data and interpret the results in a biologically meaningful manner. Efficient analysis require an efficiently designed database.
- It must allow researchers to place their query effectively and provide them with ^{all} the information they need to begin their data analysis.
- If queries cannot be performed, or if the performance is too slow, the whole system breaks down since scientists will not be inclined to use the databases.

Data Integration

- Once information has been analysed, a researcher must often associate or integrate it with the related data from the other databases.
- for example Scientists may run a series of gene expression analysis experiments and observe that a particular set of 100 genes is more highly expressed in a cancerous lung tissue than in a normal lung tissue.
- The scientist may wonder which of the genes is most likely to be truly selected to the disease.



Bioinformatics Application molecular medicine:-

- The human genome will have profound effects on the fields of biomedical research and clinical medicine. Every disease has a genetic component.
- This may be inherited as a result of the body's response to an environmental stress which causes alterations in the genome (eg - Cancer, heart disease, diabetes).
- The completion of the human genome means that we can search for the genes directly associated with different diseases and begin to understand the molecular basis of these diseases more clearly.
- This new knowledge of the molecular mechanism of disease will enable better treatment, cures and even preventive tests to be developed.



• Personalised medicine

- Clinical medicine will become more personalised with the development of the field of Pharmacogenomics.
- This is the study of how an individual's genetic inheritance affects the body's response to drug.
- At present, some drugs fail to make it to the market because a small percentage of the clinical patient population show adverse to a drug due to sequence variants in their DNA.
- As a result, potentially life saving drugs never make it to the market place.
Today's doctor have to use trial and error to the best drug to treat a particular patients as those with the same clinical symptoms can show a wide range of response to the same treatment.

Drug development

- At present all drugs on the market target only about 500 proteins.
- With an improved understanding of disease mechanisms and using computational tool to identify and validate new drug targets, more specific medicines that act on the cause, not merely the symptoms of the disease can be developed.



- These highly specific drugs promise to have fewer side effects than many of today's medicines.

Gene therapy

- In the not too distant future, the potential for using genes themselves to treat diseases may become a reality.
- Gene therapy is the approach used to treat, cure or even prevent disease by changing the expression of a person's genes.
- Currently, this field is in its infantile stage with clinical trials for many different types of cancer and other diseases ongoing.

The Reality of Bio weapon: -

- Scientists have recently built the virus poliomyelitis using entirely artificial means.
- They did this using genomic data available on the internet and materials from a mail-order chemical supply.
- The research was financed by the US Department of defence as part of bio warfare response program to provide the world the reality of bio weapon.



The researchers also hope their work will discourage officials from ever relaxing programs of immunisation. This project has been met with very mixed feelings.

Antibiotic Resistance

- Scientists have been examining the genome of *Enterococcus faecalis*, leading cause of bacterial infection among hospital patients.
- They have discovered a virulence region made up of a number of antibiotic-resistant genes that may contribute to the bacterium's transformation from a harmless gut bacteria to a menacing invader.
- The discovery of the region, known as pathogenicity island, could provide useful markers for detecting pathogenic strains and help to establish controls to prevent the spread of infection in wards.

The End