



UNIVERSITÀ
DEGLI STUDI DI MILANO-BICOCCA

SYLLABUS DEL CORSO

Analisi dei Dati

2021-1-F0901D043

Aims

Basic knowledge of the most important statistical-methodological tools of the descriptive and inferential statistics for: design of experiments, data collection and analysis, interpretation of scientific literature. Introduction to the main problems related to the computational analysis of biological sequences (DNA, RNA, proteins).

Contents

The goal of the course is to contribute to the education of the medical biotechnologist in order to be able to:

- 1) understand the principles of the experimental design in medicine and biology
- 2) understand the most important statistical techniques for data analysis
- 3) use a software for data analysis
- 4) understand the literature presenting results from statistical analysis
- 5) understand the motivations, problems and methodologies.
- 6) be introduced to NGS technologies
- 7) be able to_____
- 8) understand the main data analysis techniques: genome reconstruction and annotation; sequence comparison: global, local and multiple alignment algorithms; reconstruction of phylogenies; transcriptome analysis.

Detailed program

The module of Biostatistics is organized in three parts: descriptive statistics, inferential statistics, and interpretation of scientific literature. The first and the second part share the following characteristics:

- 1) inclusion of methodological aspects of study design and programming of experiments
- 2) are thought using motivating examples from the applied literature
- 3) involves the STATA package

Part one – Basic descriptive statistics, graphical representation of quantitative and qualitative variables, indicators of position and variability, Gaussian distribution, concepts of probability.

Part two– Basics on inferential statistics, Hypothesis testing on continuous variables, T test for paired and unpaired data, test on association between categorical variables, Chi square test, McNemar test, analysis of variance, sample size and power.

Part three – Reading, interpretation, comment of scientific papers.

The module of Bioinformatics is organized in 8 chapters.

1. Data management in life sciences
2. Basics of informatics
 - 2.1. Algorithms and programs
 - 2.2. Alphabets, word, graphs
 - 2.3. Databases
3. The NGS technology
 - 3.1. Second generation NGS platforms
 - 3.2. Third generation NGS platforms
 - 3.3. Genomic data formats
 - 3.4. Genome reconstruction and annotation
4. Basi di dati di sequenze molecolari
 - 4.1. Genomic databases (EMBL – GenBank)
 - 4.2. Protein databases (SwissProt, PDB)
 - 4.3. Database query systems
5. Sequence Analysis in molecular biology

- 5.1. Exact String matching algorithms
- 5.2. Sequence alignments
 - 5.2.1. Motivations
 - 5.2.2. Dot matrices
 - 5.2.3. Substitution matrices (PAM, BLOSUM)
 - 5.2.4. Global alignment: Needleman-Wunsch Algorithm
 - 5.2.5. Local alignment: Smith-Waterman Algorithm
 - 5.2.6. Euristic Algorithms: BLAST, Fasta, BWA
 - 5.2.7. Multiple alignment algorithms; CLUSTALW
- 6. Functional motifs finding in sequences
 - 6.1. Suffix trees
 - 6.2. Pattern discovery algorithms
- 7. Transcriptome Analysis
 - 7.1. Gene Annotation and d alternative transcripts
 - 7.2. RNA-seq data analysis
- 8. Molecular evolution: philogenetic trees reconstruction
 - 8.1. Clustering algorithms
 - 8.1.1 k-means
 - 8.1.2 Neighbor joining
 - 8.2. UPGMA
 - 8.3. Maximum parsimony methods
 - 8.4. Maximum likelihood methods

Prerequisites

The student is expected to have a basic knowledge on the use of personal computer, informatics and molecular biology

Teaching form

Audio-video lessons, On-line classes (webex in little groups), On-line quiz

Textbook and teaching resource

M.M. Triola, M.F. Triola, Fondamenti di statistica per le discipline biomediche

M. Helmer Citterich, F. Ferrè, G. Pavesi, C. Romualdi, G. Pesole, Fondamenti di bioinformatica (Zanichelli editore)

Semester

First semester

Assessment method

Written exam (Biostatistics) and Oral exam (Bioinformatics)

Office hours

To be defined with the student by email contact
