# Education and Training challenges in Bioinformatics

Pedro L. Fernandes
Instituto Gulbenkian de Ciência

# Why Bioinformatics?

Bioinformatics is a hybrid discipline that has a quite simple purpose:

To contribute to better understand Biology by using
***biological information***
and
***computational methods***

# Biological Information

***Biological information*** is a term that encapsulates a very wide range of entities, emanating from laboratory instruments, field observations, clinical activity, etc.

***Biological information*** is ideally organised into databases. Datasets can be quite large and full of **variation** and **noise**.

There is a huge field of application for computer science and statistics that results from the need to manipulate such datasets.

# Computational Methods

Bioinformatics activities require the use of custom-developed software.

Behind this, there is algorithm development and statistics, that need to be dealt-with quality concerns.

Users, on the other hand, need to be able to **critically assess** the behaviour and the performance of these tools, in order to trust the results.

# A matter of trust...



Tears Of Joy At Conference As Bioinformaticist Says "I Won't Go Into Detail Of Algorithm"

# Bioinformatics Tools, some examples

- Assessing data quality

- Comparing gene expression in large arrays

- Evaluating scores in sequence aligments

- Evaluating fitness of structural models to data

- Evaluating variation in genomic data

- Classifying gene expression patterns

- Classifying population genetics data

- Assessing association and correlation

- Inferring causality in observations and experiments

# Bioinformatics Tools, specific needs

- Fitting

- Regression

- Maximization

- Statistical tests

- Correcting for multiple testing

- Machine learning

- Design of experiments

- Testing the performance of classifiers, etc.

# Alignments and Searches

Looking for sequence similarity

Alignment

- Global

- Local (BLAST and its variants)

Searching in databases and ranking by similarity

Identification of Motifs and Domains

# BLAST statistics
## BLAST = Basic Local Alignment Search Tool

In "BLAST" by Ian Korf, Mark Yandell and Joseph Bedell, O'Reilly, 2003

# Prediction and Inference

Predicting structure from sequence

- Sequence motifs

- Domains as sets of motifs

- Fingerprints

Inferring function from sequence and structure

- By similarity (holomogy modelling)

# Acquiring confidence in results

Bioinformatics can produce results that are not **consistent**.

Often the question is pushing us back into the statistical methods to try to isolate the causes of inconsistency.

We also need to guarantee that the results are fully **reproducible**.

The need for **knowledge** in statistics is just everywhere.

# Education and Training

Bioinformatics requires knowledge in Statistics at a variety of levels, naturally different for **researchers**, **software developers** and **users**.

**Formal learning** (higher education, degrees)

**Non-formal learning**

- Training

- Continuing Professional Education  (CPE)

# Education

Revision of undergraduate curricula in Statistics

Possible enhancement via **online learning** and **flipped class** techniques

Possible enhancement via **Peer Instruction**

"A New Online Computational Biology Curriculum"  by David Searls, PLOS Comp. Bio.  **10**, 6  June 2014.

From "Computational Biology Online Course Catalog":

https://www.coursera.org/course/statistics

http://ocw.jhsph.edu  Methods in Biostatistics I

# Biostatistics Training in GTPB

Meeting the needs of:
End users
Developers
Researchers
other trainers (instructors)

# Introductory Biostatistics for Biologists

The importance of Statistics.
Quantitative observations. Accuracy and Precision. Observations with error. Chance. Probabilities. Causation.

Descriptive Statistics. Basic concepts.
Describing and summarizing data. Summary statistics and plots for univariate and bivariate data.

Review of probability theory
Probability, random variables and their properties.
Independence and conditional probability.
Distributions: discrete random variables and continuous random variables.

Statistical inference
Sampling distributions. Confidence Intervals.
Hypothesis testing (parametric tests).

Statistical inference
Hypothesis testing (non-parametric tests). Contingency tables.
Design of experiments
ANOVA: one-way, two-way, repeated measures.
Factorial design. Latin Squares.

# Advanced Biostatistics for Biologists

Significance and p-value
Multiple testing issues
Corrections for multiple testing.

Simulation modelling methodologies
Monte Carlo and Bootstrap methods
Parametric approach, Non-Parametric approach

Bayesian inference
Bayes' theorem. Principles of Bayesian methodology. Gibbs Sampling.
Statistical inference. Expectation-Maximization (EM) algorithm.

Multivariate data analysis
Organising multivariate data.
Principal component analysis.

Machine Learning in Bioinformatics.
Statistical Methods for NGS Data Analysis
Introduction. Using the EDASeq package.
Using the edgeR, DESeq packages.

# Training in GTPB

The two training courses are designed to meet strict learning objectives

The content tends to compensate deficiencies in basic knowledge

Results are quite rewarding but the effort is higher than normal

Attempts will be made to offset basic subjects by using online content

# Acknowledgements

Acknowledgements are due to several hundred collaborators, who have accompanied me in providing training to thousands of Bioinformatics users and developers, among which I have to cite the following statisticians for their very generous contributions

Maria Antónia Turkman
Lisete de Sousa
Carina Silva
Ana Luísa Papoila
Maria Fernanda Diamantino

# Thematic Session: Bioinformatics

## Thank you, for your attention