# Bioinformatics & Machine Learning

Daniel Glez-Peña

IP Leiria, June 3 2009

SING
Upgrading your knowledge ™

# Agenda

## 1. Bioinformatics

Definition, major research areas, databases

## 2. Machine Learning for bioinformatics

Algorithm types, examples in bioinformatics

## 3. DNA Microarrays

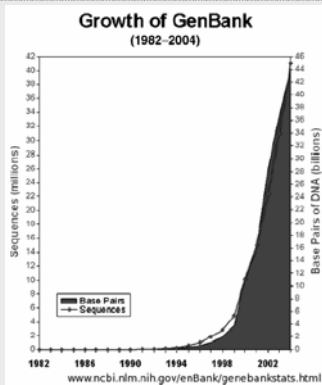Technological overview

## 4. Applications

GeneCBR and WhichGenes?

# Bioinformatics

# Bioinformatics

- "Application of the Information Technologies to the field of molecular biology"

- Creation and enhancement of:

  - Databases with biological information

  - Algorithms

  - Statistical techniques

  …to solve formal and practical problems arising from the management and analysis of biological data



Growth of GenBank
(1982–2004)

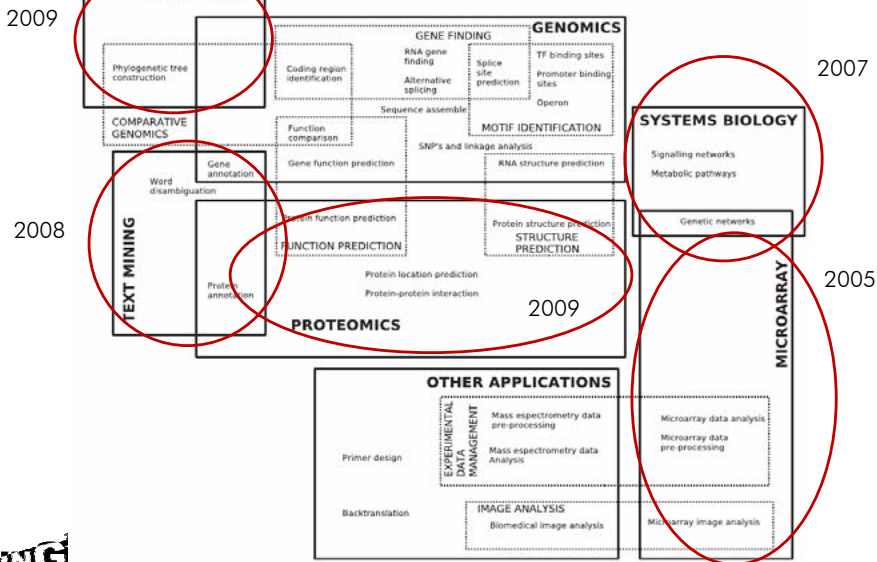www.ncbi.nlm.nih.gov/enBank/genebankstats.html
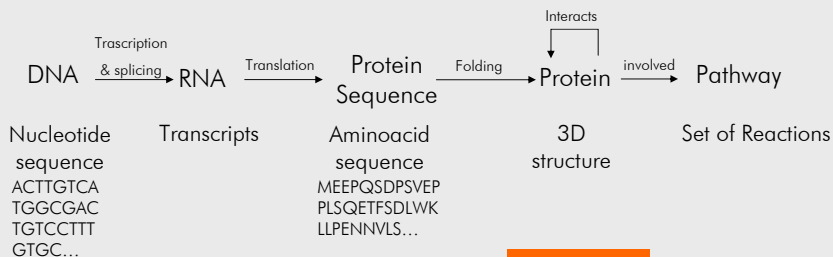
# Major research areas

- GENOMICS
  - Sequence analysis
  - Genome annotation
  - Analysis of mutations in cancer
- PROTEOMICS
  - Protein-protein docking
  - Analysis of protein expression
  - Prediction of protein structure
- MICROARRAYS
  - Analysis of gene expression
  - Genetic network induction

- TEXT MINING
  - Gene annotation
  - Protein annotation
  - Relation extraction
- EVOLUTION
  - Phylogenetic reconstruction
  - Comparative genomics
- SYSTEMS BIOLOGY
  - Modelling biological systems
- OTHER
  - Image Analysis

# Major research areas



Larrañaga *et al* (2005), Briefings in Bioinformatics 7(1):86-112

# Molecular biology dogma

DNA ──Trascription & splicing──► RNA ──Translation──► Protein Sequence ──Folding──► Protein ──involved──► Pathway

Interacts

| Nucleotide sequence | Transcripts | Aminoacid sequence | 3D structure | Set of Reactions |

Nucleotide sequence
ACTTGTCA
TGGCGAC
TGTCCTTT
GTGC…

Aminoacid sequence
MEEPQSDPSVEP
PLSQETFSDLWK
LLPENNVLS…

Interactomics

| GENOMICS | PROTEOMICS | METABOLOMICS |
|---|---|---|

Sequence analysis
Genome annotation
Analysis of mutations

**Evolution**
Phylogenetic reconstruction
Comparative genomics

Gene expression analysis [DNA microarray]

Protein expression analysis [mass spectometry]

Protein structure prediction [folding]

Protein interaction prediction [3D docking]

Modelling biological systems
Functional analylis

# Databases

## Genomics

### Sequences



### Genomes



### Cene-centric



## Proteomics

### Proteins



### Structure



### Domains



## Interactomics & Metabolomics

### Prot-Prot interactions



### Pathways



## Ontologies



## Bibliome



## Experimental data

# Machine Learning for Bioinformatics

# Machine Learning & Bioinformatics

- CLASSIFICATION (SUPERVISED LEARNING)
- CLUSTERING (UNSUPVERVISED LEARNING)
- GRAPHICAL PROBABILISTIC MODELS
- OPTIMIZATION

# ML & Bioinformatics: Classification

- Classification (supervised learning)
  - Given a set of "instances", each one with a set of measured "attributtes" and a "outcome" value we want to train a model that predicts the outcome in further problem instances
    - If the "outcome" is discrete (typical 2 o more different values) we are talking about **classification** (if not: regression)

| | $X_1$ | $\ldots$ | $X_n$ | $C$ | |
|---|---|---|---|---|---|
| $(x^{(1)}; c^{(1)})$ | $x_1^{(1)}$ | $\ldots$ | $x_n^{(1)}$ | $c^{(1)}$ | Training data |
| $(x^{(2)}; c^{(2)})$ | $x_1^{(2)}$ | $\ldots$ | $x_n^{(2)}$ | $c^{(2)}$ | |
| $(x^{(N)}; c^{(N)})$ | $x_1^{(N)}$ | $\ldots$ | $x_n^{(N)}$ | $c^{(N)}$ | |
| $x^{(N+1)}$ | $x_1^{(N+1)}$ | $\ldots$ | $x_n^{(N+1)}$ | ??? | Test data |

# ML & Bioinformatics: Classification
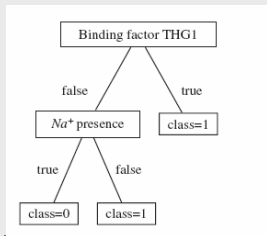
- Classification

  - Feature subset selection.

    - Are all input attributes useful?

    - Advantages: reduced cost in data adquisition, improved uniderstability of the model, faster training, and better accuracy

    - It is a search space problem ($2^n-1$), in general:

      - 1. Generate a subset
        [brute force, deterministic/not deterministic heuristic search]

      - 2. Evaluate subset
        Statistical estimation: Information Gain, X2, t-test, DFP, CFS
        Wrapper (use classifier accuracy in training set)

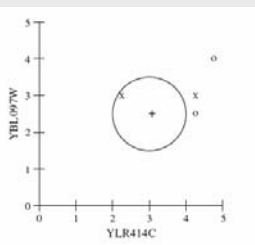      - 3. if (!halt_condition) GOTO 1

# ML & Bioinformatics: Classification

- ## Classification
  - Popular techniques
    - Logistic regression
    - Linear discriminant analysis (LDA)
    - Bayesian classifiers: Naive Bayes, semi-NB, Tree augmented NB, k dependence Bayesian…
    - Classification trees: CART, C4.5, RandomForest, J48…
    - K-Nearest Neighbours
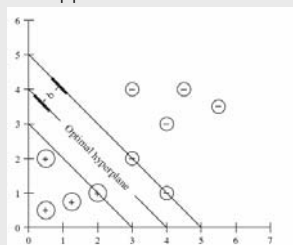    - Support Vector Machines
    - Meta: Bagging, Boosting

| Classification Tree | kNN classifier | Support Vector Machine |
|---|---|---|

# ML & Bioinformatics: Classification

- Examples of classification in Bioinformatics (I)
  - Genomics
    - Gene finding (if a sequence is a coding region)
    - Splice site prediction (if a sequence is a splice site)
    - Predict disease genes (from i.e. its sequence length?)
    - Prediction of mutation (SNP) effect
    - **Cancer prediction from gene expression (microarrays)**
  - Proteomics
    - Prediction of secondary structure (alpha-helix, beta-sheet,etc.)
    - Prediction of sub-cellular location of the protein
    - Cancer prediction from protein expression (mass spectra)

# ML & Bioinformatics: Classification

- Examples of classification in Bioinformatics (and II)

  - Systems biology
    - Predict the cell migration speed (high, low) from the phosphorilation levels of signalling proteins
    - Predict a gene regulatory level (up-regulated or down-regulated given the 'related' genes expression)

  - Text mining
    - Protein/gene recognition in biomedical literature (is this word a gene/protein given some word features: ortographic, part-of-speech, suffix, trigger words, etc…??)

# ML & Bioinformatics: Clustering

- Clustering
  - Partition a set of "instances" in several groups (clusters) given the differences between them
    - Their are based on "distances" between instances that is a problem-dependant issue
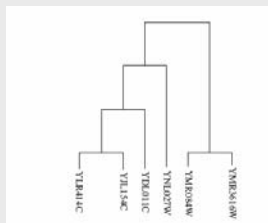      - Typical: Euclidean, Pearson, Sperman

# ML & Bioinformatics: Clustering

- Clustering
  - Popular techniques
    - Partition clustering: k-means, SOM, GCS, PAM
    - Hierarchical clustering with single-linkage, complete linkage, centroid linkage and wards-criterion
      - They produce the popular "dendograms"
    - Model-based clustering

Partition clustering

Hierarchical clustering (dendogram)

# ML & Bioinformatics: Clustering

- Clustering in Bioinformatics
  - Mainly applied to analyze gene expression data
    - Co-Expression detection (group genes with similar expression)
    - Subclass discovery (group samples given the expression of its genes)
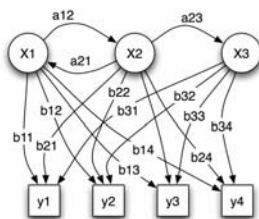    - Expression data visualization/summarization with dendograms

# ML & Bioinformatics: Probabilistic graphical models

- DAGs where nodes are random variables and links are probabilities from any kind of conditional dependence
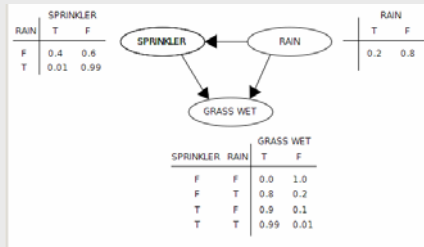
  - Examples

    - Hidden Markov Models
    - Bayesian Networks

Hidden Markov Model



Bayesian Network

# ML & Bioinformatics: Probabilistic graphical models

- Probabilistic Graph Models in Bioinformatics

  - Genomics

    - HMM to gene finding (does a gene sequence come from a coding or a non coding DNA region?)
    - Bayesian networks to detect splice sites (does a gene sequence come from a splice-site)

  - Systems Biology

    - Inference of regulatory genetic networks. Bayesian networks to expression pattern recognition (which genes cause other genes to express?)

# ML & Bioinformatics: Optimization

- Optimization
  - Search of the best solution in a huge (exponential) space.
  - Popular techniques
    - Exact optimization
      - Brute force
    - Deterministic
      - Hill climbing, local optimization
    - Stochastic
      - Monte Carlo
      - Simulated Annealing
      - Tabu search
      - Evolutionary
        - Genetic algorithms
        - Genetic Programming
        - Estimation of probability

# ML & Bioinformatics: Optimization

- Optimization techniques in Bioinformatics
    - Genomics
        - Multiple sequence alignment (used almost all optimization algorithms)
        - Splice site prediction with estimation of distribution algorithms
        - DNA sequencing
        - Cluster microarray data
    - Proteomics
        - Protein folding (predict 3D structure)
        - Protein side-chain prediction (determine the optimal set of 'angles' in the 3D structure that minimize the energy)
    - Systems Biology
        - Inference of gene networks and estimate the parameters of bioprocesses
    - Evolution
        - Inference of phylogenetic trees
        - Haplotype reconstruction
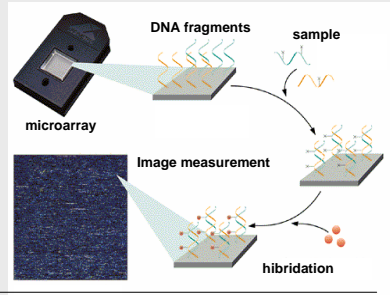
# DNA Microarrays

# DNA Microarrays

- DNA microarray. Objetive: Measure gene expression

- Description

  - Matrix with measures the expression of thousands of genes simultaneously

  - Gives a "global" vision of gene activity, and allows comparison

    - Between different individuals
    - Same individual at different times
    - Different tissues

# DNA Microarrays

- How it works
  - DNA fragments are spotted or printed in probes on the array surface
    - Each probe is a gene
  - **Hibridation** is performed with a sample putted onto the array
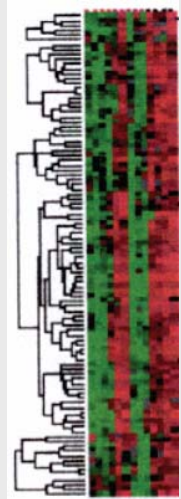  - A scanner measures the intensity in each probe



DNA fragments sample
microarray
Image measurement
hibridation



**Microarray data** → scanner → adquisition → **Dataset**

# DNA Microarrays

- Human Genome U133
  - HG U133A, HG U133B
  - 22.000 probes aprox. ($\cong$1 probe x gen)
- Human Genome U133 plus
  - 44.000 probes ($\cong$2 probes x gen)
- Exon array
  - 1.4 millions of probes ($\cong$16 probes x gen)

# DNA microarrays

- Typical analyses & ML Techniques
  - Gene-based analysis
    - Co-expression detection with clustering techniques (unsupervised)
  - Differential gene expression analysis
    - Detect which genes has a significant expression variation among samples of two or more conditions (feature selection)
  - Sample-based analysis
    - Class predicion with classification techniques (supervised)
    - Class discovery with clustering techniques (unsupervised)
  - Problems:
    - Huge number of features (thousands of genes) y low number of samples (dozens) V.S. Machine Learning
    - High false positive rate

# DNA microarrays

- Functional interpratation after data analysis

  - Typically we have a list of genes of interest (ie. differentially expressed)

  - Question: who are those genes?

  - Solution: Use the available gene annotations (Gene Ontology, Pathways, etc) and see if there is a correlation with a functional module.

    - They answer to the question: Are my genes significantly chosen from a given gene function? If so, which function?
    - On-line tools
      - List-based: FatiGO, DAVID, Pathjam
      - Gene-set based: GSEA, FatiScan

# Sample applications

# geneCBR



Translational tool for DNA microarray-based diagnostics

- www.genecbr.org

- Glez-Peña *et al.* BMC Bioinformatics 10:37 2007

- Classification guided by a clustering algorithm GCS

# WhichGenes?

**WhichGenes: a web-based tool for gathering, building, storing and exporting gene sets with application in gene set enrichment analysis**

Daniel Glez-Peña[1], Gonzalo Gómez-López[2], David G. Pisano[2] and Florentino Fdez-Riverola[1,3,*]

[1]Higher Technical School of Computer Engineering, University of Vigo, Ourense, [2]Bioinformation Unit (UBio), Structural Biology and Biocomputing Programme, Spanish National Cancer Research Centre (CNIO), Madrid and [3]Informatics Department, University of Vigo, Vigo, Pontevedra, Spain

Received January 20, 2009; Revised April 3, 2009; Accepted April 8, 2009

## ABSTRACT

WhichGenes is a web-based interactive gene set building tool offering a very simple interface to extract always-updated gene lists from multiple databases and unstructured biological data sources. While the user can specify new gene sets of interest by following a simple four-step wizard, the tool is able to run several queries in parallel. Every time a new set is generated, it is automatically added to the private gene-set cart and the user is notified by an e-mail containing a direct link to the new set stored in the server. WhichGenes provides functionalities to edit, delete and rename existing gene sets as well as the capability of generating new ones by combining previous existing sets (intersection, union and difference operators). The user can export his sets configuring the output format and selecting among multiple gene identifiers. In addition to the user-friendly environment, WhichGenes allows programmers to access its functionalities in a programmatic way through a Representational State Transfer web service. WhichGenes front-end is freely available at http://www.whichgenes.org/, WhichGenes API is accessible at http://www.whichgenes.org/api/.

## INTRODUCTION

During the past several years, bioinformatics enrichment tools have played a very important and successful role contributing to the gene functional analysis of large gene lists (ranging in size from hundreds to thousands of genes) for various high-throughput biological studies (1).

From the large amount of tools that are currently available in the community, two widely used approaches can be

identified: (i) individual gene analysis (IGA), which evaluates the significance of individual genes between two groups of samples compared, and (ii) gene set analysis (GSA), free from the problems of the 'cutoff-based' methods (2).

In recent years, GSA approach has received a great deal of attention because, from a biological perspective, functionally related genes often display a coordinated expression to accomplish their roles in the cell. In this direction, GSA methods enable the understanding of cellular processes as an intricate network of functionally related components (3). However, while extensive work has been done during last years in developing new GSA methods, little effort has been put on implementing tools that can help researchers gather, store and manage gene sets containing large 'interesting' gene lists from multiple data sources. To our knowledge, only the Gene Set Builder tool (4) shares the same fundamental concept, giving support to easily handle sets of genes. Compared with Gene Set Builder, which follows a fixed database-driven architecture to give access only to Ensembl and Gene Lynx gene catalogs, the tool presented herein follows a more flexible database-free and interactive approach allowing the integration of 14 different data sources for the *Homo sapiens* and *Mus musculus* organisms.

In this article, we present WhichGenes, an online, database-free, web-based tool for easily gathering, building, storing and exporting always-updated gene sets coming from multiple data sources. It allows researchers to elaborate custom hypotheses in the form of lists of genes in order to further use them as input in existing GSA tools. WhichGenes currently supports queries about *Homo sapiens* and *Mus musculus* organisms by retrieving up-to-date gene lists directly coming from multiple databases, currently including Ensembl, MSigDB, KEGG/Biocarta/Reactome pathway databases, GeneCards, CancerGenes, Dsigdbm, Diseases CTD, Targetscan, miRBase, Chemical CTD, AmiGO and InfAd. Generated gene sets can be

On-line geneset building tool

- Create your own genesets from multiple datasources and use them in your favourite geneset-based analysis tools like GSEA

- www.whichgenes.org

- Glez-Peña *et al*. Nucleic Acids Res (web server issue) 2009

**Bioinformatics     Machine Learning     Microarrays     Applications**

# Questions?