

## Bioinformatics and Agriculture:



Miguel Pérez-Enciso  
([miguel.perez@uab.es](mailto:miguel.perez@uab.es))

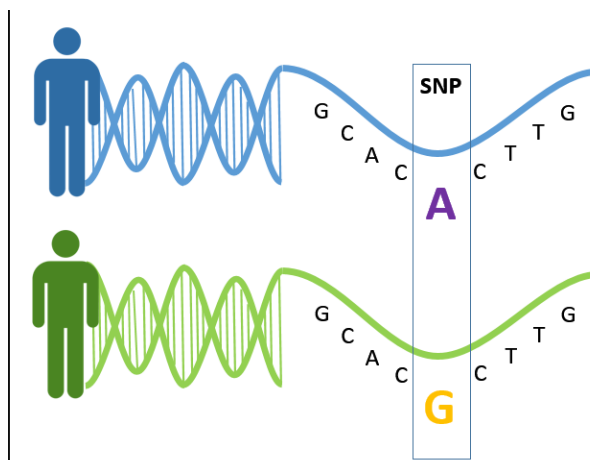


## Some Bioinformatics Applications

- ✓ Precision agriculture (Big data issues, modeling ...)
- ✓ Causal mutation discovery
- ✓ Genomic selection, machine learning
- ✓ Genome assembly
- ✓ Sequence analysis
- ✓ Functional analysis
- ✓ ...

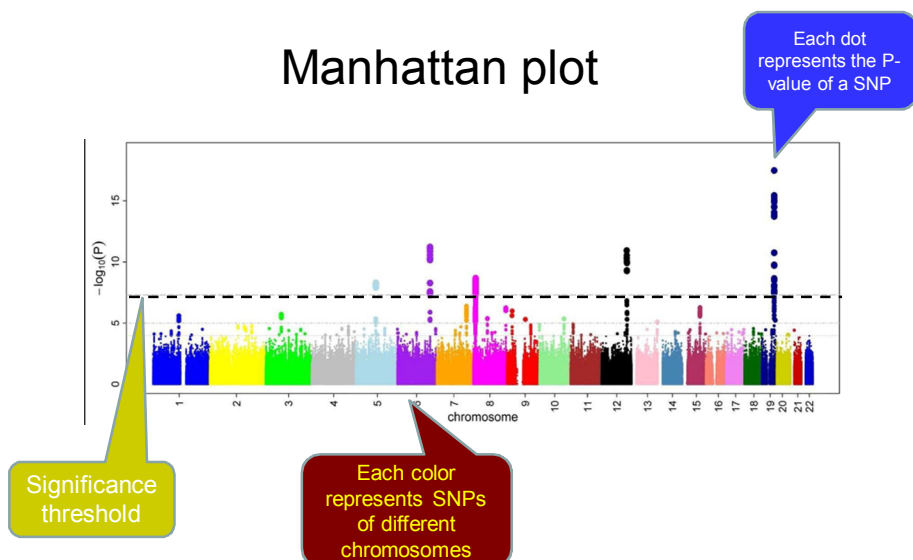
## Causal Mutation Discovery: Genome Wide Analysis Study (GWAS)

There are millions of small DNA differences between individuals called 'SNPs' (single nucleotide polymorphisms) because they affect a single base



## Causal Mutation Discovery: GWAS

### Manhattan plot



## Causal Mutation Discovery: WARNING

- Small effect
- Strong disequilibrium
- Many candidates
- Small samples

It is going to be VERY difficult to prove causality

## Machine learning, genomic selection

Machine learning is the discipline that develops algorithms allowing computer learning from data and making predictions on future data, e.g., to predict likelihood that someone suffers from a disease based on molecular information from healthy and affected people.

### Supervised learning: predicting an output variable from high-dimensional observations

- Nearest neighbor and the curse of dimensionality
- Linear model: from regression to sparsity
- Support vector machines (SVMs)

### Model selection: choosing estimators and their parameters

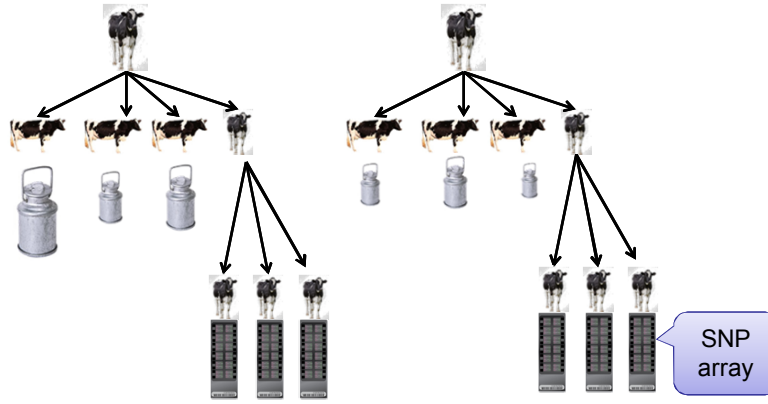
- Score, and cross-validated scores
- Cross-validation generators
- Grid-search and cross-validated estimators

### Unsupervised learning: seeking representations of the data

- Clustering: grouping observations together
- Decompositions: from a signal to components and loadings

**scikit-learn**  
Machine Learning in Python

## What is the main direction today: Genomic selection



SNP data is used to predict future performance.

## Genome Assembly



Full Paper

### Whole Genome Sequencing of Turbot (*Scophthalmus maximus*; Pleuronectiformes): A Fish Adapted to Demersal Life

Antonio Figueras<sup>1,\*</sup>, Diego Robledo<sup>2</sup>, André Corvelo<sup>3,4</sup>, Miguel Hermida<sup>5</sup>,



### The genome of melon (*Cucumis melo* L.)

Jordi Garcia-Mas<sup>1,2</sup>, Andrej Benjak<sup>3</sup>, Walter Sanseverino<sup>4</sup>, Michael Bourgeois<sup>5</sup>, Gisela Mir<sup>6</sup>, Victor M. González<sup>6</sup>, Elizabeth Hénaff<sup>6</sup>, Francisco Cámara<sup>7</sup>, Luca Cozzuto<sup>8</sup>, Ernesto Lowy<sup>9</sup>, Tyler Alloto<sup>9</sup>, Salvador Capella-Gutiérrez<sup>9</sup>, Jose Blanca<sup>9</sup>, Joaquin Caliztrone<sup>9</sup>, Pello Zariwail<sup>9</sup>, Daniel González-Iglesias<sup>9</sup>, Luis Rodríguez-Moreno<sup>9</sup>, Marcos Drogue<sup>9</sup>, Lei Du<sup>9</sup>, Miguel Alvarez-Tejado<sup>9</sup>, Belén Lorente-Galdos<sup>9</sup>, Marta Melé<sup>9</sup>, Luming Yang<sup>9</sup>, Yiqun Weng<sup>9</sup>, Arcadi Navarro<sup>10</sup>, Tomas Marques-Bonet<sup>11</sup>, Miguel A. Avendaño<sup>12</sup>, Fernando Ruiz<sup>13</sup>, Belén Pico<sup>14</sup>, Toni Gabaldón<sup>15</sup>, Guglielmo Roma<sup>16</sup>, Rodéric Guigó<sup>17</sup>, Josep M. Casacuberta<sup>18</sup>, Pere Aroca<sup>19</sup>, and Pere Puigdomènech<sup>20</sup>

<sup>1</sup>Institut de Recerca i Tecnologia Agroalimentàries, Centre for Research in Agricultural Genomics Consejo Superior de Investigaciones Científicas-Institut de Recerca i Tecnologia Agroalimentàries-Universitat Autònoma de Barcelona-Universitat de Barcelona, 08193 Barcelona, Spain; <sup>2</sup>Centre for Research in Agricultural Genomics Consejo Superior de Investigaciones Científicas-Institut de Recerca i Tecnologia Agroalimentàries-Universitat Autònoma de Barcelona

DOI: 10.1093/nar/gkz004

GigaScience

DATA NOTE

Open Access

### Genome sequence of the olive tree, *Olea europaea*

Fernando Cuccato<sup>1</sup>, Irene Julca<sup>2,3,4</sup>, Jessica Gómez-Carrión<sup>5</sup>, Daniel Loka<sup>6</sup>, Marina Marcehouben<sup>7</sup>, Emilio Cano<sup>8</sup>, Susana Galán<sup>9</sup>, Leonor Fidalgo<sup>10</sup>, Pablo Blasco<sup>11</sup>, Sophia Demaj<sup>12</sup>, Maria Gu<sup>13</sup>, Manuel Sánchez-Fernández<sup>14</sup>, Jose Luis García<sup>15</sup>, Ivo G. Gál<sup>16</sup>, Pablo Vargas<sup>17,18</sup>, Tyler S. Alloto<sup>19,20</sup> and Toni Gabaldón<sup>21,22</sup>

An active area of research, practical and evolutionary interest  
Huge variety in genome sizes and complexities:

Rice	0.4 Gb	Mammals	3.2 Gb
Maize	2.5 Gb	Chicken	1.2 Gb
Wheat	17 Gb	Fishes	0.3 -130 Gb
Conifers	25 Gb	Turbot	0.5 Gb

## Sequence Analysis

- ❑ New (> 2008) sequencing technologies have revolutionized genomics. There are already databases with thousands of human complete genomes.
- ❑ Bioinformatics is critical to profit from the deluge of data.



- Polymorphism discovery.
- Transcriptome analysis (RNAseq)
- Epigenome analysis.
- Motif, nucleosome discovery (ChIP-seq, Dnase1-seq...)
- Metagenome
- ...

Sequence can be used for numerous applications:

## Take home message

- ✓ There are no simplistic responses to complex scenarios.
- ✓ Bioinformatics can give you an initial advantage but more important is to have good questions.
- ✓ Learn as much Statistics as you can.

## (DIFFICULT) QUESTION

Do you think Big Data will help to solve causality?

