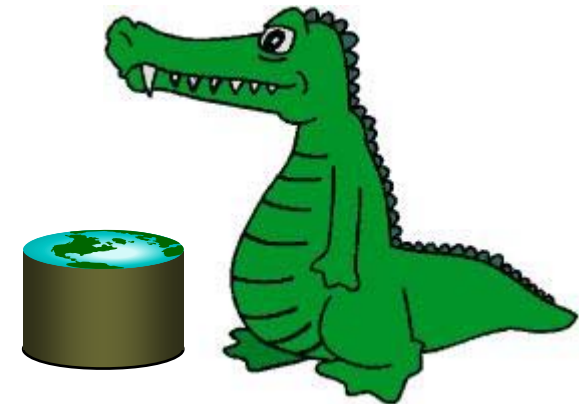# DATA SCIENCE: APLICACIONES A LA BIOLOGIA Y A LA MEDICINA CON PYTHON Y R

Toni Monleón-Getino
Section of Statistics. Fac Biology
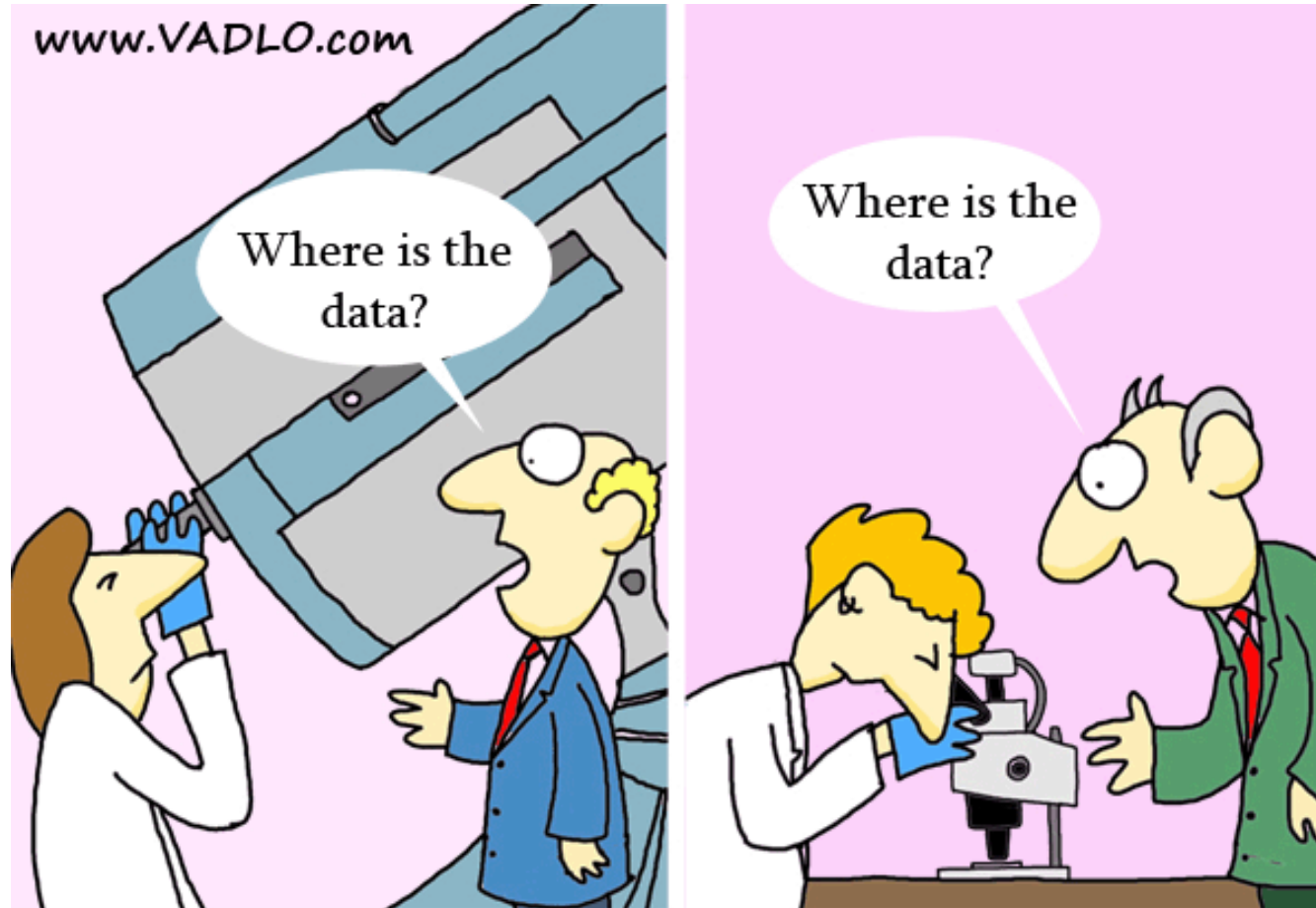
SEQUENTIA

Biost³

BIB BIOINFORMATICS BARCELONA
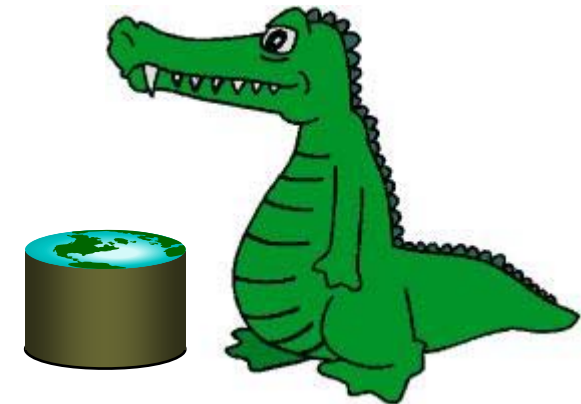
HPC Now!

GRBIO

UNIVERSITAT DE BARCELONA
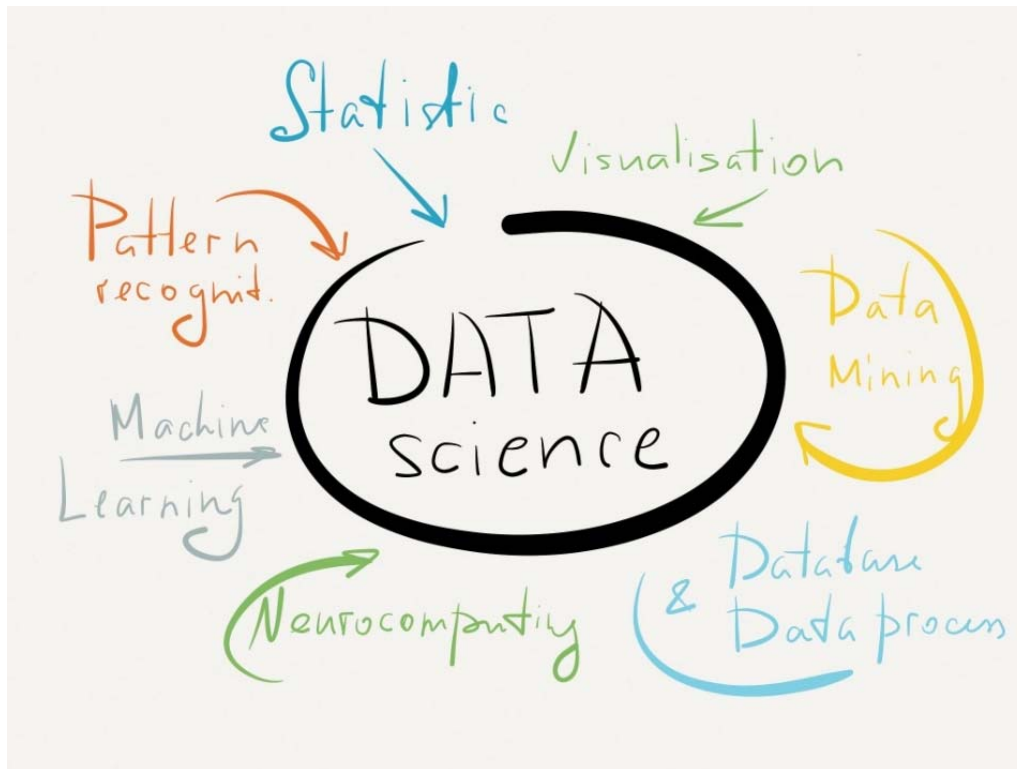
# Grand Unification of Sciences



Grand Unification of Sciences

# What is Data Science?



Why, Where, What, How, Who

# "Data Science" an Emerging Field



O'Reilly Radar report, 2011

# Data Science – A Definition

**Data Science** is the science which uses computer science, statistics and machine learning, visualization and human-computer interactions to collect, clean, integrate, analyze, visualize, interact with data to create data products and science.

# Goal of Data Science

Turn data into data products and science.

## WHY IS HEALTH DATA SCIENCE IMPORTANT?

**P** — **Personalized Medicine**

Merge and analyze data sets from from multiple sources to create personalized treatment.

**G** — **Genomics**

Inexpensive DNA sequencing and next-generation genomic technologies are changing the way health care providers do business.

**S** — **Self-Motivated Care**

It's a "patient heal thyself" world, now. Developments like personal genetic testing, online patient networks, and behavioral apps are allowing individuals to take control of their own health.

**D** — **Disease Modeling and Mapping**

One of the flashiest uses of data science in the past few years has been in tracking (and finding ways to halt or prevent) diseases.
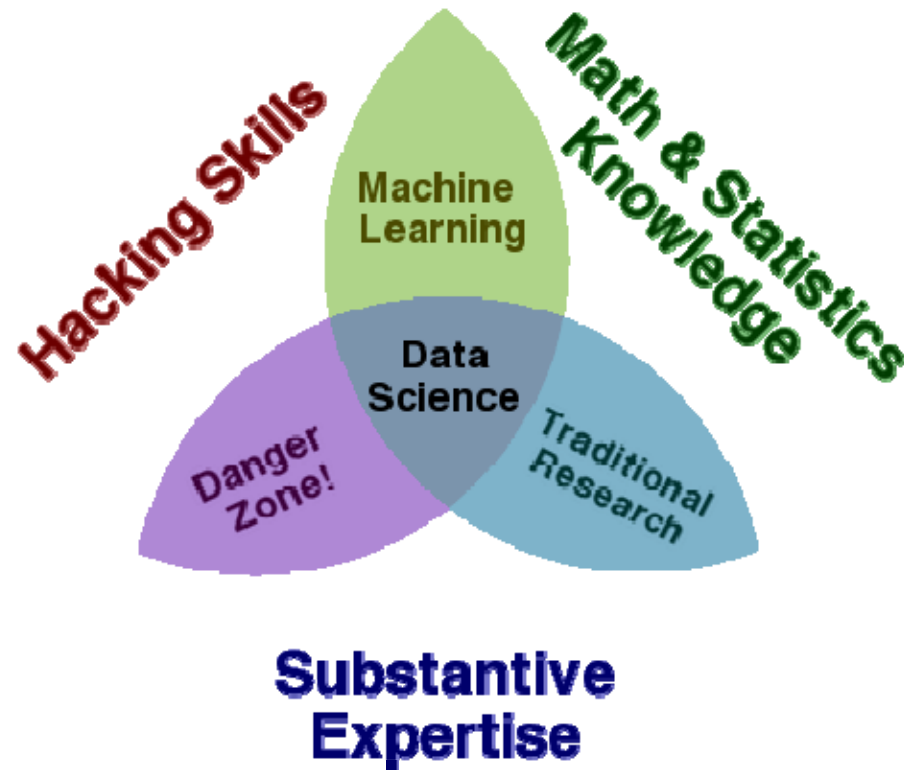
Some recent ML Competitions at https://www.kaggle.com/

NIST Pre-Pilot Data Science Evaluation – likely to be incorporated to be part of Labs/Final project



**kaggle**

Active Competitions

| | | Flight Quest 2: Flight Optimization | 33 days Coming soon $220,000 |
| GE | | Final Phase of Flight Quest 2 | |
| | | Packing Santa's Sleigh | 5.8 days 338 teams $10,000 |
| | | He's making a list, checking it twice; to fill up his sleigh, he needs your advice | |
| | Genentech | Flu Forecasting | 41 days 37 teams |
| | | Predict when, where and how strong the flu will be | |
| | GALAXY ZOO | Galaxy Zoo - The Galaxy Challenge | 2 months 160 teams $16,000 |
| | | Classify the morphologies of distant galaxies in our Universe | |
| | DENIED APPROVED | Loan Default Prediction - Imperial College Lon... | 52 days 82 teams $10,000 |
| | | Constructing an optimal portfolio of loans | |
| | | Dogs vs. Cats | 11 days 166 teams Swag |
| | | Create an algorithm to distinguish dogs from cats | |

# Data Science – A Visual Definition

# Contrast: Databases

| | **Databases** | **Data Science** |
|---|---|---|
| Data Value | "Precious" | "Cheap" |
| Data Volume | Modest | Massive |
| Examples | Bank records, Personnel records, Census, Medical records | Online clicks, GPS logs, Tweets, Building sensor readings |
| Priorities | Consistency, Error recovery, Auditability | Speed, Availability, Query richness |
| Structured | Strongly (Schema) | Weakly or none (Text) |
| Properties | Transactions, ACID* | CAP* theorem (2/3), eventual consistency |
| Realizations | SQL | NoSQL: MongoDB, CouchDB, Hbase, Cassandra, Riak, Memcached, Apache River, ... |

ACID = Atomicity, Consistency, Isolation and Durability

CAP = Consistency, Availability, Partition Tolerance

# Contrast: Machine Learning

| Machine Learning | Data Science |
|---|---|
| Develop new (individual) models | Explore many models, build and tune hybrids |
| Prove mathematical properties of models | Understand empirical properties of models |
| Improve/validate on a few, relatively clean, small datasets | Develop/use tools that can handle massive datasets |
| Publish a paper | Take action! |

# datavolution