

## COME TRASFORMARE I DATI IN INNOVAZIONE

Relatore: Alessio Passalacqua

[www.produzioneperfetta.it](http://www.produzioneperfetta.it) | [info@produzioneperfetta.it](mailto:info@produzioneperfetta.it)

Correlatore: Alessandro Tronchin

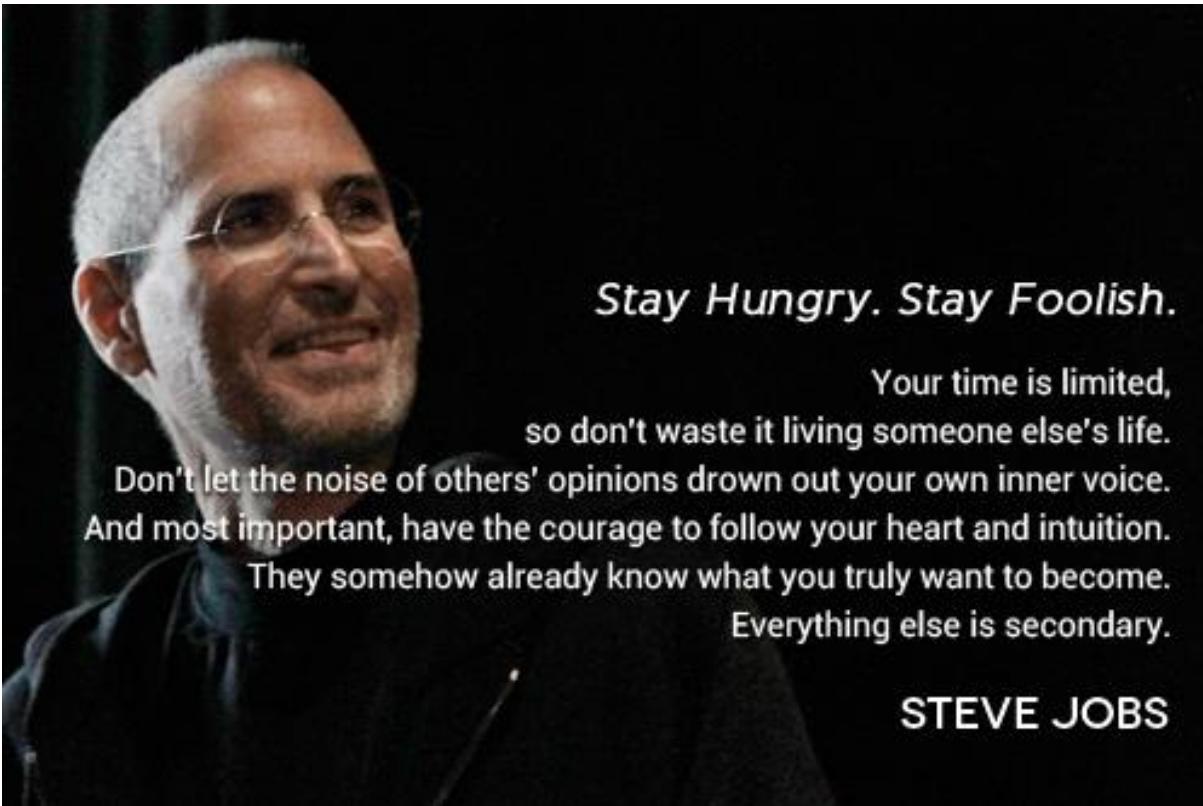
ICT Manager at Haemotronic



# *Non smettere di imparare...*

## *per innovare bisogna saper cambiare*

INTRO



*Stay Hungry. Stay Foolish.*

Your time is limited,  
so don't waste it living someone else's life.

Don't let the noise of others' opinions drown out your own inner voice.  
And most important, have the courage to follow your heart and intuition.

They somehow already know what you truly want to become.  
Everything else is secondary.

STEVE JOBS

# ANALISI DATI E' UN LAVORO DI SQUADRA

CONDIVISIONE  
OBIETTIVO



Bisogna giocare di squadra se  
si vuole arrivare ad un  
obiettivo comune

INTRO

Ribattere le responsabilità da  
un'altra parte è solo una  
perdita di tempo



«ATTRaverso il successo  
degli altri ottengo il mio  
successo»

Ing. Mazzola  
Ex Test Team Manager FERRARI

# ANALISI DATI E' UN PERCORSO A TAPPE

INIZIO FATICOSO



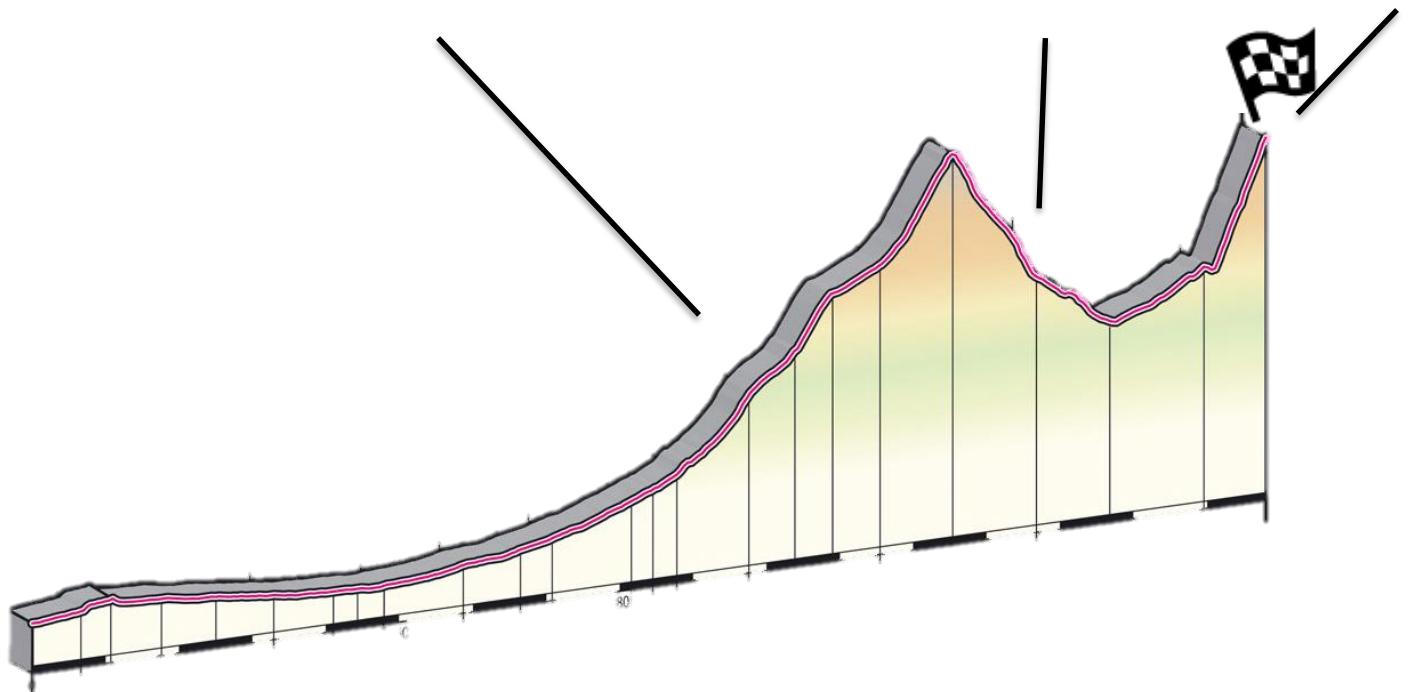
INTERMEDIO RISCHIOSO



FINALE GLORIOSO



INTRO



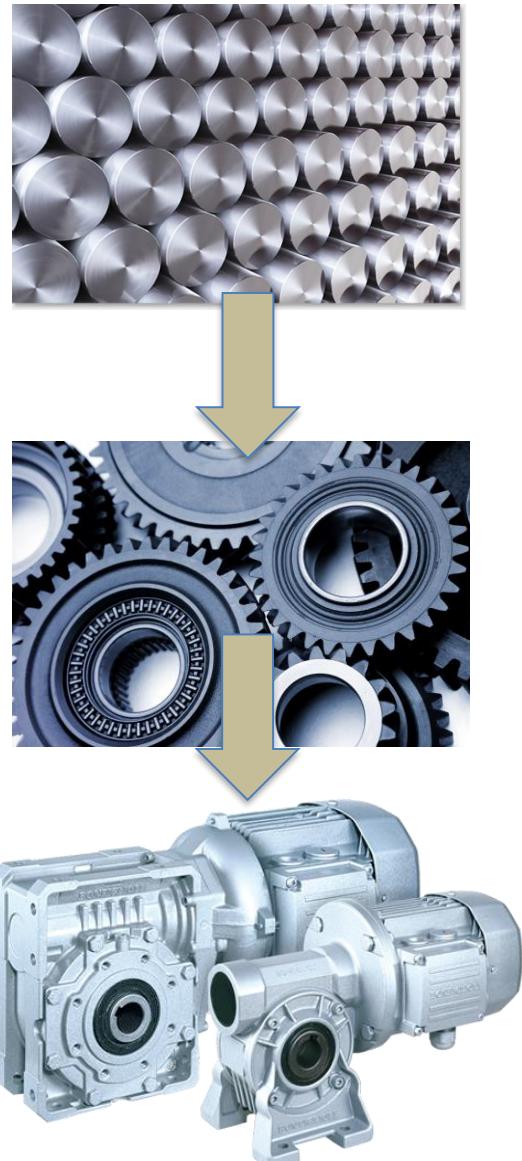
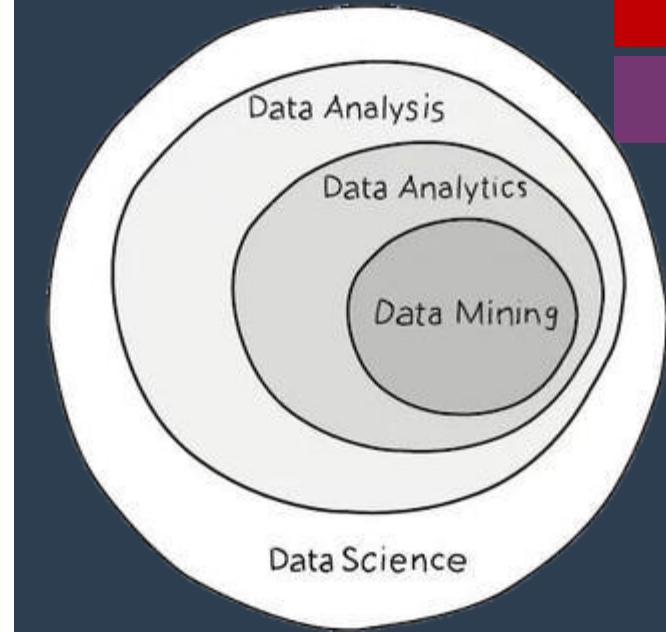
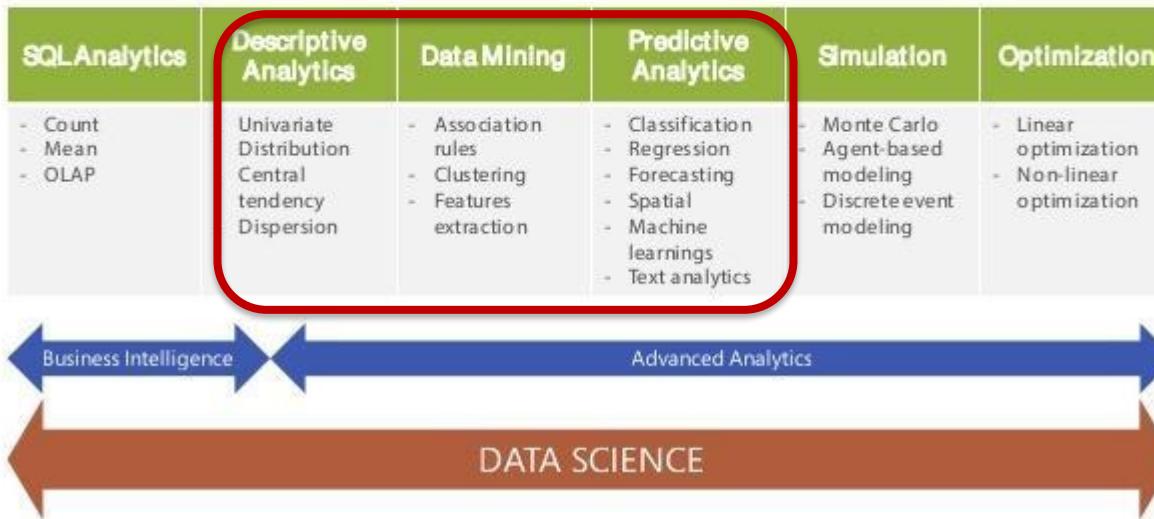
"Big data is often misunderstood by CFOs as a risk or cost rather than an asset." Dez Blanchfield, IT consultant

# ANALISI DATI: LA FILIERA

INTRO

## BI vs. Data Science

- Broader skillset and broader analytics spectrum.



# ANALISI DATI: LE PROFESSIONI

## DATA ANALYST DATA DETECTIVE

### Role

Collects, processes and performs statistical data analyses

### Mindset

Intuitive data junkie with high "figure-it-out" quotient



HIRED BY

Languages  
R, Python, HTML, Javascript, C/C++, SQL

### Skills & Talents

- ✓ Spreadsheet tools (e.g. Excel)
- ✓ Database systems (SQL and NO SQL based)
- ✓ Communication & visualization
- ✓ Math, Stats, Machine Learning

Languages  
R, SAS, Python, Matlab, SQL, Hive, Pig, Spark

### Skills & Talents

- ✓ Distributed computing
- ✓ Predictive modeling
- ✓ Story-telling and visualizing
- ✓ Math, Stats, Machine Learning

## DATA SCIENTIST AS RARE AS UNICORNS<sup>®</sup>

### Role

Cleans, massages and organizes (big) data

Mindset  
Curious data wizard



HIRED BY

## DATA ENGINEER SOFTWARE ENGINEERS BY TRADE

### Role

Develops, constructs, tests and maintains architectures (such as databases and large-scale processing systems)

### Mindset

All-purpose everyman



HIRED BY

Languages  
SQL, Hive, Pig, R, Matlab, SAS, SPSS, Python, Java, Ruby, C++, Perl

### Skills & Talents

- ✓ Database systems (SQL & NO SQL based)
- ✓ Data modeling & ETL tools
- ✓ Data APIs
- ✓ Data warehousing solutions

## DATA AND ANALYTICS MANAGER DATA SCIENCE TEAM LEADER

### Role

Manages a team of analysts and data scientists

### Mindset

Data Wizards' Cheerleader



HIRED BY

Languages  
SQL, R, SAS, Python, Matlab, Java

### Skills & Talents

- ✓ Database systems (SQL and NO SQL based)
- ✓ Leadership & project management
- ✓ Interpersonal communication
- ✓ Data mining & predictive modeling

## BUSINESS ANALYST CHANGE AGENT

### Role

Improves business processes as intermediary between business and IT

Mindset  
Resilient project juggler



HIRED BY

Languages  
SQL

- ✓ Basic tools (e.g. MS Office)
- ✓ Data visualization tools (e.g. Tableau)
- ✓ Conscious listening and storytelling
- ✓ Business Intelligence understanding
- ✓ Data modeling

## DATABASE ADMINISTRATOR DATABASE CARETAKER

### Role

Ensures that the database is available to all relevant users, is performing properly and is being kept safe

Mindset  
Master of Disaster Prevention



HIRED BY

Languages  
SQL, Java, Ruby on Rails, XML, C#, Python

### Skills & Talents

- ✓ Backup & recovery
- ✓ Data modeling and design
- ✓ Distributed Computing (Hadoop)
- ✓ Database systems (SQL and NO SQL based)
- ✓ Data security
- ✓ ERP & business knowledge

# DI COSA PARLEREMO

INTRO

- **“The Era of Data”:** L'inizio della quarta rivoluzione industriale
- **“Data-Driven Innovation”:** 4 steps per trasformare i dati in innovazione
- Confronto con i metodi di **analisi statistica Six-Sigma**
- Presentazione del **software open-source R**
- Presentazione dei metodi di **analisi statistica avanzata** e di **“Machine Learning”**
- **Case studies**

**“The Era of Data”:  
L'inizio della quarta rivoluzione industriale**

# THE DATA ERA



THE DATA ERA

## RAW MATERIALS



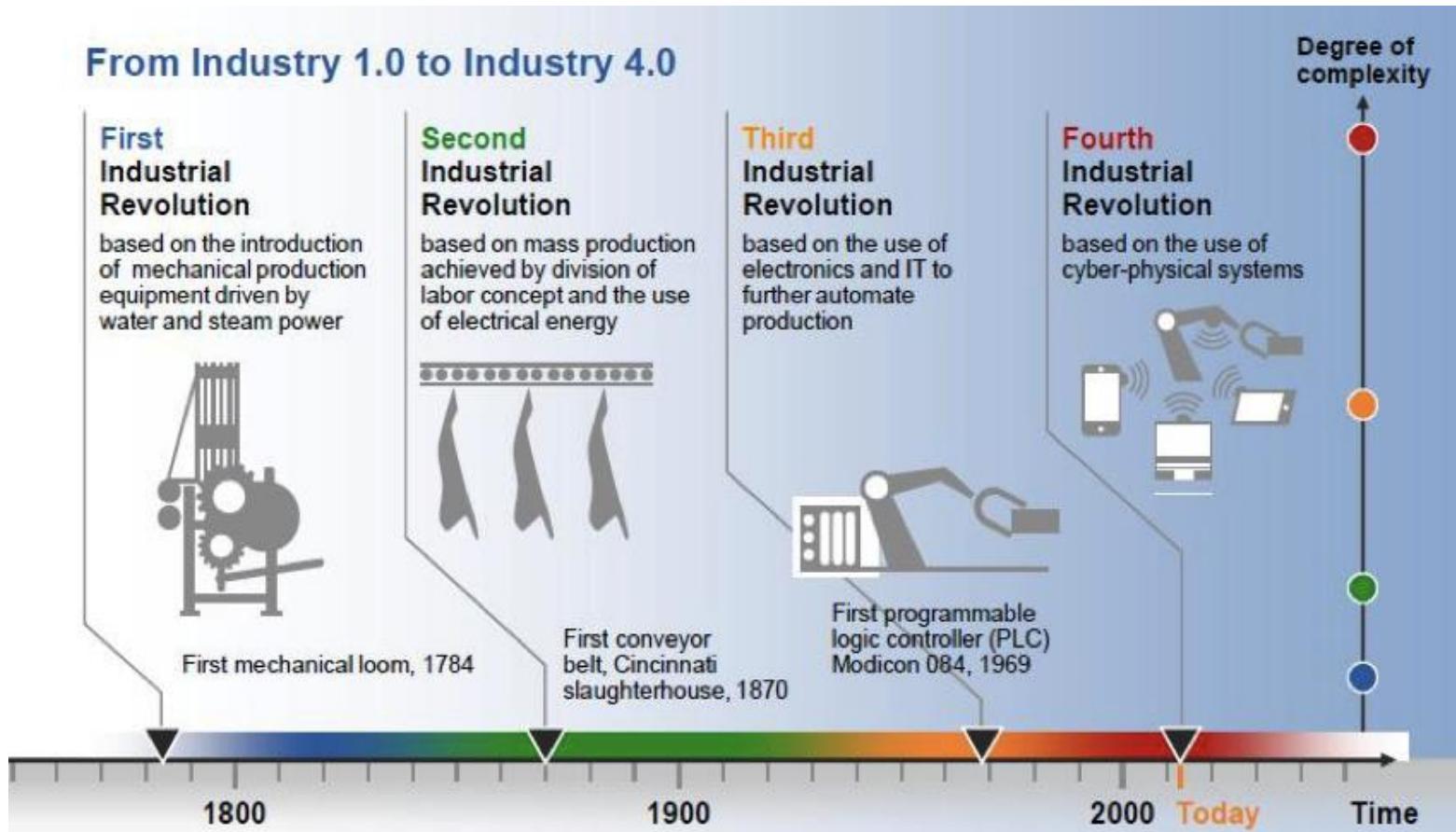
### Data

```
100100011101000000101000110110101010  
1001001110111000000111110010100100  
100001101011110101010011100001101001  
11111101000011011100101011100001011  
110011111011111110010000110101010  
01000011010011010000110000100010000  
0101111001001111011001110100010111  
00100001010110010100000100001001110  
011010011111001011101010101011100  
10001000010110010110101011000101  
0100100001001010111001100001010000  
01011000001001111010101010111010001  
0110111101011110001010001010001000  
0110100110110101000100010111001101  
00010100000110011000110010001001101  
10010101010001001110010101010111101
```

*"Information is the oil of the 21st century, and Analytics is the combustion engine."*

# INDUSTRY 4.0

## From Industry 1.0 to Industry 4.0



THE DATA ERA

# Applications of Data Analysis

- Marketing
  - Analysis of consumer behavior
  - Advertising campaigns
  - Targeted mailings
  - Segmentation of customers, stores, or products
- Finance
  - Creditworthiness of clients
  - Performance analysis of finance investments
  - Fraud detection
- Manufacturing
  - Optimization of resources
  - Optimization of manufacturing processes
  - Product design based on customer requirements
- Health Care
  - Discovering patterns in X-ray images
  - Analyzing side effects of drugs
  - Effectiveness of treatments

---

By analyzing all the data available, decision-makers can better assess competitive threats, anticipate changes in customer behavior, strengthen supply chains, improve the effectiveness of marketing campaigns, and enhance business continuity.

THE DATA ERA

According to a recent study by the MIT Sloan School of Management, organizations that use analytics are twice as likely to be top performers in their industry as those that don't.

# ANALISI DATI: COME CREARE VALORE

sapere esattamente cosa vuole il cliente (customer insight)



*(DATI+SERVIZIO)\*CLIENTE*

Creare valore al cliente / rendere la vita più semplice (make the world a better place)



Booking.com

INNOVATING NEW  
BUSINESS MODELS  
AND SERVICES

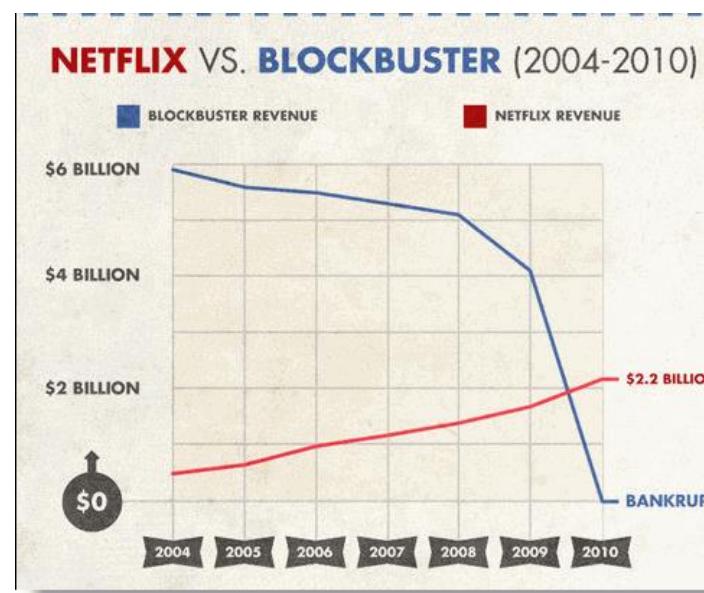
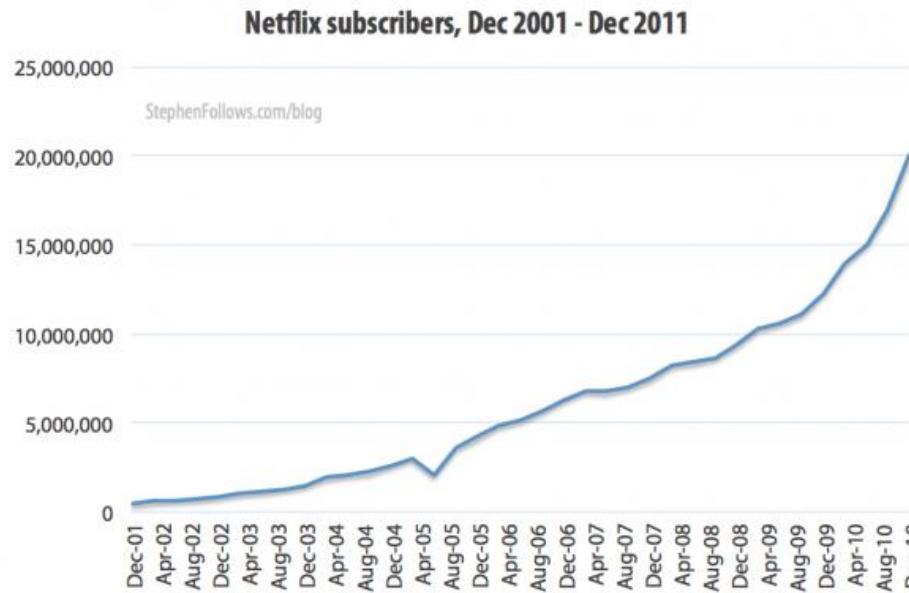
THE DATA ERA

"Consumer data will be the biggest differentiator in the next two to three years. Whoever unlocks the reams of data and uses it strategically will win."

Sophisticated analytics can substantially improve decision-making, minimise risks, and unearth valuable insights that would otherwise remain hidden.

Netflix is the world's leading Internet television network with over 75 million members in over 190 countries

- In 2000, Blockbuster refused to buy Netflix for \$50m.
- In 2013, Netflix generated \$4.4 billion in revenue.
- In just three months (July to Sept 2014) Netflix spent \$1.2 billion acquiring new content for its streaming service and a further \$16 million on new DVDs.
- In 2013 Netflix spent \$470m marketing its streaming service and just \$0.2m promoting its DVD service.



## THE DATA ERA



# (DATI+SERVIZIO)\*CLIENTE=RECOMMANDER

"WE ARE LEAVING THE AGE OF INFORMATION AND ENTERING THE AGE OF RECOMMENDATION"

- Recommandation: Estimate a utility function that automatically predicts how a user will like an item.
- As even improving recommendations a little bit could lead to dramatically increased revenue due to smaller churn
- Netflix spends \$150 million on content recommendations every year

## THE DATA ERA

*"In 2006 and 2009, Netflix announced a \$1M prize competition to advance recommendation algorithms.*



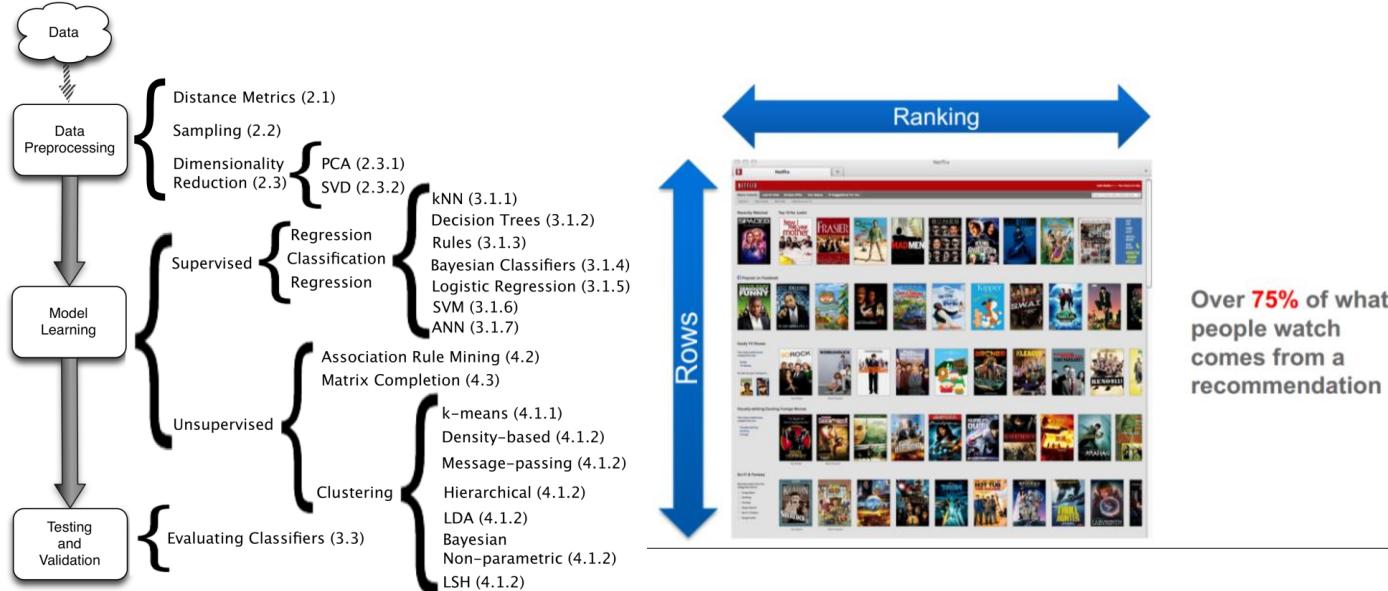
It is the combination of this data with cutting edge analytical techniques that makes Netflix a true Big Data company.

# (DATI+SERVIZIO)\*CLIENTE=RECOMMANDER

"WE ARE LEAVING THE AGE OF INFORMATION AND ENTERING THE AGE OF RECOMMENDATION"

The core of the Recommendation Engine can be assimilated to a general data mining problem

(Amatriain et al. Data Mining Methods for Recommender Systems in Recommender Systems Handbook)



## NETFLIX

Xavier Amatriain – August 2014 – KDD

## THE DATA ERA

Over 75% of what people watch comes from a recommendation

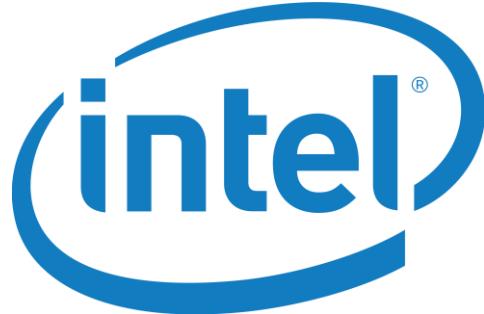


## RECOMMENDER

E-Commerce: Amazon.com, Ebay, Etsy.  
Music: Spotify, Pandora.  
Movie: Netflix, IMDB.  
News: Digg, Summly.  
Social Networks: LinkedIn, Facebook, Quora, YouTube  
Apps: Playstore, Cover

# ANALISI DATI: COME CREARE VALORE

INNOVATING PRODUCTS



*(DATI+PRODOTTO)\*INNOVARE*



THE DATA ERA

Manufacturers are using data obtained from sensors embedded in products to create innovative after-sales service offerings such as proactive maintenance to avoid failures in new products.

# ANALISI DATI: COME CREARE VALORE

INNOVATING PRODUCTS



For instance, tire manufacturer Pirelli collects data on tire pressure, temperature and wear-and-tear on trucks using sensors in its tires. This data gives Pirelli competitive advantage and help to improve Tyre quality, maintenance but it can also be monetize externally.

Car manufacturer can buy this data to understand driving patterns to improve there service.

Pirelli also offers that data as an add-on service for fleet managers and insurers, InfoWorld reports.



THE DATA ERA

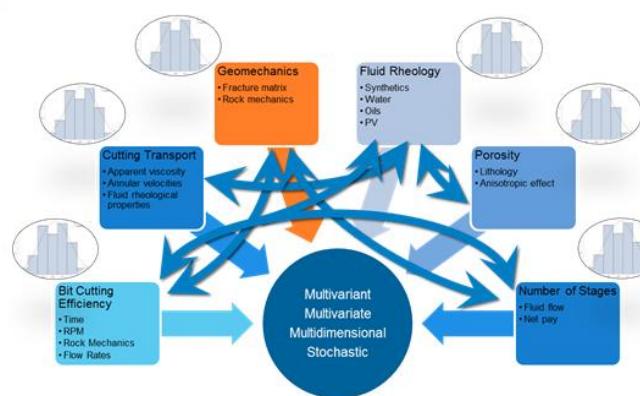
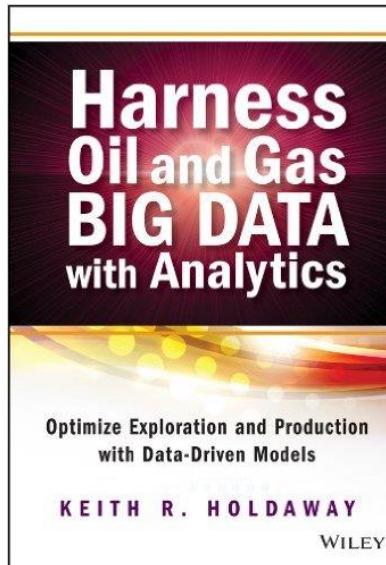
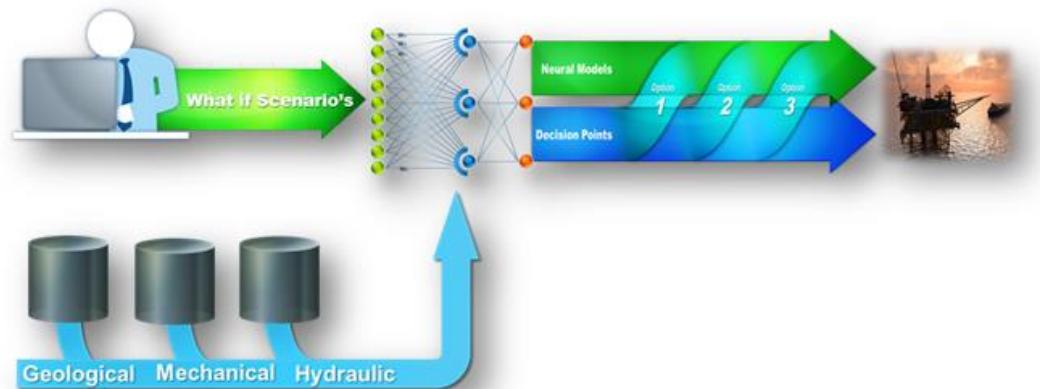


# ANALISI DATI: COME CREARE VALORE

INNOVATING PRODUCTS



Data from any prospective oil field can then be compared alongside that from thousands of others around the world, to enable geologists to make more accurate recommendations about where to drill.



Shell also uses Big Data to ensure its machines are working properly and spending as little time as possible offline due to breakdowns and failure.

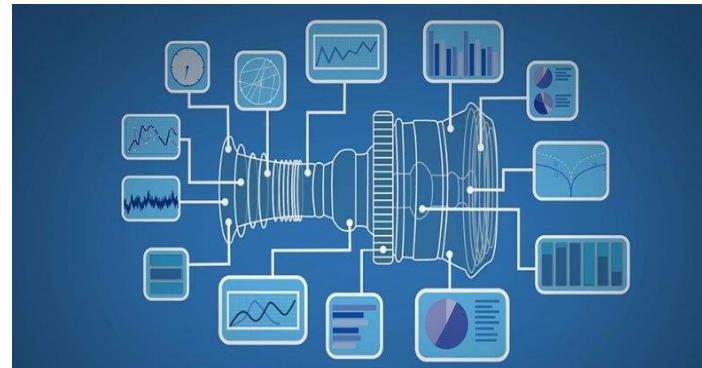
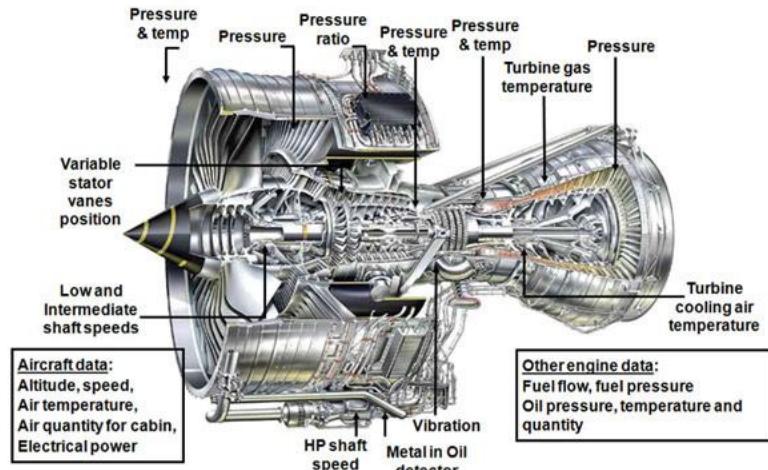
THE DATA ERA

# ANALISI DATI: COME CREARE VALORE

INNOVATING PRODUCTS

## Engine health management

Rolls Royce uses Engine Health Management (EHM) to track the health of thousands of engines operating worldwide, using onboard sensors and live satellite feeds.



200 sensors across the turbine generate 300 data points per second of performance and operation every hour.

EHM is a pro-active technique for predicting when something might go wrong and averting a potential threat before it has a chance to develop into a real problem.



THE DATA ERA

# ANALISI DATI: PROSPETTIVE FUTURE



## THE DATA ERA

### The Technology Opportunity: Monetize Big Data, Underpinned by the Cloud

#	#	Priority	2016	2015	2014	%
1	1	BI/Analytics	39%	41%	50%	▼
2	2	Infrastructure and Data Center	27%	31%	37%	▼
3	3	Cloud	25%	27%	32%	▼
4	4	ERP	21%	26%	34%	▼
5	5	Digitalization/Digital Marketing	21%	17%	11%	▲
6	6	Mobile	20%	24%	36%	▼
7	7	Security	15%	13%	11%	▲
8	8	Networking, Voice and Data Communications	10%	12%	12%	▼
9	9	Legacy Modernization	10%	7%	7%	▲
10	10	Industry-Specific Applications	9%	9%	10%	△
11	11	CRM	9%	11%	8%	▼

Note: Percentages represent the proportion of CIOs citing each priority as one of their top three areas of new IT spending.

#GartnerSYM

© 2016 Gartner, Inc. and/or its affiliates. All rights reserved. Gartner and ITXPO are registered trademarks of Gartner, Inc. or its affiliates.

Gartner  
SYMPOSIUM ITXPO 2015

"Analytics is not a technology issue. It's a strategy and operational issue." Chris Mazzei, EY

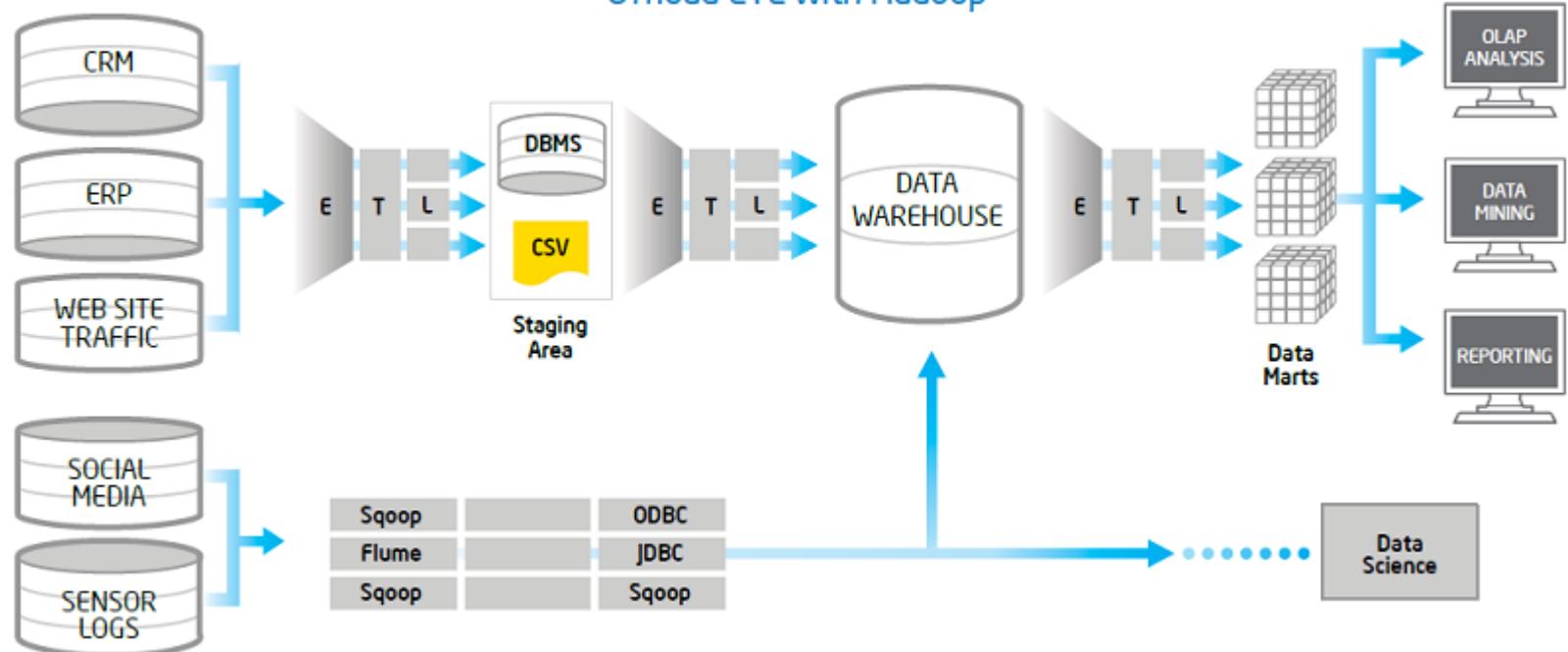
# What is big data?

Big Data is essentially a special application of data science, in which the data sets are enormous and require overcoming ***logistical challenges*** to deal with them. The primary concern is efficiently capturing, storing, extracting, processing, and analyzing information from these enormous data sets..

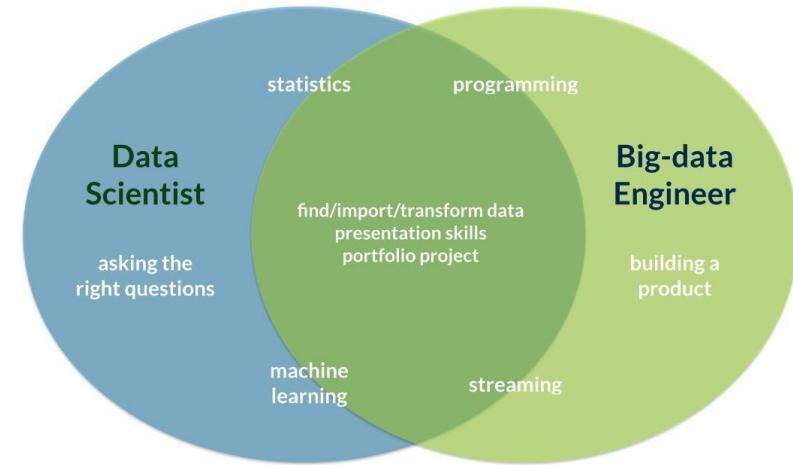
Hadoop is a technology of big data word.

Instead of using a database with structure, indexes and accelerated queries, the data is just dumped into hadoop, and when you have figured out what to do, you re-read all your data and extract the information you really need

Offload ETL with Hadoop

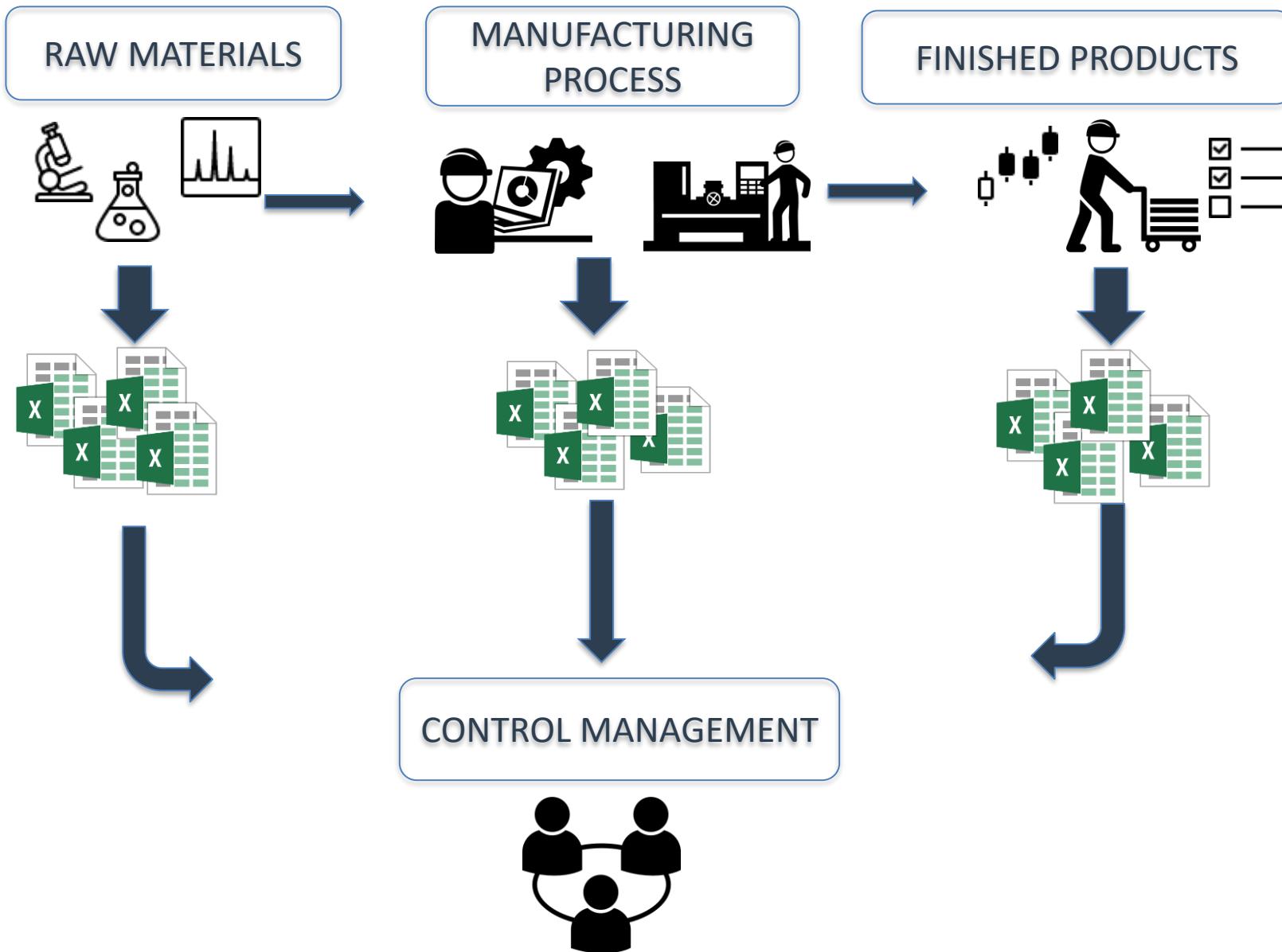


THE DATA ERA



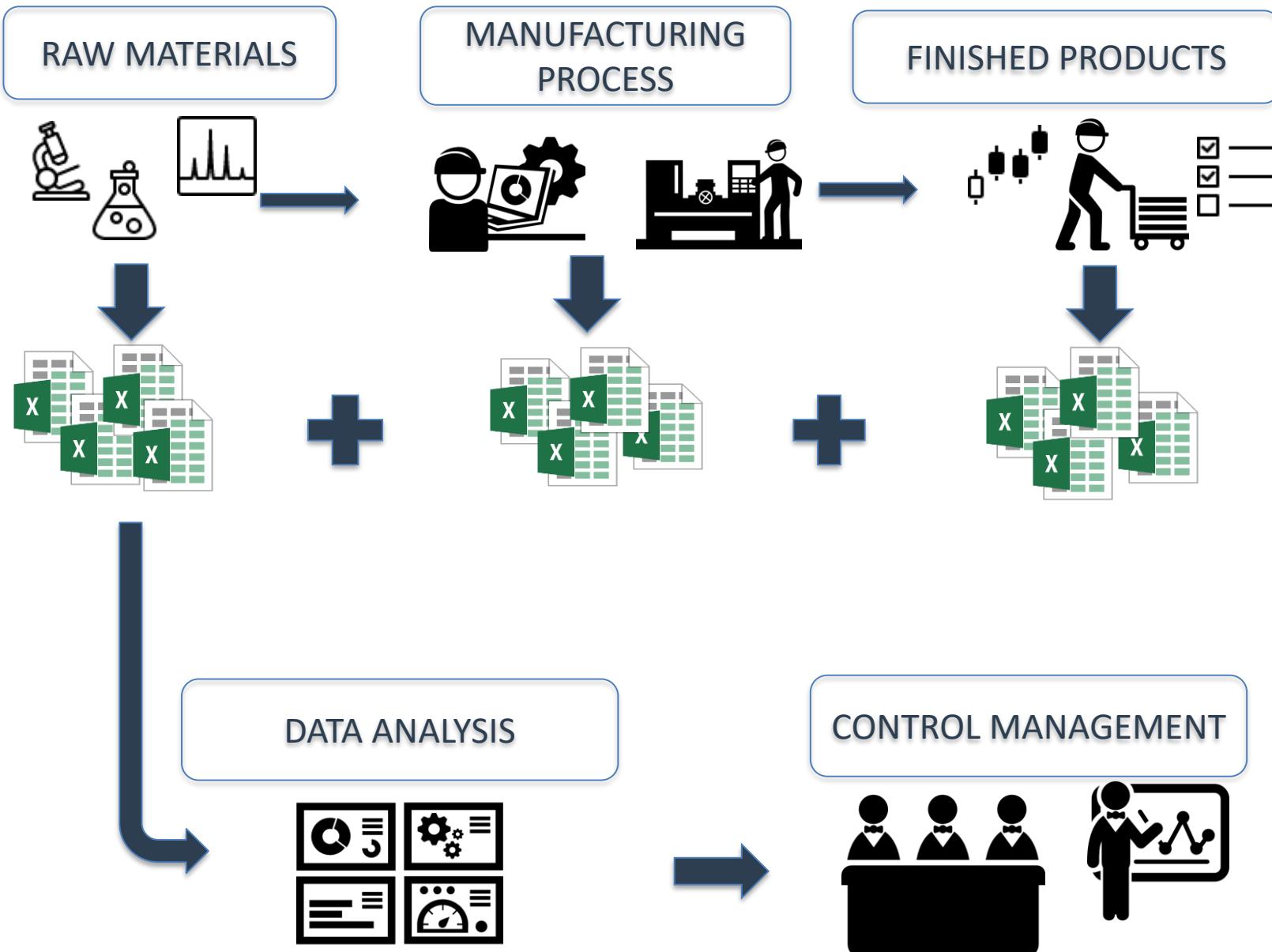
# FACTORY

TODAY



TOMORROW

# SMART FACTORY



# **“Data-Driven Innovation”: 4 steps per trasformare i dati in innovazione**

# ANALISI DATI

- Non c'è alcuna formula sicura e replicabile
- L'80% del lavoro di analisi viene svolto per raccogliere organizzare, ripulire e trasformare i dati
- L'80% dei software in commercio svolgo solo una parte dell'attività di analisi: la modellizzazione
- I risultati saranno contro-intuitivi e mai certi al 100%.

DATA DRIVEN INNOVATION



Davide Camera 2nd  
Founder & Senior Partner at Excelle

Veramente un bell'articolo, complimenti Alessio Passalacqua hai focalizzato in maniera egregia, imho, i warning di chi fa questo lavoro. hai un nuovo follower, complimenti ancora

Like · Reply · 5 days ago



# ANALISI DATI

FROM DATA TO  
BUSINESS



DATA DRIVEN INNOVATION

**Manuale Utente**

L'analisi dei dati è un processo di ispezione, pulizia, trasformazione e modellazione di dati con il fine di evidenziare informazioni che suggeriscano conclusioni e supportino le decisioni strategiche aziendali.

Da Wikipedia, l'enciclopedia libera.

# I DATI

## Data Source

### Internal



#### Structured



##### Human-Generated

- Survey ratings
- Aptitude testing
- Machine-Generated**
- Web metrics from Web logs
- Product purchase from sales Records
- Process control measures

#### Unstructured



##### Human-Generated

- Emails, letters, text messages
- Audio transcripts
- Customer comments
- Voicemails
- Corporate video/communications
- Pictures, illustrations
- Employee reviews

### External



##### Human-Generated

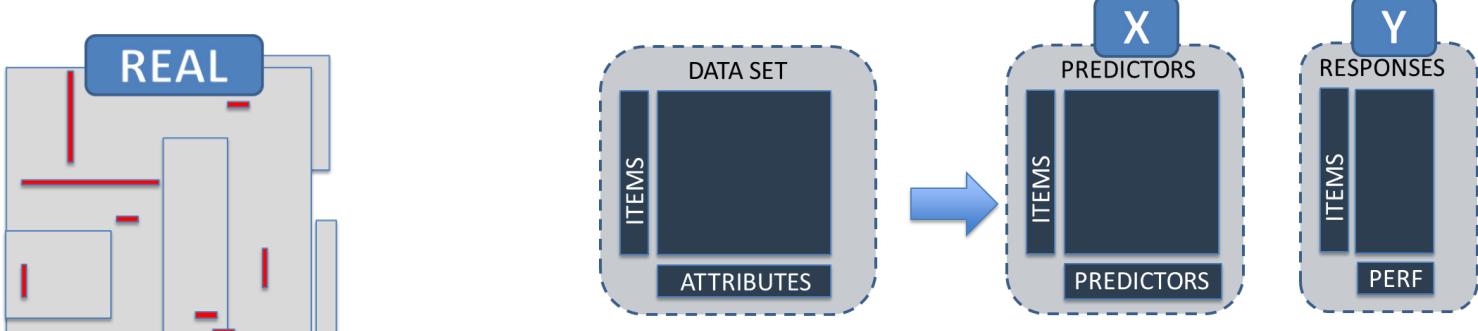
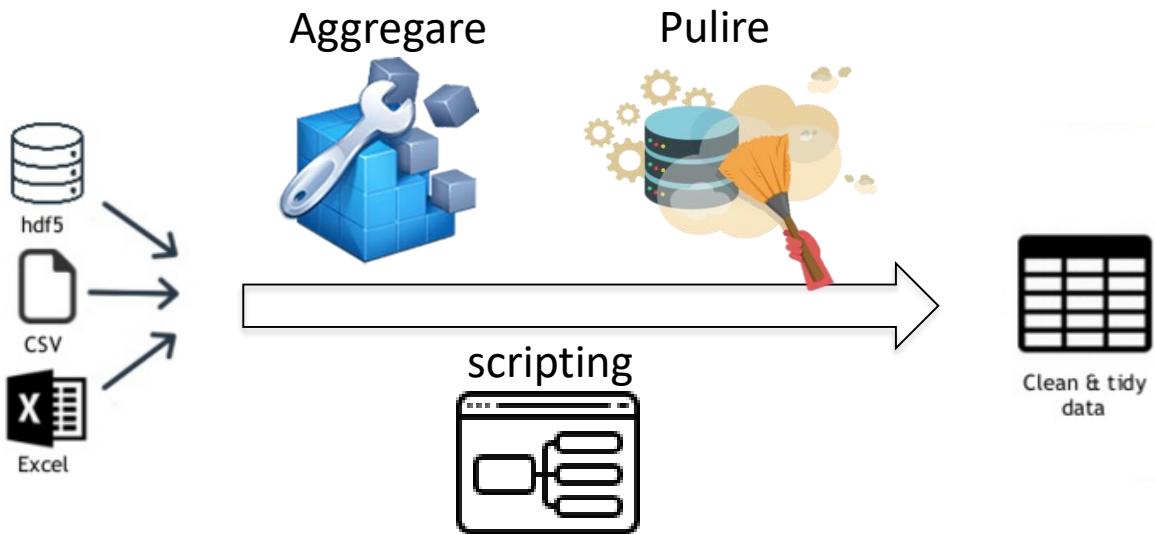
- Number of Retweets, Facebook likes, Google Plus +1s
- Ratings on Yelp
- Patient ratings ratings
- Machine-Generated**
- GPS for tweets
- Time of tweet/updates/postings

##### Human-Generated

- Content of social media updates
- Comments in online forums
- Comments on Yelp
- Video reviews
- Pinterest images
- Surveillance video

DATA DRIVEN INNOVATION

# Preparazione Dati



Selezionare  
Aggregare  
Manipolare  
Pulire

# ANALISI DATI: TOOLBOX



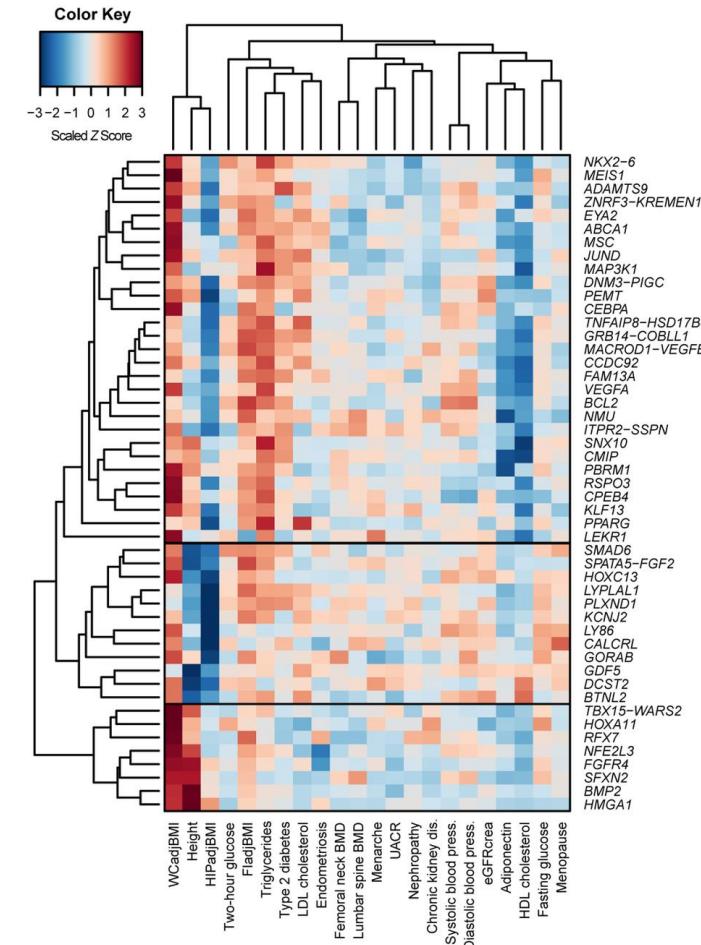
DATA DRIVEN INNOVATION

# Analisi Esplorativa

## Unsupervised

Find underlining relations in data by observing the raw data only  
(without the expected output).

- ✓ Clustering
- ✓ Dimensionality reduction



DATA DRIVEN INNOVATION

Comprendere  
Visualizzare  
Ridurre le variabili  
Gestire gli outliers

# Data clustering

**Goal:** discover natural groupings among given data points

Unsupervised learning (unlabeled data)

Exploratory analysis (without any pre-specified model/hypothesis)

DATA DRIVEN INNOVATION

## Usages

- Gain insight from the underlying structure of data (salient features, anomaly detection, etc)
- Identify degree of similarity between points (infer phylogenetic relationships)
- Data Compression (summarizing data by cluster prototypes, removing redundant patterns)

# Partitional vs. Hierarchical

## 1. Partitional algorithms (k-means)

Partition the data space

Finds all clusters simultaneously

## 2. Hierarchical algorithms

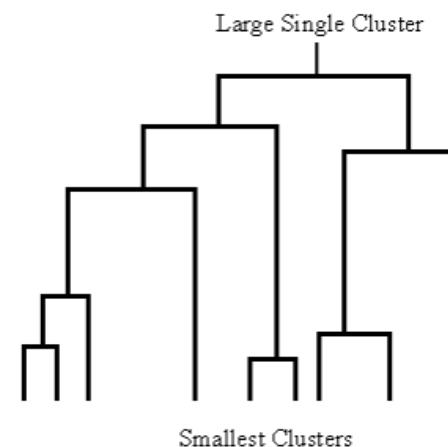
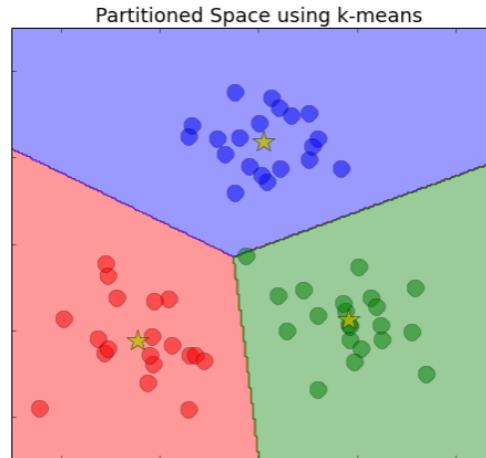
Generate nested cluster hierarchy

Agglomerative (bottom-up)

Divisive (top-down)

Distance between clusters:

Single-linkage, complete linkage,  
average-linkage

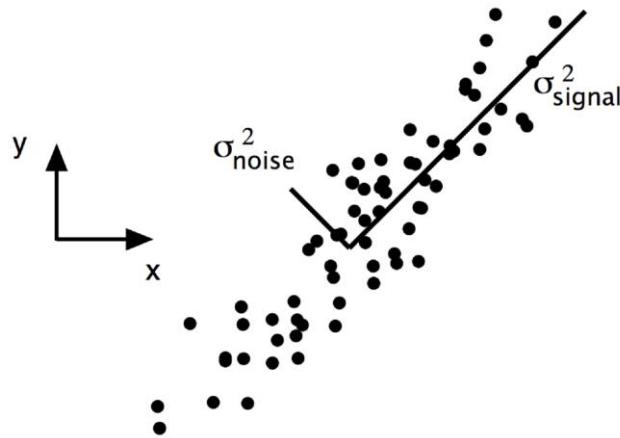


DATA DRIVEN INNOVATION

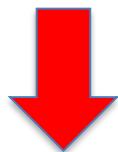
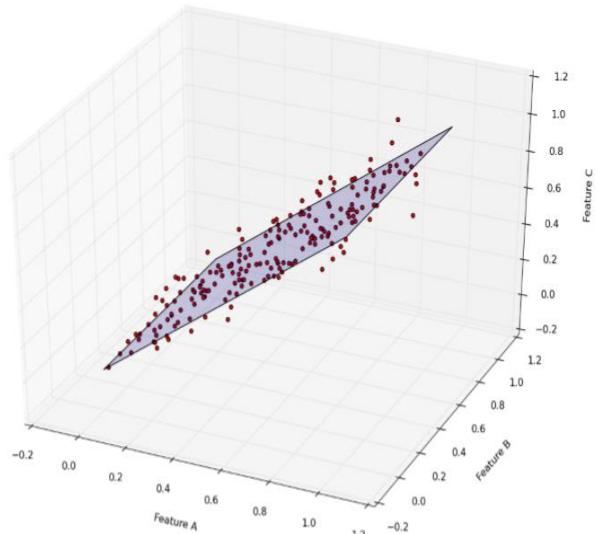
# PRINCIPAL COMPONENT ANALYSIS

LA PCA è una trasformazione geometrica. Una rotazione degli assi.

PCA ricerca delle direzioni privilegiate che massimizzano la variazione dei dati ed elimina le correlazioni



MASSIMIZZA L'INFORMAZIONE, misurata dalla varianza  
MINIMIZZA LA RIDONDANZA, misurata dalla covarianza



TOGLIE la ridondanza

FILTRA il rumore

TECNICA DI  
RIDUZIONE E  
INTERPRETAZIONE DEI  
DATI

DATA DRIVEN INNOVATION

TRASFORMAZIONE SPAZIALE

# LINEAR EXAMPLE: PCA

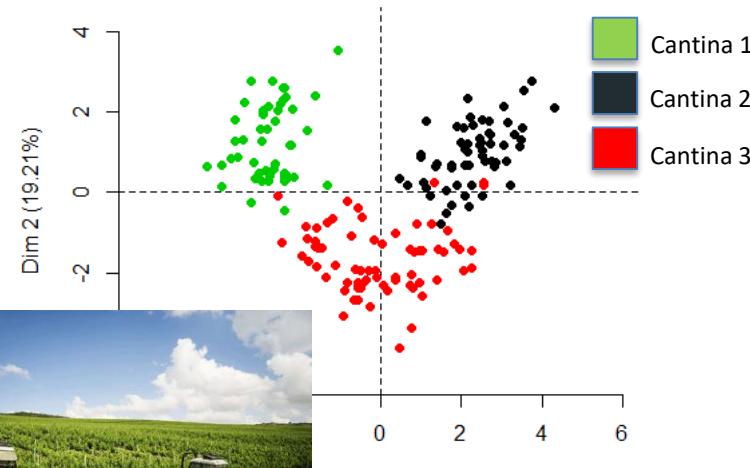
SPACE  
TRANSFORMATION

La trasformazione geometrica evidenza maggiormente la struttura dei dati

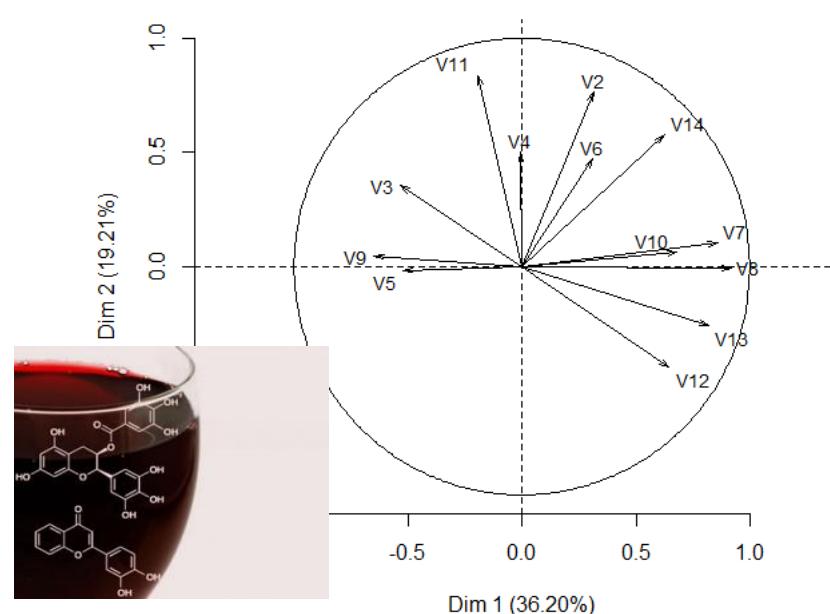
Interpretare il fenomeno attraverso il significato assunto dalle componenti principali

## WINE DATA SET

Individuals factor map (PCA)



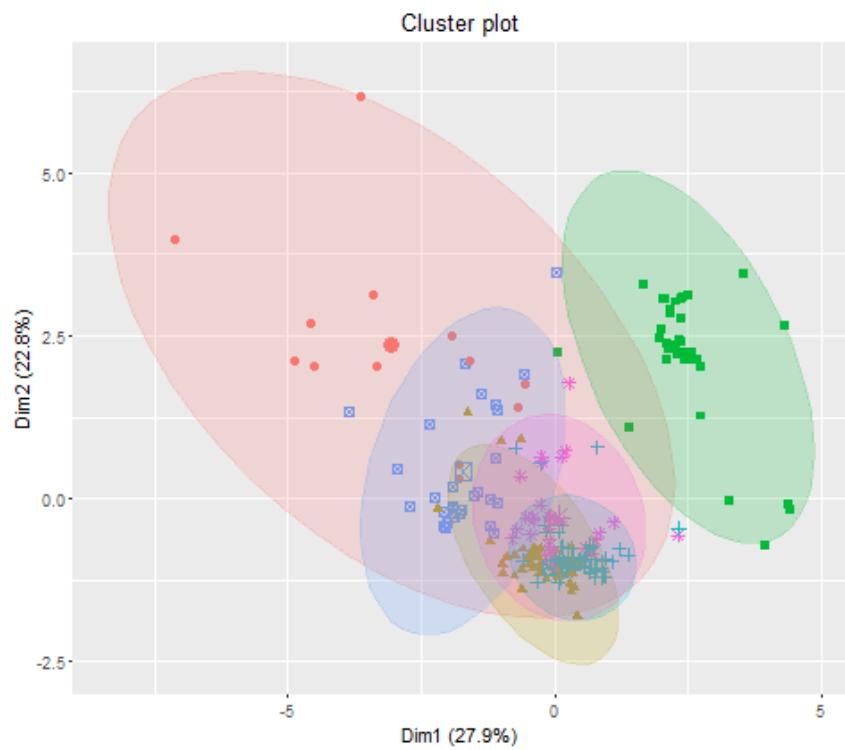
## Variables factor map (PCA)



DATA DRIVEN INNOVATION

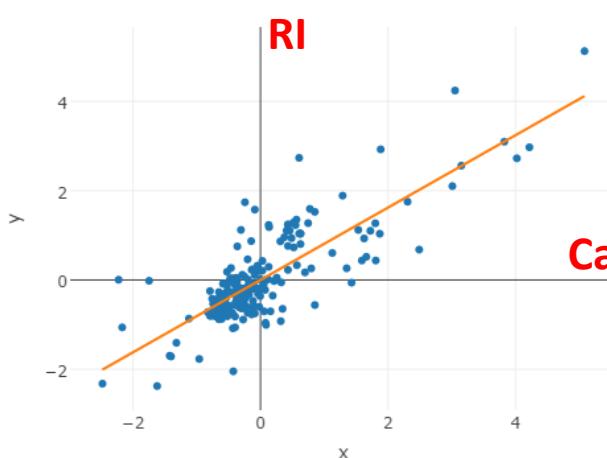
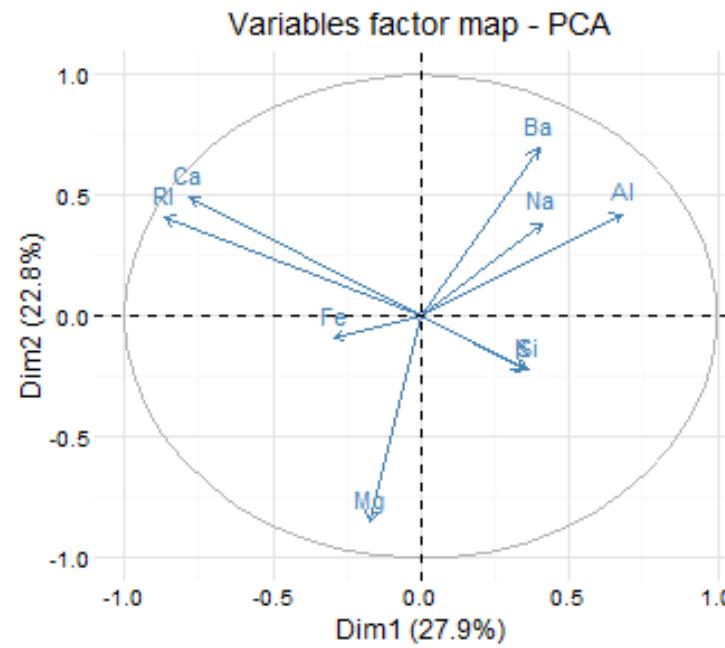
RIVELA LE CARATTERISTICHE PIU' IMPORTANTI  
E RIDUCE LE DIMENSIONI

# Analisi Esplorativa



cluster

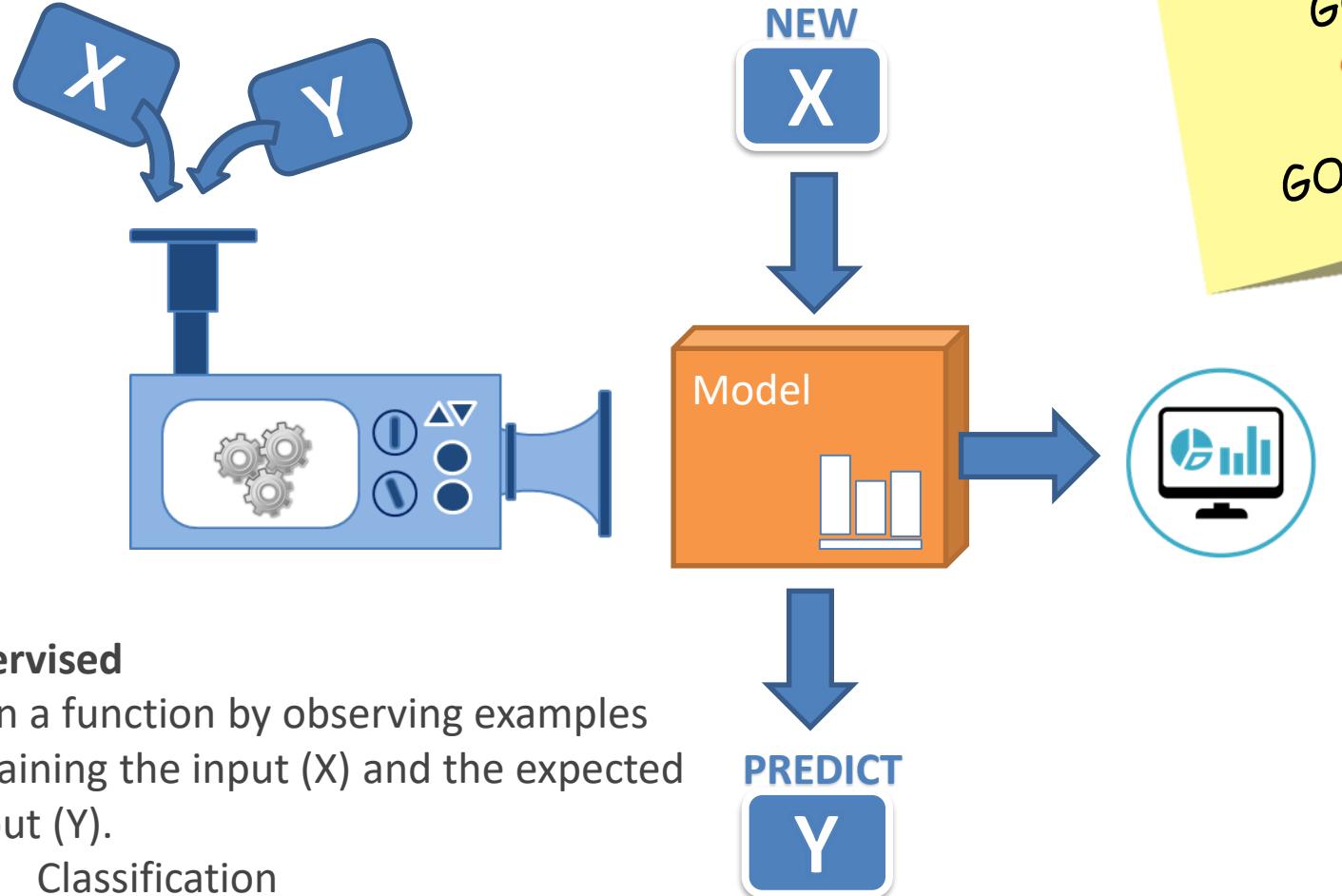
- 1
- 2
- 3
- 4
- 5
- 6



DATA DRIVEN INNOVATION

The dataset is collection of chemical analysis performed on glass samples. The main field of application of this analysis is criminological investigation, where the correct identification of the origin of a glass fragment can be very important.

# MODELING



## Supervised

Learn a function by observing examples containing the input ( $X$ ) and the expected output ( $Y$ ).

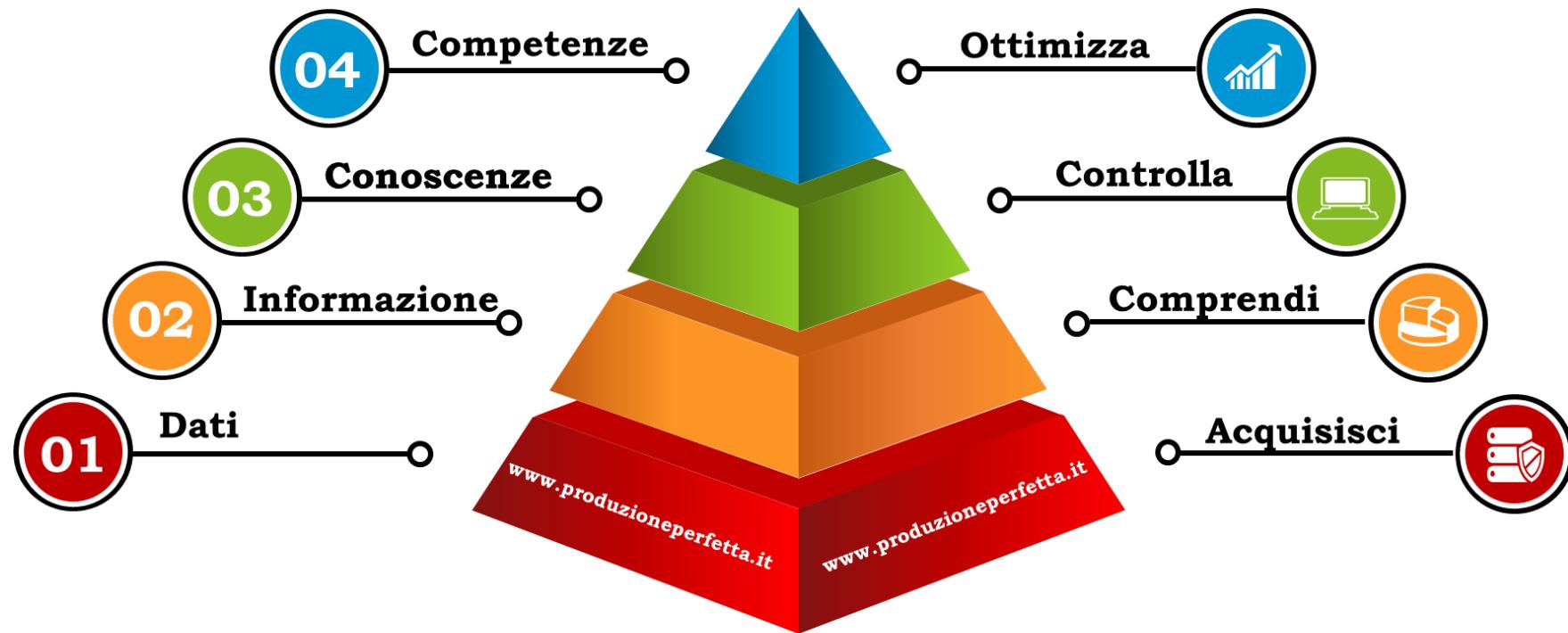
- Classification
- Regression

## DATA DRIVEN INNOVATION

il segreto del successo dell'analisi dati è avere dati di qualità. Il modello è una rappresentazione approssimata del sistema reale. Se i tuoi dati descrivono al meglio il tuo sistema allora il tuo modello descriverà al meglio la realtà.

# DATA DRIVEN INNOVATION

TRASFORMA I TUOI DATI IN  
INNOVAZIONE



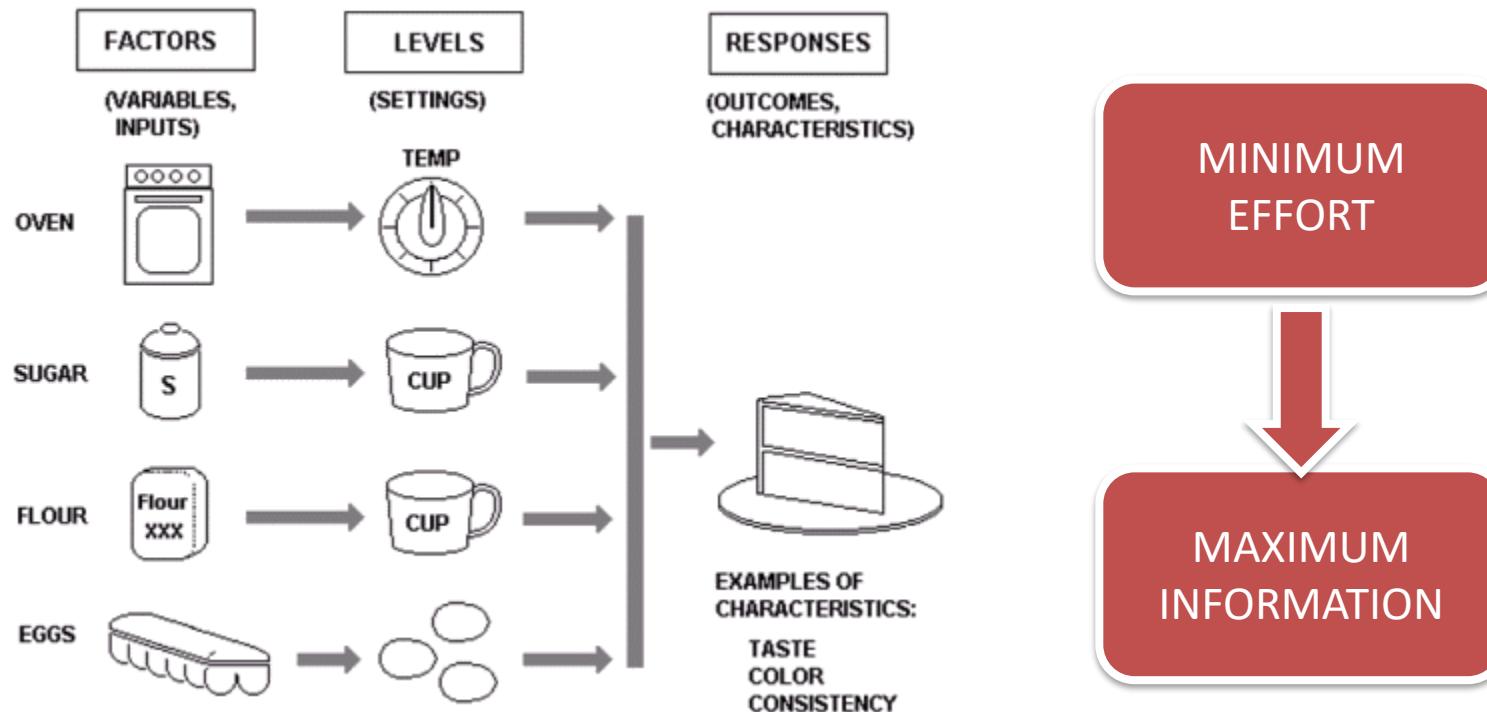
DATA DRIVEN INNOVATION

Quality is the consequence of  
knowledge

# DESIGN OF EXPERIMENTS

Experimental design can be used at the point of **greatest leverage** to reduce design costs by speeding up the design process, reducing late engineering design changes, and reducing product material and labor complexity.

## DATA DRIVEN INNOVATION



DOE in general is a useful method to problem solving, optimization, designing products, and manufacturing and engineering.

COFFEE  
**BREAK THE**  
PROJECTS



**10** MINUTES  
BREAK

# **Confronto con i metodi di analisi statistica**

## **Six-Sigma**

# Control Chart

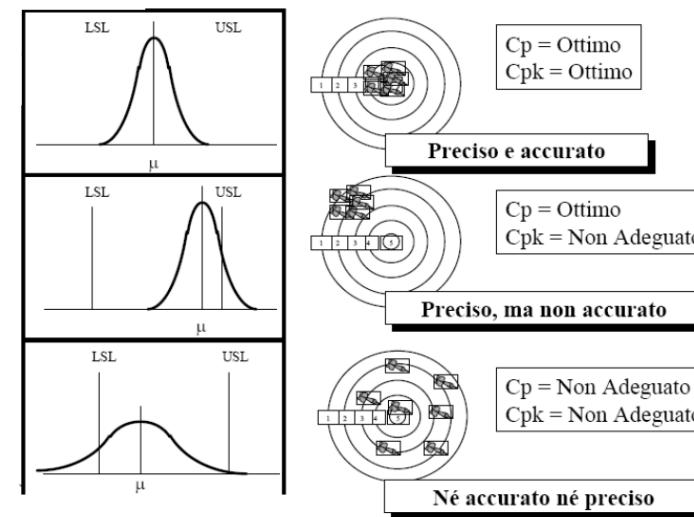
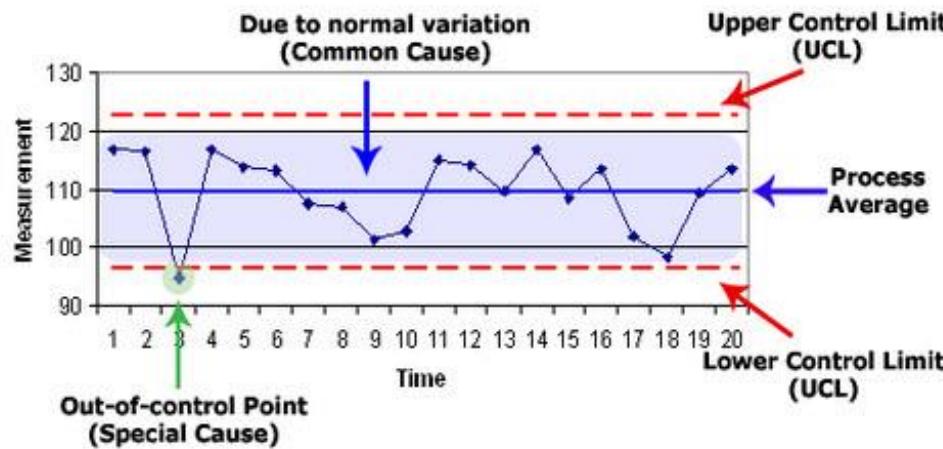
La **centerline** rappresenta la media della caratteristica di qualità corrispondente allo stato “in controllo. Ci sono poi altre due linee orizzontali che rappresentano i limiti di controllo: superiore e inferiore (**UCL = upper control limit**, **LCL = lower control limit**). I limiti di controllo sono scelti in modo tale che, se il processo è in controllo, la gran parte dei punti sta tra questi limiti (in genere si sceglie il 99.7% della variabilità naturale che, nel caso gaussiano, significa  $\pm 3\sigma$  dalla media). La carta di controllo viene usata per monitorare nel tempo eventuali modifiche della variabilità statistica che rivelano la presenza di cause speciali.

Se **Cp > 1** il processo è capace perché la **variabilità** naturale è inferiore a quella delle specifiche e questo garantisce che il processo fornisca prodotti che soddisfino le specifiche. L'indice Cp non tiene conto di dove è localizzata la media del processo rispetto alle specifiche.

Per tenere più accuratamente conto della **centratura** del processo è definito il **Cpk**, che si tratta di un indice unilaterale rispetto al limite di specifica più vicino alla media.

Mentre l'indice Cpk misura la capacità effettiva del processo, l'indice Cp misura la capacità potenziale del processo.

## STATISTICAL PROCESS CONTROL

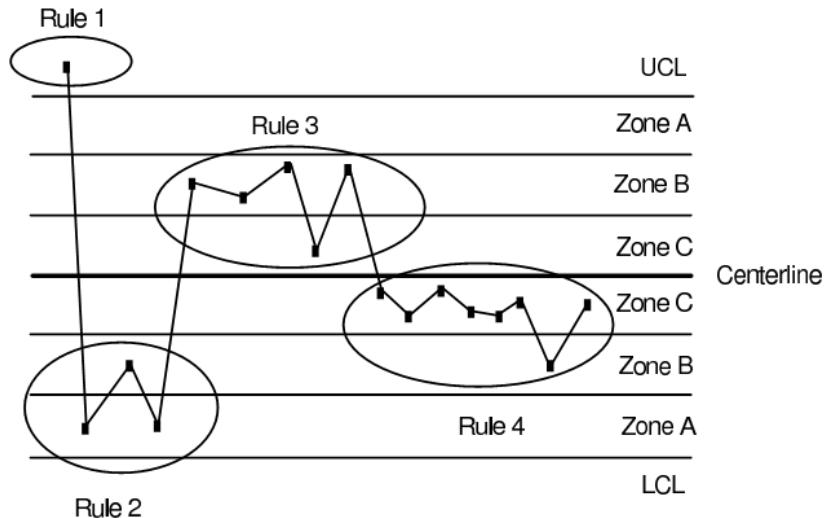


# STATISTICAL PROCESS CONTROL

Se il processo è in controllo, allora i punti rappresentati sulle carte di controllo devono avere un andamento aleatorio. Ci sono 4 regole comunemente usate per vedere se un processo è fuori controllo (Western Electric Rules, 1956):

1. uno o più punti fuori dai limiti di controllo;
2. due punti su tre consecutivi fuori dai “warning limits (banda  $\pm 2\sigma$ );
3. 4 punti su 5 consecutivi fuori dalla banda  $\pm \sigma$ ;
4. 8 punti consecutivi dallo stesso lato della centerline.

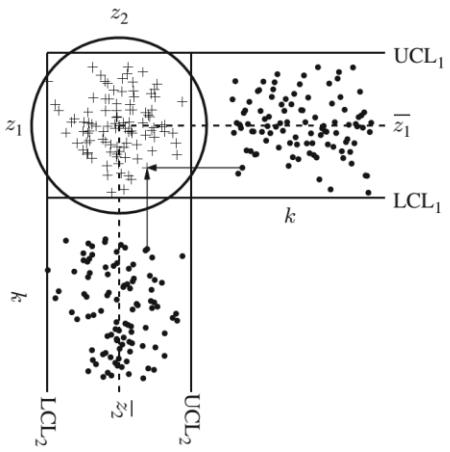
STATISTICAL PROCESS CONTROL



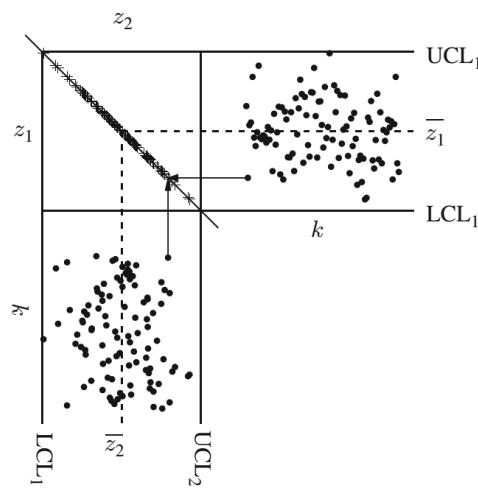
# LIMITATION

Their ease of use, implementation and interpretation make them very popular in practice. Besides the literature describing univariate SPC is very widespread. Unfortunately, **traditional (Shewhart) control charts are not suitable for monitoring many variables** mainly because of separate significance level for any chart, **not taking into account a relationship** between variables and “organizational mess” connected with using separate control charts for every important process/product variable.

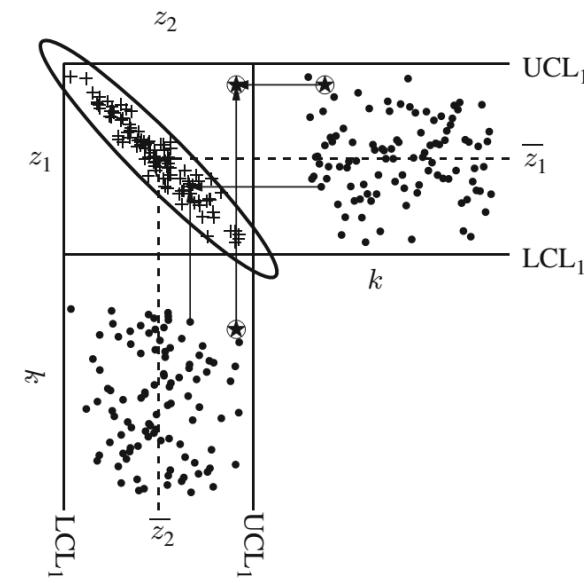
No correlation



Perfect correlated



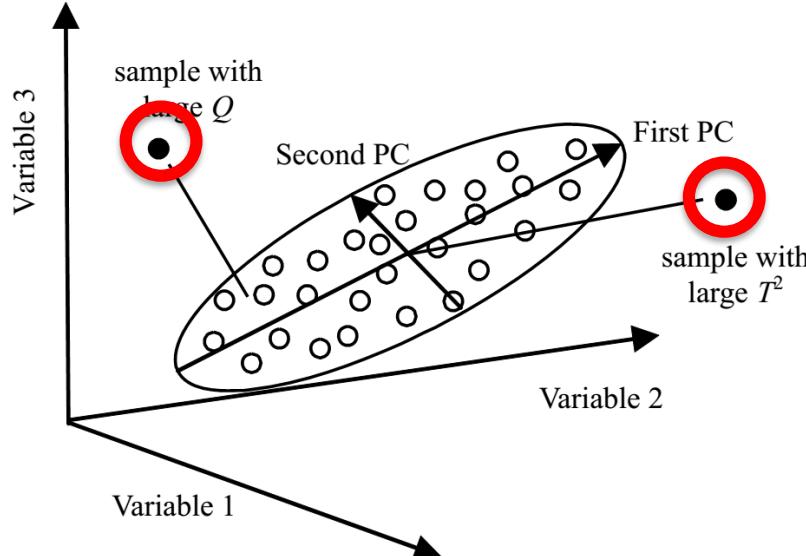
Highly correlated



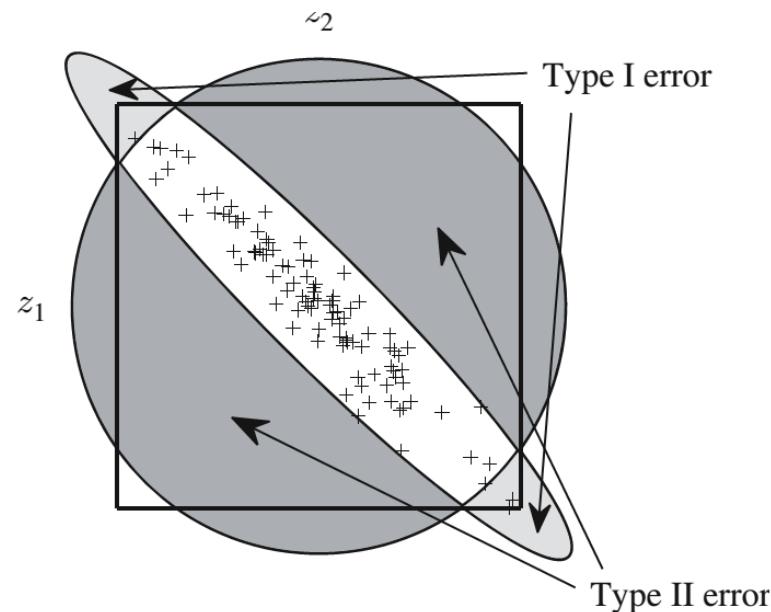
STATISTICAL PROCESS CONTROL

# MULTIVARIATE STATISTICAL PROCESS CONTROL

Changes in relationships between variables



unusual variability within the normal subspace.



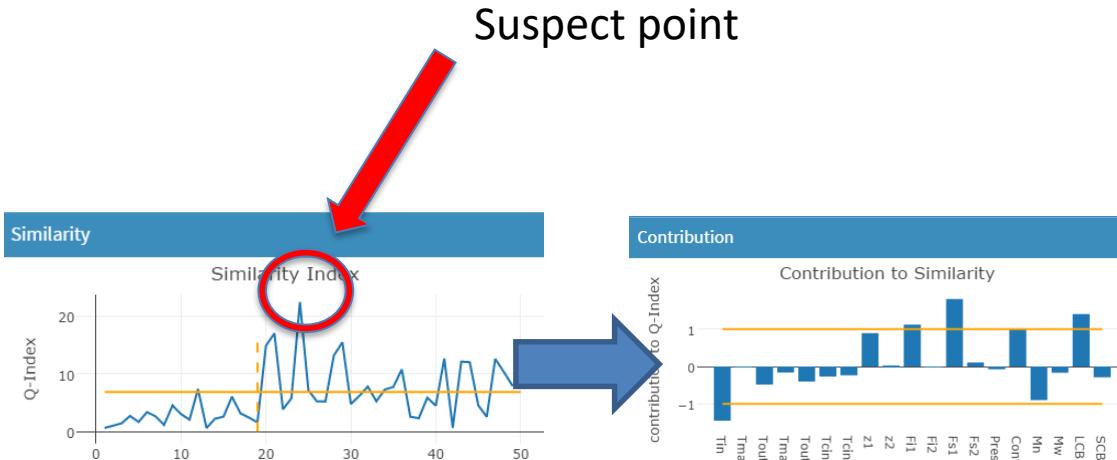
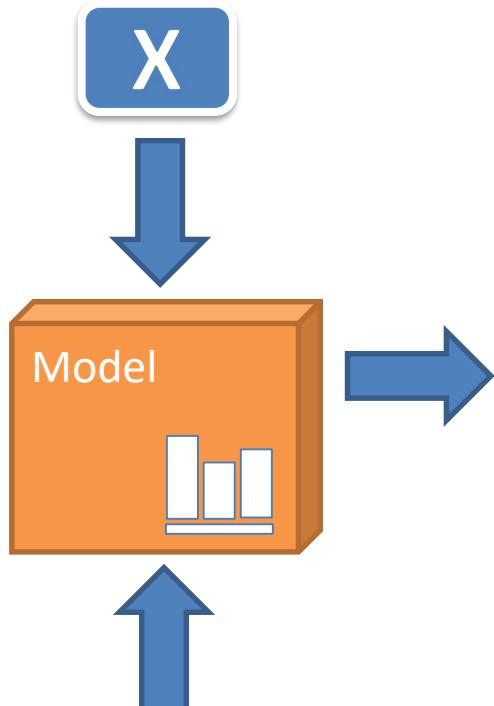
- The outliers are not detected until you look at the combination of the variables
- The information is found in the correlation pattern - not in the individual variables!

STATISTICAL PROCESS CONTROL

# MULTIVARIATE STATISTICAL PROCESS CONTROL

MONITORING

NEW Measurement



Normal Operation Condition



Root causes analysis



Operator intervention

STATISTICAL PROCESS CONTROL

procedures to monitor multivariate process variation over time, to determine whether the process is stable, and to detect and diagnose changes in a stable process.

# APP per il monitoraggio di Processo

MONITORING



STATISTICAL PROCESS CONTROL

Intuitive web-based interface  
for operators, engineers and  
managers



Interactive



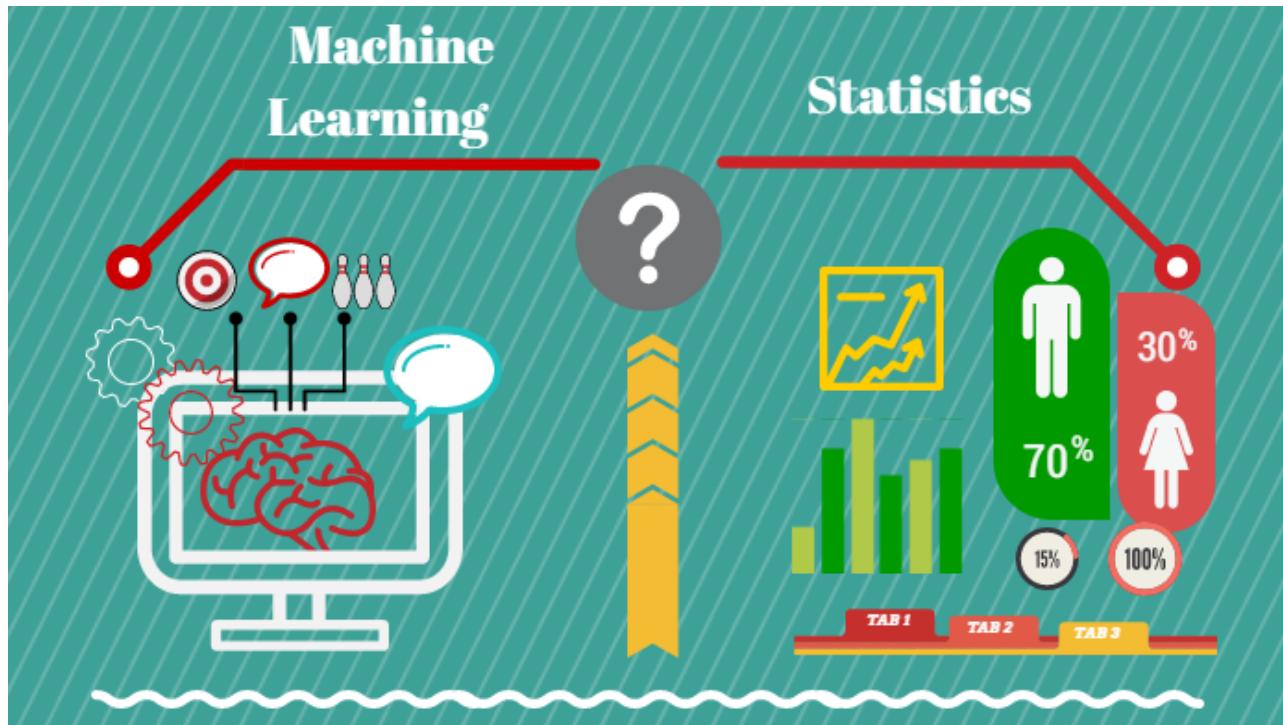
Easy to Use



Web-based

# **Confronto tra i metodi di analisi statistica avanzata e di “Machine Learning”**

# MACHINE LEARNING



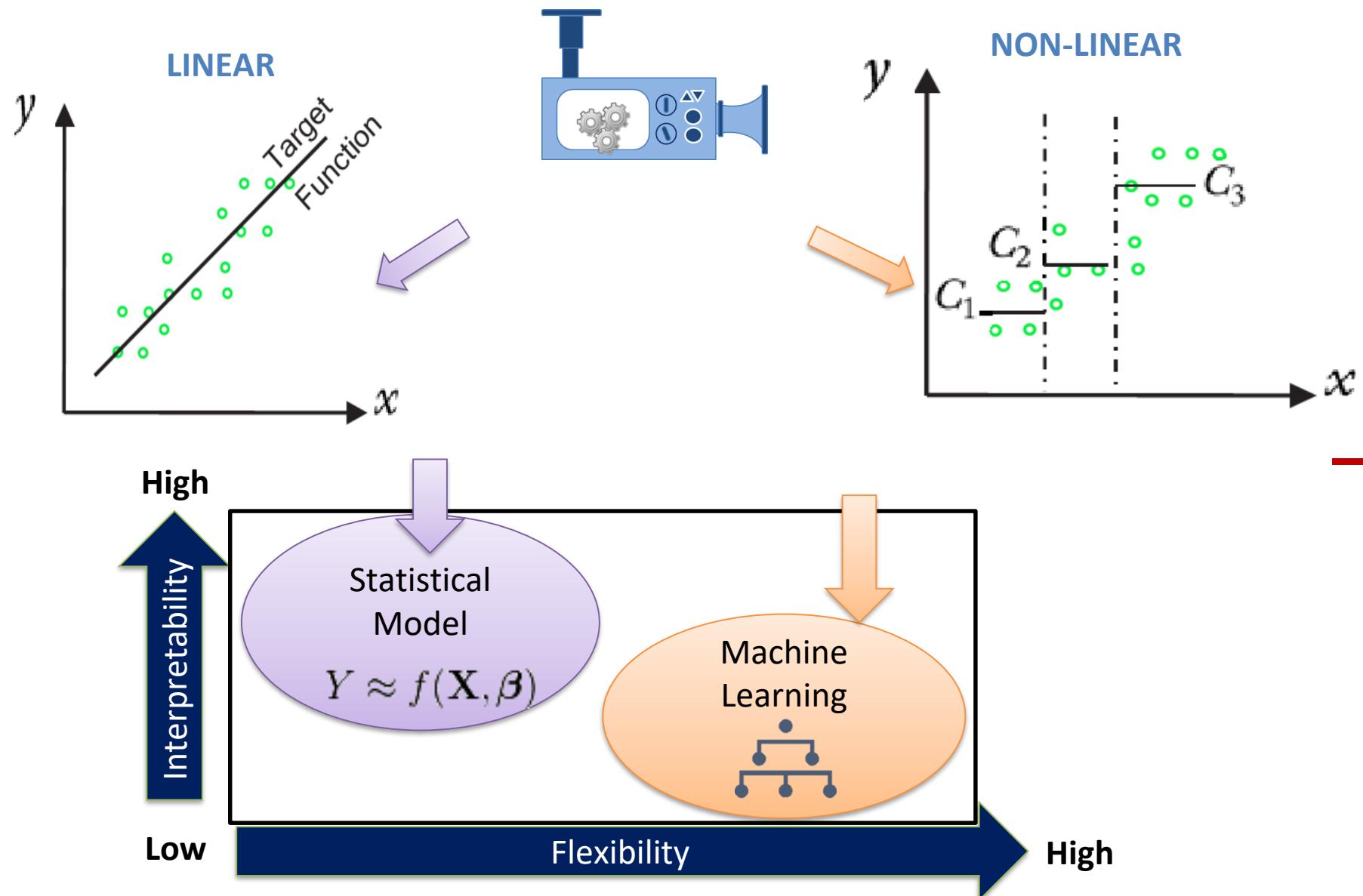
*“Learning is any process by which a system improves performance from experience.”*

*“Machine learning is a very hot topic for many key reasons, and because it provides the ability to automatically obtain deep insights, recognize unknown patterns, and create high performing predictive models from data, all without requiring explicit programming instructions.”*

STATISTICS VS MACHINE LEARNING

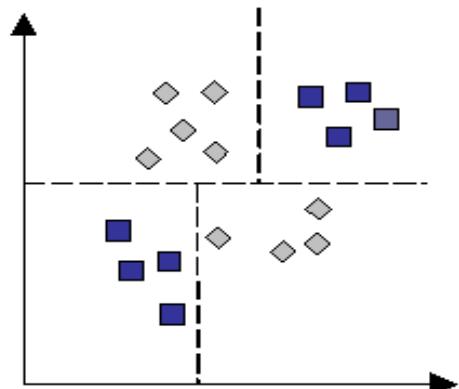
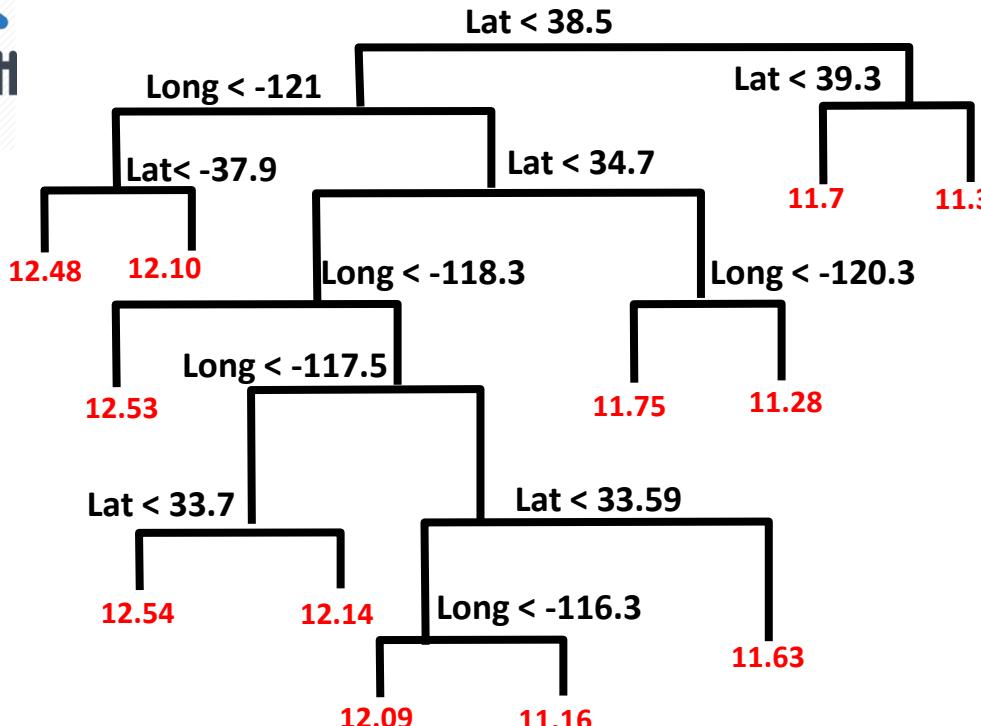
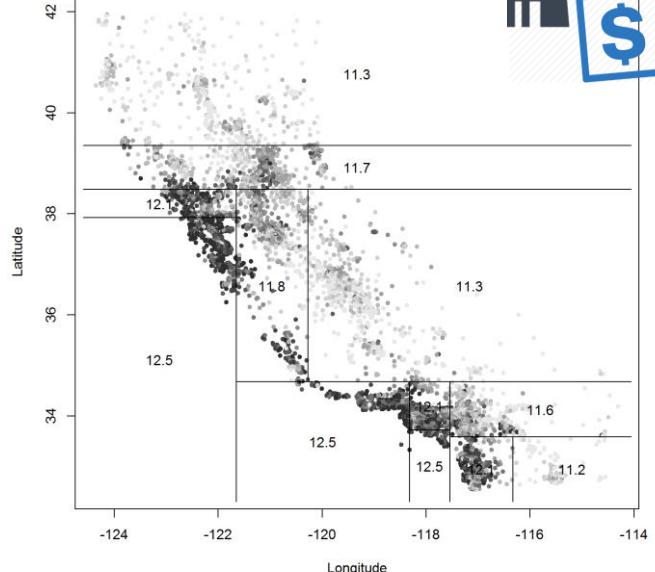
# MODELING

## STATISTICAL MODEL VS MACHINE LEARNING



# MACHINE LEARNING: DECISION TREES

HOUSE VALUE



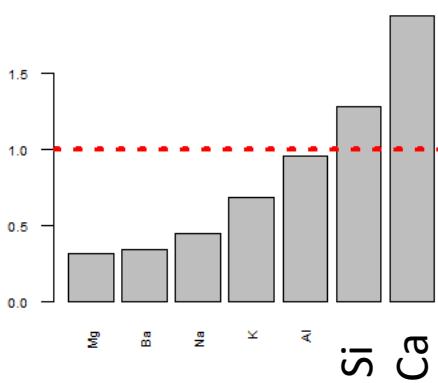
- Decision trees partition the feature space by splitting the data
- Learning the decision tree consists in finding the order and the split criterion for each node

STATISTICS VS MACHINE LEARNING

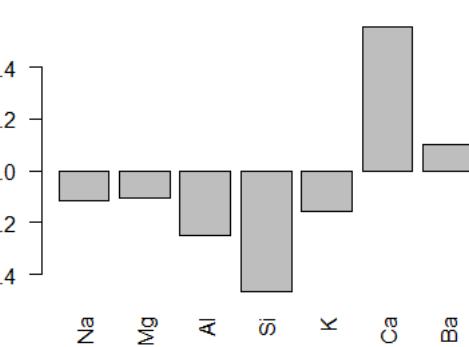
NON LINEAR

# VIRTUAL METROLOGY

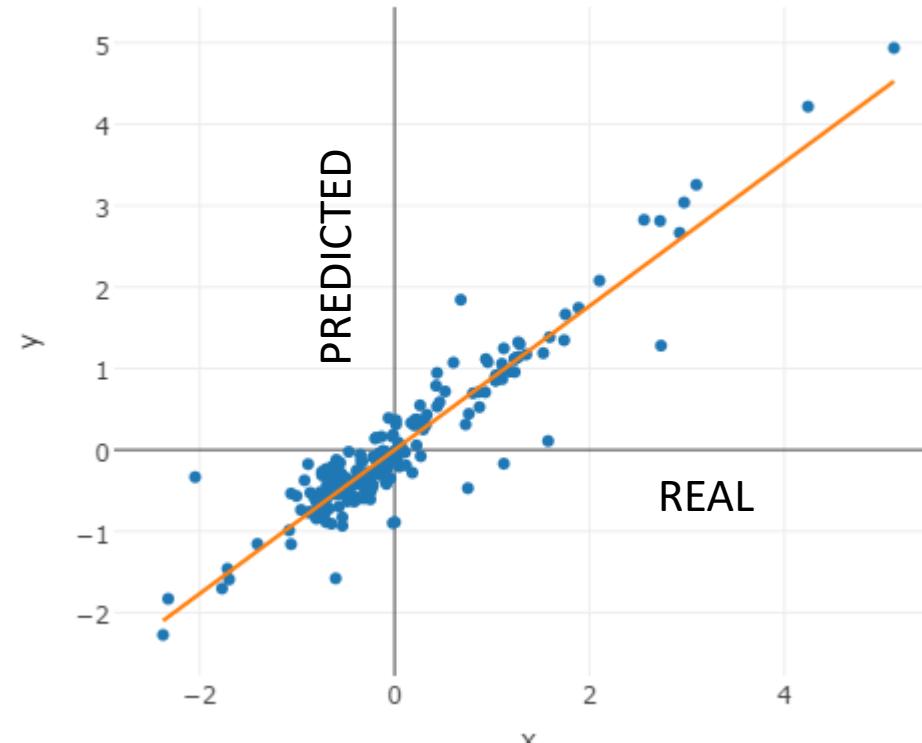
Variable Importance in the Projection



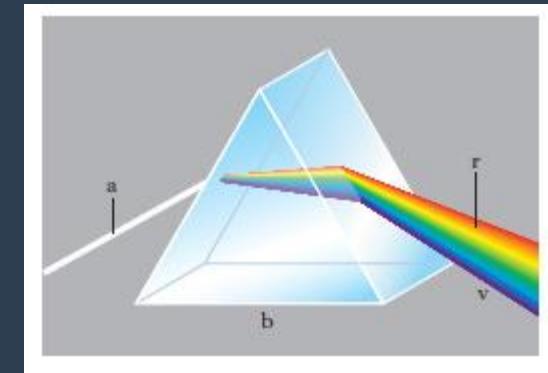
REGR. COEF (COMP=2)



REFRACTIVE INDEX



# STATISTICAL MODEL

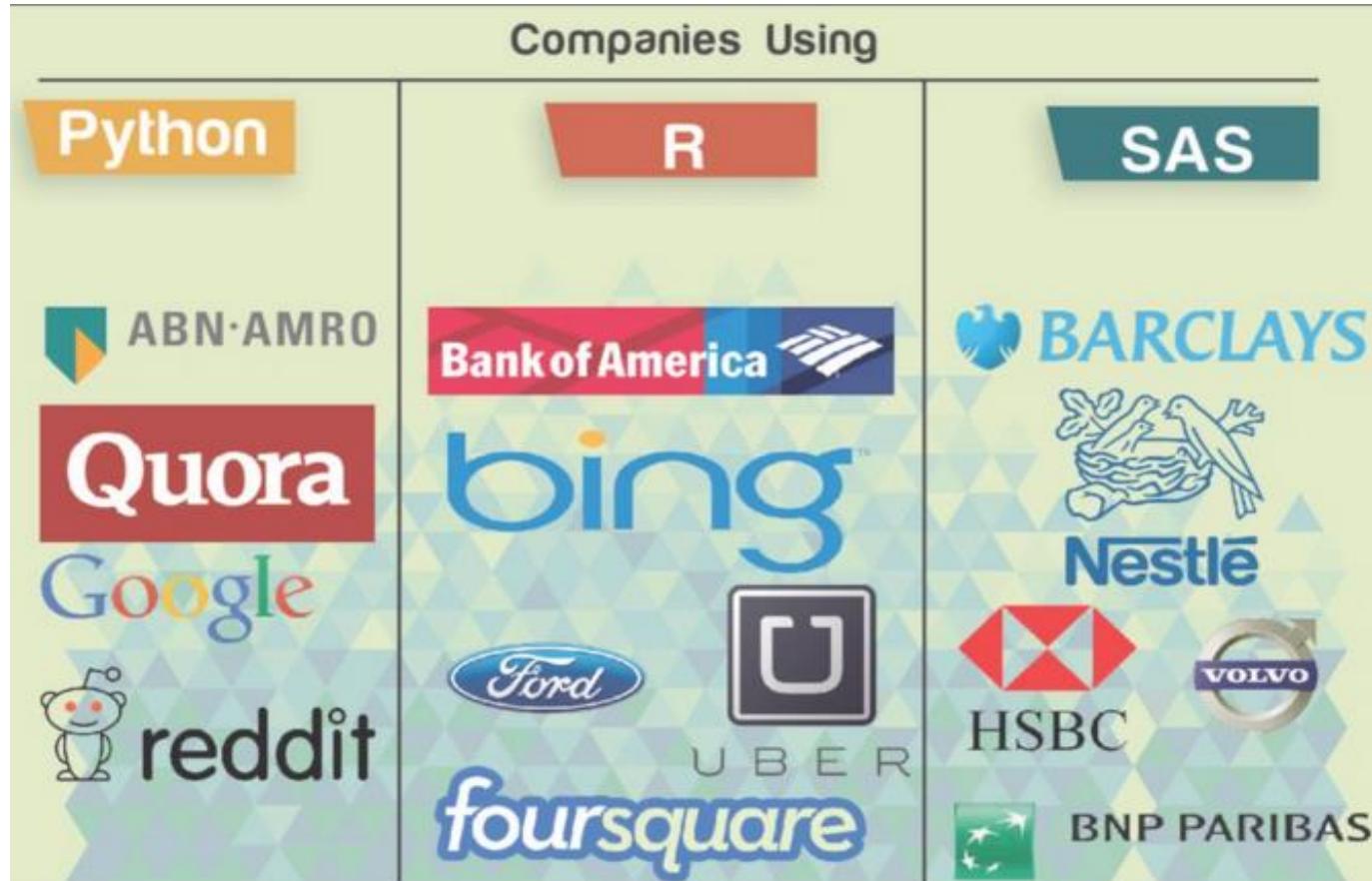


STATISTICS VS MACHINE LEARNING

PREDICT PERFORMANCE  
FROM PROCESS STATE  
VARIABLE

**Sotware open-source R, il più potente e diffuso  
software per l'analisi statistica.**

# DATA ANALISI: IL LIGUAGGIO



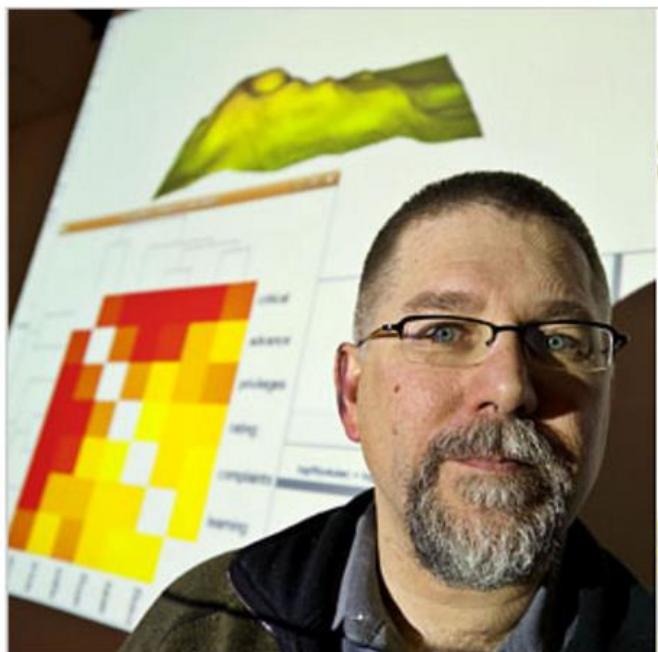
R

**The New York Times**

## Business Computing

[WORLD](#) [U.S.](#) [N.Y. / REGION](#) [BUSINESS](#) [TECHNOLOGY](#) [SCIENCE](#) [HEALTH](#) [SPORTS](#) [OPINION](#)

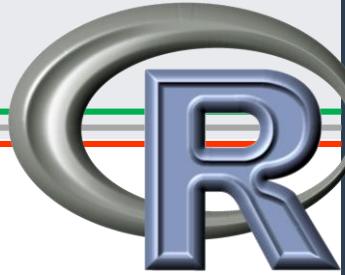
### Data Analysts Captivated by R's Power



Left, Stuart Isett for The New York Times; right, Kieran Scott for The New York Times

R first appeared in 1996, when the statistics professors Robert Gentleman, left, and Ross Ihaka released the code as a free software package.

R



# What is R?

- Most widely used data analysis software
  - Used by 2M+ data scientists, statisticians and analysts
- Most powerful statistical programming language
  - Flexible, extensible and comprehensive for productivity
- Thriving open-source community
  - Leading edge of analytics research

R

There are packages for nearly any analysis task in nearly any field. Music analysis, text analysis, image analysis, clinical trials, experimental design, astrostatistics, crop science, online marketing, bioinformatics, financial analysis are just some of the fields for which numerous packages exist. Packages have been written to facilitate other aspects of data analysis; for instance to download (scrape) data from the internet, or to download tweets for text or sentiment analysis.

# R & Friends

Because R packages are contributed by the community, traditional software vendors are having to adapt by embracing R instead of competing against it. R is now integrated throughout the enterprise analytics stack.



SQL Server 2016 will include the capability of running R directory “inside” the database

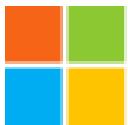


PowerBI adds support for R

R



Working with the R Foundation  
Supporting the R user community  
Continuing the growth of the R Project  
Linux Foundation collaborative project  
Non-profit trade organization



Microsoft



alteryx

Google



MANGO SOLUTIONS  
data analysis that delivers

ORACLE®

# CRAN: 7000+ add-on packages for R

Bayesian Inference  
Applied researchers interested in Bayesian statistics are increasingly attracted to R because of the ease of which one can code algorithms to sample... [\[more\]](#)

Chemometrics and Computational Physics  
Chemometrics and computational physics are concerned with the analysis of data arising in chemistry and physics experiments, as well as the simulation of... [\[more\]](#)

Clinical Trial Design, Monitoring, and Analysis  
This task view gathers information on specific R packages for design, monitoring and analysis of data from clinical trials. It focuses on including... [\[more\]](#)

Cluster Analysis & Finite Mixture Models  
This CRAN Task View contains a list of packages that can be used for finding groups in data and modelling observed cross-sectional heterogeneity. Many... [\[more\]](#)

Probability Distributions  
For most of the classical distributions, base R provides probability distribution functions (p), density functions (d), quantile functions (q), and... [\[more\]](#)

Computational Econometrics  
Base R ships with a lot of functionality useful for computational econometrics, in particular in the stats package. This functionality is complemented by many... [\[more\]](#)

Analysis of Ecological and Environmental Data  
This Task View contains information about using R to analyse ecological and environmental data... [\[more\]](#)

Design of Experiments (DoE) & Analysis of Experimental Data  
This task view collects information on R packages for experimental design and analysis of data from experiments. Please feel free to suggest enhancements... [\[more\]](#)

Empirical Finance  
This CRAN Task View contains a list of packages useful for empirical work in Finance, grouped by topic... [\[more\]](#)

Natural Language Processing  
This CRAN task-view contains a list of packages useful for natural language processing... [\[more\]](#)

Analysis of Pharmacokinetic Data  
The primary goal of pharmacokinetic (PK) analysis is to determine the relationship between the dosing regimen and the body's exposure to the drug... [\[more\]](#)

Optimization and Mathematical Programming  
This CRAN task-view contains a list of packages which offer facilities for solving optimization problems. Although every regression model in statistics... [\[more\]](#)

Phylogenetics, Especially Comparative Methods  
The history of life unfolds within a phylogenetic context. Comparative phylogenetic methods are statistical approaches for analyzing historical... [\[more\]](#)

Multivariate Statistics  
Base R contains most of the functionality for classical multivariate analysis, somewhere. There are a large number of packages on CRAN which extend this... [\[more\]](#)

Official Statistics & Survey Methodology  
Several add-on packages implement ideas and methods developed at the borderline between computer science and statistics - this field of research is usually... [\[more\]](#)

Machine Learning & Statistical Learning  
R is rich with facilities for creating and developing interesting graphics. Base R contains functionality for many plot types including cplots, mosaic... [\[more\]](#)

Graphic Displays & Dynamic Graphics & Graphic Devices & Visualization  
This CRAN task view contains a list of packages, grouped by topic, that are useful for high-performance computing (HPC) with R. In this context, we are... [\[more\]](#)

High-Performance and Parallel Computing with R

Analysis of Spatial Data  
Base R includes many functions that can be used for reading, visualizing, and analysing spatial data. The focus in this view is on "geographical" spatial... [\[more\]](#)

Time Series Analysis  
Base R ships with a lot of functionality useful for time series, in particular in the stats package. This is complemented by many packages on CRAN, which are... [\[more\]](#)

Robust Statistical Methods  
Robust (or "resistant") methods for statistical modelling have been available in S from the start, in R in package stats (e.g., median(), mean(), trim = ...). [\[more\]](#)

Survival Analysis  
Survival analysis, also called event history analysis in social science or reliability analysis in engineering, deals with time until occurrence of an... [\[more\]](#)

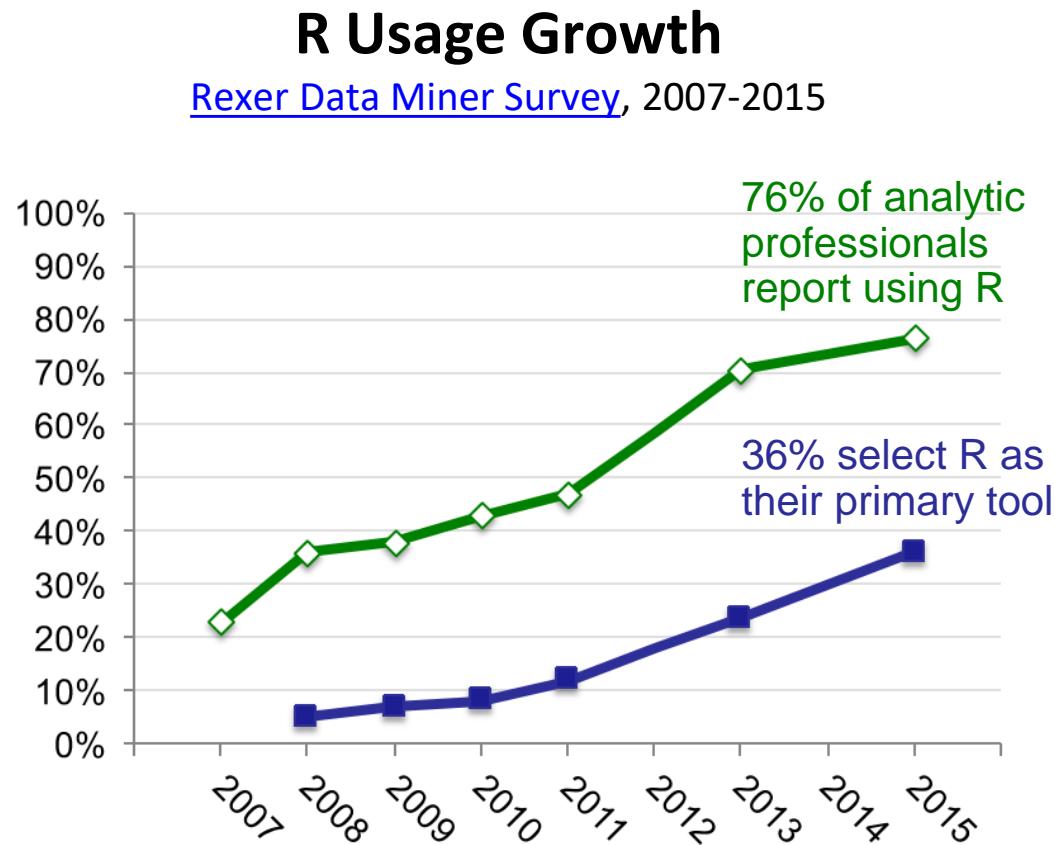
Statistics for the Social Sciences  
Social scientists use a wide range of statistical methods. To make the burden carried by this task view lighter, I have suppressed detail in some areas that... [\[more\]](#)

Graphical Models in R  
Wikipedia defines a graphical model as a graph that represents independencies among random variables by a graph in which each node is a random variable, and... [\[more\]](#)

Reproducible Research  
The goal of reproducible research is to tie specific instructions to data analysis and experimental data so that scholarship can be re-created, better. [\[more\]](#)

R

# R: The #1 software for Data Science



R

# **Case studies**

# RAPID PROTOTYPING

16 | Impresa & territori

La questione industriale. Carenza di materie prime e pezzi alle stelle mettono in ginocchio il settore della trasformazione

## Tensione nella filiera della plastica

Allarme per il fermo impianti e la cancellazione di consegne già confermate

Critico Casadel  
una carenza di materie prime plastiche e pezzi che esistono alle stelle fino a strapparle stanno mettendo in ginocchio i trasformatori di materie prime plastiche che lavorano esclusivamente per il settore del packaging e per quello dell'auto. Proprio quando le imprese sono impegnate nella rincorsa del rilancio della produzione in tutta Europa, come spiega Alexandre Drangis, managing director della European plastics converter, Italia compresa, come ribadiscono da Uniplast, l'associazione italiana dei trasformatori di materie plastiche che fa parte della Federazione della gomma plastica. Germania pure, come dicono dalla Industrievertretung Kunststoffverpackungen che denuncia ormai frequente cancellazione di consegne già confermate come anche la sorpresa, al momento della consegna, di un prezzo più alto di quello concordato.

Secondo i produttori tedeschi, in particolare, le cause di questa situazione che interessa soprattutto la componentistica auto e il packaging vanno ricordate all'uso della Forza maggiore, ossia al fermo di alcuni impianti che producono materie prime plastiche per ragioni di diverso tipo. Possono coincidere con il verificarsi di un evento eccezionale, al di fuori del controllo del contraente, che determina il decadimento dell'obbligazione a fornire una prestazione contrattualmente prevista.

Dalla Federazione della gomma plastica hanno calcolato che da gennaio 2015 in avanti, in Europa d'Occidente, è interessata dal fenomeno il ricorso allo stato di forza mag-

giore è stato registrato ben 100 volte e ha coinvolto almeno 5 impianti di polipropilene, 2 di polietilene a basso densità e 2 di polietilene ad alta densità. Questo ha portato ad un aumento straordinario dei prezzi di polietilene e polipropilene, con importanti contaminazioni anche nei mercati di altre plastiche come il poliuretano. I trasformatori esprimono il sospetto che dietro le chiusure ci sia in realtà un intento speculativo.

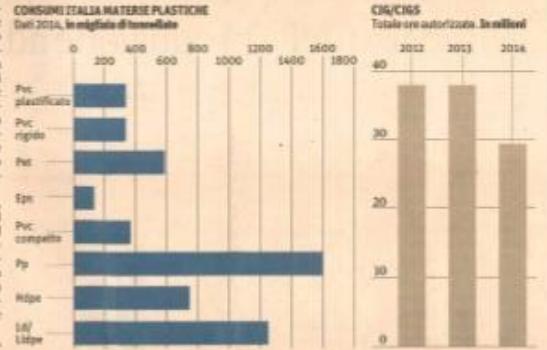
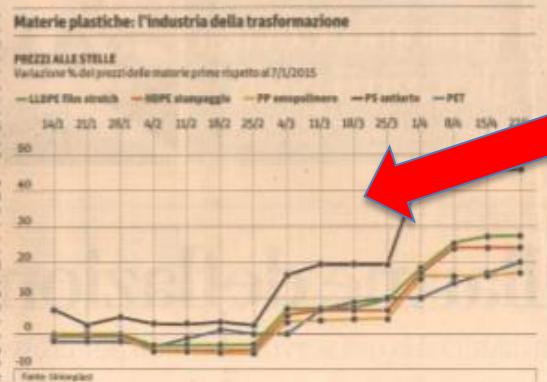
Dalla Federazione gomma-plastica spiegano che «la

DENUNCIA DI UNIPLAST.  
«La scarsità di informazioni rende difficile accettare se sostituirsi i criteri per le dichiarazioni di forza maggiore»

LA PAROLA CHIAVE  
Forza maggiore

«La "Forza Maggiore" non la molta spudore nell'ordinamento giuridico italiano. Può considerarsi come situazione appartenente a questa categoria che l'eventuale impedita la regolare esecuzione di contratti e, nella efficacia qualiasi atto dell'obbligo diretto all'elminazione. Secondo avvistamenti giurisprudenziali è caratterizzata da due elementi: strutturale e imprevedibilità.

L'inevitabilità, invece, ha natura soggettiva.



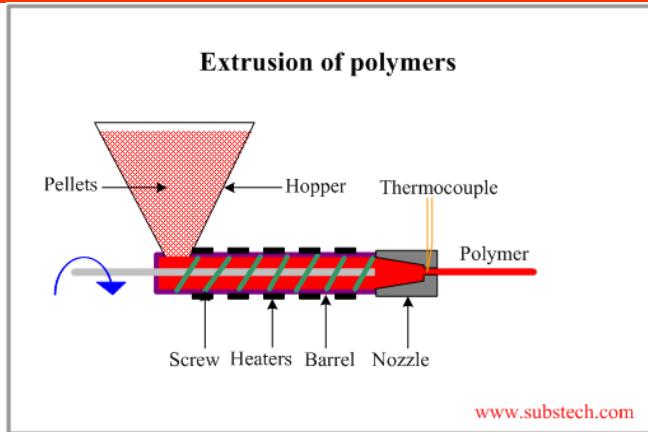
Carenze Materie prime  
=  
Prezzi alle stelle

## CASE STUDIES

# Plastic covering of mobile

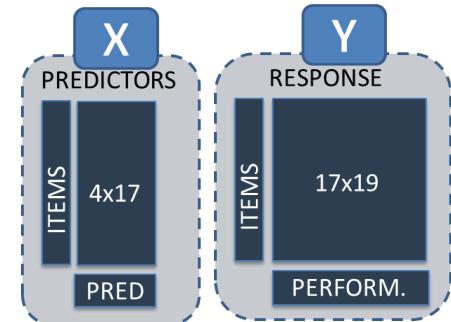
-Industria: **Plastica**

-Processo: **Estrusione+Stampaggio**



-Proprietà di interesse: **Deformazione (-) e resistenza (+)**

-Esigenze: **Ottimizzazione PROPRIETA' MULTIPLA**



-Dati: **17 formulazioni (DOE)**

4 costituenti:

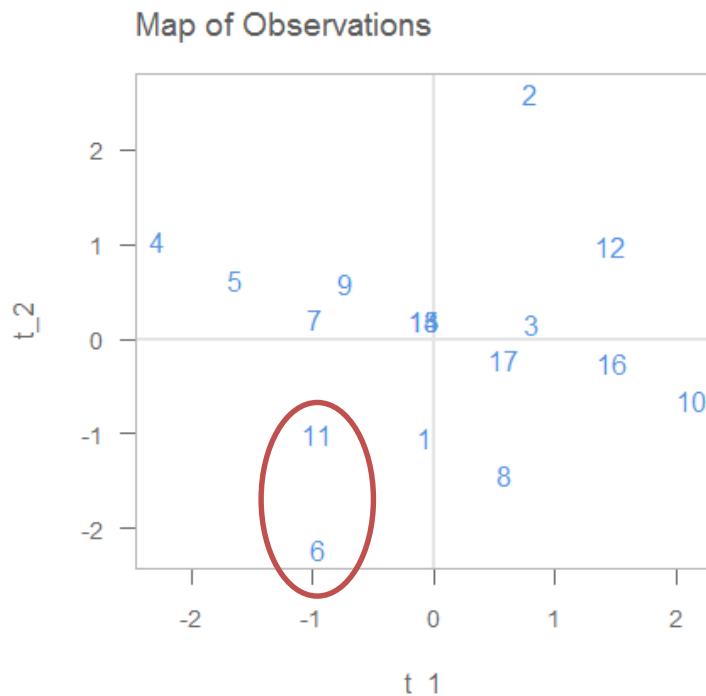
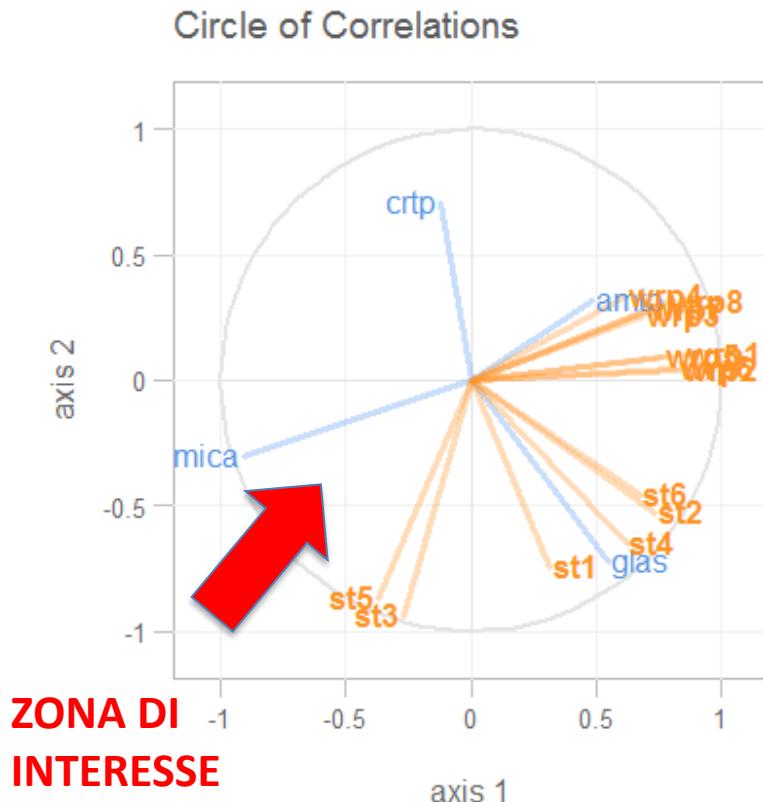
vetro/mica/crtp/amp

19 proprietà misurate

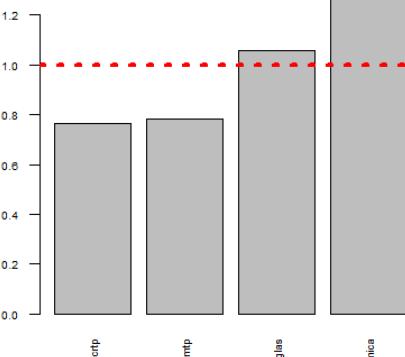


CASE STUDIES

# Plastic covering of mobile



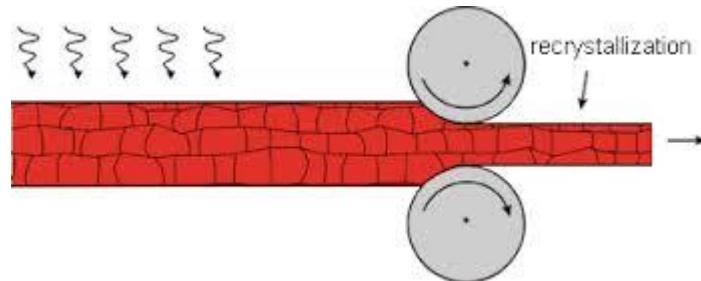
*Variable Importance in the Projection*



CASE STUDIES

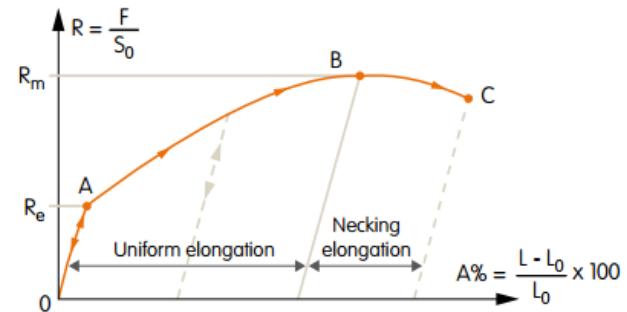
# HOT ROLLING

-Industria: **Siderurgica**



-Processo: **Hot rolling steel slabs**

-Proprietà di interesse: **Proprietà meccaniche**

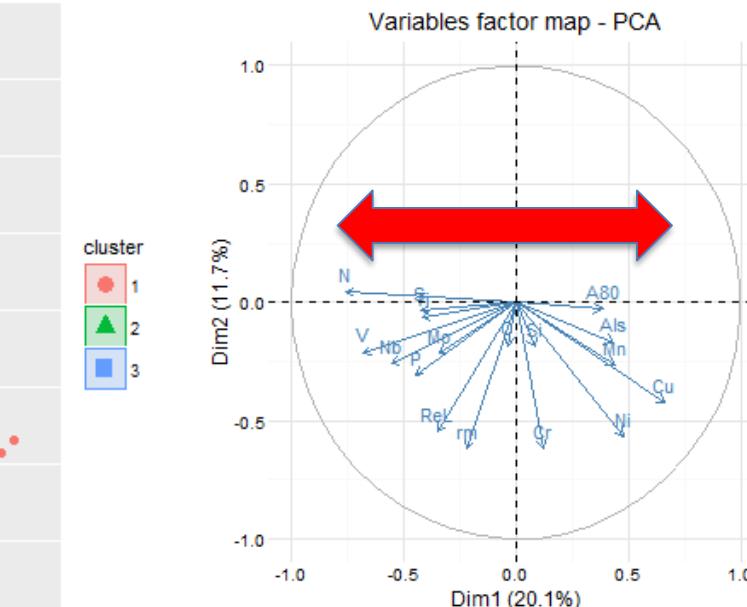
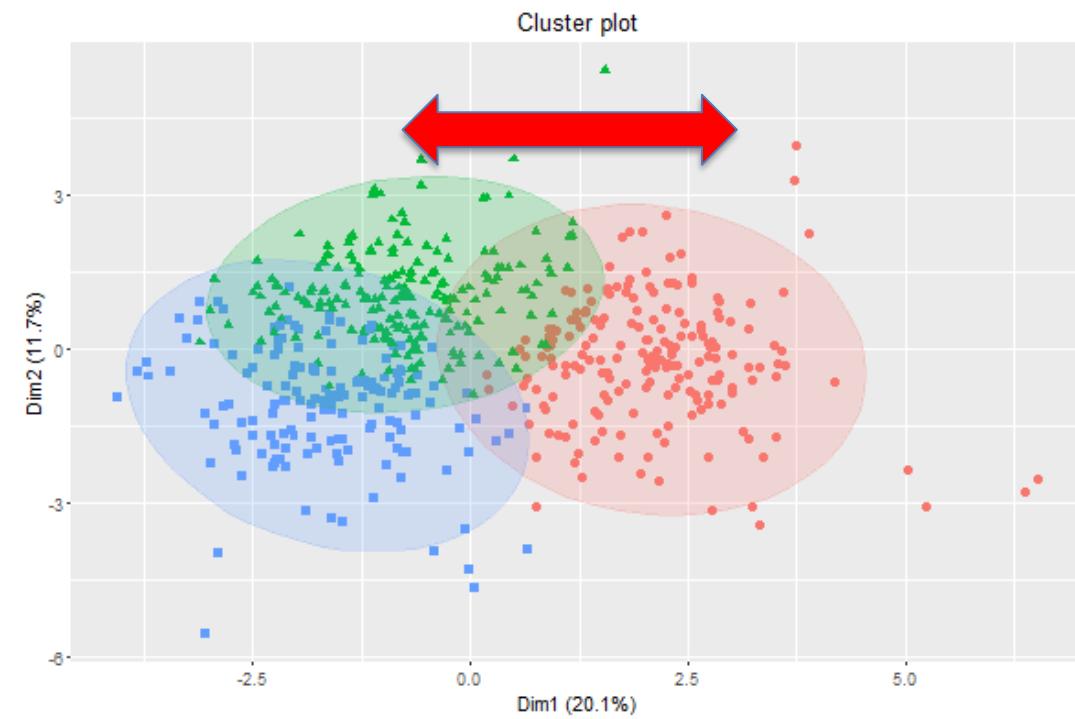
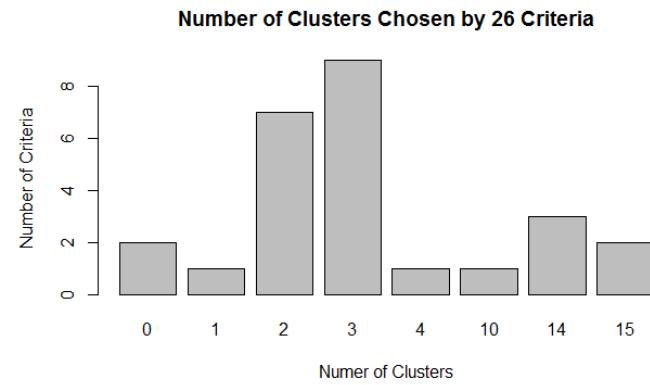
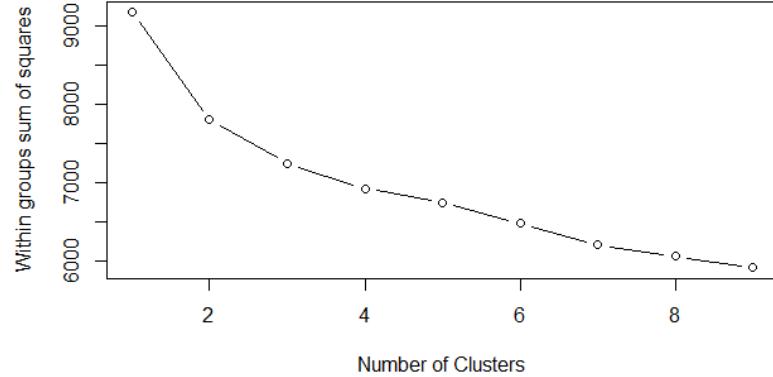


-Esigenze: Acciai dello stesso grado **SAE (SAE steel grades)** proprietà meccaniche diverse

-Dati: **Serie storiche** di processo

CASE STUDIES

# HOT ROLLING

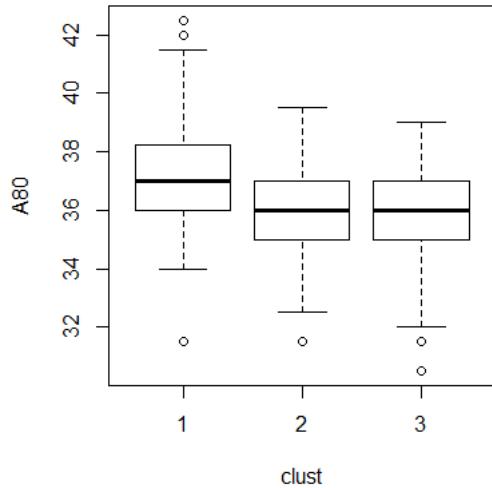


CASE STUDIES

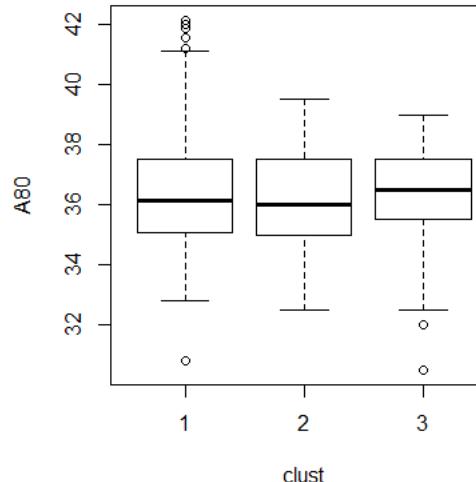
# HOT ROLLING

## NEW TEMPERATURE SETUP

PRIMA



DOPO



SCRIPT PER IL  
CALCOLO DELLA  
CLASSE

```

Sub Send_Range()
    ' In Green, must start with ":".
    ' Select the range of cells on the active worksheet.
    ActiveSheet.Range("A1:B5").Select
    ' Show the envelope on the ActiveWorkbook.
    ActiveWorkbook.EnvelopeVisible = True
    ' Set the optional introduction field that adds
    ' some header text to the email body. It also sets
    ' the To and Subject lines.
    ' www.sanuja.com
    With ActiveSheet.MailEnvelope
        .Introduction = "This is a sample work"
        .Item.To = "E-Mail_Address_Here"
        .Item.Subject = "My subject"
        .Item.Send
    End With
End Sub
End Sub

```

CICLO TERMICO  
PER CLASSE

CASE STUDIES

# LDPE POLYMERIZATION

## MULTIVARIATE STATISTICAL PROCESS CONTROLL

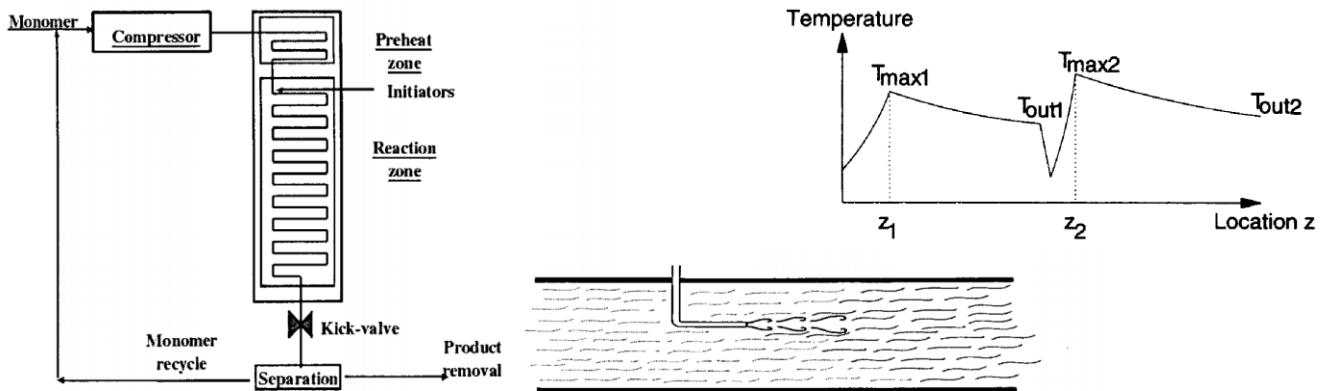
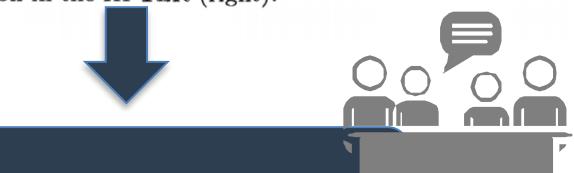


Figure 1.2 Schematic of a high-pressure tubular LDPE reactor (HPTLR) (left) and initiator injection in the HPTLR (right).



Quality Control

## PESO DEL POLIMERO

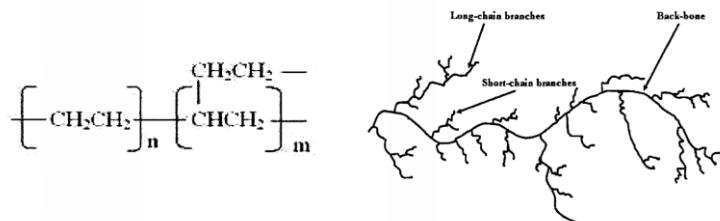
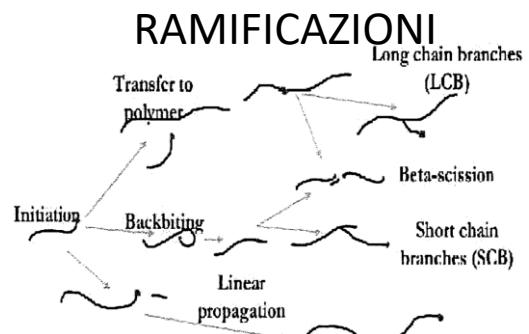


Figure 1.1 Schematic of LDPE chain structure - monomer unit (left); short- and long-chained branches (right).



## CASE STUDIES

Process Monitoring and  
Diagnosis by Multiblock PLS  
Methods  
*Mac Gregor et all,*  
*Process Systems Engineering, 1994.*

# ANALISI DATI: QUANDO SERVE



## RAPID PROTOTYPING

Per sviluppare un nuovo processo e/o un nuovo prodotto

Per supportare lo scale-up



## DATA DRIVEN INNOVATION

Per gestire processi intrinsecamente complessi



## MULTIVARIATE STATISTICAL PROCESS CONTROL

Analizzare in tempo reale i dati del tuo processo aiutando gli operatori a mantenere il processo sotto controllo.



## CLUSTERING

Problemi con materi prime



## VIRTUAL METROLOGY

Per prevedere le caratteristiche di un prodotto dai dati di processo



## PREDICTIVE MAINTENANCE

Per ridurre i costi delle riparazioni e aumentare l'efficienza



"After careful consideration of all 437 charts, graphs, and metrics,  
I've decided to throw up my hands, hit the liquor store,  
and get snockered. Who's with me?!"

"Dopo un attento esame di tutte le 437 tabelle, grafici e  
metriche, ho deciso di arrendermi, fiondarmi in un negozio di  
liquori e ubriacarmi. Chi sta con me ?!"

THANK YOU



# Dr. Alessio Passalacqua

*Data Analyst & DOE Expert*



PHONE

349 6707508



MAIL

[info@produzioneperfetta.it](mailto:info@produzioneperfetta.it)



ADDRESS

Via Chinnici 30, 41125 Modena



WEB

[www.produzioneperfetta.it](http://www.produzioneperfetta.it)