# BIG DIVE

## TECH. CUSTOM EDITION

**A project by TOP-IX**
**designed for Intesa Sanpaolo**

**#DataScience**

aizoon®
AUSTRALIA
EUROPE
USA
TECHNOLOGY CONSULTING

# aizoOn Group Profile

**aizoOn is an independent technology consulting firm.**
**Focused on innovation.**
**aizoOn has achieved record growth with offices across Australia, Europe and USA.**



USA
- New York, NY
- Troy, MI
- Lewiston, ME
- Cambridge, MA

EUROPE
- Bologna, IT
- Cuneo, IT
- Genova, IT
- Milano, IT
- Roma, IT
- Torino, IT
- Sheffield, UK

AUSTRALIA — Sydney

**Gartner.**

G00300966

**Cool Vendors in Managing Operational Technology in a Digital Business, 2016**

Published: 25 April 2016

aizoOn

AUSTRALIA
EUROPE
USA

TECHNOLOGY CONSULTING
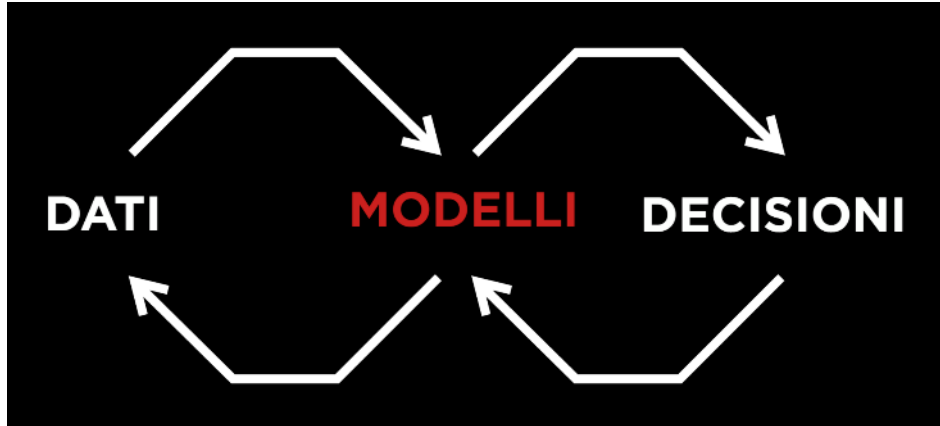
# aizoOn Organisation: Markets & Technology

**TeMI**
AUTOMATION, DURABLE GOODS, ENTERTAINMENT, MEDIA, TELCO, TRAVEL

**E2&R**
CHEMICAL, ENERGY, ENVIRONMENT, MINING, OIL&GAS, UTILITIES

**DASS**
AERONAUTICS, AEROSPACE, DEFENCE, NAVAL

**Pa_S**
CENTRAL PA, LOCAL PA, HEALTH

**DDI**
DATA DRIVEN INNOVATION

**CYSE**
CYBER SECURITY

**FIN**
INSURANCE, INVESTIMENT BANKING, REAL ESTATE, RETAIL BANKING

**F&F**
FASHION, FOOD

**LAIF**
AUTOMOTIVE, INFRASTRUCTURE, LOGISTIC, METRO, RAIL

**Scientiae**
BIOTECHNOLOGY, HEALTH CARE, LIFE SCIENCES, PHARMACEUTICAL

**SW&IT**
SOFTWARE DEVELOPMENT & ARCHITECTURES

**ITS**
INTELLIGENT THINGS & SYSTEMS

# About me

- PhD Physics of Complex Systems → useful for critical thinking

- Master in Epidemiology → useful for Stat

- Master in Data Engineer → useful for Tech (and CV)

Moved from Reasearch to Business



aizon® TECHNOLOGY CONSULTING
AUSTRALIA
EUROPE
USA

# Statistics & Data Analysis

aizOOn ®

AUSTRALIA
EUROPE
USA

TECHNOLOGY CONSULTING

# Models



- Mathematical model
- Dynamical model
- Statistical model
- Machine learning model
- Agent-based model
- Data model
- ….

# Data ..*..

"**<u>Data analysis</u>** has been generally used as a way of *explaining some phenomenon* by *extracting interesting patterns* from individual data sets with well-formulated queries.

**<u>Data science</u>**, on the other hand, aims to discover and extract *actionable knowledge* from the data, that is, knowledge that can be used to make *decisions and predictions*, not just to explain what's going on."

– I. Wladawsky-Berger

# Actionable Knowledge

1. Scoprire patterns nascosti nei dati

**The New York Times**

***What Wal-Mart Knows About Customers' Habits***

2. Automatizzare decisioni su larga scala

# Data-analytic thinking

" It is important to understand data science even if you never intend to apply it to yourself. **Data-analytic thinking enables you to evaluate proposal for data mining projects**.

[…] If a potential business application extracting knowledge from data is proposed, you should be able to assess the proposal systematically and decide whether it is sound or flawed. "

- Foster Provost

# Data Science 'complementary skills'

- Capacità di scomporre un problema di data-analytics in parti per cui esistono strumenti e tecniche note da applicare

  → non reinventare la ruota!

- Intuizione e immaginazione  per sviluppare test e analisi in grado di validare o smentire le ipotesi generate

  → Non sempre è possibile fare test in laboratorio!

- Capacità di comunicazione dei risultati e ascolto dei requisiti di business

  → Actionable Knowledge!

AIZOON
TECHNOLOGY CONSULTING
AUSTRALIA
EUROPE
USA

# Data Quality



**Garbage in, garbage out** (**GIGO**) refers to the fact that computers, since they operate by logical processes, will unquestioningly process unintended, even nonsensical, input data ("garbage in") and produce undesired, often nonsensical, output ("garbage out").

- Wikipedia

# Data Quality

Esistono molteplici definizioni e dimensioni per valutare la qualità dei dati.
Ad esempio:

# Data Quality

Data Quality

⬇

Fitness-for-use
(Wang & Strong,1996)

<u>Suggerimento</u>: verificare i dati per integrità, orizzonte temporale e spaziale, modalità di raccolta e pattern di dati mancanti

Data Exploration

Model Testing

Prediction

# Statistica descrittiva

## Dati



## Indicatori statistici

| Full-Time Statistics by Geographical Region | Average | Standard Deviation |
|---|---|---|
| Texas | $60,632 | $13,375 |
| Austin, Texas* | $52,500 | $11,529 |
| Dallas, Texas* | $60,635 | $10,041 |
| Houston, Texas* | $66,996 | $13,646 |
| Other Texas* | $49,360 | $16,406 |
| Mid-Atlantic | $54,200 | $22,451 |
| Midwest | $60,770 | $8,920 |
| Northeast | $68,069 | $15,595 |
| South | $55,220 | $11,670 |
| Southwest | $64,600 | $12,178 |
| West | $62,341 | $16,173 |
| Other | $63,588 | $5,074 |

## Visualizzare i dati

AIZOON
TECHNOLOGY CONSULTING
AUSTRALIA
EUROPE
USA

# Outliers & heavy tail distribution

- Distribuzioni con una lunga coda: **quasi tutti sono sotto la media.**

- Attenzione agli **outliers**: se Bill Gates entra in un bar, in **media** tutti gli avventori sono milionari…

  1 1 1 **2** 2 2 1000

- ….ma la **mediana** resta invariata!

AUSTRALIA
EUROPE
USA
aizoon®
TECHNOLOGY CONSULTING

# Quantili

Quantili: punti presi a intervalli regolari sulla distribuzione cumulativa

- Quartili
- Percentili
- …

# Boxplot



- Show the variation for each data bin

- More informative than averages or medians

- Useful to summarize measurements or simulated results when fluctuations are important

# Statistica descrittiva



NOTA: tutte **le variabili** $x_i$ hanno la **stessa media e stessa varianza**.
Lo stesso vale per tutte le variabili $y_i$.

# Coefficiente di Pearson

$$\rho = \frac{\sum_{i=1}^{N}(x_i - \mu_X)(y_i - \mu_Y)}{\sigma_X \sigma_Y}$$

Assunzioni:

- Correlazione lineare
- Variabili continue
- Variabili con distribuzione normale (no outliers)
- Le due variabili formano una distribuzione normale bivariata
- Omoschedasticità dei dati

credits: Wikimedia Commons

# Coefficiente di Spearman

$$\rho = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2 \sum_i (y_i - \bar{y})^2}}$$

$$\rho = 1 - \frac{6 \sum_{i=1}^{N} d_i^2}{N(N^2 - 1)}$$

Spearman correlation=1
Pearson correlation=0.88

credits: Wikimedia Commons

Assunzioni:

- Correlazione monotona
- Variabili continue o ordinali

→ Test non parametrico

# Correlazione



Tutti le variabili mostrate in figura hanno lo **stesso valore di correlazione lineare**.

(P: Pearson, S: Spearman)

# Correlazione vs. Causalità

Immaginiamo di osservare che due grandezze A e B siano correlate

A ⟶ B

'Correlation is not causation but it sure is a hint.'
– E. Tufte

B ⟶ A

CAUSALITA' INVERSA: Es. più pompieri vengono inviati e maggiori sono i danni registrati a causa dell'incendio.

C
A ↙ ↘ B
.....

A E B SONO CAUSATE DA UN FATTORE COMUNE C: Es. all'aumento del consumo di gelato, aumentano le morti per annegamento ("estate" è il fattore C).

# Confondimento e variabili proxy

C

A B

- Nello stimare la forza di associazione tra A e B, la variabile C '**confonde**' la stima. → **BAD**

    Soluzioni possibili:

    1) randomizzazione

    2) misurare C

- La variabile A (misurabile) è un **proxy** della variabile C che è associata a B. → **GOOD**

# Causality

**Hill's criteria for causation**, are a group of guidelines that can be useful for providing evidence of a causal relationship between a putative cause and an effect:

- **Strength**: A small association does not mean that there is not a causal effect, though the larger the association, the more likely that it is causal.
- **Consistency**: Consistent findings observed by different persons in different places with different samples strengthens the likelihood of an effect.
- **Temporality**: The effect has to occur after the cause (and if there is an expected delay between the cause and expected effect, then the effect must occur after that delay).
- **Plausibility**: A plausible mechanism between cause and effect is helpful
- **Analogy**: The effect of similar factors may be considered.

# Causality

**Hill's criteria for causation**, are a group of guidelines that can be useful for providing evidence of a causal relationship between a putative cause and an effect:

- **Strength + Consistency + Temporality + Plausibility + Analogy …**

- **Specificity**: Causation is likely if there is a very specific population at a specific site and disease with no other likely explanation. The more specific an association between a factor and an effect is, the bigger the probability of a causal relationship.
- **Coherence**: Coherence between inferred and laboratory findings increases the likelihood of an effect.
- **Experiment**: "Occasionally it is possible to appeal to experimental evidence".

# Feature Engeneering – Misurare l'effetto

- Valore medio            → effetti lenti e reversibili

- Valore cumulativo        → effetti cumulativi e irreversibili

- Durata dell'esposizione    → effetti cumulativi e irreversibili

- Valore di picco           → effetto reversibile e acuto

# E.g. Average Time Spent Composing 1 email

Un pregiudizio (conscio od inconscio) o un errore di disegno delle analisi possono portare a misure distorte.

# Diverse tipologie di bias

**Es.** Sbagliare la misura d'effetto può introdurre misclassificazione e solitamente causa un 'bias verso il valore nullo'.

A Catalog of Biases in Questionnaires

Bernard C.K. Choi, PhD, Anita W.P. Pak, PhD

**Es.** *Belief vs behavior:*
*1. Do you think that it is a good idea to have everyone's chest regularly checked by X-ray?*
*2. Have you ever had yours checked?*
**Es.** *Framing*
*Which operation would you prefer?*
    *[ ] An operation that has a 5% mortality.*
    *[ ] An operation in which 90% of the patients will survive.*

# Simpson's Paradox

*'il **paradosso di Simpson** indica una situazione in cui una relazione tra due fenomeni appare modificata, o perfino invertita, dai dati in possesso a causa di altri fenomeni non presi in considerazione nell'analisi (variabili nascoste).'*

- wikipedia

vudlab.com/simpsons/

AIZOON ®
TECHNOLOGY CONSULTING

AUSTRALIA
EUROPE
USA

# Simpson's Paradox

| departments | # applied | | # admitted | | % admitted | |
|---|---|---|---|---|---|---|
| | men | women | men | women | men | women |
| "Easy" | 1,372 | 294 | 856 | 234 | 62% 38% | 20% 80% |
| "Hard" | 1,319 | 1,541 | 338 | 408 | 26% 74% | 27% 74% |
| Combined | 2,691 | 1,835 | 1,194 | 642 | 44% 56% | 35% 65% |

**Simpson's paradox? yes**

# Simpson's Paradox

| departments | # applied | | # admitted | | % admitted | |
|---|---|---|---|---|---|---|
| | men | women | men | women | men | women |
| "Easy" | 1,372 | 294 | 856 | 234 | 62% / 38% | 20% / 80% |
| "Hard" | 1,319 | 1,541 | 338 | 408 | 26% / 74% | 27% / 74% |
| Combined | 2,691 | 1,835 | 1,194 | 642 | 44% / 56% | 35% / 65% |

**Simpson's paradox? yes**

$$y = \alpha + \beta X + \epsilon$$

$$y = \alpha + \beta_1 X_1 + \beta_2 X_2 + ... + \epsilon$$

# Inference of association: regression – basic concepts

$$y = \alpha + \beta_1 X_1 + \beta_2 X_2 + ... + \epsilon$$

- Assunzioni da verificare:
  - Legame lineare
  - Variabili con distribuzione normale (no outliers)
  - Le due variabili formano una distribuzione normale bivariata
  - Omoschedasticità dei dati
  - Assenza di autocorrelazione
  - Tipologia variabili
  - multicollinearità

$$y = \alpha + \beta_1 X_1 + \beta_2 X_2 + \ldots + \epsilon$$

- Valutare se il modello ha una 'buona performance' (adjusted R-squared, RMSE, F-test, .....)

- $\beta_i$ rappresenta la **forza di associazione** tra la variabile target y e la variabile indipendete $X_i$ considerando i valori delle altre variabili.
  - Verificare 1) la forza di associazione della variabile e 2) se l'associazione è statisticamente significativa

- Valutare se esistono modelli migliori → Akaike Information Criterion (AIC)

aiz**oo**n® 
TECHNOLOGY CONSULTING

AUSTRALIA
EUROPE
USA

No Association

RESPONSE

Represents the "cloud" of data around the regression line

X

AUSTRALIA
EUROPE
USA

aizoon®
TECHNOLOGY CONSULTING

Positive Association

RESPONSE

X

**RESPONSE**

**X**

No effect modification

effect modification

RESPONSE

X

# Confounding

Confounding

Testare se le differenze osservate sono 'statisticamente significative'

# Confronto fra gruppi – A/B Testing



Some see this version...

...others see this version.

Only the headlines are different.

$p_A$

$p_B$

- Gli utenti vengono divisi in due gruppi (A e B)

- Per ogni gruppo viene misurata una metrica di interesse (ad es. probabilità di conversione)

- Viene testata **l'ipotesi nulla** che le metriche misurate nei due grupppi (controllo e trattati) non siano diverse:

$$H_0 = p_A - p_B \leq 0$$

# Confronto fra gruppi – A/B Testing

- Se l'outcome di interesse è un **processo binomiale**
- L'**errore standard** è:

$$SE = \sqrt{\frac{p\,(1-p)}{N}}$$

- Con **Z-score** (i.e. numero di deviazioni standard dalla media):

$$Z = \frac{p_A - p_B}{\sqrt{SE_A^2 + SE_B^2}}$$



Normal, Bell-shaped Curve

- Fissiamo un **livello di confidenza** (es. 95%)

se Z-score > 1.65, allora possiamo rigettare l'ipotesi nulla

- Calculate the probability of an event more extreme that the observation under the "null hypothesis"



- $p < 0.05$ Moderate evidence agains null-hypothesis

- $p < 0.01$ Strong evidence against null-hypothesis

- $p < 0.001$ Very strong evidence against the null-hypothesis

- The smaller the p-value the better.

## http://abtestguide.com/calc

# P-Value (misuse)

# Kolmogorov-Smirnov test

KS è un test **non parametrico** per verificare se due gruppi differiscono significativamente.



KS–Test Comparison Cumulative Fraction Plot

ATTENZIONE: se vengono effettuati controlli multipli, la probabilità di trovare differenze dovute al caso aumenta!

Un esempio 'big data':

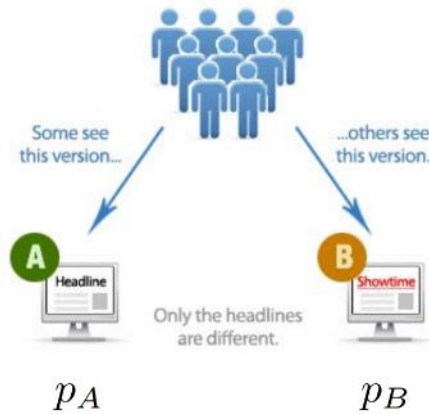# Computational approach for testing hypothesis

Keywords: bootstrap, resampling

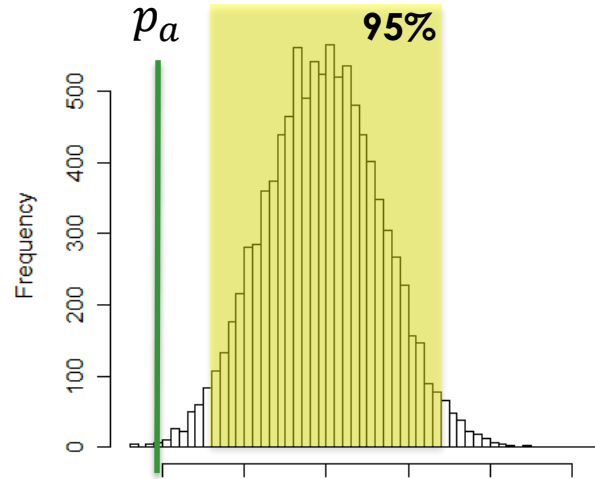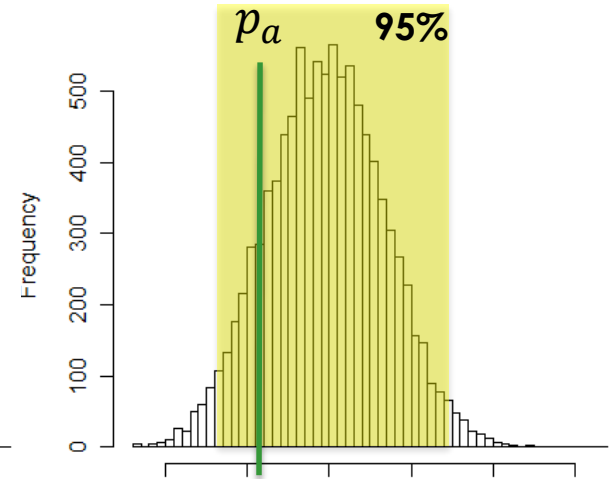# Computational approach for testing hypothesis

Si verifica se il valore misurato di $p_a$ è compatibile con valori di $p_a$ ottenuti da una **distribuzione casuale** di utenti con outcome positivo/negativo.

L'effetto di A vs. B impatta **significativamente** l'outcome

Effetto **non significativo:** esistono molti casi in cui il valore $p_a$ può essere effetto del caso

AIZOON
AUSTRALIA
EUROPE
USA
TECHNOLOGY CONSULTING

# Q & A