



# La Statistica, il ritorno.

## Dalle schede perforate ai **big data**

roberto leombruni  
università di torino e laboratorio revelli

**Analisi dei dati: cosa serve alle imprese, cosa serve alla pubblica amministrazione**

*Machine Learning, Intelligenza Artificiale, Mining, Statistica, Modelli di simulazione: cosa c'è di vecchio nel nuovo, cosa c'è di nuovo nel vecchio*

Torino  
Campus Luigi Einaudi  
10 gennaio 2019

## indice della presentazione

- ✓ i dati, dai tempi delle schede perforate ai big data...
- ✓ ...e quando non ci sono dati!
- ✓ dati sì, ma di qualità: vecchi e nuovi *epic fails*
- ✓ cosa c'è di vecchio nel nuovo: dati amministrativi e "found data"
- ✓ understanding vs forecasting



## DATO, in italiano

i dati servono per  
prendere decisioni

- ...
- (2) Ciascuno degli elementi di cui si dispone per formulare un giudizio o per risolvere un problema.
  - "raccogliere tutti i d."
  - *Dato di fatto*, elemento certo.

# DATO, in italiano, in statistica

- ...
- (2) Ciascuno degli elementi di cui si ha bisogno per un giudizio o per risolvere un problema
  - "raccogliere tutti i d."
  - *Dato di fatto*, elemento certo.
  - *Dato statistico*, la **misura** di un fenomeno collettivo risultante dalla **rilevazione** di fatti singoli della stessa specie.
  - *Dati sensibili*, vedi sensibile.

i dati sono il risultato di un  
**processo**

# DATO, in italiano, in statistica, in informatica

- ...
- (2) Ciascuno degli elementi di cui si dispone per formulare un giudizio o per risolvere un problema.
  - "raccogliere tutti i d."
  - *Dato di fatto*, elemento certo.
  - *Dato statistico*, la **misura** di un fenomeno collettivo risultante dalla **rilevazione** di fatti singoli della stessa specie.
  - *Dati sensibili*, vedi sensibile.
- (3) **Qualunque** informazione rappresentata **in modo da poter essere trattata da un calcolatore**



1492

stolisco, e tèrra per \*tersa terra, ossia la parte disseccata (cfr. Terra e Torrido). Confr. lo zendo tasta tazza e guscio, che altri però riferisce alla rad. sser. TAK-Š-grossare, comporre, fare (v. Tecnico). Vaso di terra cotta dove si coltivano le piante; Stoviglia di terra cotta, rotonda e alquanto cupa, colla quale si copre la pentola; Sorta di stoviglia di terra cotta, piatta, per uso di cuocervi sopra alcuna cosa.

Cfr. Teschio; Testa; Testaccio; Testaceo; Testicolo; Testuggine; Stoviglia.

2. fr. texte: = lat. TEXTUM da TÈXERE tessere: propr. tessuto, quindi intreccio, e traslat. discorso continuato (v. Tessere).

Ciò che è contenuto parola per parola in uno scritto; Scritto di un autore, considerato in rapporto ai commenti e alle note, che vi sono state fatte sopra.

« Testo d'Aldo » Nome di un carattere presso gli stampatori.

Deriv. Testino nome di un piccolo carattere nelle Stamperie; Testuale. Cfr. Contesto; Pre-testo.

32%



Esempi di dati: cos'hanno in comune?

i dati sono il risultato  
di un **processo**...



Esempi di dati: cos'hanno in comune?

i dati sono il risultato  
di un **processo**...



...che produce  
**numeri**

Esempi di dati: cos'hanno in comune?



# Quindi un dato è un numero, ma...

Non tutti i numeri sono dati: quale di questi due è un dato?!

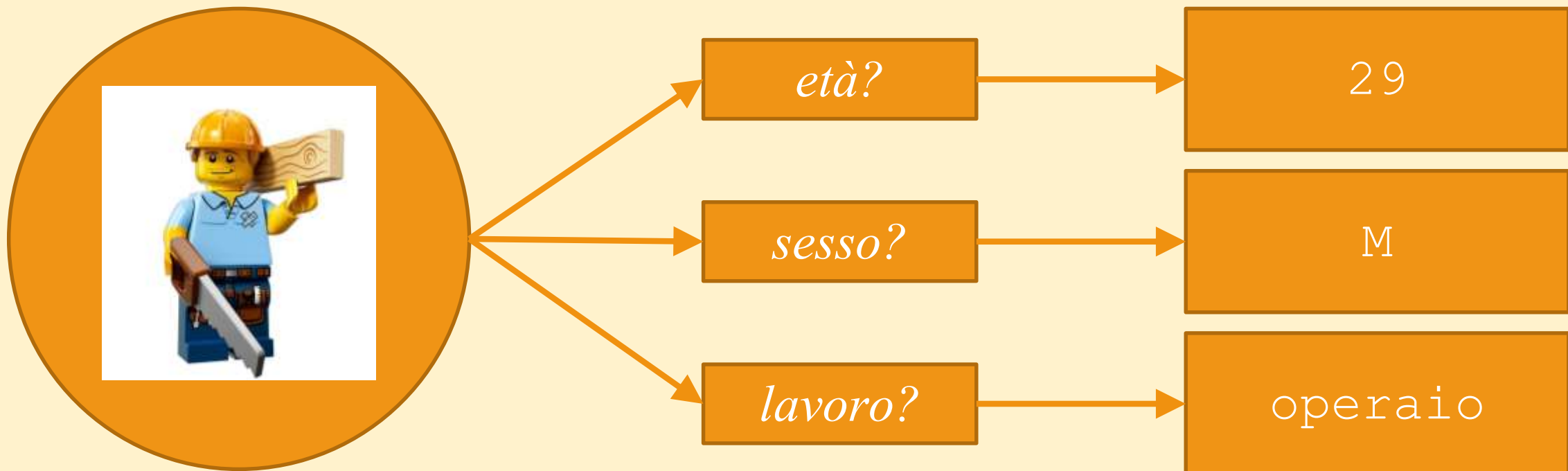
0	25	50	44	46	2D	31	2E	33	0A	25	C4	E5	E5	EB	A7	F3	A0	D0	C4	C6	0A	34	
23	20	30	20	6F	62	6A	0A	3C	3C	20	2F	4C	65	6E	67	74	68	20	35	20	30	20	52
46	20	2F	46	69	6C	74	65	72	20	2F	46	6C	61	74	65	44	65	63	6F	64	65	20	3E
69	3E	0A	73	74	72	65	61	6D	0A	78	01	ED	90	D9	8E	64	C7	71	F7	EF	CF	53	1C
92	01	16	5C	03	88	C5	B3	2F	BE	32	37	09	14	68	89	F2	B4	A0	0B	CB	10	E8	D6
115	50	1A	69	9A	14	BB	47	94	A9	87	F0	AD	AF	EC	77	D1	8B	F8	59	FC	FB	47	66
138	44	9E	AA	53	DD	03	6C	D2	02	28	D8	C0	F7	4D	33	2A	97	C8	58	FE	B1	64	56
161	E9	8B	FA	67	F5	17	75	53	2F	6D	53	DF	D4	53	9F	FE	7A	65	7F	35	F5	28	3E
184	D1	FF	FF	DB	18	F1	8B	FA	B3	F8	FB	64	74	75	FF	E8	26	D6	BD	7D	51	6B	81
207	B7	3F	7E	71	7B	FD	E2	0F	AF	FF	F8	C9	AB	FA	F6	25	FB	B7	5D	7F	EC	9A	61
230	B1	ED	86	F1	38	0F	43	BB	D4	7D	7B	9C	C6	61	6E	EB	B9	6F	8E	DD	38	2E	D5
253	F5	4D	FD	F6	87	37	6D	FD	FE	E7	B0	FD	F6	7B	77	6D	7D	7D	C7	9C	BB	EB	FA
276	ED	7F	7E	F1	EA	93	D7	2F	BF	7C	F1	DE	E7	AF	3E	BF	7D	79	F3	E2	F5	ED	CB
299	EB	BC	B8	AD	FA	56	5B	AF	F5	D2	30	E3	A6	7E	F7	AA	5A	13	71	AD	BB	F9	38
322	B7	C3	B2	D4	F3	9A	B7	BD	62	93	AB	B6	6E	EB	AB	4F	EB	7F	A9	0F	EF	BE	
345	78	56	BF	05	03	F5	E1	23	FF	E3	A5	FF	F1	9B	67	55	FA	E8	B7	4E	79	FD	AC
368	B6	B1	75	FE	F7	F9	E7	FE	C9	A7	99	E2	23	FE	E4	1F	7C	92	FE	A8	0E	B7	B1
391	93	CF	FE	B9	8F	79	FE	CE	0F	CE	16	0E	66	E2	8F	F7	9E	D5	FF	5A	5F	FD	B8
414	FE	E0	AA	92	52	DB	0B	E7	46	06	D2	E8	5B	9B	73	AF	CD	3D	E7	9E	9C	AD	FC

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
00000000	74	65	6D	70	6C	61	74	65	20	22	42	6F	6F	74	20	53
00000010	65	63	74	6F	72	20	4E	54	46	53	22	0D	0A	0D	0A	2F
00000020	2F	20	54	65	6D	70	6C	61	74	65	20	62	79	20	50	61
00000030	75	6C	20	4D	75	6C	6C	65	6E	2C	20	70	63	67	75	72
00000040	75	40	74	68	65	2D	61	6E	73	77	65	72	2E	63	6F	6D
00000050	0D	0A	2F	2F	20	6C	61	73	74	20	6D	6F	64	69	66	69
00000060	65	64	20	4A	75	6C	20	31	36	2C	20	32	30	30	30	0D
00000070	0A	0D	0A	2F	2F	20	54	6F	20	62	65	20	61	70	70	6C
00000080	69	65	64	20	74	6F	20	73	65	63	74	6F	72	20	30	20
00000090	6F	66	20	61	6E	20	4E	54	46	53	2D	66	6F	72	6D	61
000000A0	74	74	65	64	0D	0A	2F	2F	20	6C	6F	67	69	63	61	6C
000000B0	20	64	72	69	76	65	20	6F	72	20	74	6F	20	74	68	65
000000C0	20	6D	69	72	72	6F	72	20	63	6F	70	79	20	6F	66	20
000000D0	74	68	65	20	62	6F	6F	74	0D	0A	2F	2F	20	73	65	63
000000E0	74	6F	72	2C	20	77	68	69	63	68	20	77	69	6C	6C	20
000000F0	62	65	20	6C	6F	63	61	74	65	64	20	6E	65	61	72	20

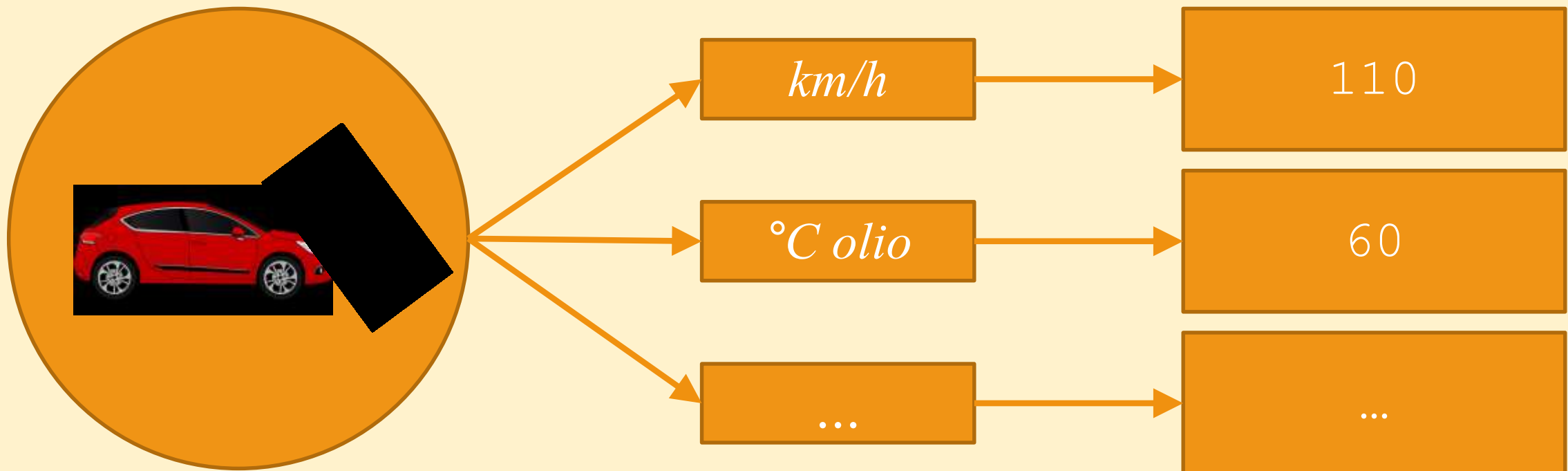
Un dato, se non conosciamo il processo che l'ha generato, non è più un dato.  
Che dato è «1492»?

Anno scoperta dell'America	1492
Prezzo bicicletta in fibra di carbonio	1492
Articolo codice civile sugli effetti della garanzia di un prodotto	1492
...	

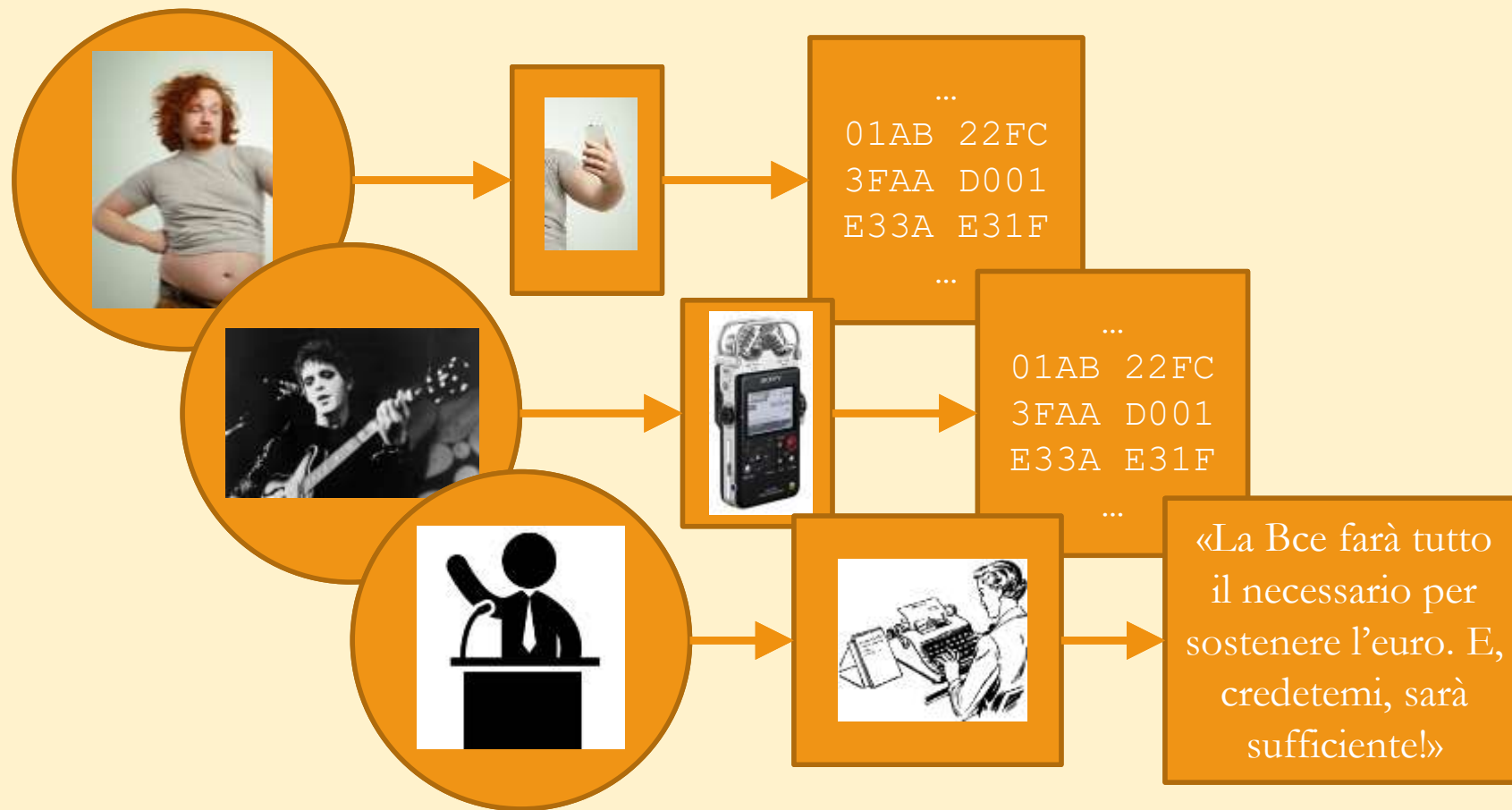
# Esempi di PROCESSO: una rilevazione statistica...



...un *device* della Internet of things...



# ...o qualunque altro mezzo di codifica!



Non tutti i processi generano numeri



# Non tutte le evidenze empiriche hanno la forma di un dato!

- Nella ricerca empirica, il ricercatore “classico” sa che a volte gli serve una rilevazione statistica, altre volte una analisi qualitativa.
- La sovrabbondanza di dati può portare al pregiudizio che qualunque decisione **può** essere supportata soprattutto da analisi quantitative.
- O addirittura che **deve** essere basata sulla elaborazione automatica dei \*byte di dati immagazzinati.

data driven *vs* evidence driven

## “The Cost of Missing Something”

- TED Conference di Tricia Wang, consulente di NOKIA nel 2009.



- In ambito big data si usano spesso metafore tipo l'ago nel pagliaio, o il granello di sabbia nella spiaggia.
- E gli analytics in NOKIA dicevano che gli smartphone non avevano mercato: troppo pesanti e costosi.



你比你想象的更强大

你拥有一种力量, 来开创, 来塑造, 来分享自己的生活,  
这种力量就在你的手中, 背包中, 或口袋中, 这就是 iPhone 5s.

观看影片“强大”

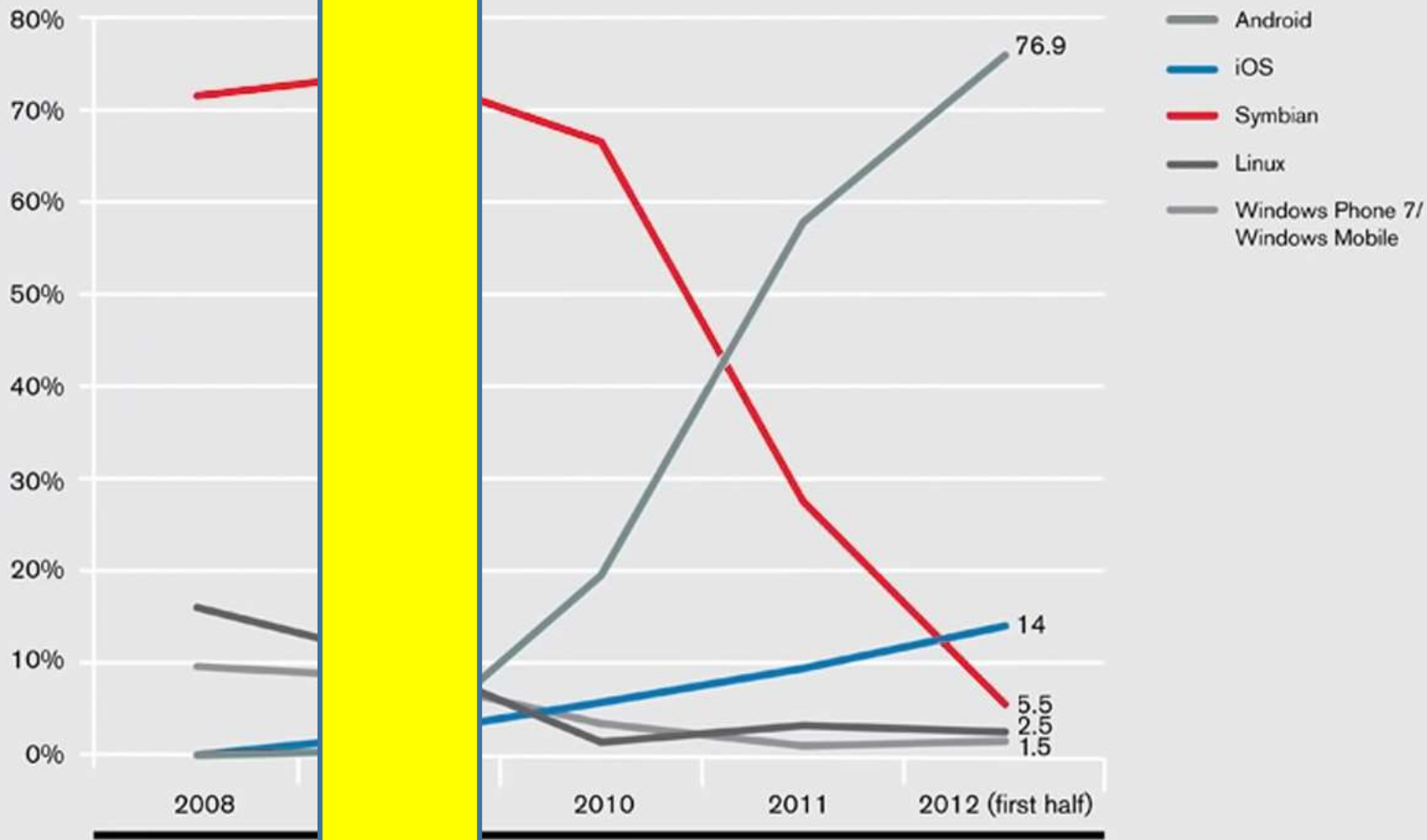


l'analisi qualitativa in action



## Upward Mobility

China smartphone market (percentage of total shipments) by operating system



Source: IDC Asia/Pacific Phone Tracker, August 2012

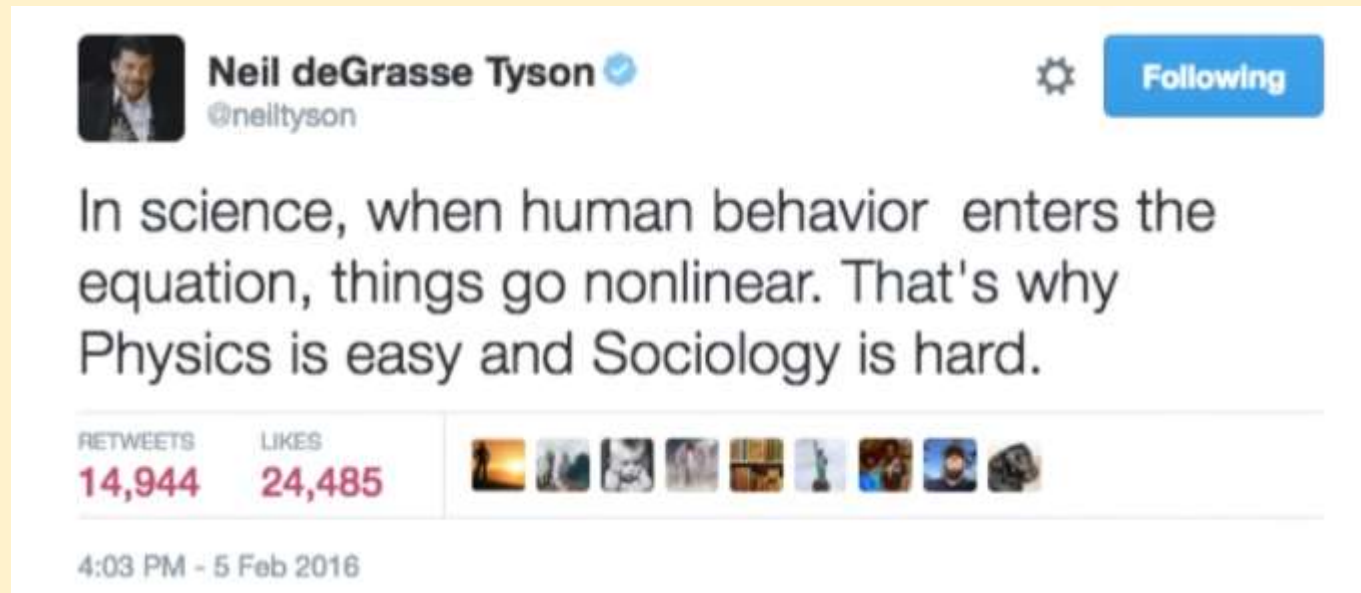
NOKIA prima e dopo la *strategic decision*



Il pagliaio allargando lo sguardo...

## altre due conclusioni

1. Dati e analisi qualitativi sono fondamentali. Quelli che, con termine “nuovo”, Patricia Wang chiama “THICK DATA”.
2. Sono tanto più importanti quanto più si ha a che fare con le persone e le loro scelte.



## Tornando ai dati quantitativi

- La caratteristica fondamentale che determina la qualità di un dato statistico risiede tutta nel **processo** che ha portato a quel dato.
- Processo che in ambito statistico si concretizza in una **raccolta sistematica** di misure del fenomeno.

## Cosa si intende per raccolta sistematica

Adrian Smith, *Mad Cows and Ecstasy: Chance and Choice in an Evidence-Based Society*, ci propone una lista dei passaggi tipici che sono necessari per produrre una SOLIDA evidenza statistica.

- the framing of questions;
- design of experiments or surveys;
- drawing up protocols for data collection
- collection of data
- monitoring compliance with protocols
- monitoring data quality
- data storage, summarization, presentation
- stochastic modelling
- statistical analysis
- model criticism and assumptions assessment
- inference reporting and the use of results for prediction, decision-making or hypothesis generation

raccolta dati

## Un esempio di successo in ambito BD-A

- Il progetto ARTEMIS della University of Ontario, per il monitoraggio e la prevenzione delle infezioni in ospedale nelle unità di terapia intensiva per nati prematuri.
- Con un “Big Data approach”, è stata implementata una tecnologia per rilevare e registrare in real time più di 1,000 misure al secondo relative a elettrocardiogrammi, battiti cardiaci, ritmi respiratori, saturazione dell’ossigeno nel sangue...
- Sono stati in grado di identificare “early warnings” di infezione 24 ore prima rispetto all’approccio tradizionale.

# Un esempio di raccolta non-sistematica: Google flu trends

nature

Vol 457|19 February 2009|doi:10.1038

## LETTERS

### Detecting influenza epidemics using search engine query data

Jeremy Ginsberg<sup>1</sup>, Matthew H. Mohebbi<sup>1</sup>, Rajan S. Patel<sup>1</sup>, Lynnette Brammer<sup>2</sup>, Mark S

Seasonal influenza epidemics are a major public health concern, causing tens of millions of respiratory illnesses and 250,000 to 500,000 deaths worldwide each year<sup>1</sup>. In addition to seasonal influenza, a new strain of influenza virus against which no previous immunity exists and that demonstrates human-to-human transmission could result in a pandemic with millions of fatalities<sup>2</sup>. Early detection of disease activity, when followed by a rapid response, can reduce the impact of both seasonal and pandemic influenza<sup>3,4</sup>. One way to improve early detection is to monitor

By aggregating historical logs of online web search queries submitted between 2003 and 2008, we computed a time series of weekly counts for 50 million of the most common search queries in the United States. Separate aggregate weekly counts were kept for every query in each state. No information about the identity of any user was retained. Each time series was normalized by dividing the count for each query in a particular week by the total number of online search queries submitted in that location during the week (Supplementary Fig. 1).

#### How Google Flu Trends Works



By aggregating historical logs of online web search queries submitted between 2003 and 2008, we computed a time series of weekly counts for 50 million of the most common search queries in the United States. Separate aggregate weekly counts were kept for every query in each state. No information about the identity of any user was retained. Each time series was normalized by dividing the count for each query in a particular week by the total number of online search queries submitted in that location during the week, resulting in a query fraction (Supplementary Fig. 1).

## Un esempio di raccolta non-sistematica: Google flu trends



The screenshot shows the top portion of a web page from the journal Nature. The header features the 'nature' logo in white on a dark red background, with the tagline 'International weekly journal of science' below it. A navigation bar contains links for Home, News & Comment, Research, Careers & Jobs, Current Issue, Archive, and Audio & Video. Below this is a breadcrumb trail: Archive > Volume 494 > Issue 7436 > News > Article. The main content area has a sub-header 'NATURE | NEWS' and a share icon. The article title is 'When Google got flu wrong' in a large, dark font. Below the title is a subtitle: 'US outbreak foxes a leading web-based method for tracking seasonal flu.' The author's name, 'Declan Butler', is listed in blue. The date '13 February 2013' is at the bottom left.

“When influenza hit early and hard in the United States this year, it quietly claimed an unacknowledged victim: one of the cutting-edge techniques being used to monitor the outbreak.”



## La differenza tra Artemis e GFT sta nella raccolta dati

- Nel caso di Artemis, è stata seguita tutta la checklist dello “statistician cook-book”.
  - *L’aspetto nuovo rispetto a una indagine tradizionale sta nella mole dei dati da processare in real time, che richiede una **infrastruttura IT** dedicata e modelli in grado di **automatizzare le analisi**.*
- Nel caso di GFT i dati non sono stati prodotti per gli scopi dell’analista. Sono “**found data**”, che a posteriori hanno rivelato una correlazione con il fenomeno di interesse.
  - *Non era il modello statistico ad essere sbagliato: più semplicemente, la variabile proxy utilizzata per predire la diffusione dell’epidemia era una proxy di **cattiva qualità***



## Quality is KEY

- Non è perché i dati sono **BIG**, o il computer che usiamo è **POWERFUL**, o perché il posto dove registriamo i dati è addirittura una *NUVOLA*, che le evidenze che produciamo possono essere utilizzate a supporto di decisioni strategiche.
- L'ingrediente fondamentale è la **qualità dei dati di partenza**.

# Non è perché i dati sono BIG

Un altro esempio di “smart forecast”, che nonostante la mole **BIG** di dati si è convertita in un *epic fail*.

Indagine su 10 milioni di cittadini americani per prevedere il futuro presidente degli Stati Uniti.

**La previsione, “Landon in a Landslide”:**

*Landon, 57.1%, Roosevelt, 42.9%*

**I risultati veri, oops:**

*Roosevelt 60.8%, Landon 36.5%*

## Messy Matters

Bring your own data

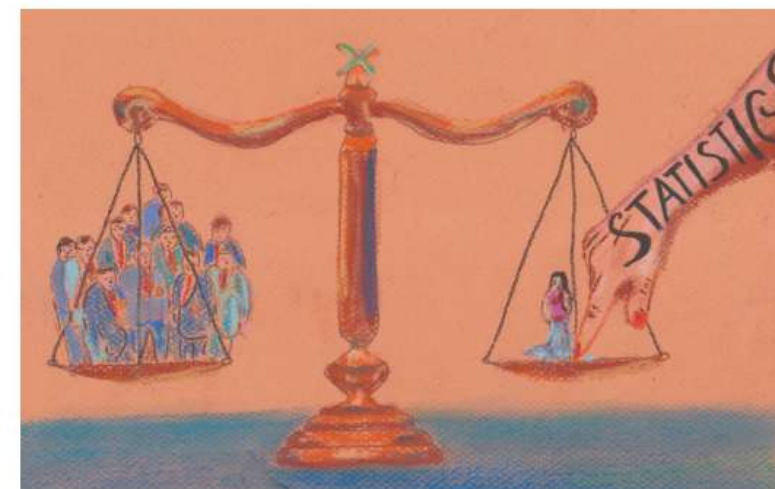
« Google Glass(es)

Yootles Postmortem »

### Forecasting Elections with Dirty Data

Monday, September 30, 2013

By Sharad Goel



During the 1936 U.S. presidential campaign, the popular magazine *Literary Digest* conducted a mail-in election poll that attracted over two million responses, a huge sample even by today's standards. Unfortunately for them, size isn't the only thing that matters. *Literary Digest* notoriously and erroneously predicted a landslide victory for Republican candidate Alf Landon. In reality, the incumbent Franklin D. Roosevelt decisively won the

## Made Data vs Found data (I)

«The “big data” that interests many companies is what we might call “found data”, the digital exhaust of web searches, credit card payments and mobiles pinging the nearest phone mast...»

*Tim Harford, cit.*

Le differenze (e i pericoli) rispetto a una indagine statistica tradizionale risiedono esattamente qua: I dati non provengono da un protocollo rigoroso di data procurement, ma sono stati “trovati” da qualche parte, li ha prodotti qualcun altro.

→ *c'è qualcosa di “già visto” anche in questo caso...*

## Made Data vs Found data (II)

Anche nell'ambito della statistica ufficiale c'è un progressivo abbandono della raccolta dati con indagini campionarie (*made data*), a favore dell'utilizzo di **dati amministrativi**, cioè a dati raccolti da amministrazioni pubbliche nella loro normale attività gestionale (dati fiscali, contributivi, assicurativi...).

Il punto in comune a tutti i *found data*, che vale per i dati amministrativi così come per un tweet o una traccia GPS, è che sono generati/raccolti per scopi che non sono quelli dell'analista che li ha "trovati".

si  
può  
fare!



and the answer is...

# Forecasting vs ...

- 1) Forecasting vs understanding
- 2) Prediction vs explanation
  - black-box models vs  $y=f(x)$
  - modelli strutturali vs forma ridotta
  - microsimulazioni statiche vs dinamiche
- 3) Correlation vs causation
  - modelli quasi sperimentali
  - social random experiments... *already on the go!*
    - Esperimento sociale randomizzato di Facebook nel 2012
    - Google per lo sviluppo di nuovi algoritmi di ricerca



# Experimental evidence of massive-scale emotional contagion through social networks

Adam D. I. Kramer<sup>a,1</sup>, Jamie E. Guillory<sup>b,2</sup>, and Jeffrey T. Hancock<sup>b,c</sup>

<sup>a</sup>Core Data Science Team, Facebook, Menlo Park, CA 94025; and <sup>b</sup>Departments of <sup>c</sup>Communication and <sup>d</sup>Information Science, Cornell University, Ithaca, NY 14853

Edited by Susan T. Fiske, Princeton University

Emotional states can be transmitted through social networks, leading people to experience emotions without their awareness. Emotional contagion was tested in laboratory experiments, with participants exposed to negative emotions to others. A social network, collected over a 20-year period, was used to study moods (e.g., depression, happiness) in laboratory experiments, with participants exposed to negative emotions to others. In a social network, collected over a 20-year period, we tested whether exposure to negative emotions outside of in-person interactions affected the amount of emotional content expressed in posts. When expressions were reduced, people posted more negative posts; when expressions were increased, the opposite pattern of emotions expressed by others was observed, constituting experimental evidence of emotional contagion via social networks. These findings contrast to prevailing assumptions that verbal cues are not strictly necessary for emotional contagion and that the observation of other people's emotions is a positive experience for people.

computer-mediated communication

## Editorial Expression of Concern and Correction

### PSYCHOLOGICAL AND COGNITIVE SCIENCES

PNAS is publishing an Editorial Expression of Concern regarding the following article: “Experimental evidence of massive-scale emotional contagion through social networks,” by Adam D. I. Kramer, Jamie E. Guillory, and Jeffrey T. Hancock, which appeared in issue 24, June 17, 2014, of *Proc Natl Acad Sci USA* (111:8788–8790; first published June 2, 2014; 10.1073/pnas.1320040111). This paper represents an important and emerging area of social science research that needs to be approached with sensitivity and with vigilance regarding personal privacy issues.

Questions have been raised about the principles of informed consent and opportunity to opt out in connection with the research in this paper. The authors noted in their paper, “[The work] was consistent with Facebook’s Data Use Policy, to which all users agree prior to creating an account on Facebook, constituting informed consent for this research.” When the authors prepared their paper for publication in PNAS, they stated that: “Because this experiment was conducted by Facebook, Inc. for internal purposes, the Cornell University IRB [Institutional Review Board] determined that the project did not fall under Cornell’s Human Research Protection Program.” This statement has since been [confirmed by Cornell University](#).

### PSYCHOLOGICAL AND COGNITIVE SCIENCES

Correction for “Experimental evidence of massive-scale emotional contagion through social networks,” by Adam D. I. Kramer, Jamie E. Guillory, and Jeffrey T. Hancock, which appeared in issue 24, June 17, 2014, of *Proc Natl Acad Sci USA* (111:8788–8790; first published June 2, 2014; 10.1073/pnas.1320040111).

The authors note that, “At the time of the study, the middle author, Jamie E. Guillory, was a graduate student at Cornell University under the tutelage of senior author Jeffrey T. Hancock, also of Cornell University (Guillory is now a postdoctoral fellow at Center for Tobacco Control Research and Education, University of California, San Francisco, CA 94143).” The author and affiliation lines have been updated to reflect the above changes and a present address footnote has been added. The online version has been corrected.

The corrected author and affiliation lines appear below.

**Adam D. I. Kramer<sup>a,1</sup>, Jamie E. Guillory<sup>b,2</sup>, and Jeffrey T. Hancock<sup>b,c</sup>**

<sup>a</sup>Core Data Science Team, Facebook, Inc., Menlo Park, CA 94025; and Departments of <sup>b</sup>Communication and <sup>c</sup>Information Science, Cornell University, Ithaca, NY 14853



