

BIG DATA: INTRODUCTION, RISKS, AND OPPORTUNITIES



SAPIENZA
UNIVERSITÀ DI ROMA



Master in Data Intelligence e Strategie Decisionali
Dipartimento di Scienze Statistiche

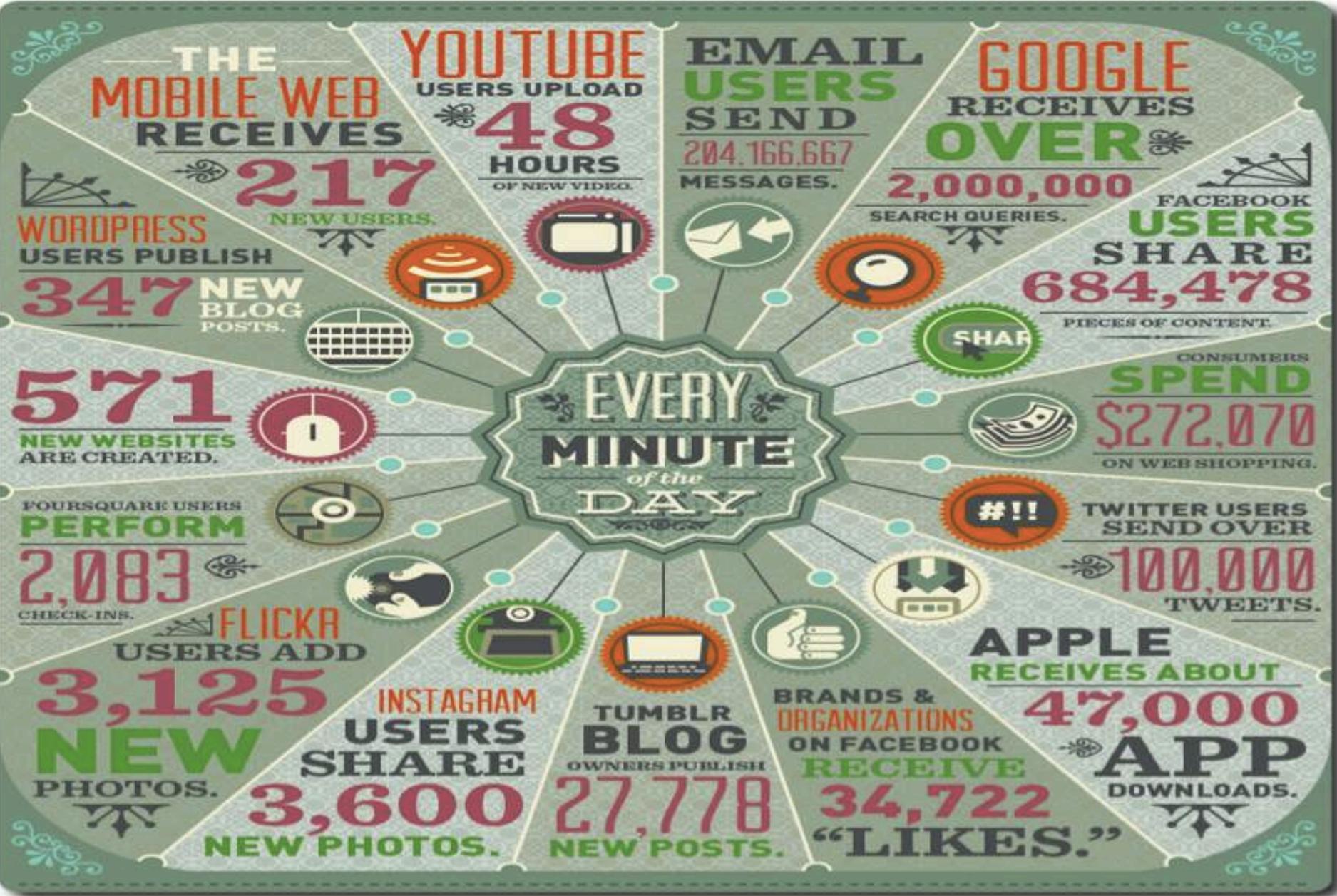
Paolo Dell' Olmo

Organizzazione delle lezioni

- 1. Introduzione**
- 2. Tecniche ed esperienze**
- 3. Tecnologie**
- 4. Data Driven Decision Making**
- 5. Machine Learning**
- 6. Marketing and BI**
- 7. Cognitive Computing**
- 8. Progetti**

Outline

- 1. Introduzione ai Big Data e agli Analytics**
- 2. Small, Open e Big Data**
- 3. Rischi**
- 4. Opportunità**
- 5. Esempi and Takeaways**



Competenze

Quali sono le difficoltà percepite nel mondo Big Data

Abbiamo chiesto ai partecipanti al test di ammissione al Master in Data Intelligence e Strategie Decisionali di esprimere un grado di difficoltà (10 massimo, 1 minimo) su alcuni punti o criticità dei Big Data e Analytics sulla base della propria esperienza.

I punti attenzionati sono riportati nelle slide successive

Questionari 1/4

- Qualità dei dati
- Definizione-strutturazione del problema
- Adeguatezza dei modelli di analisi
- Dati insufficienti per lo specifico problema
- Troppi dati da gestire
- Troppi dati da analizzare

Questionari 2/4

- Funzionalità di Strumenti-Piattaforme utilizzati
- Correttezza approccio metodologico
- Competenze delle risorse umane
- Processi interni all'organizzazione
- Rappresentazione (visualizzazione) dei risultati

Questionari 3/4

- Troppe variabili descrivono il problema
- Mancanza di collegamenti espliciti tra i dati
- Troppe soluzioni da valutare
- Modelli matematici non adeguati al problema
- Modelli matematici alimentati con dati non adeguati o corretti
- Complessità computazionale per trovare/valutare le soluzioni

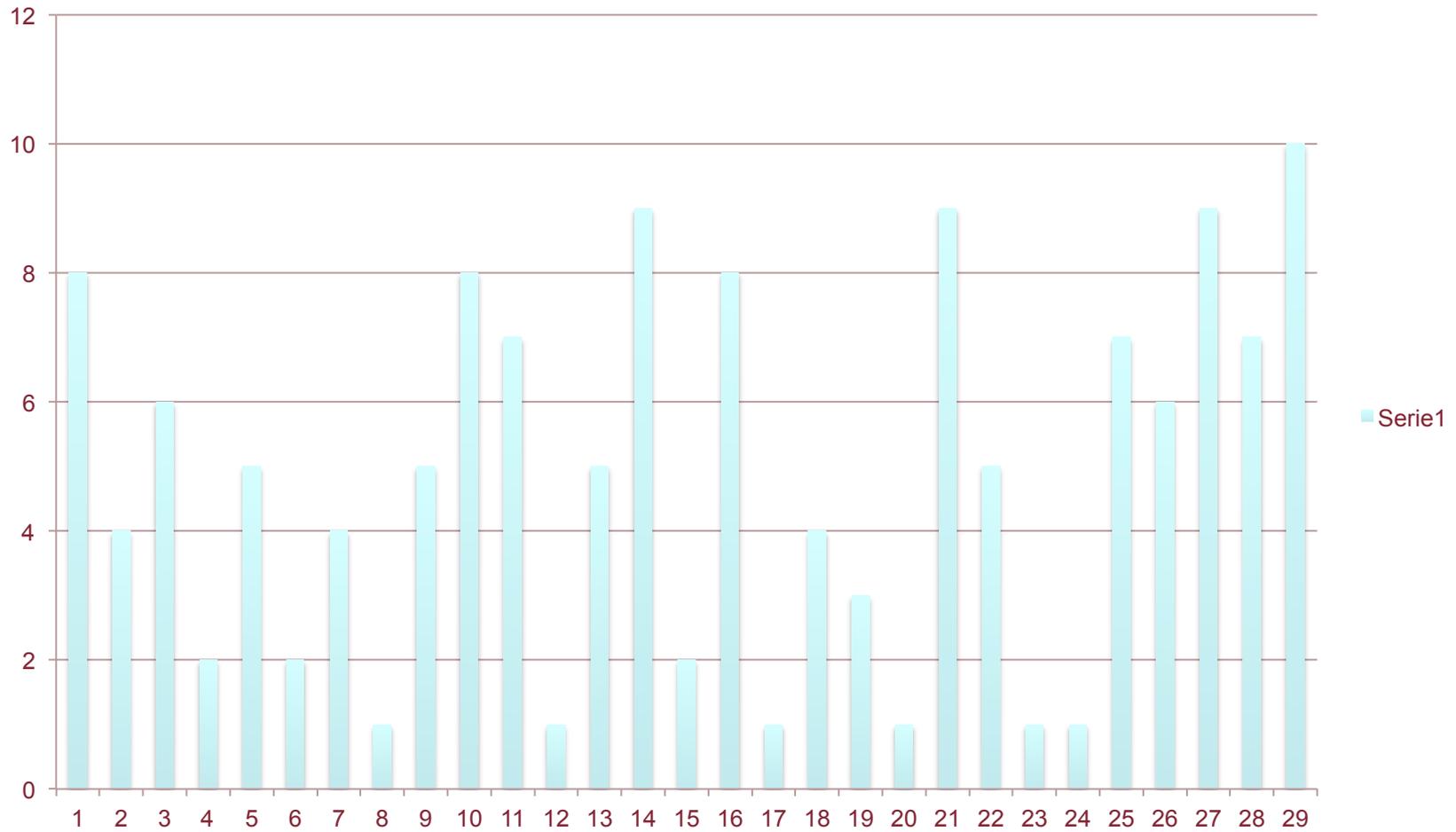
Questionari 4/4

- Uso non corretto degli strumenti
- Comunicazioni errate o poco chiare all'interno del team
- Comunicazioni errate o poco chiare con i livelli superiori
- Comunicazioni errate o poco chiare con i livelli operativi

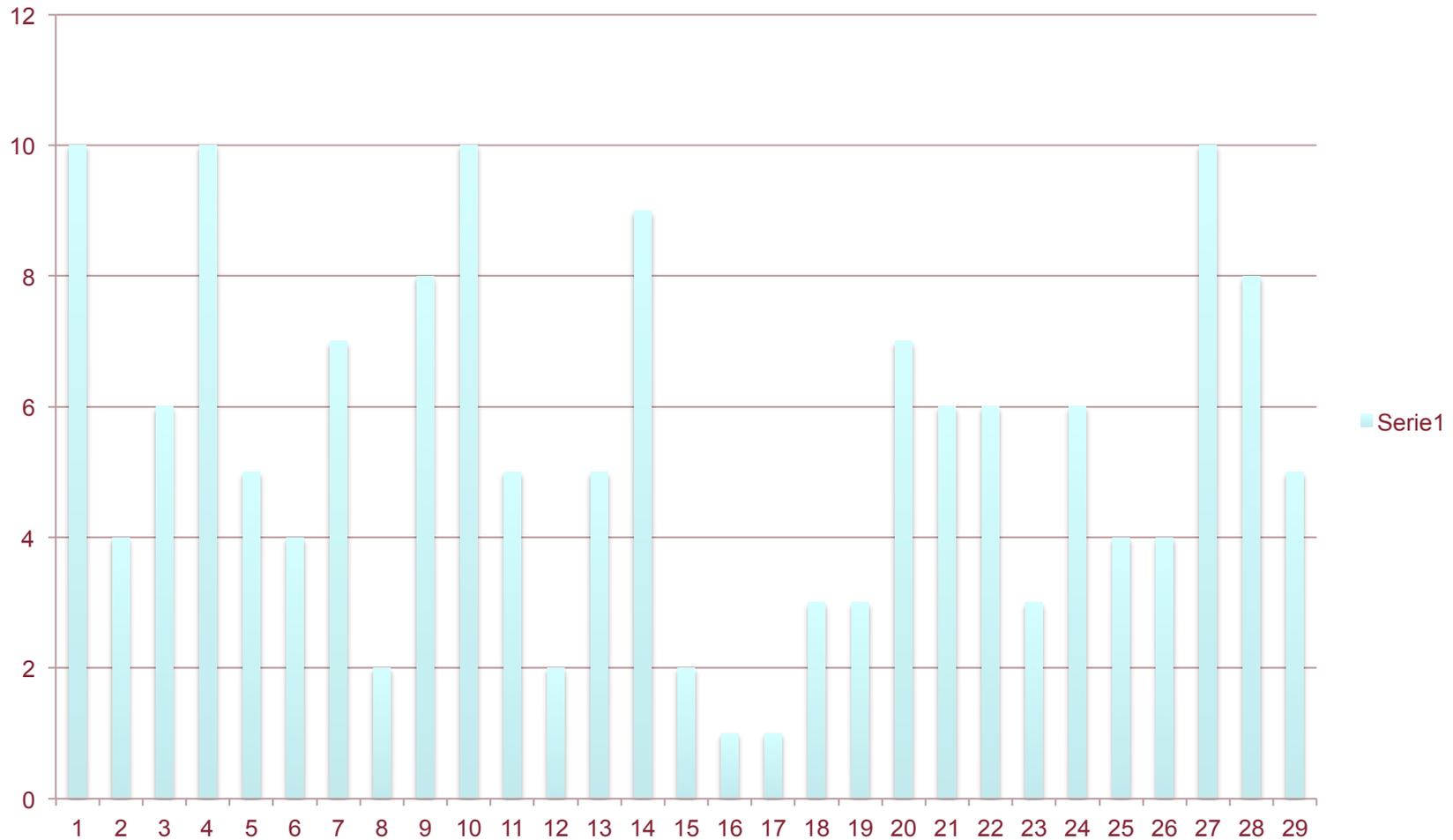
Risposte

DOMANDA	MEDIA	MEDIANA	VAR
Qualità dei dati	7,90	8	3,74
Definizione-strutturazione del problema	8,69	9	3,94
Adeguatezza dei modelli di analisi	8,72	9	1,71
Dati insufficienti per lo specifico problema	7,31	8	3,36
Troppi dati da gestire	4,86	5	8,69
Troppi dati da analizzare	4,59	5	9,11
Funzionalità di Strumenti-Piattaforme utilizzati	6,52	7	5,97
Correttezza approccio metodologico	8,41	9	2,25
Competenze delle risorse umane	7,79	8	3,38
Processi interni all' organizzazione	7,52	8	3,47
Rappresentazione (visualizzazione) dei risultati	7,14	8	4,77
Troppe variabili descrivono il problema	5,69	6	3,51
Mancanza di collegamenti espliciti tra i dati	5,72	6	5,28
Troppe soluzioni da valutare	5,38	5	7,74
Modelli matematici non adeguati al problema	6,66	7	5,95
Modelli matematici alimentati con dati non adeguati o corretti	7,41	8	5,82
Complessità computazionale per trovare/valutare le soluzioni	6,52	7	5,83
Uso non corretto degli strumenti	7,61	7,5	4,69
Comunicazioni errate o poco chiare all' interno del team	8,21	9	4,47
Comunicazioni errate o poco chiare con i livelli superiori	8,07	8	3,99
Comunicazioni errate o poco chiare con i livelli operativi	8,19	8	2,72

Troppi dati da gestire



Troppe soluzioni da valutare



Tecnologie

Progetto Grandi Attezzature TeraDrive - Sapienza

TeraDrive, la prima piattaforma disponibile in Sapienza per la gestione e l'analisi di Big Data

- Spazio di storage dell'ordine dei 256 Terabyte, espandibile sino ad oltre 1 Petabyte.
- Infiniband basata su fibra ottica a bassa latenza.

Progetto Grandi Attezzature TeraStat - Sapienza

- 1. Prevalenza e incidenza di sclerosi multipla ad elevata risoluzione spaziale**
- 2. Sistemi di raccomandazione su larga scala**
- 3. Analisi della mobilità delle persone tramite Big Floating Car Data**
- 4. Multi Objective Optimization e Big Data per la Sicurezza nei Trasporti Pubblici Locali**
- 5. Gestione del Rischio nei mercati finanziari: nuove sfide**
- 6. Small area estimation: measurement error, benchmarking and record linkage**

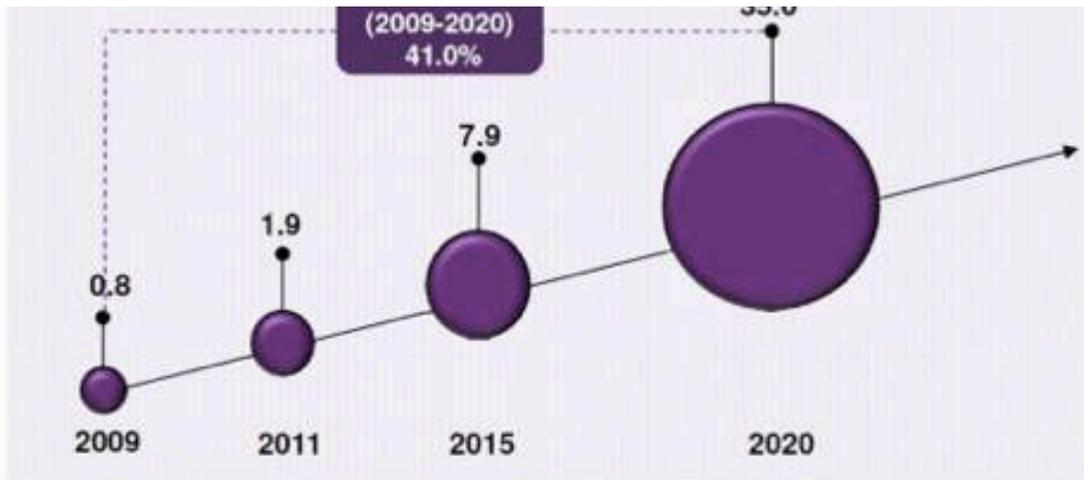
Big and Small

Se non BIG ... Analytics

- La prevalenza dei problemi comunemente denotati come “Big Data” può essere convenientemente risolto utilizzando strumenti tradizionali



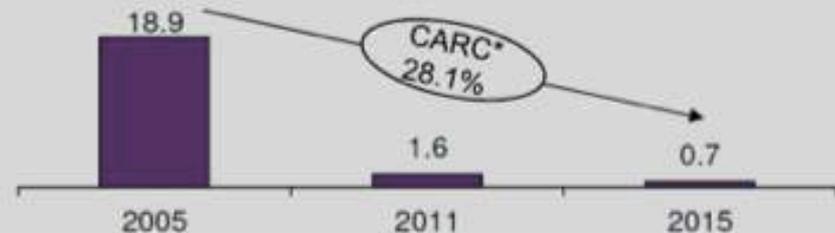
Stock di informazioni



- Need for large storage capacity
- Need for quick retrieval of data
- Enable informed decision-making effectively, leveraging large data sets. For e.g.,:
 - Turn 12 TB of Tweets created each day into improved product sentiment analysis
 - Convert 350 billion annual meter readings to better predict power consumption

Key drivers for deployment of larger datasets

Total storage costs, 2005-2015
USD / gigabyte



- Creation of data from multiple sources and touch-points
- Distributed storage techniques and cloud computing enabling organizations to store large amount of data at lower costs
- Emergence of innovative open source software and architecture such as Hadoop Distributed File System and MapReduce
- Roll-out of 100G Ethernet cables for fast information retrieval

Stock di informazioni

- A typical PC might have had 10 gigabytes of storage in 2000.
- Today, Facebook ingests 500 terabytes of new data every day.
- Boeing 737 will generate 240 terabytes of flight data during a single flight across the US.
- The smart phones, the data they create and consume; sensors embedded into everyday objects will soon result in billions of new, constantly-updated data feeds containing environmental, location, and other information, including video.

Cambio di prospettive

- Stiamo quindi passando da un'analisi basata sulla ricerca del “**Perché**”, della causa di un evento, ad un'analisi basata sulla rilevazione del “**Cosa**” stia accadendo (dalla causalità alla correlazione).
- Un altro cambio di prospettiva (Mayer-Schonberger V. & Cukier K.), è riferito ai profili professionali: nel mondo dei Big Data si attenua l'importanza degli specialisti, non c'è più la rigida ripartizione per competenze. Oggi si richiedono **conoscenze diversificate** ed accumulo di esperienze su tematiche eterogenee.

Small Data vs. Big Data

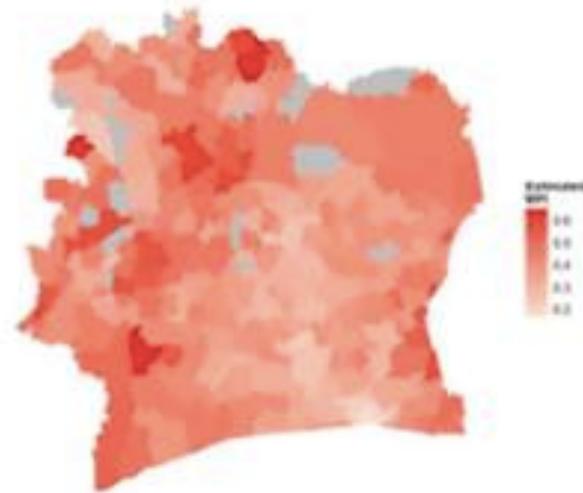
In un mondo “**Small Data**” le decisioni venivano prese sulla base di informazioni che erano:

- Limitate
- Esatte
- Di natura causale

In un mondo “**Big Data**”

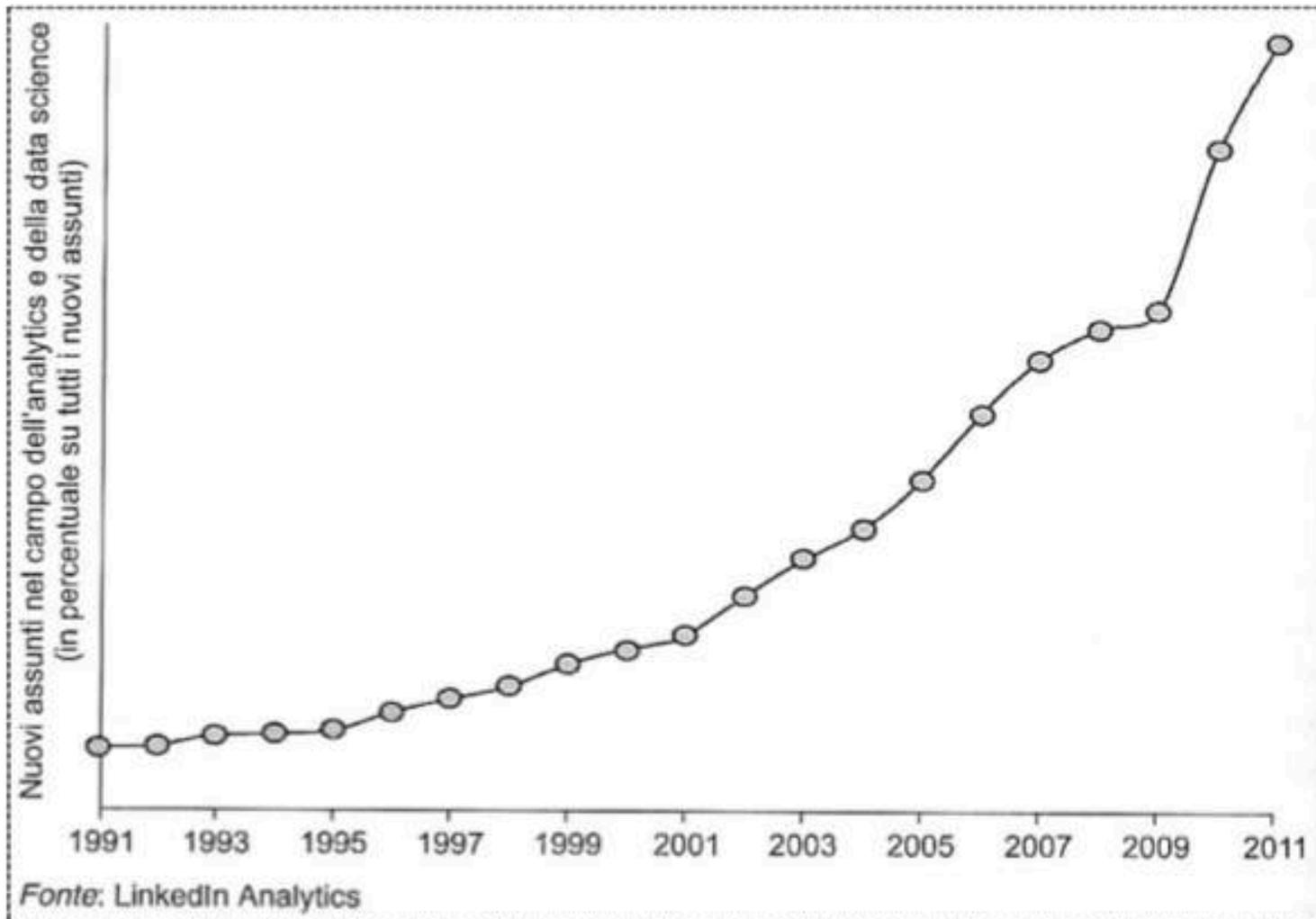
- Poco campionamento
- Rinunciare all’ esattezza per la comprensione del livello macro del fenomeno
- No causalità ed analizzare correlazioni (sono più facili e soprattutto più rapide da scoprire e certamente più economiche).

Small Data vs. Big Data



Poverty map in Côte d'Ivoire
estimated by mobile phone data

Richieste di professionalità



Hype Cycle 2015

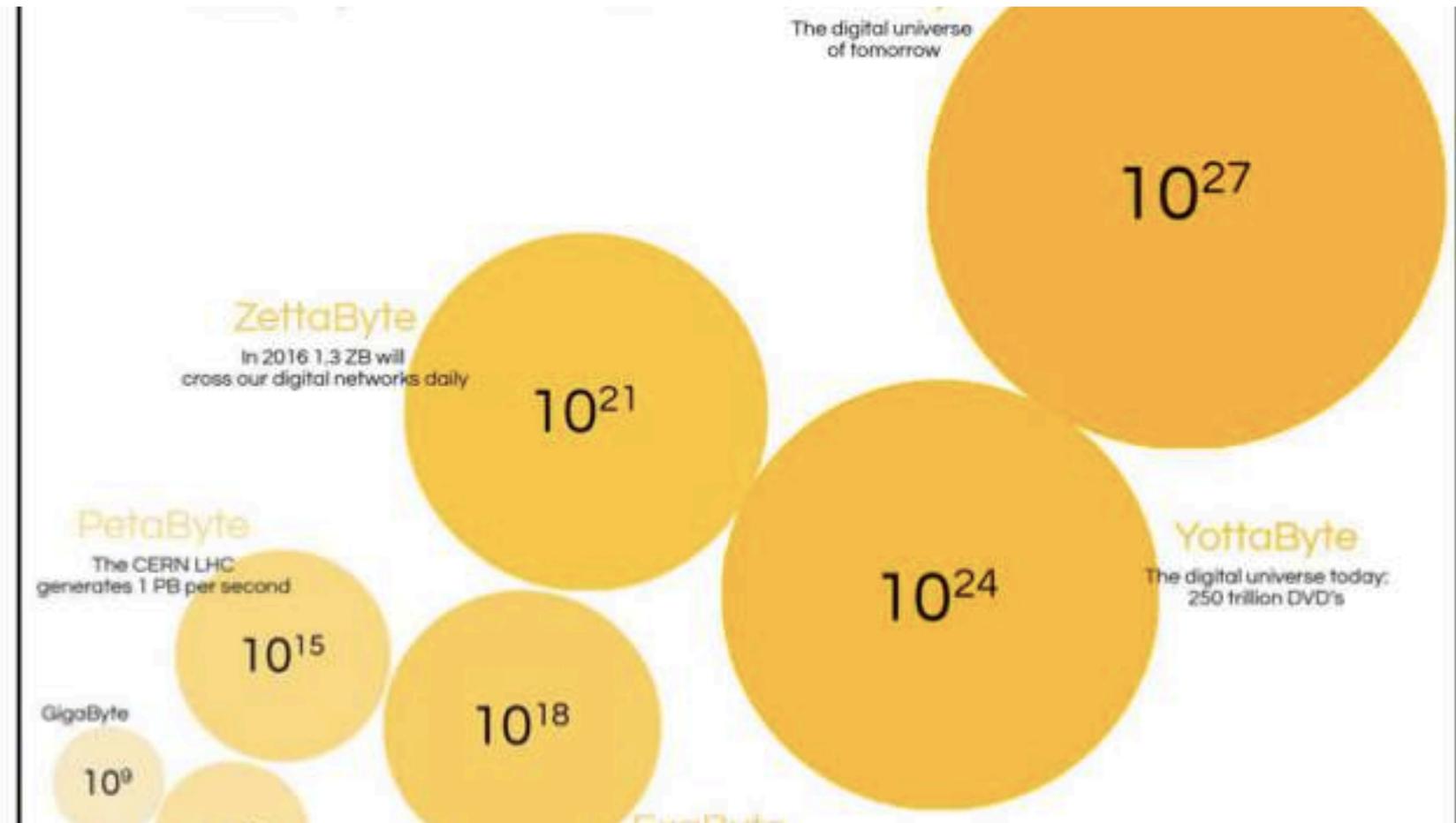


Cosa si intende per Big Data (2011-14)

- **Roger Magoulas della O' Reilly Media**
- **McKinsey Global Institute (2011)**
- **Commissione Europea**
- **Gartner**
- **SDA Bocconi**

Caratteristiche

Volume (esempio): Internet of Things



Caratteristiche - Varietà

- **Dati Strutturati (es. tabelle) 5%**
- **Dati non Strutturati (es. audio e video)**
- **Dati Semi strutturati (es. XML)**

Tipi di dati

Structured Data

- Resides in formal data stores – RDBMS and Data Warehouse; grouped in the form of rows or columns
- Accounts for ~10% of the total data existing currently

RDBMS (e.g., ERP and CRM)



Data Warehousing



Microsoft Project Plan File



Tipi di dati

Unstructured Data

- Comprises data formats which cannot be stored in row/column format like audio files, video, clickstream data,
- Accounts for ~80% of the total data existing currently

Video



Audio



Text message



Blogs



Weather patterns



Location co-ordinates



Web logs & clickstreams



Sensor data/
M2M



Email



Social media



Geospatial data



Tipi di dati

Semi-Structured Data

A form of structured data that does not conform with the formal structure of data models

Accounts for ~10% of the total data existing currently



Tipi di dati

	Video	Image	Audio	Text/ numbers	
Banking	High	High	High	High	
Insurance	Low	Low	Low	High	
Securities and investment services	Low	Low	Low	High	
Discrete manufacturing	Medium	Medium	Low	High	
Process manufacturing	Medium	Medium	Low	High	
Retail	Medium	Low	Low	High	
Wholesale	Low	Low	Low	High	
Professional services	Medium	Medium	Medium	High	
Consumer and recreational services	Medium	Low	Medium	Medium	
Health care	Low	High	Low	High	
Transportation	Medium	Medium	Low	High	
Communications and media ²	High	Medium	High	High	
Utilities	Medium	Medium	Low	High	
Construction	Low	High	Low	Medium	
Resource industries	Medium	Medium	Low	High	
Government	High	Medium	High	High	
Education	High	Medium	High	Medium	

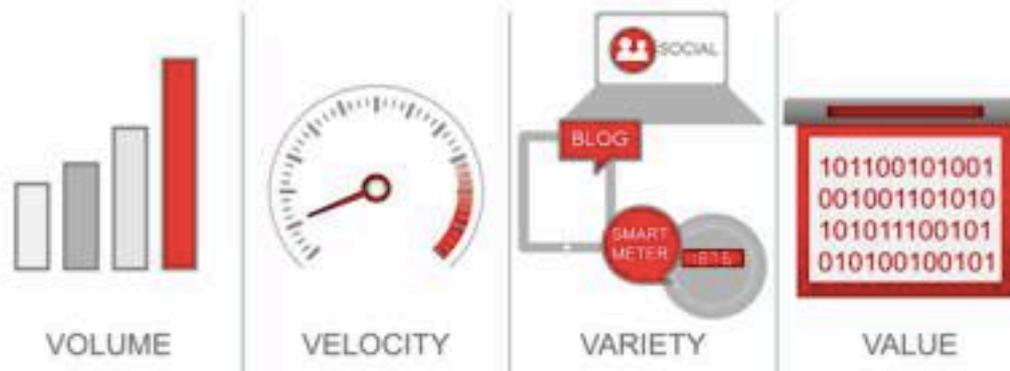
Penetration

- High
- Medium
- Low

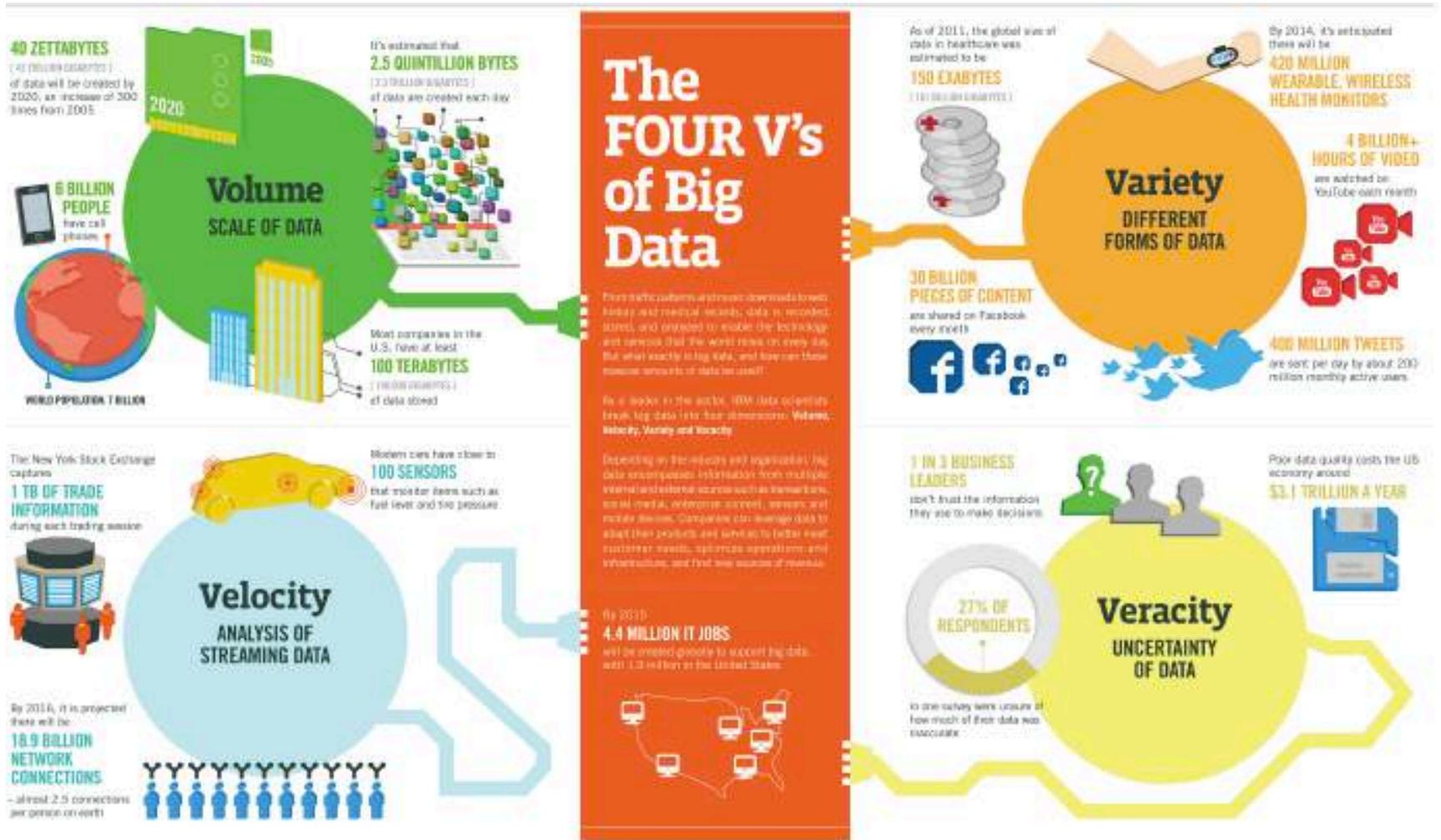
Big Data: V3 + valore

Big Data: V³+VALUE

- **Volume:** Gigabyte(10^9), Terabyte(10^{12}), Petabyte(10^{15}), Exabyte(10^{18}), Zettabyte (10^{21})
- **Varietà:** Structured, semi-structured, unstructured; Text, image, audio, video, record
- **Velocità:** Periodic, Near Real Time, Real Time
- **Valore:** Può generare grossi vantaggi competitivi!



Le "tre" V



Sources: McKinsey Global Institute, Terabit, Cisco, Gartner, IDC, SAS, IBM, HETPCC, GAO



Caratteristiche - Veracity

Veracity: genuinità e bontà del dataset. Questo elemento è stato introdotto da IBM, in un suo report del 2012, per porre l'accento sulla necessità di selezionare e trattare dati affidabili (esempio: i dati generati da Social Network potrebbero essere in parte inaffidabili)

Fattori abilitanti

- **Eccezionale crescita del volume dei dati accumulati**
- **Nuove “affordance” tecnologiche offerte da**
- **Affermazione di nuovi approcci teorici**

Le opportunità

Start up israeliana

The API to display what's in the air

Include air quality in your products and technology easily using our API. Enrich the user experience, educate and engage with your consumers, by offering intuitive and relevant air quality information, color-coded heat maps, health recommendations and more.

The screenshot displays the BreezoMeter API interface. At the top left, there is a search bar with the placeholder text "Type street address". Below the search bar, a yellow pop-up box shows the current air quality: "54 Moderate air quality BreezoMeter AQI 0 to 100". The main area is a map of Paris, France, with a yellow pin indicating the current location and an AQI of 54. To the right of the map is a sidebar with several sections: 1. Pollutant selection: A row of circular buttons for CO, NO₂, O₃, PM₁₀ (highlighted in yellow), PM_{2.5}, and SO₂. Below this, it shows "PM₁₀ | 57.48 ug/m3 Dominant". 2. Health Sensitivities: A section with icons for a house, a tree, a stethoscope, a horse, and a person running. The text reads: "Health Sensitivities Exposure to air hazards is dangerous for people with health sensitivities, so it is important to monitor air quality at this time". 3. Forecast and History: A section with two tabs, "Forecast" (active) and "History". Below the tabs is a line graph showing air quality over time, with a y-axis ranging from 70 to 90.

Start up israeliana (2)

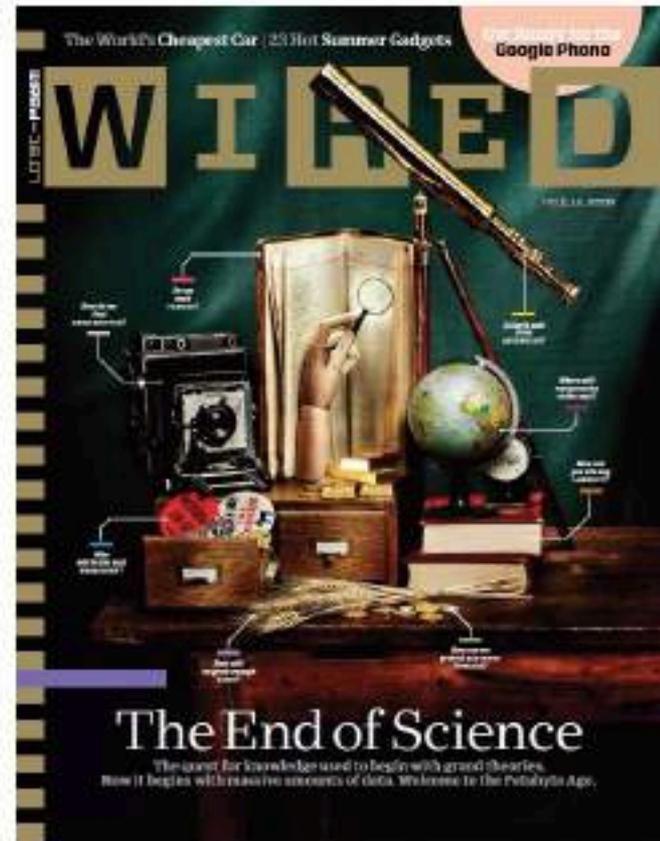


AGF

Rischi

New Paradigm

“This is a world where massive amounts of data and applied mathematics replace every other tool that might be brought to bear. Out with every theory of human behavior, from linguistics to sociology. Forget taxonomy, ontology, and psychology. Who knows why people do what they do? The point is they do it, and we can track and measure it with unprecedented fidelity. With enough data, the numbers speak for themselves.” – Chris Anderson (Wired 2008)



Google Flu

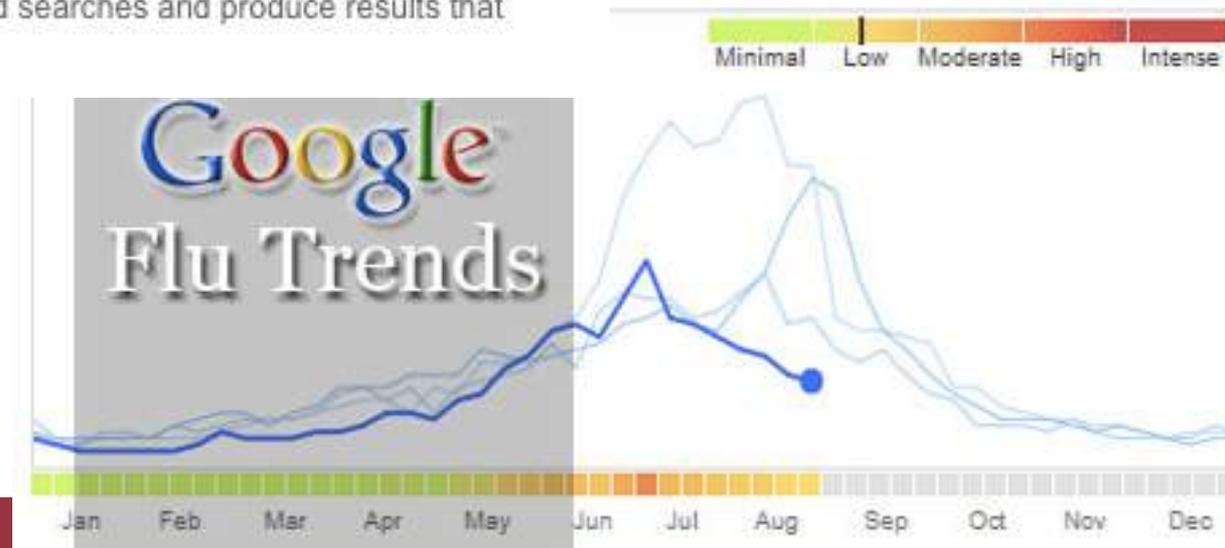
Flu outbreaks could be monitored by internet searches

Agree **2** Disagree **0**

By Hicbd

Mon Nov 26 2012 11:08 am

The Centers for Disease Control and Prevention (CDC) can take up to two weeks to collect and release survey data about where and when the flu is spreading around the country. As shown by Google's Flu Trends, internet search tools could use real-time data to identify regions around the world experiencing a high number of flu-related searches and produce results that correlate closely with CDC data.



Failing of Big Data



GEORGES GOBET/AFP/Getty Images

BIG DATA

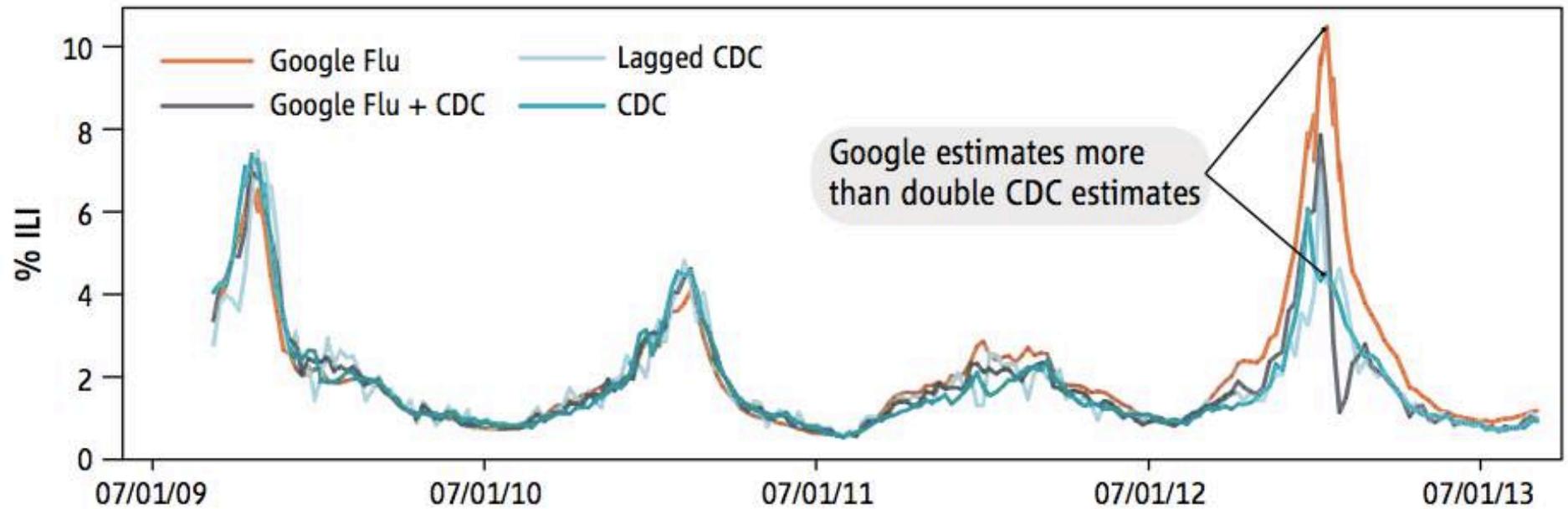
Google's Flu Project Shows the Failings of Big Data



DIGITAL ACCESS TO
SCHOLARSHIP AT HARVARD

The Parable of Google Flu: Traps in Big Data Analysis

Overestimation



Grazie per l'attenzione !

paolo.delloolmo@uniroma1.it