



Indice

1. Big Data

- *Evoluzione dei dati e delle tecniche di analisi*
- *Le 5V dei Big Data*
- *Teorema CAP*
- *Pipeline dell'analisi dei Big Data*



2. Principali Contesti Applicativi

3. Criticità e Rischi dei Big Data

Evoluzione dei dati e delle tecniche di analisi

- Partendo dai **dati** i processi di analisi vogliono **trasformarli** in **informazioni utilizzabili per supportare i processi decisionali** (in contesti aziendali e non).
- Negli **anni '60** i dati erano immagazzinati su **dischi e supporti magnetici**. Si svolgevano **analisi statiche e limitate** (es. il numero di vendite dell'ultimo semestre...)
- Negli **anni '80** i **database Relazionali e SQL** (*Structured Query Language*) permette di realizzare **analisi più dinamiche**.
- Analisi svolte su DB **operazionali**, su cui è registrata ad esempio l'attività giornaliera di un'azienda.

Evoluzione dei dati e delle tecniche di analisi

Problemi delle basi di dati **operazionali**

- L'analisi è svolta da applicativi differenti:
 - Gestione degli ordini.
 - Gestione delle anagrafiche.
 - Contabilità e fatturazione.
- Applicazioni differenti non garantiscono l'uniformità e la coerenza dei dati:
 - **Dati replicati** e manipolati in sw differenti.
 - **Possibili differenze di formato.**
 - **Aggiornamenti dei dati non garantiti.**

Evoluzione dei dati e delle tecniche di analisi

Problemi delle basi di dati **operazionali**

- Sono di tipo **OLTP** (*On Line Transaction Processing*), e presentano un modello dati fortemente **normalizzato**.
- **(+)** La normalizzazione favorisce **inserimenti, cancellazioni e modifiche dei dati** (attività transazionali).
- **(-)** Non è però adatta alle letture.
- **(-)** Incremento notevole del numero di tabelle.
- **(-)** Molte operazioni di JOIN per *denormalizzare* (ricostruire la forma tabellare) e quindi estrazione dei dati complessa.
- **(-)** Mancanza di una profondità storica dei dati.

Evoluzione dei dati e delle tecniche di analisi

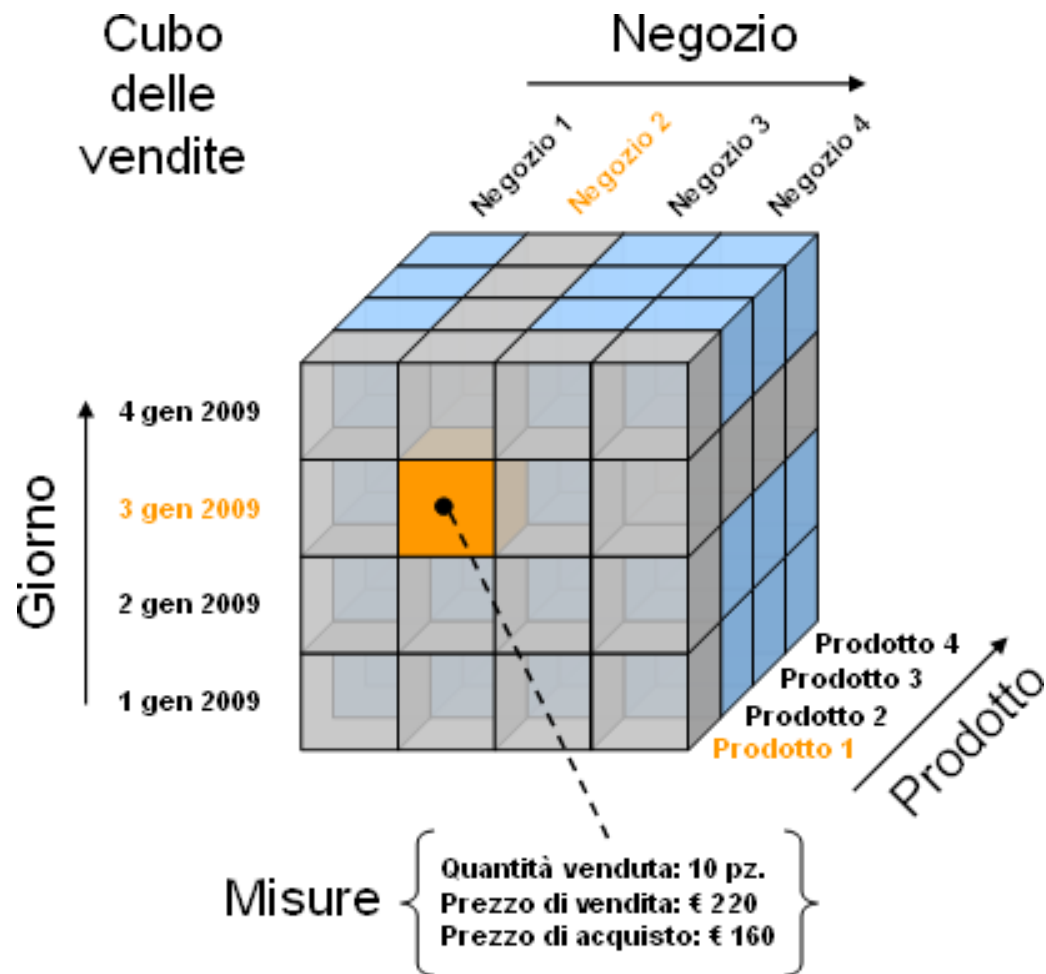
- A causa di questi limiti a partire dagli anni '90 si inizia a parlare di ***data warehouse***, cioè db che integrano dati provenienti da diversi sistemi operazionali.
- I dati sono **integrati, certificati e consistenti** ossia il punto di partenza perfetto per le attività di analisi dei sistemi di **BI**.
- **BI** (*Business Intelligence*) = è un insieme di metodi, modelli, processi, persone e strumenti che permettono una raccolta dati regolare e organizzata.
- **Dati** che possono poi essere **elaborati, aggregati, analizzati e trasformati (valorizzati) in informazioni**, che vengono conservate e rese accessibili in modo semplice e flessibile.

Evoluzione dei dati e delle tecniche di analisi

Evoluzione dei sistemi di BI

- Possibilità di analisi su data warehouse, con query SQL su basi di **dati multidimensionali** (dati e metadati insieme).
- Questi db sono sistemi di tipo **OLAP** (*On Line Analytical Processing*). Hanno una struttura multidimensionale, chiamata **Ipercubo** (spesso semplificata in tre dimensioni).
- (-) Questi sistemi offrono comunque una visione storica:
 - Valutazioni di ciò che è accaduto o che sta accadendo.
 - Valutazione statica.

Evoluzione dei dati e delle tecniche di analisi



Navigazione dei dati più semplice grazie ad operazioni di:

- Drill down
- Drill-up
- Slicing
- Dicing

Evoluzione dei dati e delle tecniche di analisi

- Dai primi anni Duemila viene fuori la necessità di **un'analisi dei dati in grado fare previsioni e dare suggerimenti per anticipare gli eventi.**
- Si inizia a parlare di ***data mining***, termine che identifica un insieme di tecniche in grado di **“scavare”** nei dati per **estrarre nuove informazioni e significati**, non evidenti immediatamente.
- Queste tecniche portano spesso alla definizione ***pattern*** (cioè un modello di rappresentazioni di alcune informazioni) e **relazioni tra i dati.**
- *Numerose applicazioni:* la segmentazione della clientela, market basket analysis, campagne pubblicitarie mirate, previsioni etc.

Evoluzione dei dati e delle tecniche di analisi

Dal **2010** le principali tendenze evolutive nell'analisi dei dati e BI sono:

- **Sviluppo di strumenti di business analytics.**

Tecnologie e applicazioni che fanno uso di modelli matematici e statistici per operazioni di data analysis e data mining.

Solitamente offrono **funzionalità per migliorare la visualizzazione dei dati e favorirne la navigazione**, e strumenti di **ottimizzazione nella gestione dei processi** (suddivisione carico di lavoro).

- **Collaboration e information sharing**

La collaborazione e la condivisione delle informazioni (report, documenti, modelli, valutazioni e analisi già svolte) è un requisito sempre più importante, soprattutto in un contesto aziendale.

Es. *Microsoft Share Point*, portale web per pubblicare informazioni.

Evoluzione dei dati e delle tecniche di analisi

Dal **2010** le principali tendenze evolutive nell'analisi dei dati e BI sono:

- **Cloud Computing**

- Risorse HW e SW disponibili come servizi su internet.
- Accesso alle risorse da diversi luoghi e con diversi dispositivi.
- Basti costi iniziali di investimento (determinabili a priori).
- Architettura scalabile.
- Gestione e manutenzione piattaforma (aggiornamenti sw, backup, fault-tolerance..) sono a carico del provider.



Come si è arrivati ai Big Data

Tra le principali fonti di dati, che nel tempo hanno contribuito allo sviluppo del fenomeno dei Big Data troviamo:

- **Fonti Operazionali**
- **Sensori, DCS (Distributed Control System) e strumenti scientifici**
- **Dati non-strutturati e semi-strutturati**

Basi di dati Operazionali

Sono quei **dati** che hanno a che fare con **l'attività giornaliera di un'azienda** (industrie, banche o GDO). Alcuni esempi sono:

- *Applicativi di gestione della produzione* (materie prime, consumi...)
- *Applicativi di gestione degli acquisti* (prodotti, ordini, magazzino...)
- *Applicativi di contabilità* (fatture, saldo, movimenti...)
- *Applicativi di gestione del personale* (anagrafica, premi, malattie...)
- *Applicativi di gestione del cliente* (abitudini, marketing mirato...)

In alcuni casi i **dati operazionali arrivano a creare dei volumi rilevanti**. Esempio consideriamo una banca di grandi dimensioni:

A diagram illustrating the calculation of data volume. It consists of three blue rounded rectangular boxes on the left, followed by a multiplication symbol (X), then another blue rounded rectangular box, followed by an equals sign (=), and finally an orange rounded rectangular box on the right. The first blue box contains the text "10 Mln di Clienti". The second blue box contains the text "250 gg lavorativi". The orange box contains the text "2,5 Mld record/anno".

$$\begin{array}{|c|} \hline 10 \text{ Mln} \\ \hline \text{di Clienti} \\ \hline \end{array} \times \begin{array}{|c|} \hline 250 \text{ gg} \\ \hline \text{lavorativi} \\ \hline \end{array} = \begin{array}{|c|} \hline 2,5 \text{ Mld} \\ \hline \text{record/anno} \\ \hline \end{array}$$

Basi di dati Operazionali

- Le basi di dati operazionali in genere fanno riferimento ai **database relazionali** o **RDBMS** (*Relational Data Base Management System*), tra i più famosi ci sono **MySQL, IBM DB2, Oracle, Microsoft SQL Server**.
- **Aumento dei dati** = Gestione e storicizzazione complessa e onerosa in termini di risorse.
- Gli RDBMS mettono a disposizione alcune tecniche di ottimizzazione:
 - **Indicizzazione**
 - **Compressione**
 - **Partizionamento**

Basi di dati Operazionali

- **Indicizzazione:** Utilizzo di Indici (strutture ordinate).
 - (+) Recupero rapido di informazioni.
 - (-) Scritture lente e aumento dello spazio occupato dal DB.
- **Compressione:** Applicazione di algoritmi di compressione.
 - (+) Meno spazio per il salvataggio dei dati.
 - (-) Tempo di esecuzione degli algoritmi e decompressione dei dati dopo averli recuperati.
- **Partizionamento:** Suddivisione di una tabella in più parti sulla base di uno specifico criterio.
 - (+) Query limitate ad una parte limitata del DB (es. tutti i record da una certa data in poi).
 - (-) I vantaggi si perdono se le query impattano più partizioni.

Dati non-strutturati e semi-strutturati

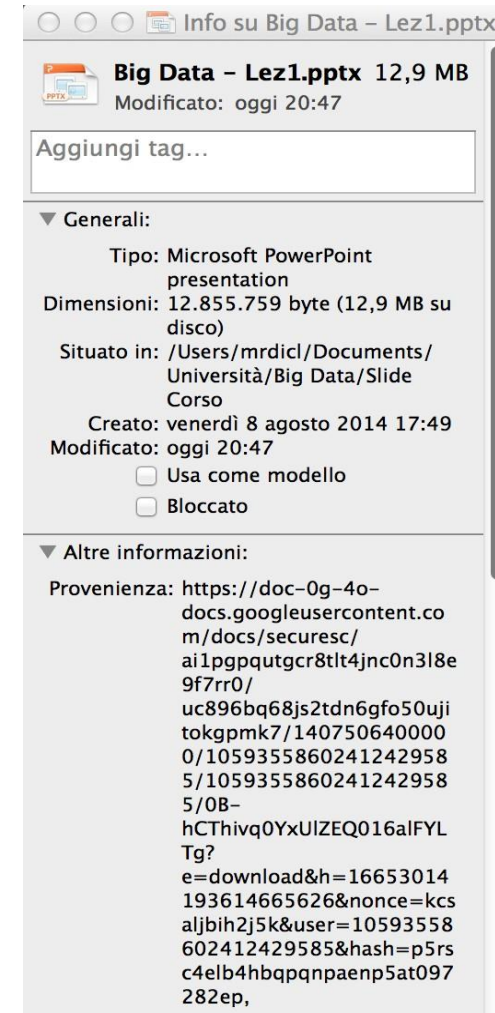
Sono dei **dati che non presentano una struttura predefinita** e che quindi **non** si prestano ad essere gestiti con uno **schema tabellare**. Alcuni esempi presenti in un contesto aziendale sono:

- **Documenti di varia tipologia** (PDF, Word, Excel, PowerPoint etc.)
- **E-mail**
- **Immagini in vari formati** (JPEG, TIFF, GIF, RAW etc.)
- **Strumenti Web 2.0** (Forum e Wiki)

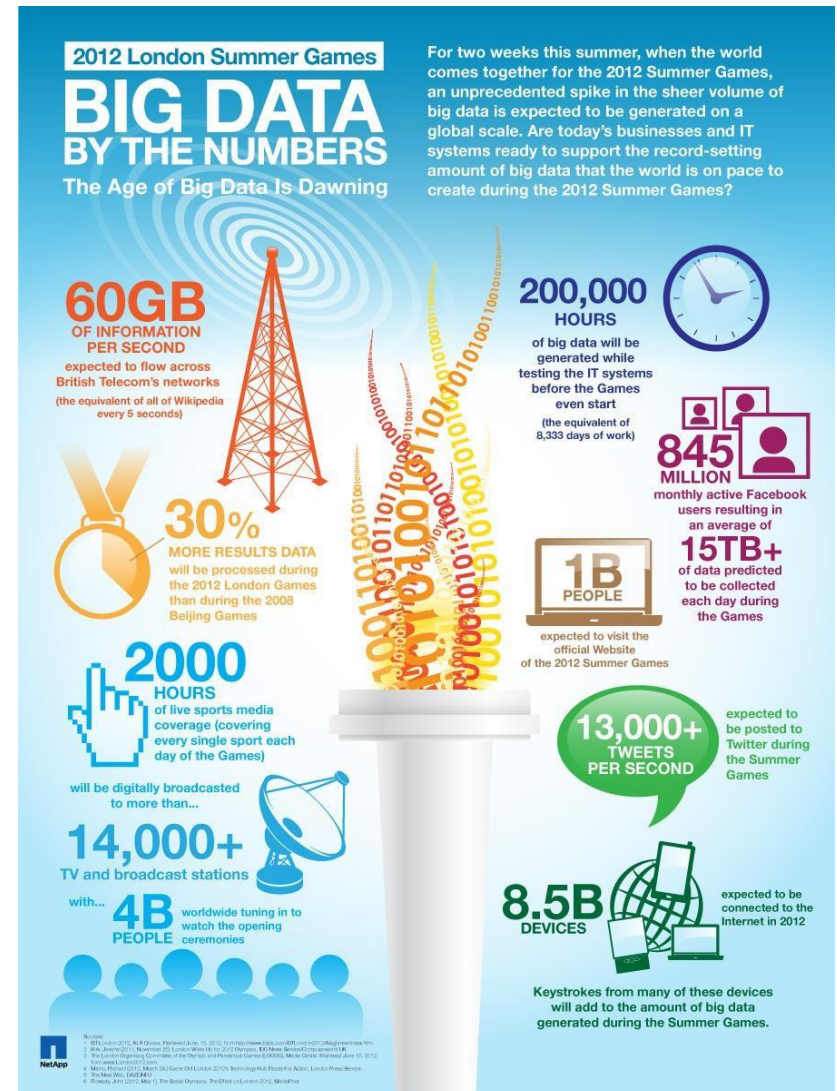
Alcuni di questi documenti in realtà non sono del tutto privi di struttura, si potrebbero definire **semi-strutturati**, per la presenza di informazioni aggiuntive rappresentabili in tabella, **i metadati**.

Dati non-strutturati e semi-strutturati

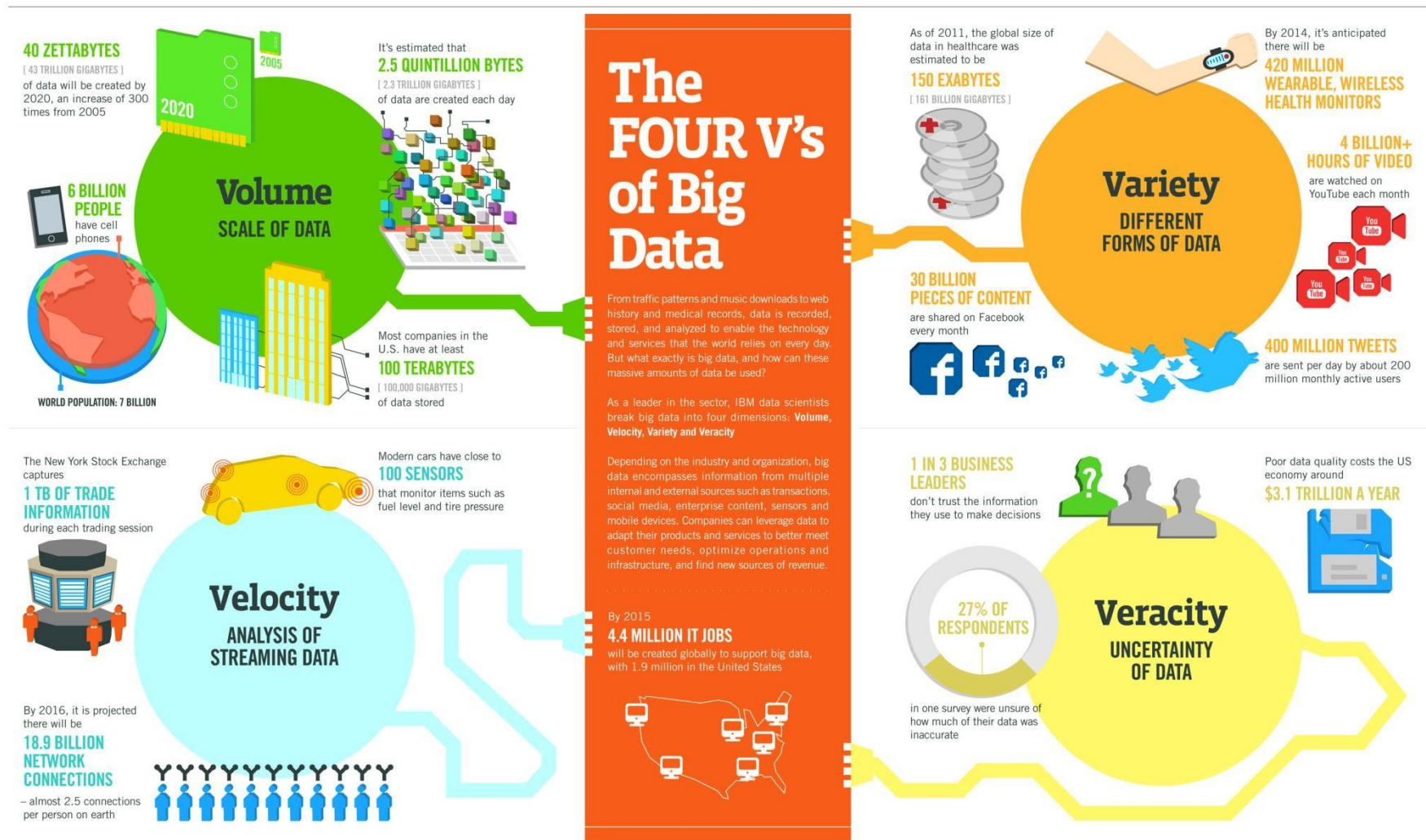
- **Metadati:** sono dei dati utilizzati per descrivere altri dati.
- Sono **facilmente estraibili** dai documenti che descrivono e messi in tabelle.
- Possono essere utilizzati per operare delle ricerche di e sui documenti.



Big Data



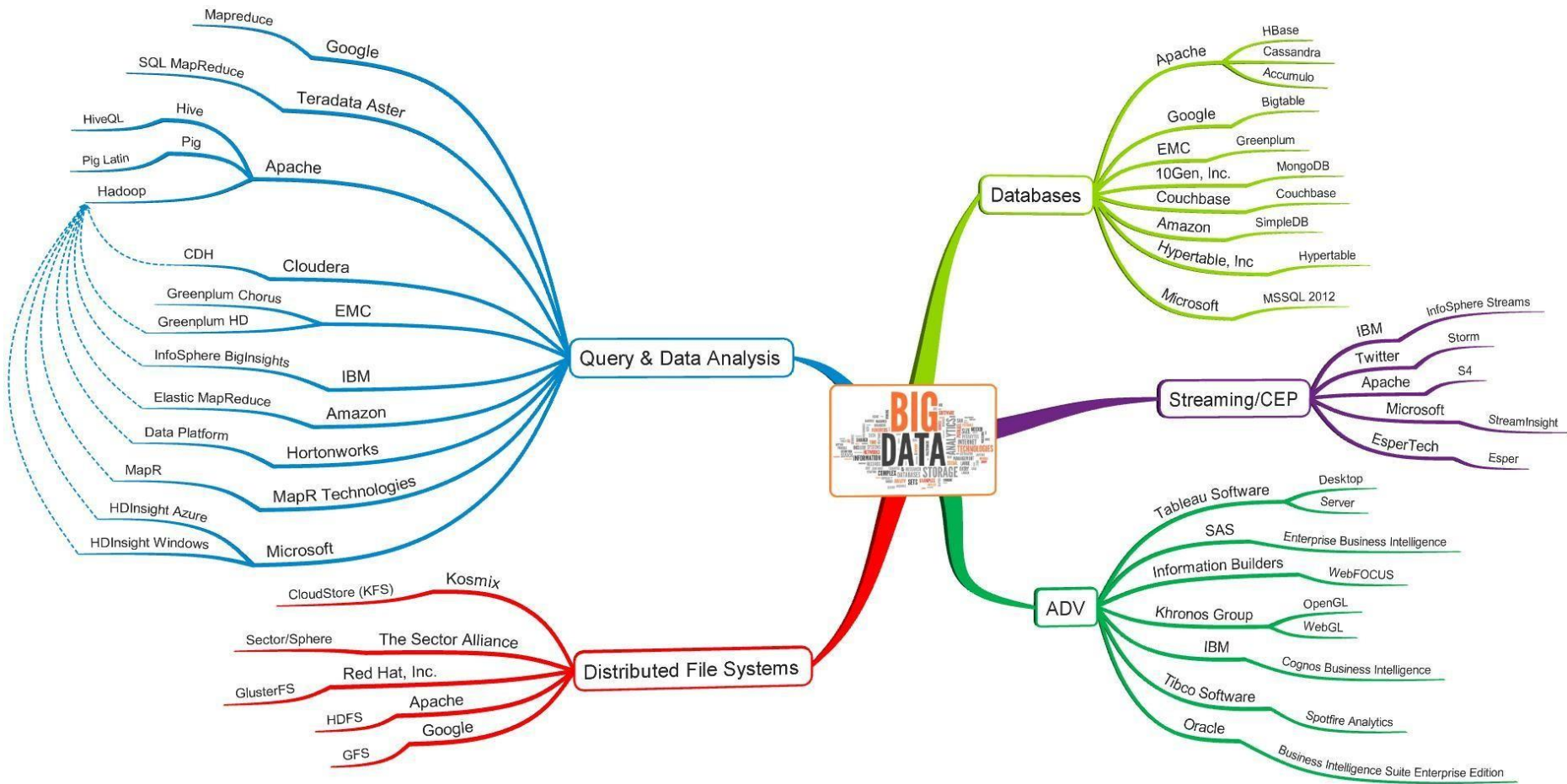
Evoluzione dei dati e delle tecniche di analisi



Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTec, QAS

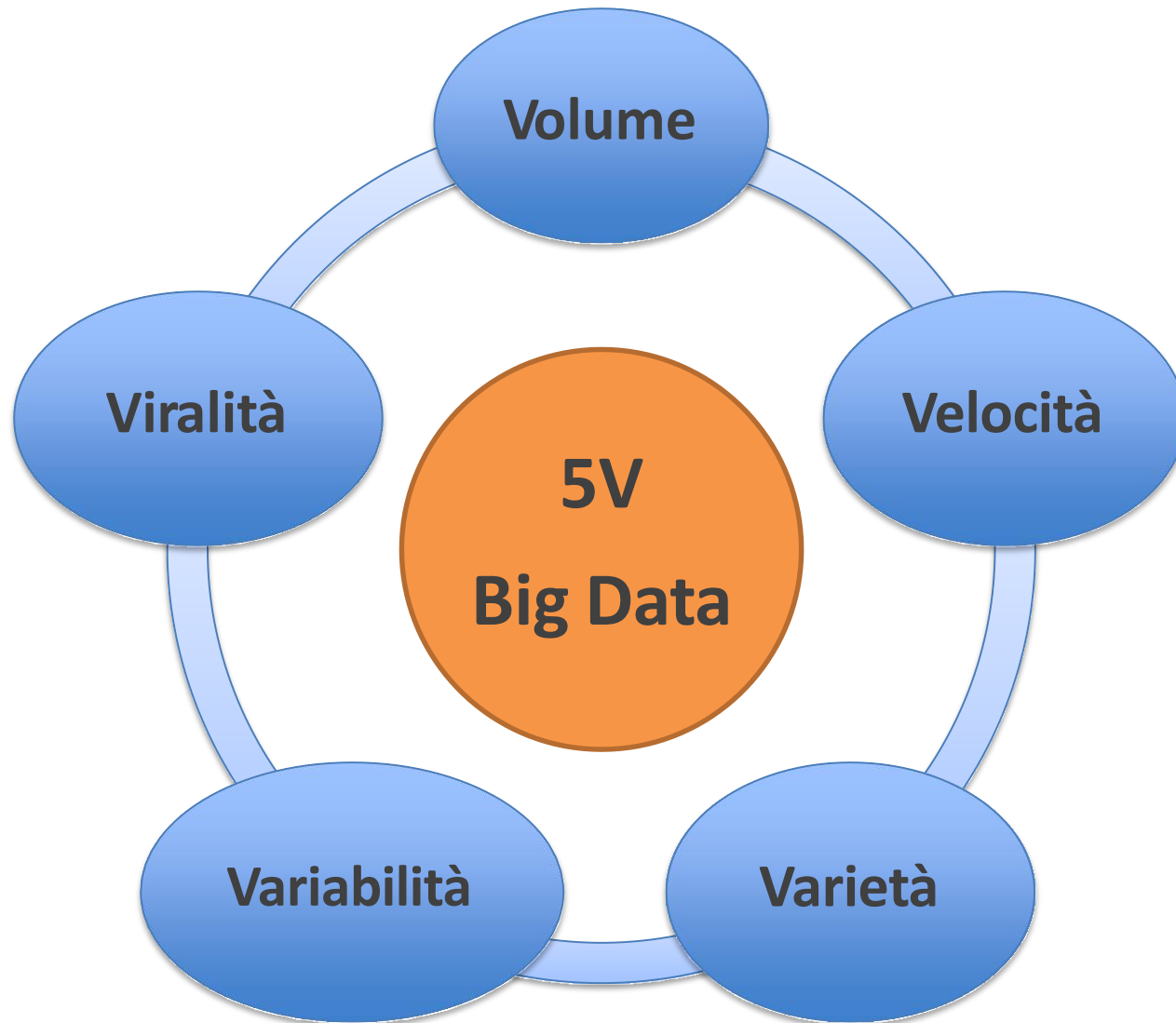
IBM

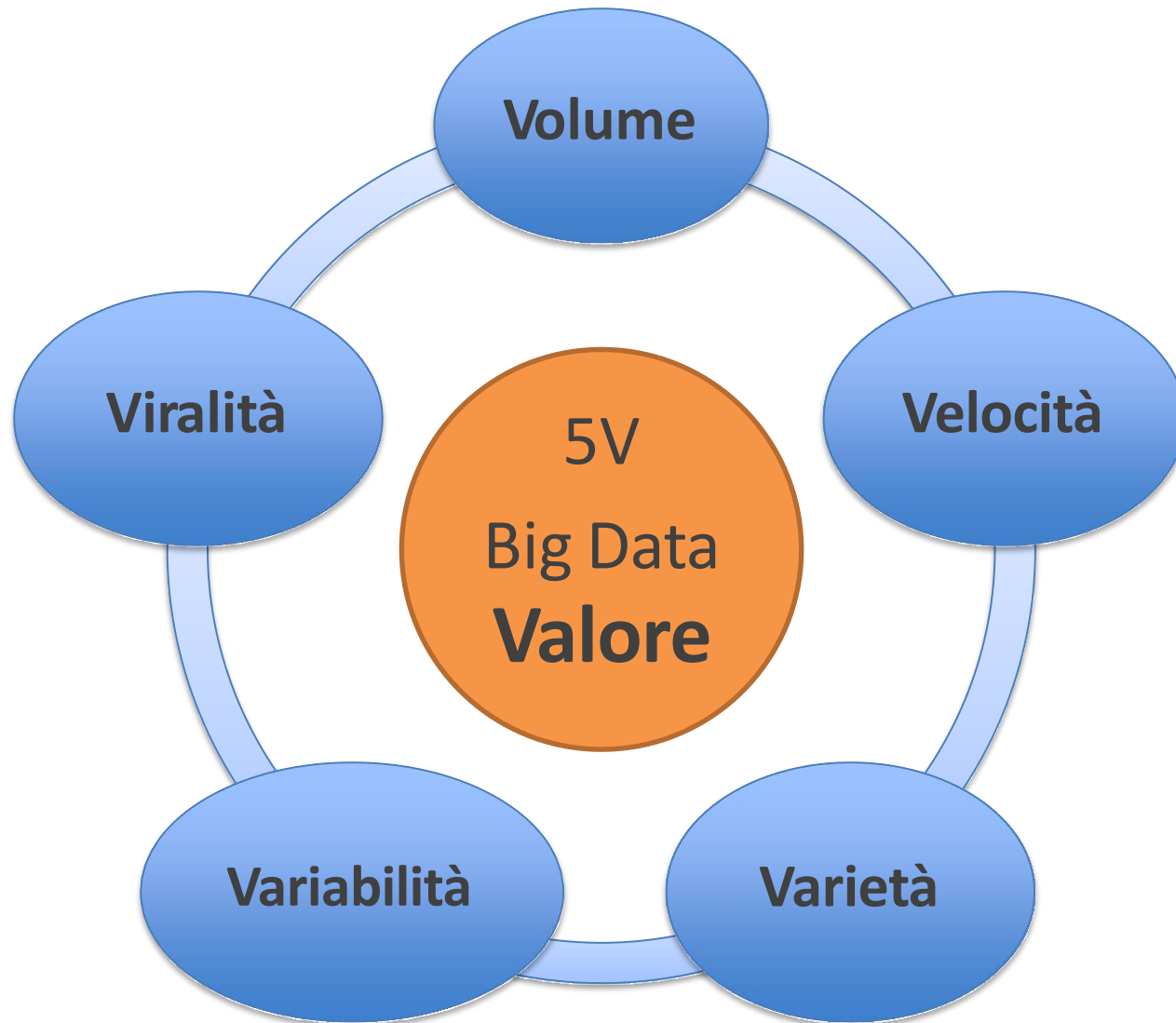
Evoluzione dei dati e delle tecniche di analisi



I Big Data sono dati che superano i limiti degli strumenti tradizionali.

- Sono ***dati*** solitamente disponibili in ***grandi volumi***, che si presentano in ***differenti formati*** (spesso privi di struttura) e con ***caratteristiche eterogenee***, prodotti e diffusi generalmente con una ***elevata frequenza***, e che ***cambiano spesso nel tempo***.
- Per questo motivo sono identificati con le **5V (+1)**.





Big Data: le 5 V - Volume

Volume: forse la caratteristica più immediata, dal momento che si tratta di dati presenti in **grandi quantità**. In **1 minuto** infatti:

- **100 mila tweet** trasmessi nel mondo.
 - **35 mila "Like" FB** a siti ufficiali di organizzazioni.
 - **160 milioni (circa) di email** inviate.
 - **2 mila check-in su 4square** effettuati.
-
- Ciò va aggiunto alle restanti **"attività digitali"**, generando una enorme mole di dati e informazioni a loro volta incrociabili.
 - Aziende, martketeers, analisti (ma anche la politica) sono le figure più ingolosite dalle potenzialità di tutto ciò.

Big Data: le 5 V - Volume

- Alcune tipologie di Big Data sono **transitorie**:
 - Dati generati da sensori.
 - Log dei web server.
 - Documenti e pagine web.
- Il **primo passo** quando si opera con i Big Data é allora l'**immagazzinamento**. L'**analisi (e la pulizia)** avvengono in una fase successiva (per evitare di perdere potenziali informazioni).
- Ciò richiede **importanti investimenti in termini di storage** e di **capacità di calcolo** adatta all'analisi di grandi moli di dati.
- Tecnologia open source più diffusa e utilizzata: **Apache Hadoop**.

Big Data: le 5 V - Velocità



Big Data: le 5 V - Velocità

Velocità: è una caratteristica che ha più di un significato.

- Si riferisce in primis alla **elevata frequenza con cui i dati vengono generati** – si ripercuote sulla quantità (Volume).
- Il secondo aspetto riguarda la **velocità con cui le nuove tecnologie** permettono di **accedere e di analizzare questi dati**.

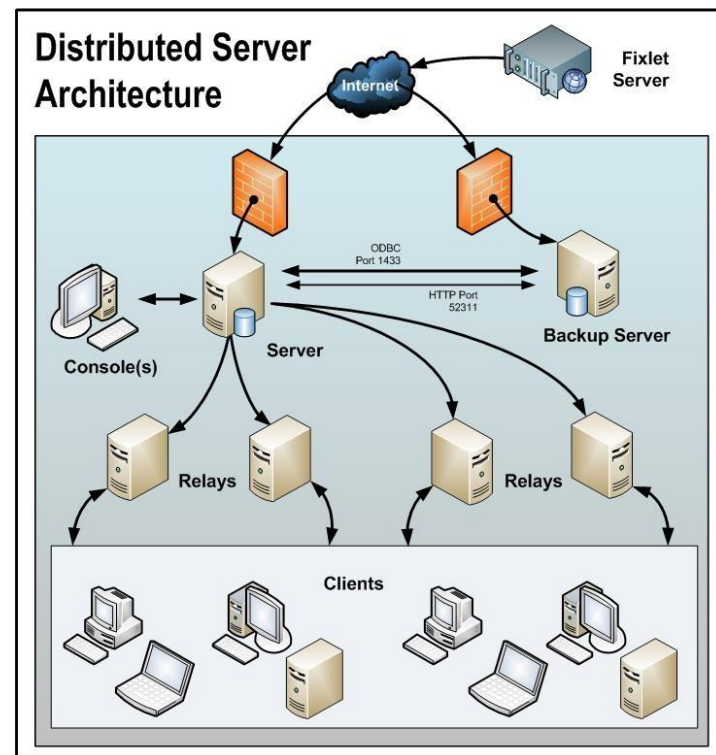
Maggiore è la velocità di accesso ai dati
Maggiore sarà la velocità in un processo decisionale
Maggiore/migliore competitività sui diversi panorami del mercato

Quali tecnologie?!

Big Data: le 5 V - Velocità

Velocità: è una caratteristica che ha più di un significato.

- Particolarmente adatte sono le **architetture distribuite**.
- Gestione di strutture dati anche complesse.
- Accesso ai dati in tempo reale.
- Velocità di elaborazione grazie a tecniche di calcolo distribuito.
- Database non relazionali come i *column DB* e *key/value DB* (NoSQL).



Big Data: le 5 V - Varietà

Varietà: caratteristica che ha a che fare con la forma in cui i dati si presentano.

- Nel contesto **Big Data** le **informazioni** da trattare sono dati non-strutturati (o semi-strutturati). Non adatti ad essere lavorati con le tecniche tradizionali dei database relazionali.
- Dati come **email, immagini, video, audio, stringhe di testo** a cui dare un significato non si possono memorizzare in una tabella.
- Per la gestione e il salvataggio di questi dati si ricorre spesso ai **database NoSQL**. Non impongono uno schema rigido per organizzare i dati (*schemaless database*).

Big Data: le 5 V - Variabilità

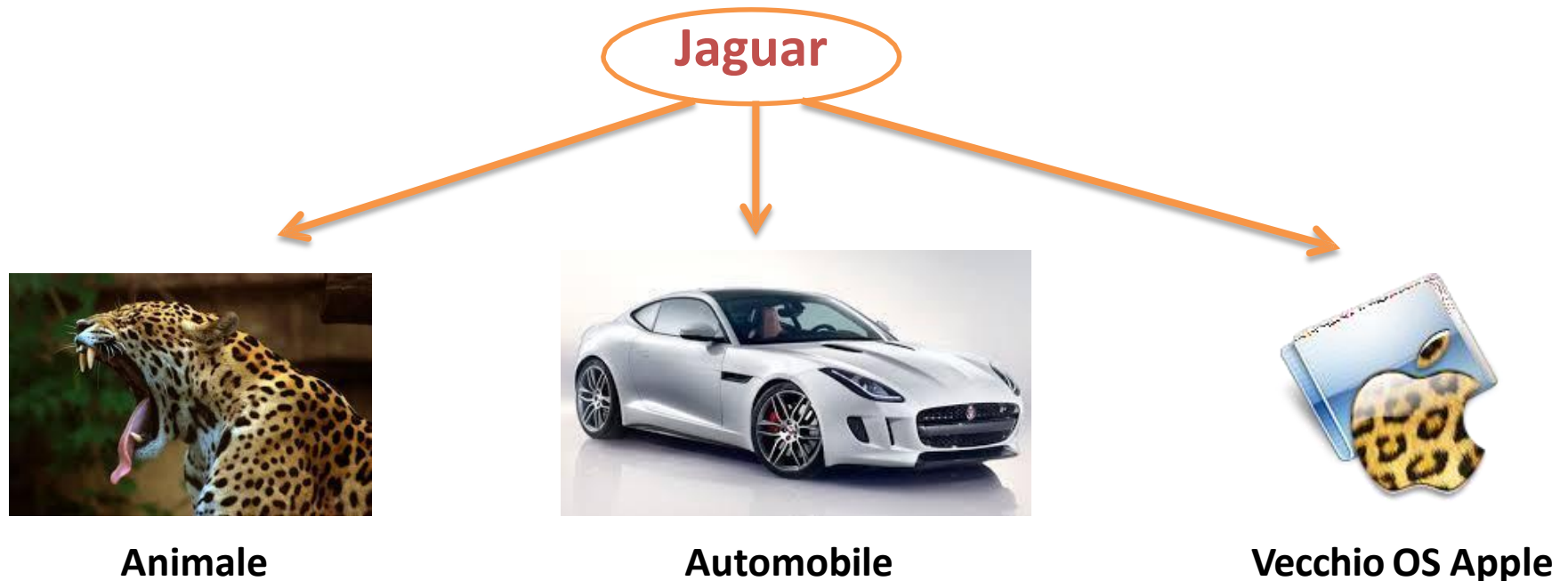
Variabilità: caratteristica relativa alla **contestualizzazione** di un dato.

- Il significato o l'**interpretazione di uno stesso dato** può **variare** in base al contesto in cui esso viene raccolto e analizzato.
- Esempio la frase "*leggete il libro*", essa avrà un **significato positivo** in un blog che parla di letteratura, mentre avrà una **connotazione negativa** in un blog per appassionati di cinema.
- Il **significato** di un dato **può essere differente** anche in base al **momento in cui viene fatta l'analisi**, spesso è fondamentale l'analisi in tempo reale (Velocità).

Big Data: le 5 V - Variabilità

Variabilità: caratteristica relativa alla **contestualizzazione** di un dato.

E' importante trovare dei meccanismi che riescano a dare una **semantica ai dati** in base al contesto in cui sono espressi.



Big Data: le 5 V - Viralità

Viraltà: caratteristica che ha a che fare su **quanto e come i dati si diffondono** (Propagazione dei dati).

- La **grande quantità** di dati (spesso correlati tra loro) e l'**alta velocità con cui sono prodotti** implica una **diffusione virale** delle informazioni.
- **Esempio: una notizia o un evento** diffusi tra diversi canali. **Diffusione amplificata** con i collegamenti nei vari **social network**.



Big Data: le 5 V - Viralità

Istituzioni e alcune organizzazioni sfruttano questa caratteristica/potenzialità per migliore attività di pronto intervento.



Big Data: le 5 V - Viralità

Virale è anche la **crescita del Volume** dei dati generati dalle attività digitali dell'uomo (user-generated content):

- Nel **2010** è stata stimata una produzione di **1,2 zettabyte** di dati (**1ZB corrisponde a mille miliardi di GB**).
- Nel **2011** è cresciuta a **1,8ZB**.
- Nel **2013** si è arrivati a **2,7ZB**.
- La proiezione per il **2015** parla di **8ZB**.

Big Data: le 5 V – Viralità (Curiosità)

Buzzsumo, una società di analisi dati, ha analizzato recentemente *milioni di contenuti in rete* per capire quali sono le **caratteristiche che rendono un contenuto virale**.

Dimensioni

- **Maggiore è la lunghezza** dei contenuti, **maggiori saranno le condivisioni**. Contenuti lunghi e ricchi di informazioni (**> 2000 parole**) ottengono più share rispetto ai contenuti brevi (**< 2000 parole**).

Emozioni

- Un contenuto deve *generare emozioni*, **le persone amano condividere elementi che facciano ridere e stupiscano i lettori** (42% dei contenuti studiati). Di contro, le emozioni **meno gradite** sono la **tristezza e la paura**, che arrivano al 7%.

Big Data: le 5 V – Viralità (Curiosità)

Buzzsumo, una società di analisi dati, ha analizzato recentemente *milioni di contenuti in rete* per capire quali sono le **caratteristiche che rendono un contenuto virale**.

Immagini

- I contenuti visivi attirano l'attenzione degli utenti, favoriscono una **comprensione immediata** e quindi tendono ad avere **maggiori interazioni**. Per la loro natura, **le immagini aumentano le condivisioni sui social**.
- Il **65%** delle persone usa **Facebook** per condividere **post che contengano almeno un'immagine**. Questi post sono quelli con cui poi si interagisce maggiormente.
- **Più del 20%** degli utenti su **Twitter** preferisce **pubblicare contenuti in cui sia presente un'immagine**.

Big Data: le 5 V – Viralità (Curiosità)

Buzzsumo, una società di analisi dati, ha analizzato recentemente *milioni di contenuti in rete* per capire quali sono le **caratteristiche che rendono un contenuto virale**.

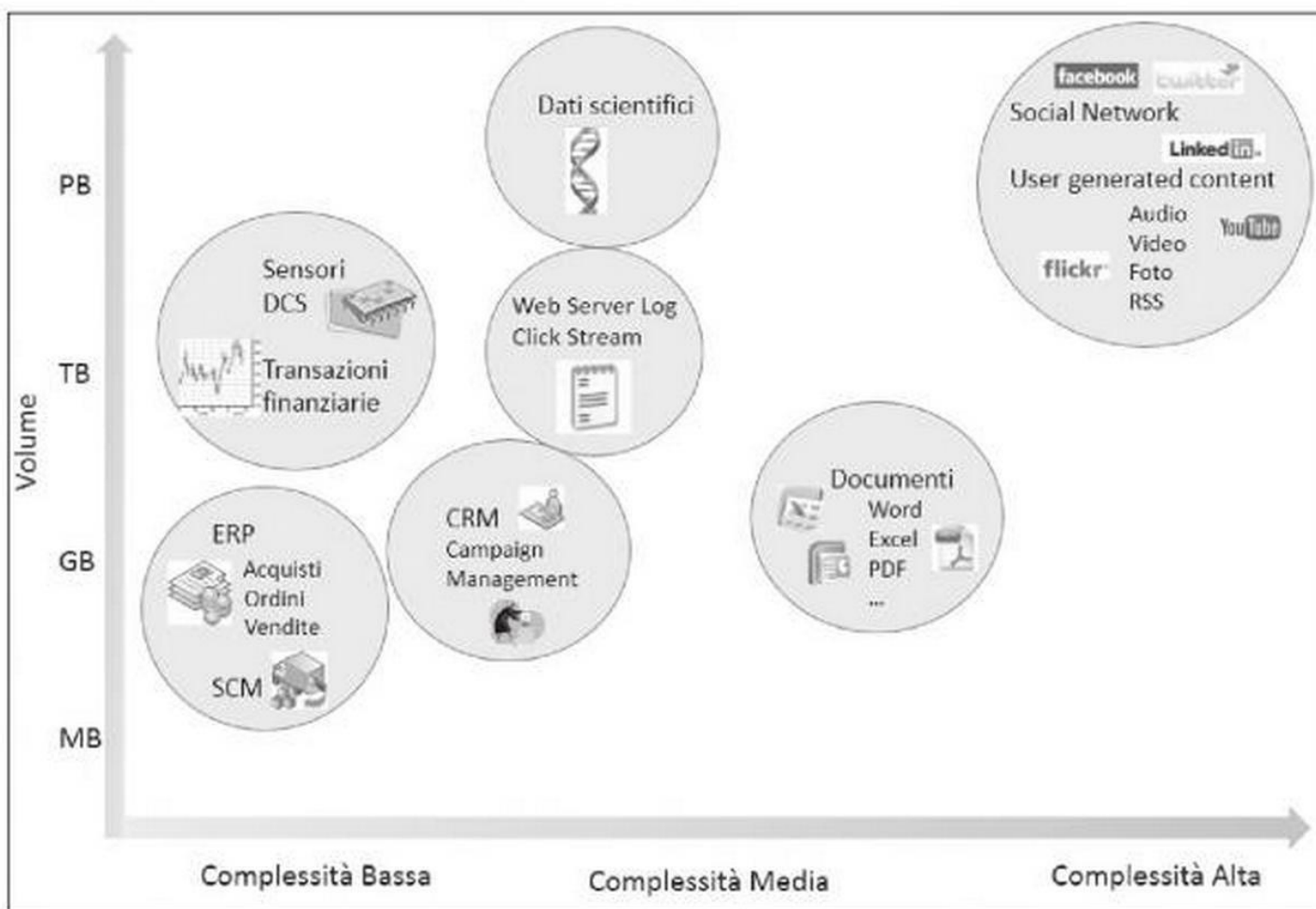
Elenchi puntati

- **Lettori e utenti web amano gli elenchi puntati, le infografiche e gli how to.** Ciò dipende dal fatto che questi contenuti permettono di **sintetizzare in forma visiva gli aspetti salienti di un post**, facilitando la comprensione.

Influencer

- **Contenuti condivisi da persone, organizzazioni e aziende ritenuti “esperti” in uno specifico settore** raggiunge un maggior numero di **utenti “targettizzati” e interessati a determinate informazioni**.

Big Data: le 5 V



Classificazione dei dati per volume e complessità

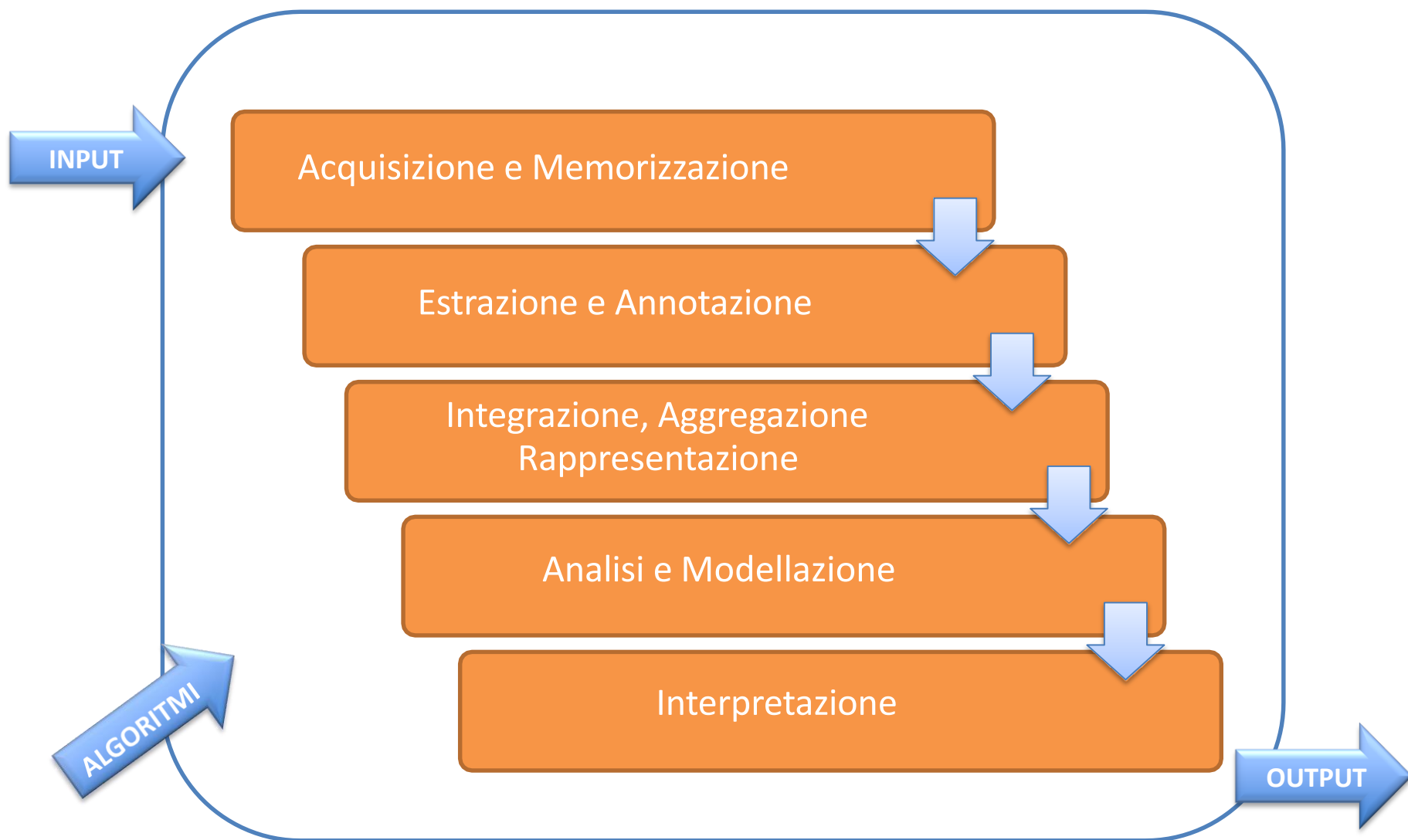
Valore: è necessario comprendere e gestire in modo adeguato i dati e tutti questi aspetti ad essi legati in modo da riuscire ad **estrarre il potenziale informativo**.

- I Big Data **nascondono un grande valore**. Al primo utilizzo di solito se ne estrae soltanto una parte, il *valore rimanente rimane “dormiente”* fino ad un successivo utilizzo.
- E' quindi importante adottare metodologie e tecnologie che permettano la **continua integrazione di nuove informazioni**, in seguito ad un utilizzo reiterato, con l'obiettivo di **costruire una base di conoscenza sempre più ampia**.

Problematiche

- **Elevato numero di campi applicativi** diversi tra loro.
 - I **differenti canali** attraverso i quali i dati vengono raccolti.
 - Identificare una possibile **architettura adattabile** a tutte le aree.
 - Come è possibile scoprire il “**Valore**” dei Big Data?
- **Utilizzo di complesse analisi e processi di modellazione.**
 - **Formulazione di ipotesi -> implementazione di modelli semantici, visuali e statistici -> validazione.**

Pipeline di Analisi dei Big Data



Pipeline: Acquisizione dei Dati e Memorizzazione

- Grandi quantità di dati possono essere **filtrati e compressi** a diversi ordini di grandezza.
 - Sfida: **Definire dei filtri opportuni in modo che non vadano perse informazioni di interesse.**
- **Dettagli** inerenti a **condizioni sperimentali e procedure** possono essere richiesti per interpretare i risultati correttamente.
 - Sfida: **Generazione automatica dei metadata corretti.**
- Possibilità di ricerca sia all'interno dei metadata che nei dati di sistema.
 - Sfida: **Creare e utilizzare delle strutture dati ottimizzate che consentano le ricerche in tempi accettabili.**

Pipeline: Estrazione delle Informazioni e Pulizia

- Le **informazioni** raccolte spesso **non** sono in un **formato pronto per l'analisi** (es. immagini di sorveglianza VS immagini scattate da fotografi)

Sfida: Realizzare un processo di estrazione delle informazioni che le fornisca in un formato adatto alla fase di analisi.

- I Big Data sono **incompleti** a causa di **errori** commessi durante la fase di acquisizione.

Sfida: Definire dei vincoli e modelli per la gestione e correzione automatica di errori in diversi domini Big Data.

- **I Dati sono eterogenei** e può non essere abbastanza raccogliarli all'interno di repository.

Sfida: Creare delle strutture dati di memorizzazione che siano in grado di adattarsi alle differenze nei dettagli sperimentali.

- **I modi di memorizzare dati sono diversi**, alcuni modelli hanno dei vantaggi rispetto ad altri per determinati scopi.

Sfida: Creare dei tool di supporto al processo di progettazione dei database e alle tecniche di sviluppo tenendo conto del contesto applicativo e d'uso dei dati.

Pipeline: Query, Modellazione Dati e Analisi

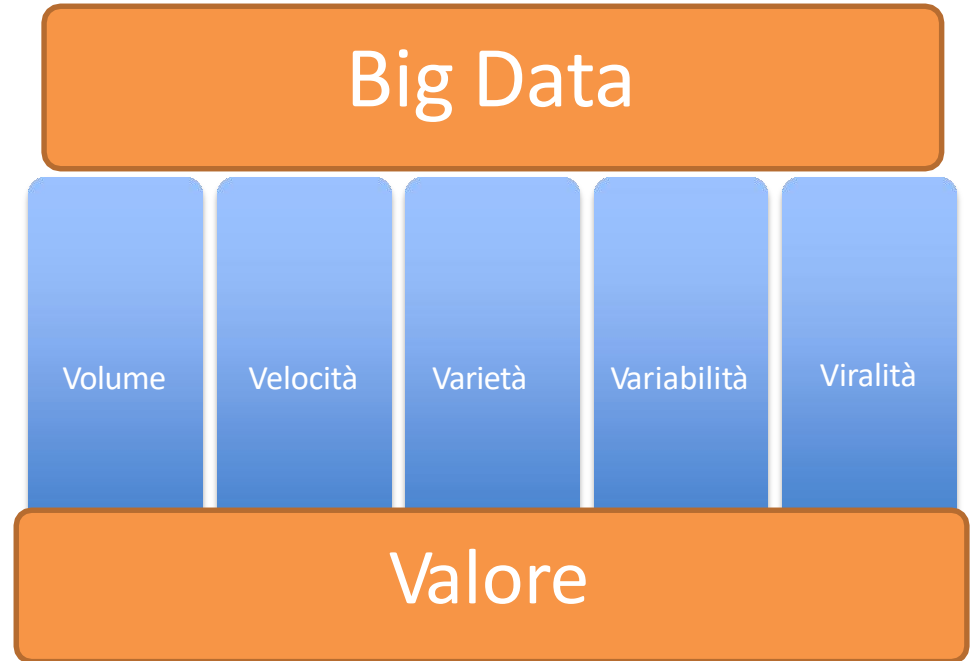
- I **metodi per investigare e interrogare** i Big Data sono **differenti** dalle tradizionali analisi statistiche.
Sfida: Creare delle tecniche per l'elaborazione di query complesse e scalabili (sull'ordine dei TeraByte), considerando delle risposte interattive nel tempo.
- **Big Data interconnessi** formano delle **reti di dati eterogenee**, in cui la **ridondanza dei dati** può essere sfruttata per compensare l'assenza di alcune informazioni, per verificare situazioni di conflitto e evitare che ci siano relazioni nascoste.
Sfida: Rendere coordinati i sistemi DB e le interrogazioni SQL, con i tool di analisi che realizzano diverse forme di elaborazione non-SQL (data mining, analisi statistica).

Datification

- Prendere **informazioni su qualsiasi cosa e trasformarle in un qualsiasi formato dati** in modo da renderle **quantificabili**.
- **Utilizzare queste informazioni** in un nuovo modo con l'obiettivo di **tirar fuori il loro valore implicito e nascosto**.
- *Quando i dati sono pochi è desiderabile che siano accurati (campionamento random). I Big Data hanno cambiato il concetto di aspettativa della precisione: Trattare queste grandi quantità di dati spesso imprecise e imperfette permette di fare delle previsioni superiori (Analisi Predittiva).*

Campi di Applicazione


- Il Problema dei Big Data si riferisce alla combinazione di un **grande volume di dati** che deve essere **trattato in tempi abbastanza rapidi**.



- Sono molte aree applicative in cui i Big Data sono attualmente utilizzati con risultati interessanti ed **eccellenti prospettive future** per affrontare le **principali sfide** come **Analisi dei Dati**, **Modellazione**, **Organizzazione** e **Ricerca** (Data Retrieval).

Campi di Applicazione

Investimenti crescenti nei Big Data possono portare fondazioni, enti e organizzazioni di nuova generazione a interessanti **scoperte** *in campo scientifico, nella medicina, vantaggi e guadagni nel settore ICT e in contesti Business, nuovi servizi e opportunità per cittadini digitali e utenti web.*

- Sanità e Medicina
- Ricerca Scientifica (Analisi dei Dati)
- Istruzione
- Settore Energetico e dei Trasporti
- **Social Network – Servizi Internet – Web Data** 
- Finanza/Business – Marketing
- Sicurezza

Social Network Big Data

2012



▪ **Facebook:** più di 10 milioni di foto caricate ogni ora, 3 miliardi di “like” e commenti ogni giorno.



▪ **Google Youtube:** 800 milioni di utenti caricano ~1h di video ogni secondo.



▪ **Twitter:** più di 400 milioni di tweet ogni giorno.



▪ **Instagram:** 7.3 milioni di utenti univoci ogni giorno.

2009

▪ **Facebook:** 3 milioni di foto caricate ogni mese, implementazione del pulsante “like”

▪ **Google Youtube:** tutti gli utenti caricavano 24h di video ogni minuto.

▪ **Twitter:** 50 milioni di tweet ogni giorno.

▪ **Instagram:** è stato creato nel 2010.

Social Network – Servizi Internet – Web Data

- Il **volume di dati** generati dai **servizi internet**, **siti web**, **applicazioni mobili** e **social network** è grande, la velocità di produzione è invece variabile, a causa del **fattore umano**.
- Da queste grandi quantità di dati raccolti in particolare attraverso i social network, aziende e ricercatori cercano di **prevedere il comportamento collettivo** e analizzare i **trend topic**.
- Ad esempio attraverso il **monitoraggio degli hashtag (#)** di Twitter è possibile identificare dei **modelli di influenza**.
- In senso più ampio da tutte queste informazioni è possibile **estrarre conoscenza** e evidenziare le **relazioni tra i dati**, in modo da migliorare l'**attività di query-answering**.

Alcuni ricercatori hanno proposto un **utilizzo alternativo** di tali dati per creare una **nuova forma di vivibilità urbana**, in un'iniziativa/progetto chiamato ConnectiCity.

Gli aspetti chiave sono:

- Creare un **set di tool** per catturare **in real-time** differenti forme di **contenuti** rilevanti **generati dai cittadini/utenti**, provenienti da diverse tipologie di sorgenti:
 - *Social network*
 - *Siti web*
 - *Applicazioni mobile*

Alcuni ricercatori hanno proposto un **utilizzo alternativo** di tali dati per creare una **nuova forma di vivibilità urbana**, in un'iniziativa/progetto chiamato ConnectiCity.

Gli aspetti chiave sono:

- **Mettere in relazione questi contenuti al territorio** utilizzando tecniche di *Geo-Referencing*, *Geo-Parsing* e di *Geo-Coding*. **Analizzarli e classificarli** utilizzando tecniche di *Natural Language Processing* per identificare:
 - *Topic di interesse*
 - *Espressioni emozionali e sentimenti*
 - *Analisi della rete per capire la propagazione delle informazioni e i modelli di comunicazione*

Alcuni ricercatori hanno proposto un **utilizzo alternativo** di tali dati per creare una **nuova forma di vivibilità urbana**, in un'iniziativa/progetto chiamato ConnectiCity.

Gli aspetti chiave sono:

- Rendere queste **informazioni disponibili e accessibili** sia a livello centrale e periferico, per consentire la creazione di **nuove forme di processi decisionali**, nonché di sperimentare modelli innovativi di partecipazione, peer to peer, iniziative generate dai cittadini/utenti.

Project link - <http://www.connecticity.net/>
<http://www.opendata.comunefi.it>

Social Network – Curiosità

- L'università di Cambridge ha pubblicato uno studio nel 2013 dove emerge come **è possibile descrivere i tratti della personalità dalla semplice analisi dei Like degli utenti di Facebook.**

<http://www.pnas.org/content/early/2013/03/06/1218772110.full.pdf>

- Tra gli **attributi analizzabili** troviamo:
 - *Orientamento Politico*
 - *Orientamento Sessuale*
 - *Orientamento Religioso*
 - *Aspetti caratteriali*
 - *Livello di soddisfazione della propria vita*


Social Network – Curiosità

- Il modello proposto può essere applicato a qualsiasi **insieme di dati** in grado di **esprimere una preferenza dell'utente**.
- I dati di Facebook sono pubblici e con **Facebook Connect**, sono facilmente ottenibili, sempre previa autorizzazione da parte dell'utente.
- L'algoritmo è stato sviluppato implementato da una start up Italiana, **Cube You** e può essere testato all'indirizzo:

<http://youarewhatyoulike.com/>

Campi di Applicazione

Investimenti crescenti nei Big Data possono portare fondazioni, enti e organizzazioni di nuova generazione a interessanti **scoperte** *in campo scientifico, nella medicina, vantaggi e guadagni nel settore ICT e in contesti Business, nuovi servizi e opportunità per cittadini digitali e utenti web.*

- Sanità e Medicina
- Ricerca Scientifica (Analisi dei Dati)
- Istruzione
- Settore Energetico e dei Trasporti
- Social Network – Servizi Internet – Web Data
- **Finanza/Business – Marketing** 
- Sicurezza

Finance/Business e Marketing

- Il compito di **trovare modelli nei dati aziendali** non è nuovo. Tradizionalmente gli analisti di business usano tecniche statistiche.
- Oggi l'uso diffuso di **PC e tecnologie di rete** ha creato grandi **repository elettronici** che memorizzano numerose transazioni commerciali.
 - La grandezza di questi dati varia tra **50-200 PBs** al giorno
 - Gli **accessi ad Internet** in Europa sono circa **381 milioni di visitatori unici**.
 - **40% dei cittadini Europei** fa shopping online.
- Questi dati possono essere analizzati per definire:
 - **Previsioni sul comportamento degli utenti.**
 - **Identificare modelli di acquisto di clienti individuali o gruppi.**
 - **Fornire nuovi servizi personalizzati.**

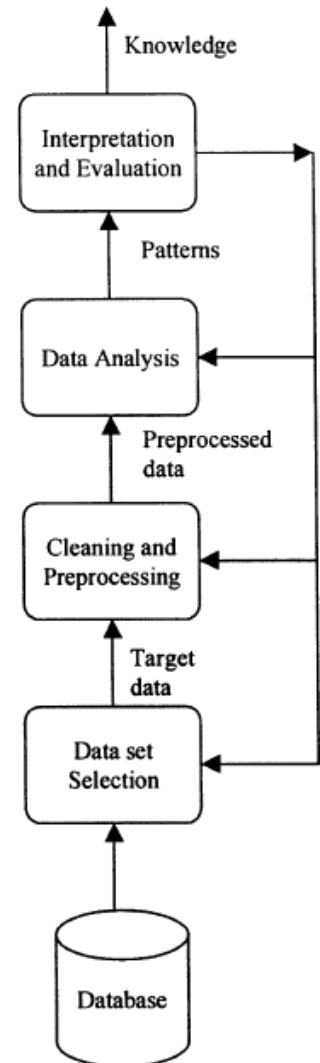


Con l'uso di tecnologie di data warehousing e tecniche di apprendimento automatico mature.

Finance/Business e Marketing

In campo finanziario, invece, si possono creare **piani di investimento e di business grazie a modelli predittivi** ottenuti con tecniche di ragionamento o per scoprire modelli interessanti e significativi dai dati aziendali.

1. **Selezione dei dati per l'analisi** (da una rete di DB).
2. **Operazioni di raffinamento** per rimuovere discrepanze e inconsistenze.
3. I **dati sono analizzati** per identificare dei **Pattern** (modelli che mostrano la relazione tra i dati).
4. Dovrebbe essere possibile la **traduzione del modello in business plan praticabile**, che aiuti l'azienda a raggiungere il suo obiettivo.
5. Modelli/Pattern che soddisfano queste condizioni diventano **business knowledge**.




▪ **GoodData**, società di San Francisco, che fornisce una *piattaforma con un insieme di tool di BI*. L'obiettivo è di supportare le aziende ad analizzare la loro **enorme mole di dati** (indagini di mercato, resoconto vendite, costi etc.) per *favorire il processo decisionale*.



Campi di Applicazione

Investimenti crescenti nei Big Data possono portare fondazioni, enti e organizzazioni di nuova generazione a interessanti **scoperte** *in campo scientifico, nella medicina, vantaggi e guadagni nel settore ICT e in contesti Business, nuovi servizi e opportunità per cittadini digitali e utenti web.*

- Sanità e Medicina
- Ricerca Scientifica (Analisi dei Dati)
- Istruzione
- Settore Energetico e dei Trasporti
- Social Network – Servizi Internet – Web Data
- Finanza/Business – Marketing
- **Sicurezza** 

- **Intelligence, Sorveglianza, e Recognition (ISR)** definiscono argomenti che sono adatti per analisi computazionali di tipo data-centrico.
 - Dimensione vicina ad uno **zettabyte (10^{21} bytes o un miliardo di TeraByte)** di dati digitali sono generati ogni anno.
- Importanti fonti di dati per i sistemi di intelligence sono
 - **Immagini satellitari e aeree (da veicoli UAV).**
 - **Comunicazioni intercettate:** civili e militari, tra cui voce, e-mail, documenti, i registri delle transazioni (log) e altri dati elettronici – **5 miliardi di telefoni cellulari in uso in tutto il mondo.**
 - **Dati di tracciamento radar.**

- Importanti fonti di dati per i sistemi di intelligence sono
 - **Sorgenti di dominio pubblico** (siti web, blog, tweet e altri dati Internet, televisione, carta stampata e radio).
 - **Dati di Sensori** (dai meteorologici, oceanografici, riprese di telecamere di sicurezza).
 - **Dati Biometrici** (Immagini facciali, DNA, impronte digitali, scansioni dell'occhio, registrazione del portamento).
 - **Informazioni strutturate e semi-strutturate** fornite da aziende e organizzazioni: *log delle compagnie aeree, carte di credito e transazioni bancarie, registrazioni telefoniche, elenco del personale dipendente, cartelle cliniche elettroniche, rapporti investigativi e dati nei registri di polizia.*

- La sfida per i servizi segreti è quello di **trovare, combinare e definire modelli** e tendenze nelle tracce di informazioni ritenute importanti.
- Occorre trovare **modelli di evoluzione significativi in modo tempestivo** tra diverse informazioni potenzialmente offuscate provenienti da fonti multiple. Necessità di metodi sofisticati per individuare modelli accurati, **senza generare un gran numero di falsi positivi** in modo che non emergano cospirazioni o allarmi dove non esistono.
- *Un esempio di come sfruttare efficacemente fonti che producono dati su larga scala, sono le principali aziende in ambito web, come Google, Yahoo e Facebook.*

- All'interno del mondo dei servizi di Intelligence le **tecnologie informatiche** e le consolidate **tecniche di apprendimento automatico** devono essere considerate come un elemento per aumentare le capacità degli analisti piuttosto che come un modo per sostituirli.
- L'idea chiave dell'apprendimento automatico è:
 - **Applicare** ad un certo dataset **prima un'analisi statistica strutturata** al fine di generare un ***modello predittivo***.
 - **Poi applicare questo modello a diversi flussi di dati** per **supportare diverse forme di analisi** e ottenere nuovi risultati.

- Nel dicembre 2012 il comune di **Philadelphia** ha rilasciato un **dataset con l'elenco dei crimini** dal 1° gennaio 2006.
- Ogni **crimine** (*furto, rapina, omicidio...*) è **taggato nella posizione esatta** in cui è stato commesso.
- Con questi **dati** è possibile la **creazione di tool e statistiche** utili sia al cittadino che alla pubblica amministrazione.



Criticità e rischi dei Big Data

Come ogni “nuova tecnologia” i Big Data offrono grandi prospettive e potenzialità, ma non presentano esclusivamente caratteristiche positive. Vi sono alcuni aspetti critici che è bene prendere in considerazione:

- Problematiche legate alla **qualità e all'affidabilità dei dati**.
- Problematiche relative alla **privacy e alla proprietà dei dati**.

Criticità e rischi dei Big Data – Qualità dei dati

La **qualità dei dati** è determinata da un insieme di caratteristiche:

- **Completezza:** la presenza di **tutte le informazioni necessarie** a descrivere un oggetto, entità o evento (es. anagrafica).
- **Consistenza:** i dati **non devono essere in contraddizione**. Ad esempio il saldo totale e movimenti, disponibilità di un prodotto richiesto da soggetti differenti, etc.
- **Accuratezza:** i dati devono essere corretti, cioè **conformi a dei valori reali**. Ad esempio un indirizzo mail non deve essere solo ben formattato **nome@dominio.it**, ma deve essere anche valido e funzionante.

Criticità e rischi dei Big Data – Qualità dei dati

La **qualità dei dati** è determinata da un insieme di caratteristiche:

- **Assenza di duplicazione:** Tabelle, record, campi dovrebbero essere memorizzati **una sola volta**, evitando la presenza di copie. Le informazioni duplicate comportano una doppia manutenzione e possono portare problemi di sincronia (consistenza).
- **Integrità:** è un concetto legato ai database relazionali, in cui sono presenti degli strumenti che permettono di implementare dei **vicoli di integrità**. Esempio un **controllo sui tipi di dato** (presente in una colonna), o sulle chiavi identificative (impedire la presenza di due righe uguali).

Criticità e rischi dei Big Data – Qualità dei dati

Nei contesti applicativi che coinvolgono l'uso di database tradizionali, la **qualità complessiva dei dati** può essere minata da:

- **Errori nelle operazioni di data entry** (campi e informazioni mancanti, errati o malformati).
- **Errori nei software di gestione dei dati** (query e procedure errate).
- **Errori nella progettazione delle basi di dati** (errori logici e concettuali).

Criticità e rischi dei Big Data – Qualità dei dati

Nel mondo Big Data invece:

- **Dati operazionali:** i problemi relativi alla qualità sono conosciuti e esistono diversi strumenti per realizzare in modo automatico la pulizia dei dati.
- **Dati generati automaticamente:** i dati scientifici o provenienti da sensori sono privi di errori di immissione. Spesso però sono “deboli” a livello di contenuto informativo, c’è la necessità di integrarli con dati provenienti da altri sistemi per poi analizzarli.
- **Dati del Web:** Social network, forum, blog generano dati semistrutturati. La parte più affidabile sono i metadati (se presenti), il testo invece è soggetto a errori, abbreviazioni, etc

Criticità e rischi dei Big Data – Qualità dei dati

Nel mondo Big Data invece:

- **Disambiguare le informazioni:** Uno stesso dato può avere significati differenti (es. calcio). La sfida è cerca di trovare quello più attinente al contesto in esame. Un aiuto sono i **tag**, etichettando i dati si cerca di evidenziare l'ambito di pertinenza.
- **Veridicità:** Notizie, affermazioni, documenti non sempre veri o corrispondenti alla realtà.

OSS. *La qualità dei dati è però legata anche al contesto in cui essi sono analizzati. Operazioni di filtraggio e pulizia devono essere fatte procedendo per gradi per evitare di eliminare dati potenzialmente utili.*

Criticità e rischi dei Big Data – Privacy

Il tema Big Data si apre a problemi di Privacy, proprietà e utilizzo dei dati da parte di terzi.

- **Dati del Web:** gli *user-generated-content* sono condivisi accessibili a tutti. E' etico il loro utilizzo?
- **Dati sensibili:** i dati presenti nei DB degli ospedali relativi alla storia clinica dei pazienti sono opportunamente protetti?
- **Dati di posizione:** l'uso di smartphone, GPS, sistemi di pagamento elettronico, ma anche social network lasciano delle tracce da cui è possibile ricavare gli spostamenti degli utenti.