

Big Data

We have big data today. But we can't simply store and analyze it.

Reasons: Large amount of storage and access speed.

Storing data in multiple drives equally distributed and using them in parallel can make it easy, fast and possible.

Types of Data:

~~Structured Data~~ Data can be stored in

Structured Data	Semi-structured	Unstructured
<ul style="list-style-type: none">• They have relational key.• Matured transaction control• eg: Relational database• Arranged in rows and cols• Flexible: Schema dependent but less flexible	<ul style="list-style-type: none">• They can be made into relational after some but not always.• Transaction is adopted from DBMS but not mature• XML files• There is no conventional format• Schema dependent ^{less} but more flexible than structured	<ul style="list-style-type: none">• They are not relational.• No transactions.• Word, PDF, Media, etc.• Contains text or binary. No format.• Schema less and most flexible.

<ul style="list-style-type: none"> Scalability: Difficult to scale DB because of schema. 	<ul style="list-style-type: none"> More scalable than structured data 	<ul style="list-style-type: none"> More scalable
<ul style="list-style-type: none"> Queries: Allows complex joins 	<ul style="list-style-type: none"> Queries using nodes 	<ul style="list-style-type: none"> Only text queries
<ul style="list-style-type: none"> Very robust 	<ul style="list-style-type: none"> not very spread 	

Vs of Big Data.

Velocity: Data is generating with high speed and need to be handled on time with speed. Late decisions leads to missing opportunities. e.g: 3.5B+ searches on google daily.

Variety: Data is of different types. It can be structured, unstructured or semi-structured.

Volume: Volume refers to the amount of data. We have to deal with huge amount of data. Eg: Mobile traffic in 2016 was 6.2 exabytes ^{per month}. It will be 40,000 exabyte exabyte of data in 2020.

Veracity: Data can be inconsistent, incomplete. Eg: Larger data set may be inconsistent can create ambiguity while less amount of data will be incomplete information or half info.

Value: Data itself has no value. We have to make it valuable.

Variability: How Fast the data changes or its shape changes. Eg: you are drinking the same drink^{but} but its taste changes.

Venue: Location of data.

Vocabulary: Semantics.

Harnessing Big Data.

OLTP

OLAP

- | | |
|---|---|
| o Online Transaction processing. | o online Analytical processing. |
| o Provides transactions based applications. | o Use different databases from different DBS to get insights. |
| o ATM transaction | o Netflix Recommendation system |
| o It provide security restrictions. | o Better security features. |
| o It cant be used for decision making. | o It cant be used for decision making. |
| o Read, write, delete operations. | o Keep data consistent. |
| o Data consistency, integrity. | o Handle large data. |

Real Time Analytical Processing \rightarrow RTAP

It process Large amount of data with less reaction time.

Predictive Analytics and Data Mining

Business Value

- More of a real time analysis.
- Large Datasets.
- Any type of data.
- Complex statistical analysis.
- Deals in Data mining techniques.
- Small to mid sized.
- Structured Data.
- Ad-hoc queries and reporting.

Application of BD:

- Recommendation system.
- SNA

Recommendation system:

- Track each visit of the user.
- Track content and metadata about the visited pages.
- Track ^{meta} data of the user.

eg: News Recommendation

- Can be based on current article
- User history.
- ~~CTR~~ Clickstream.

Clickstream Analyze: Analyze the collected data to get insights.

Clickstream Analytics: Collect data from visited pages and their content.

- In general, combination of text stream (Article) with Click stream.
- To make a good model that describes the user.

Why is it Hard?

- Cold start Problem
- Providing Accurate recommendations.

Cold start:

No existing data of the user to predict recommendations.

System Network Analysis: SNA

observe social and communication phenomena at a planetary scale.

o Small World problem:

- Average distance between two random users is 6.6 degrees.
- Facebook reported in 2016, it is 3.5 degrees of separation.