

Big data formats  
structured, unstructured, semi-structured

front end which gives data and show the data, and retrieve the data  
back end which store data e.g. Database  
column refers as fields  
rows refers as records

Page No.	7	6	19
Date	7	6	19

## Unit - I

### chapter 1. Big Data

- Data is raw facts which can be stored, processed, transmitted.
- Information is processed data or collection of organic data.
- Database is collection of related data, it is storage area, which can be stored storage related type data.
- Big data is large size, which can be complex type data.

- \* Data :-
- The quantities characters or symbols on which operations are performed by a computer which may be stored and transmitted in the form of electric signal and recorded on magnetic, optical or mechanical recording area.

- \* Big data :-
- Big data is also data but with a Huge size big data is term Huge to describe collection of data that is Huge in size and yet growing exponential with time such a data is also large and complex that none of the traditional data management able to store or process it effectively.

Structure & Unstructure data - fixed size, type.  
Unstructure data - not fixed size, type.

Database create nth no. of database you can create in system  
RDBMS not follows Big data  
nth no. of table we can create in database

Page No. \_\_\_\_\_  
Date 18 6 19

## \* Sources of Big Data :-

- Generation of Big Data → sharing of Big Data → utilization of Big Data

example : Social ad and digital search over internet  
Social media, industries, education, entertainment,  
business, etc.

11/6/2019

## \* Categories of Big Data :-

Big data can be found in three basic forms :-

(i) Structured Data

(ii) Unstructured Data

(iii) Semi-Structured Data

(i) Structured data :-

- Any data that can be stored, access and process in the form of fixed format is term as structured data.
- Over the period of time science have achieved greater success in developing techniques for working such type of data where the format well known in advance.
- It also help to derive value out of it.

example of structured data :-

Any Relational Database

for example :-

employee table database or student table database.

- However nowadays we are facing issues when size of such data grows to a huge extent where data can not be stored in a traditional fixed pattern.

semi-structured data - fixed location, position  
but not fixed size, type.

e.g. Google History, log files

Mongod is server and Mongo is client

Page No.

Date

- Advantage of structure data is easy retrieval, manipulation and storage of data.

### (ii) Unstructured Data :-

- Any data with unknown form or the structure is classified as unstructured data. e.g. Google search.
- In addition to the size to be huge unstructured data poses multiple size challenges in terms of processing or deriving value out of it.
- Typical example of unstructured data is heterogenous data source containing a combination of simple file, text, files, image, video, etc.
- Nowadays organization have a wealth of data available with them but unfortunately they don't know how to derive the value out of it since the data is in raw or unstructured format.

### (iii) semi-structured data :-

- semi-structured data can contain both the form of data we can see semi-structured data as a structure in form but it is actually not defined.
- example of semi-structured data :  
A data growth over years.

12/06/2019

imp \*

characteristics of Big Data :-

(three V's of Big Data)

(i) Volume → Large size

- (ii) Variety - complex type
- (iii) Velocity - growing rapidly

#### (i) Volume :-

- The name big data itself is related to a size which a enormous. Hence crucial factor is size of data plays very important role in determining value of data.
- Also whether the particular data can actually be considered as big data or not is dependent upon volume of data.
- Hence volume is one characteristic we needs to be consider while dealing with big data.

#### (ii) Variety :-

- The next aspect of big data is its variety. Variety refers to heterogenous sources and nature of data both structured and unstructured data previously spread sheet and databases where the only source of data considered by most of the applications.
- Nowadays data in form of image, photos, pdf, audio, etc. is also been consider in the analysis application.
- The variety of unstructure data passes certain issues of storage mining and analyzing data.

Segmentation and customization e.g. if we need to show the direction of particular city/area then you give a particular area direction like Google map

Page No.	
Date	

### (iii) Velocity

- The term velocity refers to the speed of generation of data.
- How fast data is generated and processed to meet the demand determines the real potential to in the data.
- big data velocity deals with the speed of data flow in form of sources like business processing application logs, network and social media site, sensor, etc.

The flow of data is massive and continuous.

13/06/2019

### \* Benefits of Big Data processing

- (i) Businesses can utilize outside intelligence while taking decision.
- (ii) improved customer service.
- (iii) an early identification of risk to the producer
- (iv) Better operational efficiency

14/06/2019

### \* Usages of Big Data

- (i) Visibility
- (ii) Segmentation and customization
- (iii) Innovation
- (iv) Aiding decision making
- (v) discover and analysing analyse information.

## \* Big data challenges :-

(i) policy and procedures

(ii) Access to data

(iii) technology and techniques

(iv)

(v) policy and proce

## \* Legacy system and Big Data :-

(i) structure of Big data

(ii) storage of Big data

(iii) Data processing

## (\* Big data Technologies :-

(i) new storage and processing technologies designed specifically for large unstructured data.

(ii) Parallel Processing

(iii) clustering (creating small-small parts)

(iv) Large Grid Environment

(v) High connectivity

(vi) cloud computing (external storage of data)

e.g. drive

Nosql is a term for unstructured data. Nosql follows BASE. Sql → Structured query language fixed data sql follows ACID. Commit means to save the data.

Page No.	
Date	22/06/19

## Q. Explain CAP Theorem

### NoSQL

- NoSQL follows BASE protocol of consistency
- BASE (Basically Available soft state Eventually consistent)
- BASE derived by theorem CAP theorem.
- Imp - CAP theorem stands for consistency, Availability, partition Tolerance
- Big data follows CAP Theorem but follows mysql follows BASE

### \* Advantages of NOSQL :-

- High scalability
- Manageability and Administration
- Low cost (we can use small servers to store data)
- flexible data model
- 

### \* Disadvantage Of NOSQL :-

- Maturity (it is not fully developed)
- Support (Not create in high company. Create by small ones)
- limited query capability (only complex query not handle)
- Administration (it is open source, so we can run individual)
- Expertise (They are not expert people, few people have knowledge about NOSQL)

24/06/2019

- J. ~~Op~~
- Read and write operation :-  
Read means you can view the data we can not edit the data.  
Write means you can modify, view, edit the data.

## Write Operation :-

$N \rightarrow$  No. of data copies

$R \rightarrow$  No. of data copies we Read

$W \rightarrow$  No. of data copies we write

### Example :-

If  $N \neq W$

(i)  $N = W$

- if we have 30 students data in your database but we can add one more student data in your database then here all the data is update and Refresh, if the data store starting then all data is update and refresh, if the data store in middle then previous data is refresh and remaining data refresh update, if the data store in bottom then all data is refresh and same is time consuming. Write operation is time consuming.

(ii)  $W = 1, N = R$

- Here, write one data and Read all the data; therefore, Read operation is time consuming.

(iii)  $N > W > R$

- We can update few data in particular time otherwise it gives lots of time. Not all updated at a same time. It is use in large amount of data.

## Read Operation :-

(i)  $R = 1$

- you are not sure data is first, last, second data any one data we can read.

(ii)  $R > 1$

- Now, we can read more than one data but not all.

Hence, Read the new and old data. It is time consuming.

Scale In	W+R means updated data / newly updated data means maximum it generates size of data	NOSQL supports tabular, column, etc key value, etc. format
Scale Out		
means	(iii) $N < W + R$	
data - store in Groups or small part.	Reading only that data which is new / update new data. Here only reads the new data it is not time consuming. It gives only required data to be Read or show.	
*	SQL vs. NOSQL :-	
Scale In		
means	parameter	SQL and NOSQL
increases size of database	(i) type	SQL supports all types of standard & type of data whereas
then store the data		types of standard & type of data whether format support.
but in an History	(ii) Development	SQL developing of NOSQL development
Scale Out		
it divides data into parts	(iii) Example	MYSQL, SQLDB, mongoDB, HBASE, Server, SQL, cassandra, marklogic, graphicdb, Oracle, SQLplus.
then store the data	(iv) Data storage model	It stores Rows, It stores particular and column type of data format.
	(v) Schemas	fixed dynamic
	(vi) Scalability	It is used in form It is used out in form scale.
	(vii) Support and transaction	ACID will be According CAP and BASE Format

(viii) Consistency → Strong (soft) Normal  
Centralized or distributed or hybrid

(ix) Support → High support less support  
Sufficient & high support

(x) Maturity → more less  
Simpler needs, either → high entry point

(xi) Query capability → All query statements → Only simple  
apply

(xii) Expertise → More less

26/06/2019

\* categories of NOSQL Database →

(i) Document Based → MongoDB technique

- If you use MongoDB then data store in collections into documents.
- Documents contain for example, if you want to store student info. the each students info makes single document to store student info.

(ii) XML Based database → Marklogic technique

for example.

<xml>

<student>

<name> ABC </name>

<id> 01 </id>

</student>

</student> </xml> → at the end

} One student record

- Here, if you want to store data in XML format then we can use above pattern. like in above e.g. only store the one student info, if you want multiple info then open <student> after all data entered then XML close.

Page No.	
Date	

(iii) Graph Based database → GraphDB technique

- if you want to information store in Graph format then we use this ~~to~~ at Graph Based database by using GraphDB technique.

(iv) Key value store → Redis, cassandra technique

- MongoDB is store unstructure data then we can store data in key:value format.

if you give cap then compulsory to give size of saturation point value

only one collection can accept

Page No.	1	7	19
Date			

\* creation of collection :-

- Three ways to create collection as follows ;

(i) db.createCollection(name)

e.g.

db.createCollection("TYITA")

(ii) db.createCollection(name,[option])

e.g. Options are

(i) capped boolean value (Yes/No)

(ii) size bytes (saturation point value)

(iii) max (no. of documents)

(iv) autoIndexId boolean (Yes/No)

e.g.

db.createCollection({ "student", {capped : true, size : 1045, max : 500, autoIndexId : true}})

Here, student is db collection name and we give to option to capped for cap collection.

(iv) if you want to create directly document then automatically create collection but it is not idle way to create.

db e.g.

db.customer.insert({ "name" : "abc", "roll no." : 01 })

\* if you want to see the list of collections :-

syntax :-

show collections

\* Drop collection :-

db.collection-name.drop()

e.g. db.TYITA.drop()

mongodb  
is case  
sensitive.

square bracket indicates we can insert multiple value.

Page No.	5	3
Date	5	07

## \* Inserting Data & creating documents

Syntax : db.collection-name.insert( {key:value, key:value} )

example :

- (i) db.student.insert( { "name": "ABC", "id": 01 } )
- (ii) db.student.insert( { "RothkoId": 02, "name": "PQR" } )
- (iii) db.student.insert( { "id": 03, "name": "XYZ", "Gender": "M", "Address": "HF" } )
- (iv) db.student.insert( { "name": "LMN", "id": "IT03", "gender": "m", "phoneno": [ 99000, 98000 ] } )

multiple value

- (v) db.student.insert( { "name": { "firstname": "XYZ", "Lastname": "IJK" }, "id": 04, "gender": "F", "Address": { "city": "panvel", "state": "Maharashtra" } } )

## \*

It Read & Insert Data

- It is also one way inserting data.

example :

```
user1 = { "name": "ABC", "id": 02, "gender": "F" }
```

```
db.student.insert( { user1 } )
```

Here, we insert document by using variable.

## Inserting Data by Loop :

Ex for e.g. :-

```
for( var i=1 ; i<=20 ; i++ ) db.student.insert( { "name": "username"+i, "Subject": "NGT", "age": 20 } )
```

If the parameter is same and particular key

is increment / decrement then -

- \* **Update Data :-**
- for example :- `db.student.update({ "gender": "F"}, { $set: { "country": "India" } })`
- Here, we can update only one row condition can be change.
- (E) if you want to all conditions are present then update then use `{ multi: true/false }` by default is false.
- e. `db.student.update({ "gender": "F"}, { $set: { "country": "India" } }, { multi: True })`
- if you all data can be update then `db.student.update( {}, { $set: { "country": "India" } }, { multi: True })`
- if you want to remove key then we can use `$unset`.
- `db.student.update({ "gender": "F"}, { $unset: { "country": "India" } }, { multi: True })`

- \* **Delete Data :-** document OR collection
- if you want to remove the record but collection remains same.
- `db.student.remove({ })`

8/07/2019

- (\*) query document :- (i) selector (According what you selecting data, cond<sup>n</sup>)
- (ii) projector (what is display), without cond<sup>n</sup>)
- for example :-
- (E) `select * from student where address = "Panvel" ;`
- projector
- selector

when MongoDB select or Retrieve by using find.

e.g. db.collection\_name.find()

db.collection\_name.find().pretty()  
→ it gives proper structure

db.collection\_name.find().pretty()

Q1

Syntax :

selector db.collection\_name.find({selector}, {projector})

gives value example : db.collection("student").find({})

in key : (i) db.collection\_name.find({})

Value - it gives all data without cond.

projector db.collection\_name.find({},{})

gives value (ii) db.student.find({},{ "name": 1, "Rollno": 1 })

- it gives name and rollno column without  
true false any cond & it means gives all data according to

default (name and rollno)

O. (iii) db.student.find({ "address": "panvel" })

- it gives all numbers data whose address is -

"panvel"

(iv) limit :

limit is (iv) db.student.find({ "address": "panvel" },

used to { "name": 1, "Rollno": 1 })

display -

the Q. Write Mongodb query to display Age, name, Rollno,

particular where gender is female and should not display

no. of identifiers or - id and gender.

records.

db.student.find({ "gender": "female" }, { "name": 1,  
"age": 1, "Rollno": 1, "-id": 0, "gender": 0 })

\* (v) db.student.find({ "gender": "F" }).limit(10)

Here, Only display 1st 10 female data with all records.

(vi) db.student.find({ "gender": "F" }, { "name": 1 })

.limit(10)

(iv) skip:

skip is used to display the first 10 female data with only name column.

(vii) db.Student.find({ "gender": "F"}, { "name": 1 }).limit(5).skip(5)

or last Here, skip 1st 5 records of female and after 5 records display the name.

(v)

sort (viii) db.Student.find({ "gender": "F"}, { "name": 1 }).sort({ "name": 1 })

is used to display Here, display only female name according to ascending particular Order.

order for

ascending (ix) db.Student.find({ "gender": "F"}, { "name": 1 }).sort("name" -1 )

for desending Here, display only those female name according to Name is ascending and Rollno desending condn is

satisfied. e.g. ABC 01, ABC 03, ABC 02

for e.g.

ABC 01, ABC 03

ABC 02, output, ABC 02

ABC 03, ABC 01

(x) findOne():

if only one record want to display.

db.Student.findOne({ "gender": "F" })

Here, Only display the one record whose gender is female.

- \* Conditional Operators :-
- (i) less than (\$lt)
  - (ii) greater than (\$gt)
  - (iii) not equal (\$ne)
  - (iv) less than or equal (\$lte)
  - (v) greater than or equal (\$gte)
  - (vi) equal to (==)
  - (vii) in (\$in) → inside
  - (viii) not in (\$nin) → outside

for examples:

- (i) db.student.find({ "id": 05 })
- (ii) db.student.find({ "id": { \$lt: 05 } })
- (iii) db.student.find({ "id": { \$gt: 05 } })
- (iv) db.student.find({ "id": { \$lte: 02 } })
- (v) db.student.find({ "id": { \$gte: 05 } })
- (vi) db.student.find({ "id": { \$ne: 08 } })
- (vii) db.student.find({ "class": { \$in: [A, C] } })

01 ABC A  
02 PQR B       $\rightarrow$  01 ABC A  
03 XYZ C

- (viii) db.student.find({ "class": { \$nin: [A, C] } })

01 ABC A  
02 PQR B      ~~not in~~  $\rightarrow$  02 PQR B  
03 XYZ C

\* AND and OR Operators :-

- (i) AND Operators :-

01 ABC	A	F	India	AND	01 ABC A F India
02 PQR	B	F	USA		
03 XYZ	C	M	India		

**Syntax :-**

`db.collection_name.find({key: value}, {key: value})`

`db.collection_name.find({$and: [{key: value}, {"key": value}]})`

**for example :-**

(i) Write MongoDB query to retrieve data where gender is female AND country is india.

`db.collection_name.find({{"gender": "F"}, {"country": "India"})`

(ii) Write MongoDB query to display data where id is less than equals to 4 and city is pune.

`db.student.find({{"id": {$lte: 4}}, {"city": "Pune"})`

(iii) OR Operators :-

**Syntax :-**

`db.collection_name.find({$or: [{key: value}, {key: value}]})`

**for example :-**

(i) Write a mongodb query where id is equal to One OR id is greater than 5

`db.student.find({$or: [{"id": 1}, {"id": {$gt: 5}}]})`

\* **AND OR Operator :-**

`condn1 AND [condn2 OR condn3]`

**Syntax :-**

`db.collection_name.find({"key": value, $or: [{key: value}, {key: value}]}))`

for example :

Write MongoDB query to find documents where name = ABC AND id = 1 OR 3

```
db.student.find({ "name": "ABC", $or: [ { "id": 1 }, { "id": 3 } ] })
```

07/2013

### \* Using indexes:-

`ensureIndex()` - for creating index on a field

Syntax :- `db.collection.ensureIndex( [key], options )`

for collection applying index

`db.collection-name.ensureIndex( { "key": value } )`

- Here, value contained 1, -1

1 for applying index in ascending order

-1 for applying index in descending order

It is known as single key index.

compound Key Index :-

In this key index, we have two keys to apply

index value because we apply the key in name

but we have one more document in collection

they have no name in document then we can

apply compound key index to collection.

for example :-

ABC 01 Pune

PQR 02 Mumbai

JKL 03 Mumbai

MNO 04 Mumbai

Namespace allow the area for allocating the data.

Page No.	10
Date	20/07/2019

Syntax:

```
db.collection_name.ensureIndex({ "Key": value }, {  
    "Key": value } )  
(1,-1) (1,-1)
```

Unique:

- It is used apply for you can have unique data / records in your document. It will not allow to same data again.

Syntax:

```
db.collection_name.ensureIndex({ "Key": value ,  
    "Key": value }, { unique: true } )  
(1,-1) (1,-1)
```

dropDups:

- In this, we can use unique keyword in after entering data then we use dropDups.
- dropDups scan all the record if the record is duplicate then second record will be delete.

Syntax:

```
db.collection_name.ensureIndex({ "Key": value ,  
    "Key": value }, { unique: true }, { dropDups: true } )  
(1,-1) (1,-1)
```

3/07/2019 getIndexes(): "name" { }

- if you create the index, then indexes data store in System.index
- if you want to retrieve the data by using getIndexes then it's only display the information of index.

```
db.collection_name.getIndexes()
```

\$regex also define regular expression

e.g.

name : \$regex/(ma)/i

Page No.

Date

13 7 19

dropIndex():

Here, drop the index which is previously created  
indexes are drop/delete/remove.

(db.collection\_name.dropIndex())

- we can use with option, sometime we create  
data document to make index then if you  
want to drop particular index then used

db.collection\_name.dropIndex({ "name": 1 })

ReIndex : ( { name: 1 } )

db.collection\_name.reIndex()

- Here, we undo the index if you want to dro  
back to drop index for then we use reIndex.

\*2 Regular Expression :

- if you find the particular Name whose start  
with ma but here only display the \$ small ma.

db.collection("student").find({ "name": /(ma)\* })

if you find the particular Name whose start  
with ma but you not remember the ma is  
small / capital then you can use

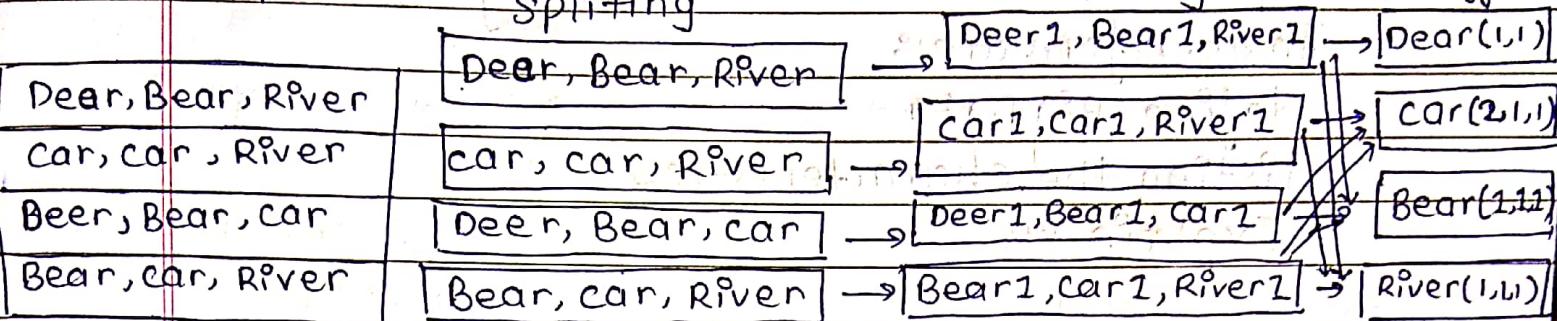
db.student.find({ "name": /(ma)\* /i })

- if you want to find two letter then more letter  
in your collection then

db.student.find({ "name": /(Ra/ch)\* /i })

Here it will display the option whose name  
start with 'Ra' and 'ch' record.

## \* Mapreduce :- Data Processing and Mapping / Shuffling



Example :-

Map	Total Deer(2)	Reduce
	Total car(4)	
	Total Bear(3)	
	Total River(3)	

- Here, your data first splitting in different parts then mapping into the splitting data then using shuffle phase then reduce and count the data. it is working of mapreduce.

(i) var map = function() {emit(this.Gender, 1);};

- Here, Gender are divided into two. Here, your documents / data are divided by gender means male and female separated.

(ii) var reduce = function(key, value) {return Array.sum(value);}

(iii) db.Student.mapReduce(map, reduce) {out: "mapreduce", count: 1}

(iv) db.mapreduceCount().find()

- Mapreduce is very complicated or more steps are required.
- Aggregation is directly apply and it is done in only one steps.

Page No.

Date

**Q. Explain the characteristics of Big data?**

Answer format:

- define Big Data

- listing the characteristics

- Explain V3

**Q. Explain the challenges faced by an organization to manage big data using legacy system.**

- define big data, legacy system

9/07/19

(ii) Directly apply Aggregation

- To perform particular / specific Group it is known as aggregation.

- This first is divided into particular group and then apply the aggregate operation in group.

Syntax:

db.collection-name.aggregate(aggregate  
aggregate Operation)

Example:

(i) Write Mongodb query, to display the total quantity of the products.

product name	quantity	price
pen	10	₹05
Eraser	15	₹3
pencil	20	₹15

If you want to count then use `1` for particular col<sup>n</sup>.

If count apply the particular field.

Page No.		
Date		

for Addition :

for storing data

```
db.stock.aggregate([{$group: {_id: "$productname", "total": {$sum: "$quantity"} }}])
```

(i) Write MongoDB query to display the sum total no. of male and female in the collection.

for counting :

```
db.Student.aggregate([{$group: {_id: "$gender", "totalCount": {$sum: 1}} }])
```

for using some column  
to count of male  
and female

(iii) Write MongoDB query to display the total no. of product in the collection.

```
db.stock.aggregate([{$group: {_id: "$productname", "totalCount": {$sum: 1}} }])
```

Q. Write a short note on Segmentation, Customisation and Innovations (Innovation - 1/2 M, Seg. Customisation - 1/2 M)

Q. Justify technique and technology in the challenge challenges for big data.

20/07/19

- Aggregation not display the details of data. It will display the sum, count, average, minimum, maximum data.

Aggregation function are different operation as follows :-

(i) sum (addition of value and count of data)

(ii) avg

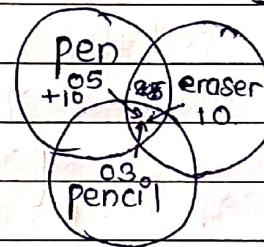
(iii) min & max (minimum & maximum of data)

(iv) max (maximum of data)

Example of sum Operation :-

Group by product

product	quantity
pencil	05
pencil	03
eraser	10
pen	10



for addition :-

```
db.stock.aggregate([{$group: {_id: "$productname",
    "total": { $sum: "$quantity" } }}])
```

for counting :-

```
db.stock.aggregate([{$group: {_id: "$productname",
    "totalCount": { $sum: 1 } }}])
```

Example of average Operation :-

Q. Write mongo db query to find the average quantity of product.

```
db.stock.aggregate([{$group: {_id: "$productname",
    "averageQuantity": { $avg: "$quantity" } }}])
```

Q. Write a mongodb query to find average score subject wise.

$\downarrow$   
 $db \cdot \text{student} \cdot \text{aggregate}(\{\$group: \{\_id: "\$subject", "AverageScore": \{\$avg: "\$score"\}\}\})$

	ch	B <sup>0</sup>	M
01	ABC	20	19
02	PQR	19	20
03	XYZ	25	25

Example of minimum Operation :-

Q. Write mongodb query to find minimum marks of particular subject.

$db \cdot \text{student} \cdot \text{aggregate}(\{\$group: \{\_id: "\$subject", "Minimum Marks": \{\$min: "\$marks"\}\}\})$

Example of maximum Operation :-

Q. Write a mongodb query to find maximum quantity of product.

$db \cdot \text{stock} \cdot \text{aggregate}(\{\$group: \{\_id: "productName", "Maximum Quantity": \{\$max: "\$quantity"\}\}\})$

Q. Write a short note on visibility ?

(ii) Adding and Aiding Decisions Making

Q. Variety is characteristic of big data justify ?

## Opinions

- \* Designing an application data model & Embedding - putting an info/ data someone else's database E.g - FB, book.

Referencing Referencing Referring info/data from one database to another database.

E.g - Book referring

- \* Design of Data Model

(i) Manage life cycle of data there are two ways are following

- (a) setting timing of group documents
- (b) Capped Collection

(iii) Indexes by default created

(iv) Sharding: division of data into customers that is distributing data among servers

(v) collection of group documents

(vi) Growth of document

Q. Write a short note on SQL and give reason why traditional SQL is not useful for big data

Q. Differentiate between SQL and NoSQL

Q. Explain brewer's theorem