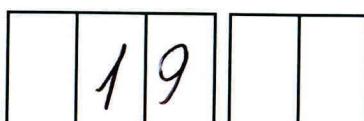


NOM BOUKHELOUA

Prénom Rémi

Promo 2020

Date



BOUKHELOUA Rémi  
M1-2018-TD BD

## MATIÈRE Data Structuring

1. Untidy data examples:

- if there are variables in column ~~and~~ names.  
for example ~~dates~~ observations ~~from~~ from different dates.

It should be

temp	year	obs
Paris	2015	3
Paris	2016	4

(We do it with stack.)

- if you have a same type of observation over several tables :

table A

Paris 2015 3

table B

Paris 2016 4.

It should be in the same table; we use concatenation.

- if you have different type of observation ~~in~~ in a single table. For example, temperatures and pressure. We should normalize it to be in two different tables.

- if you have several variables in the same column. For example: city and country. ~~They~~ They should be in two ~~as~~ different columns (thanks to split).

2. ACID stands for Atomicity, Consistency, Isolation, Durability. It means that a transaction is atomic, so unless it's performed, nothing can access the considered as one point in time. We should not have While it's performed, nothing else can modify the part it modifies (Isolation).

4. MongoDB hierarchy is Database - Collection - Document.

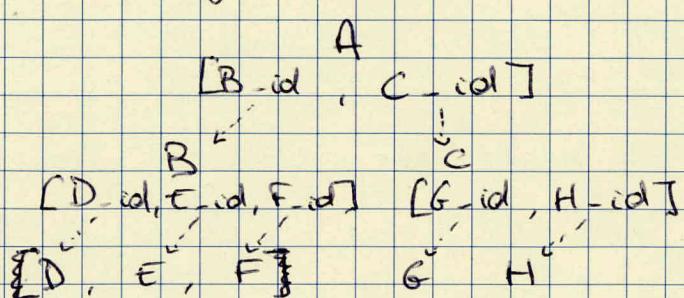
5. In order to model a "Foreign Key" in MongoDB, we use the ids of document we create. How we model it is dependent of the relation we have.

If it's a 1 to 1: We can store the id of the first one in the second document. If we are sure that ~~the document~~

If it's a 1 to Many : We can either store the id of the 1 in each document of the Many or store all ids of the many in a map / array in the 1. To extend In the 2<sup>nd</sup> choice, in order to add one more relation to the one we will only push the new document's id in the array.

If it's a Many to Many: we may do ~~it~~ the same thing that we used in 1 to Many but as much time as we ~~do~~ have documents in one of the Many side of the relation.

To store hierarchical data such as list of employees with management hierarchy, we will use the 1 to Many in order to form a "tree".



Other wise we could share information of B and C in an array in A and with B information we would chose D, E and F information and so on.

```
{ ... , 'manage': [ { ... , 'manage': [ { ... , '...', '...', '...', '...', '...', '...', '...', '...'] } ] }
```

~~But this is not considered right~~ This solution is more complex and heavy but makes queries quicker. Overall I prefer the first one even though joins may not be the best.

6. We use MongoDB through PyMongo where db.coll is the collection showing the given data.

a. db.coll.find({ "keywords": "python" }) will give us a cursor to go through in order to fetch the elements.

b. db.coll.find({ "author-name": { "\$eq": "A.BC" } })

c. db.coll.update({ "author": { "\$eq": "DEF.GG" }, { "time": 100 } })

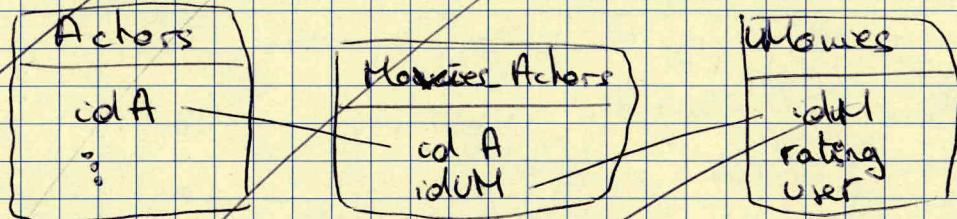
We may want to add a third parameter { "mult": "true" } if there are several documents with DEF Author.

7. One of the ~~big~~ difference between XML and JSON.

• That XML can store informations such as variable names or versions in the (values in ~~the~~ French) meanwhile JSON can't.

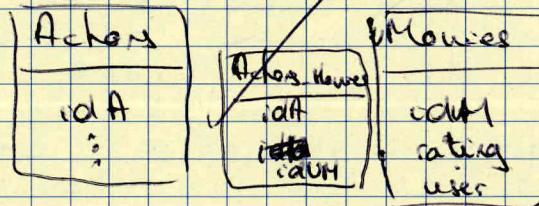
8. We have several possibilities depending on what we want to do with it and how large is the dataset:

- We could use ~~HDFS~~ SQL database if the dataset is not too big, we would store it like that.



It would be the easiest ~~way~~ way to do it and would allow us to do every business questions we need of we ~~want~~ don't ~~only~~ focus ~~only~~ on Movies, Actors or the links.

- If the number of record is too import we could decide to store it in HBase through HDFS and use HiveQL to query it.



NOM BOUKHELOUAT

Prénom Rémi

Promo 2020

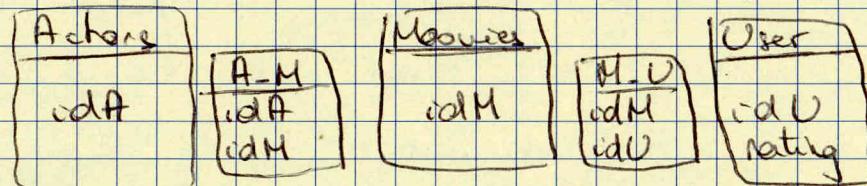
Date 2

## MATIÈRE Data Structuring (2).

8. Model will depend on how large the dataset is and what we want to focus on.

- If we want to have an easy and simple view of the dataset, we would use SQL with the following:

scheme:  
if the user can  
rate several  
times a movie



- if the dataset is too large, we would rather

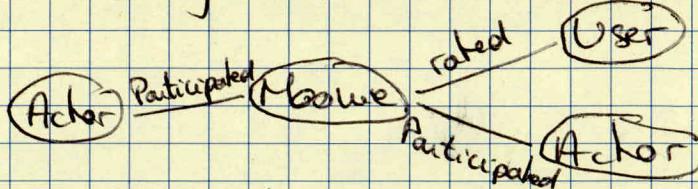
use HBase and process query it through Hive.  
Both are ~~easy~~ easy to join and look for lists of ~~other~~ movies and counting their reviews or finding mean of rates given by a single user.

- We could also use MongoDB or any document-based database. That would be more interesting if we wanted to focus on user ratings or on actors. We would use same implementations as described on Q5.

either simulating "Foreign Key" and then joining or integrating ~~in another document~~ all ~~user~~ an actor's document all movies he participated in and the its.

their ratings. It would be easier to search about ~~how~~ ~~we~~ defining the actor's domain and how well / bad he ~~is~~ regarding user's ratings.

- Finally we could use a Graph database such as Neo4J to focus mainly on links. It would be something like that:



we could easily figure a relation between two ~~for~~ movies, actors  
or users. ~~like searching in~~ <sup>We could</sup> to ~~for~~ <sup>movies</sup> to ~~for~~ <sup>users</sup> one user based on the ones that other user did.

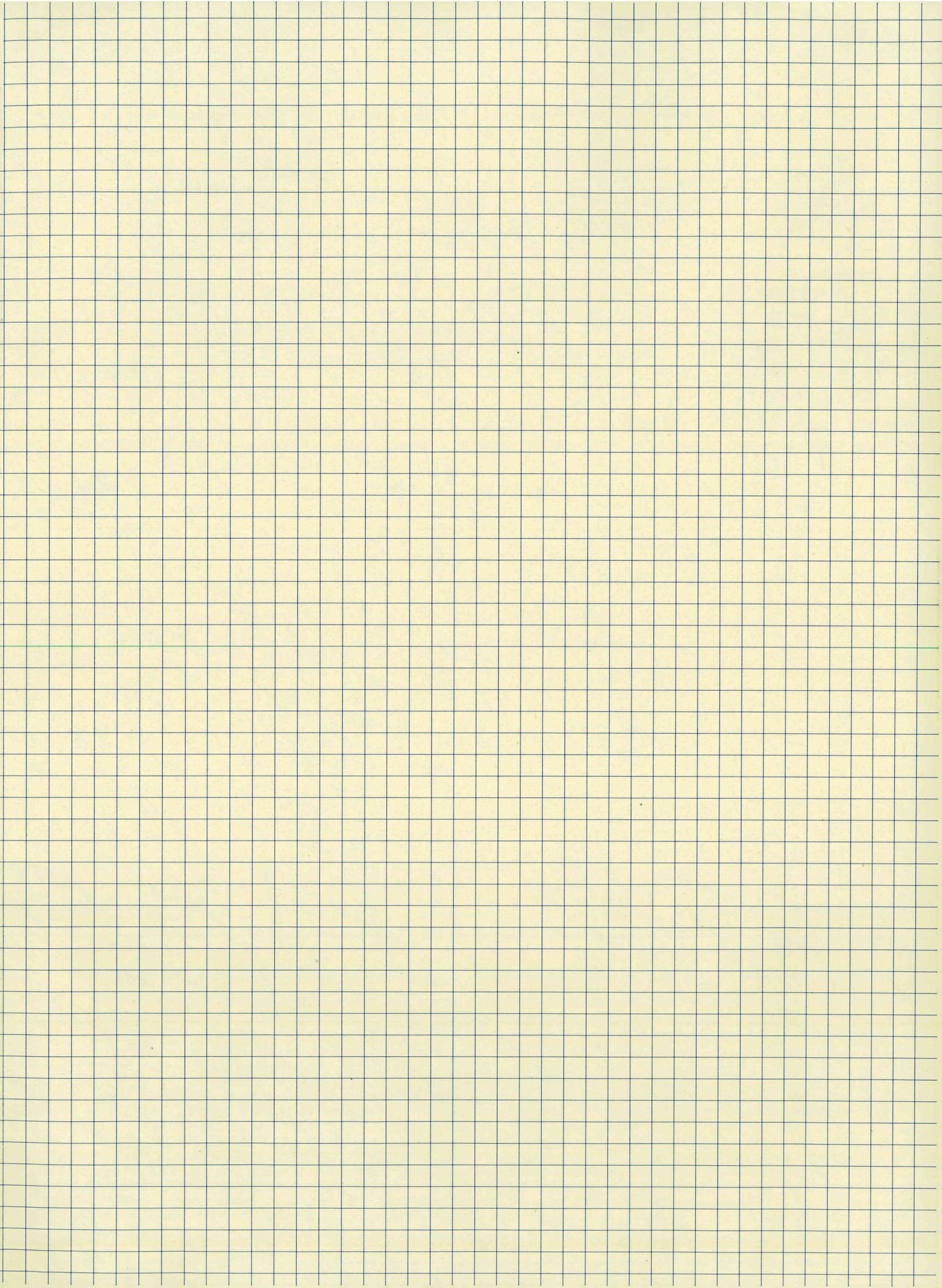
Q. Wikipedia web pages are ~~structured~~ ~~at least~~ semistructured.  
As it's ~~an~~ important amount of data, we may not be able to store ~~and~~ in SQL, ~~and~~ however it would not fit very well as elements are not always the same.

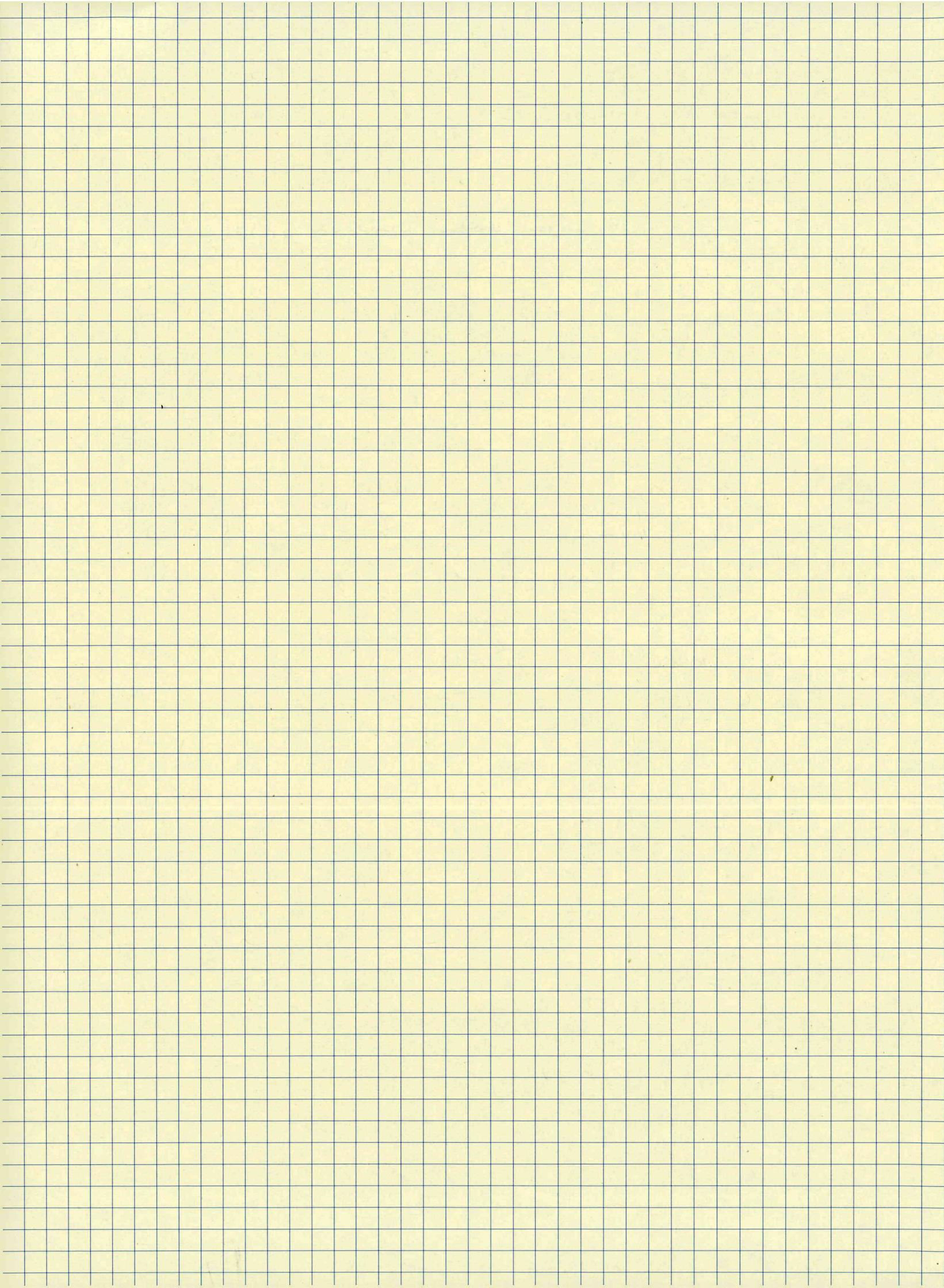
So, we should go either to MongoDB storing a web page as a document : { 'title': ..., '-id': ..., 'link': ..., 'topic': ... }

and the references to other web pages would be stored ~~as~~ ~~to~~ thanks to '-id' and images will ~~be stored~~ <sup>put</sup> share their ~~address~~ ~~of~~ connection in the system.

Otherwise, ~~as~~ ~~as~~ as Wikipedia web pages are HTML, or through XML for a dictionary, we could ~~put~~ <sup>put</sup> them in a File System such as Hadoop which handles semistructured data. Then we will be able to use Hive ~~or~~ or ~~Python~~ Spark to process it and find out easily what we need.

HDFS would fit well because it has been designed for handling such type of data and ~~large~~ large datasets.







Prénom ..... *Rémi* .....

Ne rien inscrire dans ce cadre

Nom ..... *Boukhelouf* .....

Promotion ..... *2020* .....

Groupe ..... *Big Data* .....

### M1

#### Data structuring and NoSQL databases

ST2DST

DE - 1h45 min

Date   Horaire

#### Sujet proposé par :

Calculatrice autorisée :  OUI       NON

Documents autorisés :  OUI       NON      Type de documents :

Ordinateur portable autorisé :  OUI       NON

Internet :  OUI       NON

Traducteur électronique, dictionnaire :  OUI       NON

#### **Consigne :**

Merci de restituer uniquement : **les copies quadrillées à rendre accompagnées de l'annexe**

#### **Rappel :**

- Tous les appareils électroniques (téléphones portables, ordinateurs, tablettes, montres connectées ...) doivent être éteints et rangés.
- Il est interdit de communiquer.
- Toute fraude ou tentative de fraude fera l'objet d'un rapport de la part du surveillant et sera sanctionnée par la note zéro, assortie d'une convocation devant le conseil de discipline. Aucune contestation ne sera possible. Tous les documents et supports utilisés frauduleusement devront être remis au surveillant.
- Aucune sortie de la salle d'examen ne sera autorisée avant la moitié de la durée de l'épreuve.

1. Provide 4 different examples of « Un-tidy » data (you can sketch them if necessary)
2. Explain the meaning of ACID letters for SQL transactions
3. Provide examples of XML queries and explain what they do
4. What is MongoDB schema hierarchy? (like Database/Table/Row for SQL)
5. How to model a "Foreign Key" in MongoDB? What are different modelling options? For example, how to store hierarchical data, like a list of employees with management hierarchy?
6. For a data

```
{'_id': ObjectId('5ba207f488811f06f433e7f7'), 'author': 'ABC', 'author_name': 'A. B C', 'icon_filename': 'icon1.png', 'text': 'Some text', 'time': 1234487, 'keywords': ['test', 'python', 'mongodb'], 'responses': [{ 'author': 'DEF', 'author_name': 'D. E F', 'icon_filename': 'icon2.png', 'text': 'thanks a lot!'}, { 'author': 'ABC', 'author_name': 'A. BC', 'icon_filename': 'icon1.png', 'text': 'cheers!', 'responses': [{ 'author': 'ABC', 'author_name': 'D. EF', 'text': 'Thanx again'}]}]}
```

write the Mongo query to

- a. Fetch only the elements with a keyword 'python'
- b. Fetch only the elements with "author\_name" == 'A.BC'
- c. Increase time by 100 for author "DEF"

7. What is a difference between XML and JSON?

8. You have data on

- Users watching and rating movies
- Actors playing in the movies

What are your different options to model these data and store it in a single SQL or NoSQL database (sketch it if necessary)? What databases can you use? Provide examples of "business" questions that are better adapted to each storage/modeling option.

9. You have downloaded a copy of Wikipedia (or of an encyclopedia), how would you store it? Why? What are other options and how organize in this case (SQL, NoSQL, or something else)?



