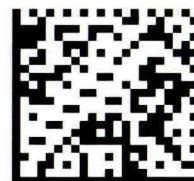


NOM BOUKHELOUA  
Prénom Rémi  
Promo 2e de  
Date 16.04.



20150257: BOUKHELOUA Rémi  
M1:  
ST2FDP-DE (16/04/2019)  
Amphi orange

## MATIÈRE Functional Data Programming (1/2)

Q1:

- a. Immutability protects ~~the~~ variables from side effect.  
This is when a variable is changed somewhere it shouldn't ~~be~~ and when it gets a value it shouldn't have. Immutability prevent from modifying its value ~~so~~ from side effect. ✓
- b. A high order function is a function that takes another function as a parameter or/and return a function. It is possible in Scala as we use the paradigm "everything is a function". ✓✓
- c. Map(f) ~~create~~ return a new list containing the image ~~of~~ by f of initial element. ✓
- d. FlatMap is a fusion of map and flatten where flatten means transforming List of List into a simple List. For example, List(List("a", "b"), List("c")).flatten() returns List("a", "b", "c"). ✓

Q2: There is no Q2 in subject.

Q3:

- a. A RDD is a Resilient ~~flex~~ Distributed Dataset. Resilient means fault tolerant, distributed means divided into partitions and dataset means partitioned data. It supports two types of manipulation: actions (from RDD to not RDD) and transformations (from RDD to RDD). It can be created from various sources such as file system (HDFS for e), distributed source (KAFKA, Cassandra), data from other or another RDD.

b. A DataFrame is a distributed collection with columns structure (unlike RDD)

✓ It's ~~represented~~ by similar to MySQL tables and called table in SPARQL  
It can be created from various sources such as filesystem, other RDDs,  
external DB and Hive tables (little better than RDD sources). Finally it's at  
~~least~~ a higher level than RDD and has great performance.

c. Given "lines: RDD[String]"

lines.map(p => p.split(" ")).  
\* flatMap(w => (w, 1))  
\* reduceByKey(- + -) X

It will divide line into  
bitwords its composed of  
flat and create tuples  
(word, 1)  
Sum tuples of same word

d. A Spark computation is triggered by any evaluation that requires a result.

X

Q4:

a. 1. A broker is in charge of receiving message from producers and persist them. It can handle several millions of message each second.

a. 2. Zookeeper is the centralization tool. Intern in Kafka, it is in charge of managing task performed in Kafka. ✓

b. A ~~topic~~ partition is where partitons are stored depending on their key. It works as a append only and guarantee the order. Partitions are dispatched on a cluster. We could resume it as following:



✓

It is important because if we have a consumer group formed by 3 consumers, ~~Kafka~~ will provide a third of the partitions to it. For example: 3 consumers A, B, C and 6 partitions 1, 2, 3, 4, 5, 6.

A will get 1; 2, B will get 3; 4 and C will get 5 and 6.

It allows distributed work.

But from different groups will have every partitions!   
~~the consumer~~ each

c. As one broker is in charge of ~~other~~ partitions, if it fails,  $\checkmark$  partitions will be lost. That's why they introduced replications in order to ~~store partitions~~ on partition ~~on~~ on different brokers. (estimatively) So if partition #1 is stored on broker A and replicated on B, and if A stops, the copy on B will become the partition 1. It's even more important that Kafka is a distributed system that runs over cheap ~~on~~ materials (Big Data horizontal scaling) so those material can may more frequently crashing. So we need replication to make sure not to loose data!

d. The offset is an integer that corresponds to the position of last element read by a given consumer group in a partition. It is incremented after the consumer gets ~~a~~ a message and committed.



Q6: a. def fact(nb: Int, acc: Int = 1): Int = nb match {

case 0 => return acc

case \_ => return fact(nb - 1), acc \* nb)

}

puntln(fact(5))

✓ ✓

b. def boys(students: ~~RDD~~[Student]: RDD[Student] = {

students.filter(\_.\_isMale)

}

✓

c. def boysAge(students: RDD[Student]): RDD[Int] = {

~~age~~ boys(students).map(s => age(s.birthday))

{

✓

d. def youngest(students: RDD[Student]): Int = {

~~students.foldLeft(age(students(0).birthday))~~

f(acc, e) => if acc < e

students.map(s => age(s.birthday))

X

• foldLeft(age(students(0).birthday)) f(acc, e) => if(acc < e) acc

else e {{}}

{

e. def hasHomonyms(@Student: Student, students: RDD[Student]): Boolean = {

val c = students.filter(s => s.name == @Student.name)

• count()

✓

~~count~~ c >= 1

{

NOM BOURHFLOUA  
Prénom Rémi  
Promo 2020  
Date 16.04

## MATIÈRE Functional Data Programming (2/2)

Q5: Regarding the customer page, I think we should ~~be record~~ record when a customer click on a model to have more information, when he add or delete an item from his shopping cart and when he confirm his payment. With following we would get info about model viewed, hesitation / annotations and payment. These three message would be composed of:

Click msg:

login  
time  
model~~list~~  
cart (is it already in cart or not?)

Add/Delete msg:

login  
time  
model~~list~~  
action ("Add" or "Delete")

Payment msg:

login  
time  
cart: [ , ]  
totalPrice

As this is going to be used for statistical purpose, ~~so~~ I think that we want to take it as a DF in SPARK ~~so~~ Evolutionary in order to ~~have~~ handle it easily in Spark Notebook (with Spark SQL). But also, we want ~~the~~ a consumer to store it into some kind of HDFS with ~~one~~ XLS files. For those two reasons, we want a consumer to handle only one ~~type~~ of format of message. That's why I ~~would~~ create a Topic Customer ~~is stored over three boxes~~ partitions (~~so~~ force at least one for each).

I want to create a Topic ~~for each message type~~ for Customer's Action where every message will fetch with a key depending on its format ( "click", "modify", "pay"). So consumer ~~will~~ will be able to get one partition.

Regarding the Administration page, we have to handle three different actions : add a new shirt, set stocks and change description

These actions would produce following messages:

Add msg:

model

size: [ ]

stock: [ ]

price

description

time.

Set msg:

model

size: [ ]

stock: [ ]

time.

Des~~Change~~ msg:

description

model

time.



As for Customer interface, I think it would be interesting to have a topic reserved to Administration partitions regarding message event type. We won't ~~need~~ need real time view but to show and modify DB ~~as~~ it would be easier to ~~select~~ each have a different consumer in each partitions.

~~Those~~ ~~the~~

Those consumers would be linked to HDFS/~~real~~ time view (for Customer) and to a Available and Partition Tolerant DB such as Cassandra or CouchDB to keep a track of stocks on the website. (for Administrator). We don't mind of consistency ~~because~~ because we prefer the ~~system~~ to be available at anytime and fault tolerant to ~~keep~~ keep any informations.

~~It's~~ The Architecture would look like:

