

→ spark RDDs are lazily evaluated, & sometimes we may wish to use RDD multiple times. Spark will recompute RDD & all its dependencies each time we call action RDD which is expensive for iterative algorithm which took at data many times. To avoid computing RDD multiple times we can use spark to persist data. It will recompute lost partitions of data when needed. We can replicate our data on multiple nodes if we want to be able to handle node failure without shutdown.

Level	space used	optime	Memory	disk	comment
MEMORY_ONLY	H	L	Y	Y	
MEMORY_ONLY	L	H	Y	N	
MEMORY_AND_DISK	H	M	some	some	split to disk if there is too much data to fit in memory
MEMORY_AND_DISK_2	L	H	some	some	split to disk if there is too much data to fit memory
DISK_ONLY	LOW	HIGH	N	Y	

spark has many layers of persistency to choose from based on what goals are. In python we always serialize the data persist stories. we write data out disk or heap storage, data is always serialized.

## part-B

Q25

Import java.util.Scanner;

object checkerking {

def test(x1: Int, y1: Int, x2: Int, y2: Int):

Boolean =

{

if ((x1 - x2) &gt;= -1 &amp;&amp; (x1 - x2) &lt;= 1 &amp;&amp;

(y1 - y2) &gt;= -1 &amp;&amp; (y1 - y2) &lt;= 1)

return true;

return false;

};

def main(args: Array[String]): Unit = {

var scanner = new Scanner(System.in);

parent("First cell row &amp; column :");

var x1 = scanner.nextInt();

var y1 = scanner.nextInt();

parent("Second cell row and column :");

var x2 = scanner.nextInt();

var y2 = scanner.nextInt();

if (test(x1, y1, x2, y2))

println("Yes");

else

println("No");

{}

264

```
import scala.collection.mutable.ArrayBuffer;  
object fruits extends App {  
    def main(args: Array[String]) {
```

```
        var fruits = ArrayBuffer[String]()  
        ("Apple", "Banana", "Mango", "Orange") ::
```

```
        fruits += ("strawberry", "pineapple");  
        scala.util.Sorting.quickSort(fruits);  
        fruits = fruits.reverse();
```

```
        fruits.remove(fruits.length - 1, fruits.length - 2);  
        fruits.toArray;
```

}

3

1BM18CS035.  
Harshavardhanak

26) spark-shell

```
→ val sqlContext = new org.apache.spark.  
→          SQLContext(spark)  
→ val df = sqlContext.read.json("cars.json")  
→ df.select("carName", "carprice").orderBy("carprice")  
      .show(10)  
→ df.select(df.columns(10)).show()
```

3bs

Tuple: is basically grouping data in of different types:

- \* To create tuple we can do val tup = (4, 3, "James")
- \* To access members of tuple tup[-1]
- \* To copy & change some values:  
tup: copy (-1=0)
- \* To pretty print them tup.tostring.
- \* To swap values tup.swap.

Maps: are useful data structure to make associations between keys and values.

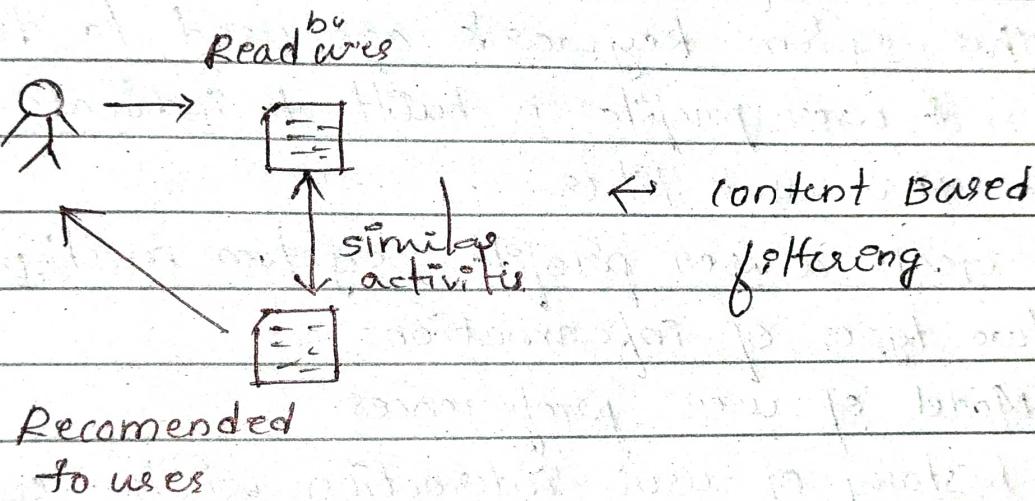
- \* We can create map like val map = Map("Hans" → 234, "Njas" → 79)
  - \* Check values exist in map map.contains("Joh")
  - \* Access values in map map("Hans")
  - \* Add new pairing to map by creating new one val anothermap = map + ("Bob", 45)
- ⇒ program:

```
def getCapital(county: string,emap:  
Map [string, string]): string = {  
    if (emap.contains(county)) {  
        return emap.get(county);  
    }  
    return emap.get("India");  
}
```

part-c.

qas) two types are:

- \* content based filtering: focus on properties of items, similarity of items is determined by measuring similarities in their properties.
- \* collaborative based filtering: focus on relationship b/w users & items. Similarity of items is determined by similarity of rating of those items by users who have rated both items.

Read by both users

```

graph TD
    User1((User 1)) <--> User2((User 2))
    User1 --> ReadBoth[Read by both users]
    User2 --> ReadBoth
    ReadBoth --> Item[Item]
    Item --> Recommended[Recommended to him]
    
```

collaborative filtering

Similar users

Recommended to him

Read by her  
recommended to him

## Content based recommendations / filtering -

- \* Based on description of item & profile of user's preference.
- \* This method are best suited to situations where there is known data on item but not on user.
- \* They treat recommendation as user-specific classification problem & learn from classifier for user's like & dislikes based on item's

In this system keywords are used to describe items, & user profile is built to indicate type of item this user likes.

To create a user profile, systems mostly focus on two types of information:

- \* Model of user preferences
- \* history of user interaction with recommender system.

Consider an example of recommending new actions to user. Let's say we have 100 articles & vocabulary of size N. We first compare tf-idf score for each of words for every article. Then, we construct a vector:

If Item vector : this is a vector of length N. It contains 1 for words that have high tf-idf score in that article otherwise 0.

	bacon	beans	beautiful	blue	breakfast	brown	dog	0.00
1	0.00	0.00	0.60	0.53	0.00	0.00	0.00	0.00

1	0.00	0.00	0.49	0.43	0.00	0.00	0.00	0.00
---	------	------	------	------	------	------	------	------

2) User vector: Again  $l \times n$  vectors, we store probability of word occurring in articles that user has consumed. Once we have these profiles, we compute similarity between users and items. 1) user has highest similarity 2) has highest similarity with other items that user has used.

2 common methods: We simply take cosine similarity between user & item vector to obtain user-item similarity. To recommend items that are most similar to items that user bought, we compute the cosine similarity.

$$\text{cosine}(x, y) = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$$

Jaccard similarity: intersection over union, formula

$$J(x, y) = \frac{|x \cap y|}{|x \cup y|}$$

This is used for item similarity. We compare item vectors with each other & return items that are most similar.

Jaccard similarity is useful only when the vectors contain binary values. If they have ranking or rating that can take on multiple values, Jaccard similarity is not applicable.