# ST446 – Distributed Computing for Big Data

## Seminar 3

Dr. Marcos Ennes Barreto

Assistant Professorial Lecturer

4 February 2021

# Dr. Marcos Ennes Barreto

Web: https://www.lse.ac.uk/Statistics/People/Dr-Marcos-Barreto

Contact: m.e.barreto@lse.ac.uk

Room: COL 5.04

Office hours (book through LSE Student Hub):

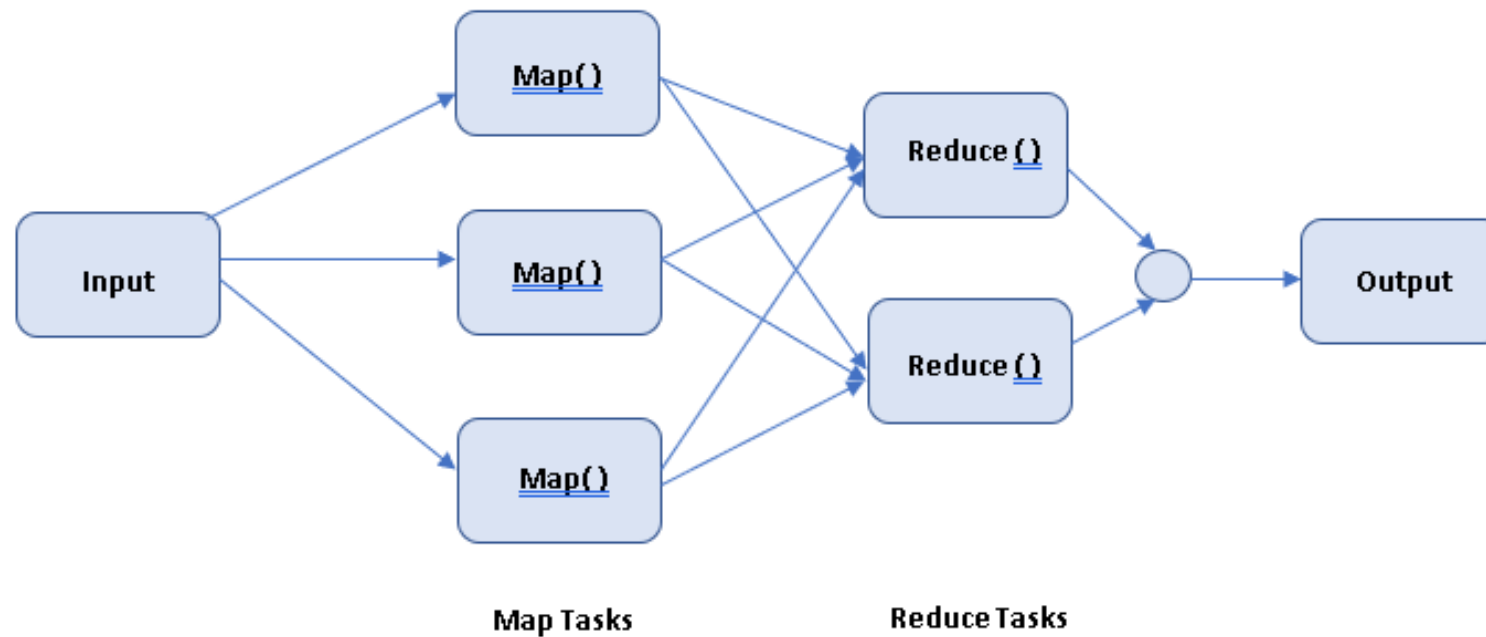✓ Tuesday – 09:00 – 10:00 / 14:00 pm – 15:00 pm

**slack**

**ST446 Q&A**

**st446-lt2021-qa.slack.com**

4  February 2021

## MapReduce

A **programming framework** that allows us to perform **distributed** and **parallel** processing on large data sets in a distributed environment.
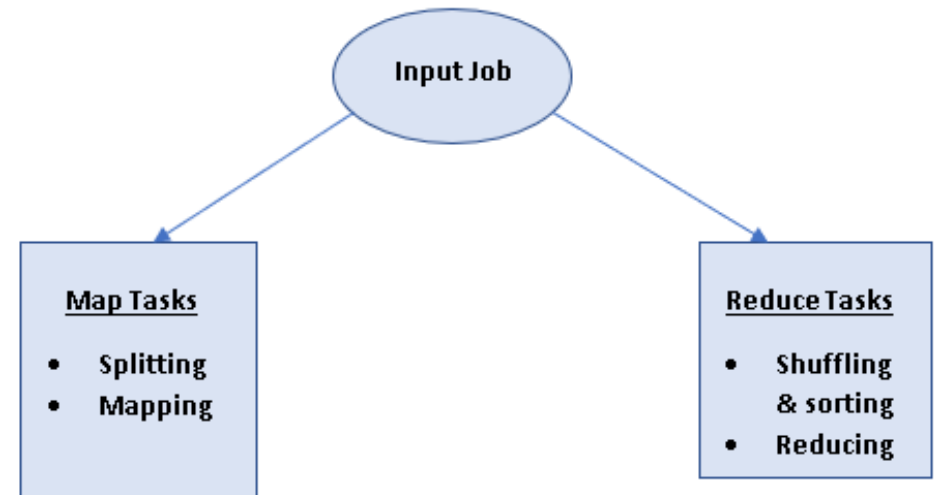


https://www.datasciencecentral.com/profiles/blogs/hadoop-for-beginners-part-2

# How MapReduce works?



## Map Tasks

• **Splitting:** Input data is divided into fixed-size chunks called input splits.
• **Mapping:** In this phase each input split is passed to a mapping function which divides the split into *List (Key, Value).*

## Reduce Tasks

• **Shuffling and Sorting:** Reduce tasks are the combination of shuffle/sort and reduce. This phase consumes output of the Mapping phase. Its main task is to club together the relevant record in sorting manner from the output of mapping phase. The output is in the form of *Key, List (Value).*

• **Reducing:** In this phase, output from shuffling and sorting are aggregated and returns single *(Key, Value)* output value. This final output value is then written in the output file in HDFS.
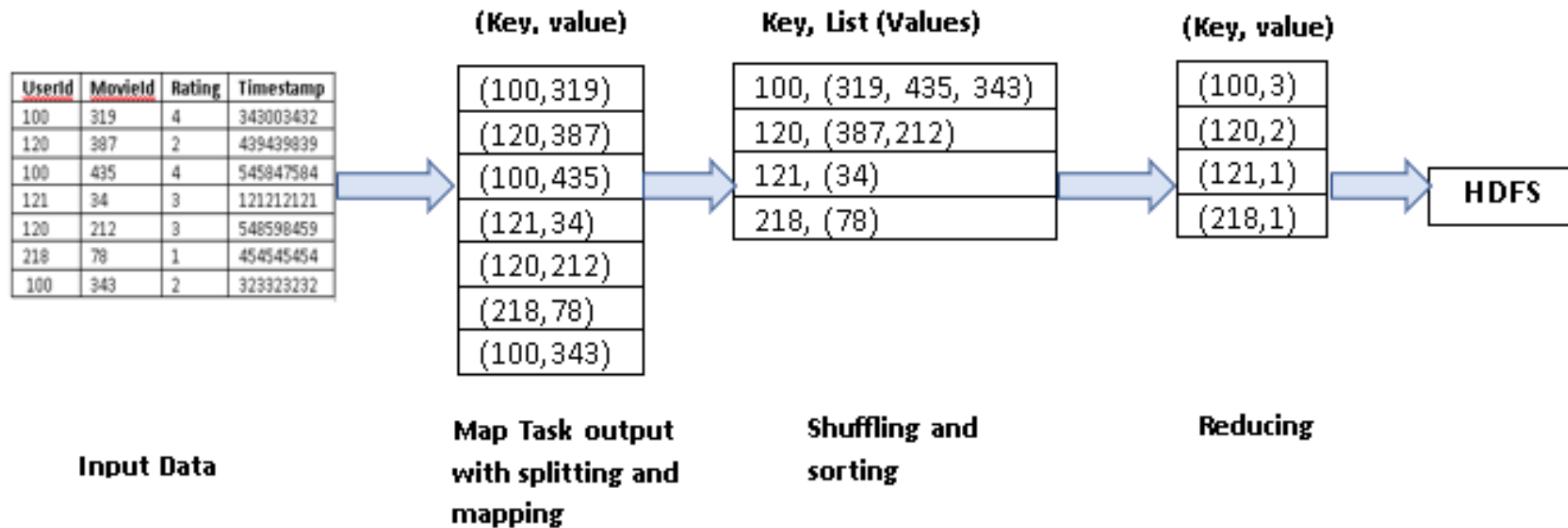
https://www.datasciencecentral.com/profiles/blogs/hadoop-for-beginners-part-2

**Working example**

- **Task** – How many movies did each user rate in the Movie data set?
- **Sample Dataset (Input File)-**

| UserId | MovieId | Rating | Timestamp |
|--------|---------|--------|-----------|
| 100 | 319 | 4 | 343003432 |
| 120 | 387 | 2 | 439439839 |
| 100 | 435 | 4 | 545847584 |
| 121 | 34 | 3 | 121212121 |
| 120 | 212 | 3 | 548598459 |
| 218 | 78 | 1 | 454545454 |
| 100 | 343 | 2 | 323323232 |

https://www.datasciencecentral.com/profiles/blogs/hadoop-for-beginners-part-2

# Working example

| UserId | MovieId | Rating | Timestamp |
|--------|---------|--------|-----------|
| 100 | 319 | 4 | 343003432 |
| 120 | 387 | 2 | 439439839 |
| 100 | 435 | 4 | 545847584 |
| 121 | 34 | 3 | 121212121 |
| 120 | 212 | 3 | 548598459 |
| 218 | 78 | 1 | 454545454 |
| 100 | 343 | 2 | 323323232 |

**(Key, value)**

(100,319)
(120,387)
(100,435)
(121,34)
(120,212)
(218,78)
(100,343)

**Key, List (Values)**

100, (319, 435, 343)
120, (387,212)
121, (34)
218, (78)

**(Key, value)**

(100,3)
(120,2)
(121,1)
(218,1)

HDFS

**Input Data**

**Map Task output with splitting and mapping**

**Shuffling and sorting**

**Reducing**

https://www.datasciencecentral.com/profiles/blogs/hadoop-for-beginners-part-2

**Another working example**



MAPREDUCE ABSTRACTION

DATA

Key          Value

MAP → Hypertension, 1
      Diabetes, 1

MAP → Heart Disease, 1
      Hypertension, 1

```
Map(patientRecord)
{
    disease_list = find_disease_mentions(patientRecord)
    for(disease in disease_list)
    {
        emit(disease,1)
    }
}
```

https://www.udacity.com/course/big-data-analytics-in-healthcare--ud758

# MAPREDUCE ABSTRACTION

**Another working example**



| | Hypertension: 1, 1 | | Hypertension: 2 |
| Key    Value | Diabetes: 1 | REDUCE | Diabetes: 1 |
| Hypertension, 1 | Heart Disease: 1 | | Heart Disease: 1 |
| Diabetes, 1 | | | |

Shuffle

Heart Disease, 1
Hypertension, 1

```
Reduce(disease, counts)
{
emit(disease, sum(counts))
}
```

https://www.udacity.com/course/big-data-analytics-in-healthcare--ud758

Hadoop is based on *acyclic data flow* from stable storage to stable storage.

**When to use Hadoop?**

Input

Map

Map

Map

Reduce

Reduce

Output

Hadoop

=

Distributed
Storage

+

Distributed
Computation

+

Fault Tolerance

https://www.udacity.com/course/big-data-analytics-in-healthcare--ud758

**When not to use Hadoop?**



## ITERATION IN MAP-REDUCE

Initial Model — Map — Reduce — Learned Model

$w^{(0)}$

Training Data

$w^{(1)}$

$w^{(2)}$

$w^{(3)}$

Hadoop is inefficient for applications that repeatedly *reuse* a working set of data:

**Iterative Algorithms**
- Machine learning
- Graph analysis
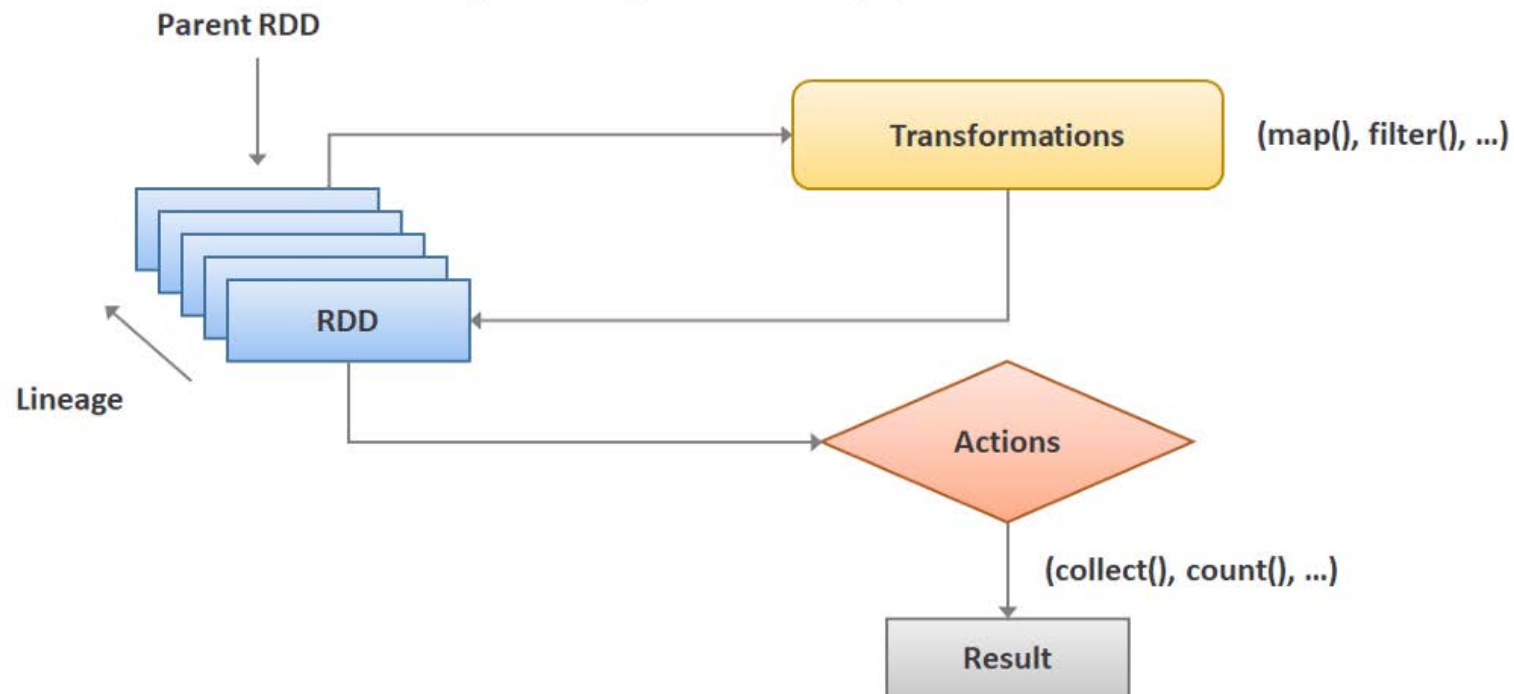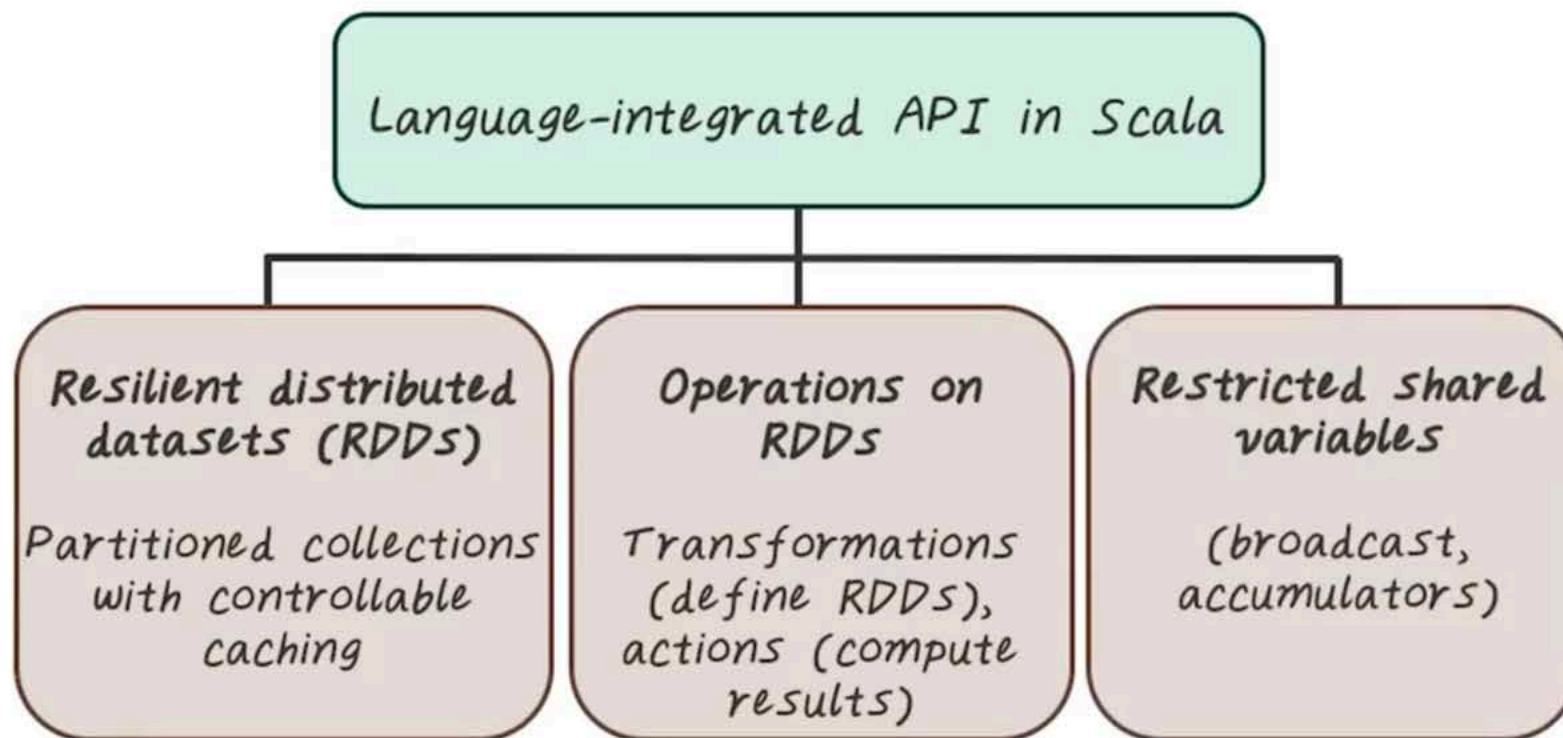
**Interactive Data Mining Tools**
- R
- Python

https://www.udacity.com/course/big-data-analytics-in-healthcare--ud758

Initial Model $w^{(0)}$ — Training Data — Map — READ 1, READ 2, READ 3 — **Repeatedly** load same data — Reduce — Learned Model $w^{(1)}$, $w^{(2)}$, $w^{(3)}$

Training Data — Map — Reduce — **Redundantly** save output between stages — Learned Model $w^{(1)}$, $w^{(2)}$, $w^{(3)}$

https://www.udacity.com/course/big-data-analytics-in-healthcare--ud758

Spark RDD (Unstructured) Operations

https://medium.com/analytics-vidhya/spark-rdd-low-level-api-basics-using-pyspark-a9a322b58f6

# SPARK OPERATIONS
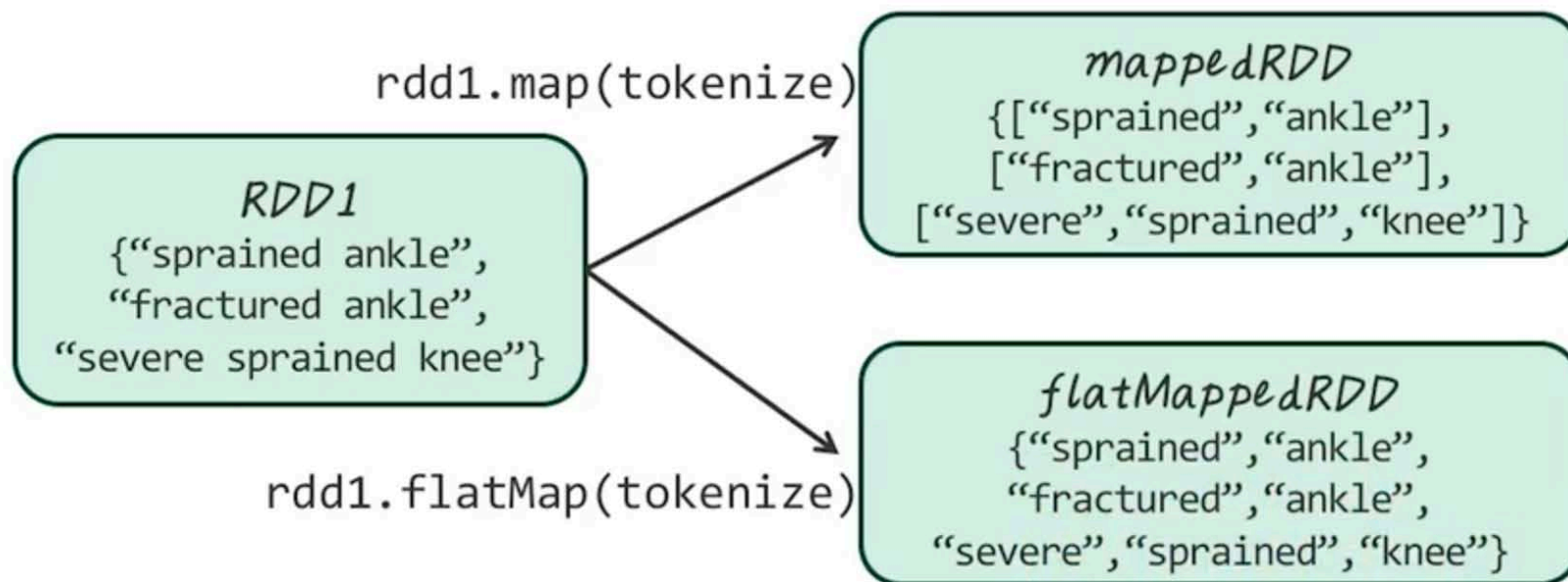
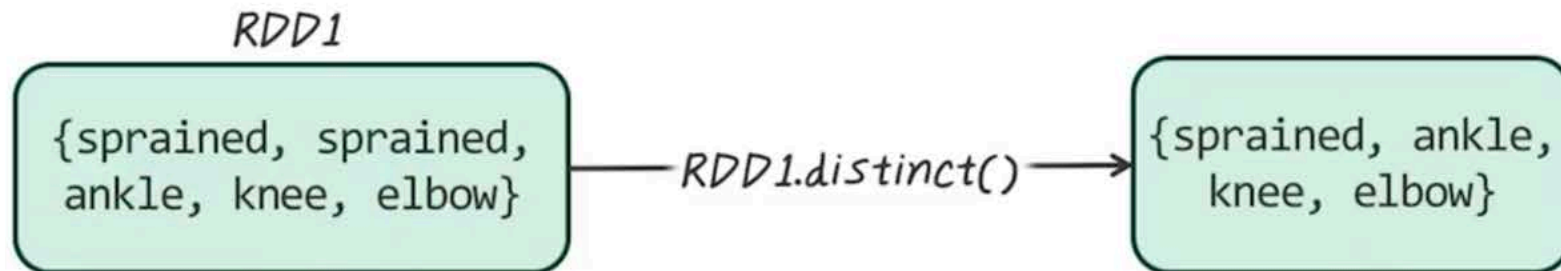| | | |
|---|---|---|
| **Transformations** (define a new RDD) | map<br>filter<br>sample<br>groupByKey<br>reduceByKey<br>sortByKey | flatMap<br>union<br>join<br>cogroup<br>Cross<br>mapValues |
| **Actions** (return a result to driver program) | collect<br>reduce<br>Count<br>save<br>lookupKey | |

https://www.udacity.com/course/big-data-analytics-in-healthcare--ud758

RDD TRANSFORMATIONS

map() vs flatmap()

tokenize("sprained ankle")=List("sprained","ankle")

rdd1.map(tokenize)

RDD1
{"sprained ankle",
"fractured ankle",
"severe sprained knee"}

mappedRDD
{["sprained","ankle"],
["fractured","ankle"],
["severe","sprained","knee"]}

rdd1.flatMap(tokenize)

flatMappedRDD
{"sprained","ankle",
"fractured","ankle",
"severe","sprained","knee"}

https://www.udacity.com/course/big-data-analytics-in-healthcare--ud758

https://www.udacity.com/course/big-data-analytics-in-healthcare--ud758

RDD TRANSFORMATIONS
Operation: Subtract()

RDD1
{sprained, sprained, ankle, knee, elbow}

RDD2
{sprained, knee, fracture}

RDD1.subtract(RDD2) → {ankle, elbow}

Cost: Expensive