

# Scalable machine learning and distributed systems

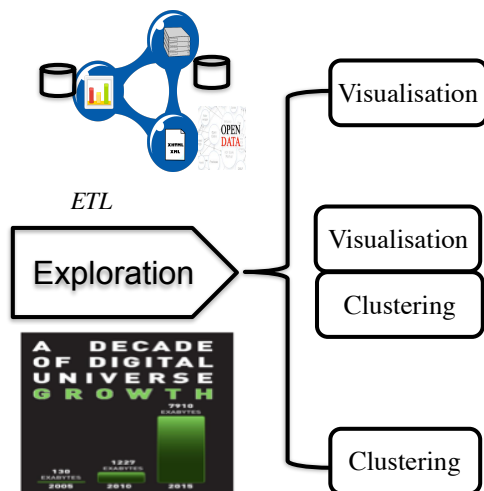
Mustapha LEBBAH  
LIPN - Univ. Paris 13  
Machine Learning and Applications (A3)

CAEN, SBDS 2017 9 June 2017

H. Azzag, D. Bouthinon, T. Sarazin, M. Ghesmoune, N. Doan, T. Duong, A. Chaibi,

Lebbah et al

## Context



### Difficulties

- Existence of a structure
- Similarity measure ?
- Number of groups (Combinatory)
- Validation (no labels)
- Nature of the variables



Tutorial, IEEE IEEE BigData 2014

Lebbah et al

3

## Outline

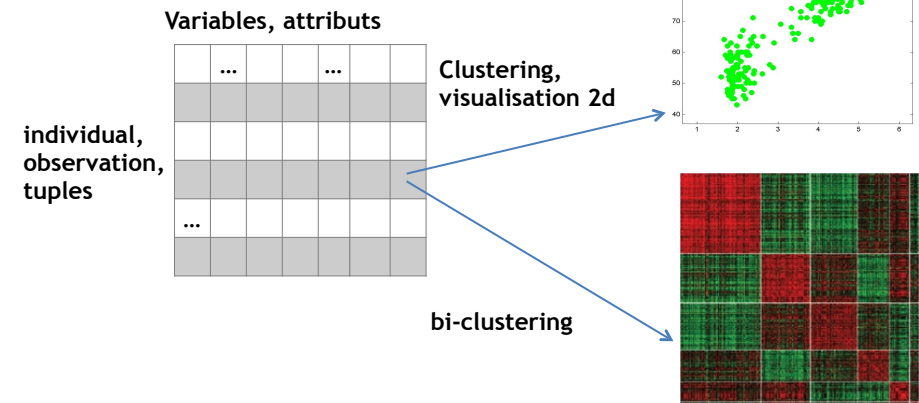
- Context and Issues
- Clustering and new paradigm
  - K-means
  - Topological model (SOM)
  - Mean-shift
- Conclusion

Lebbah et al

2

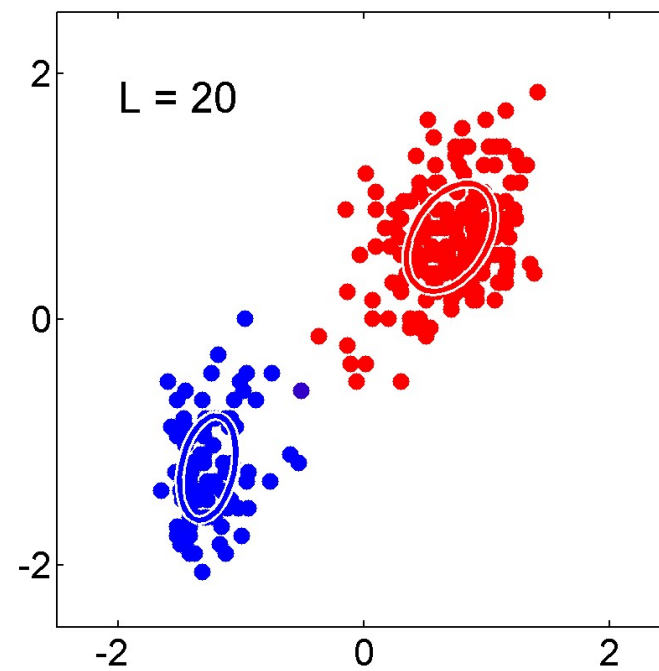
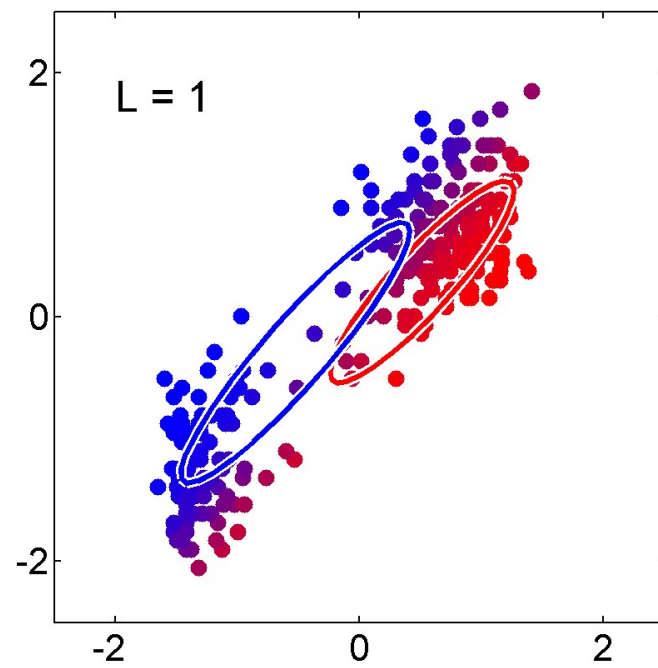
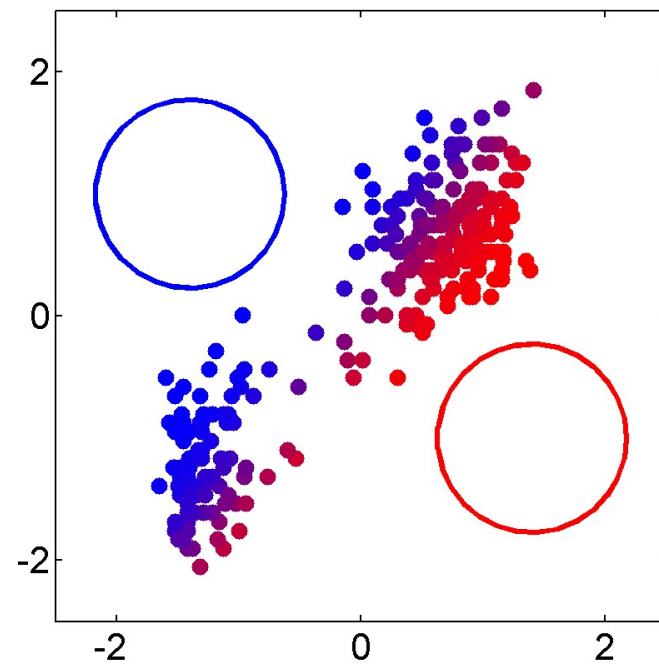
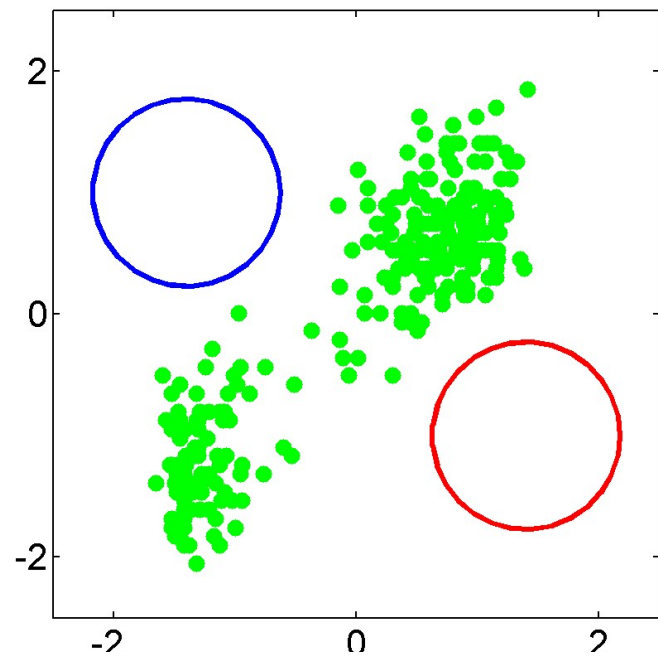
## Unsupervised learning

- > Hierarchical approaches
- > Partitioning approaches
- > Self-organizing approaches
- > Probabilistic approaches



Lebbah et al

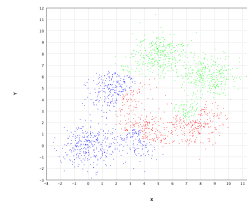
4



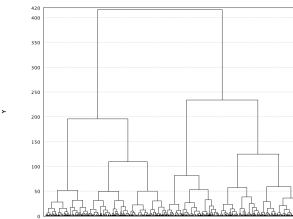
## Challenge - scaling up

Parallelism, distributed system, intermediate, collaborative (adaptable)

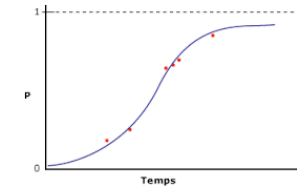
## Challenge - scaling



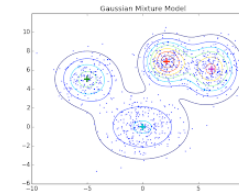
K-means



Hierarchical clustering

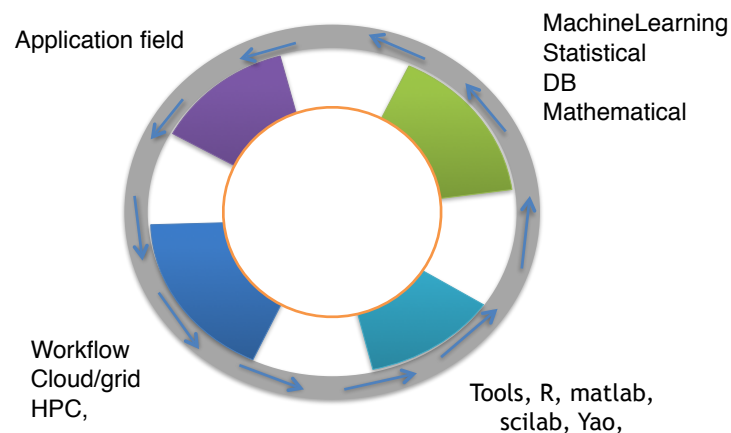


Regression



- Collaborative methods
- Distributed methods
- Ensemble methods

## Challenge - Bringing scientific communities together



## Challenge - New Systems

### Why do we need new management and analysis systems?

#### From an analyst point of view

- Fewer iterations to converge
- Data are always available

#### From a system point of view

- More iterations per second
- The algorithm is seen as a black box
- Complex data management (different types of data)

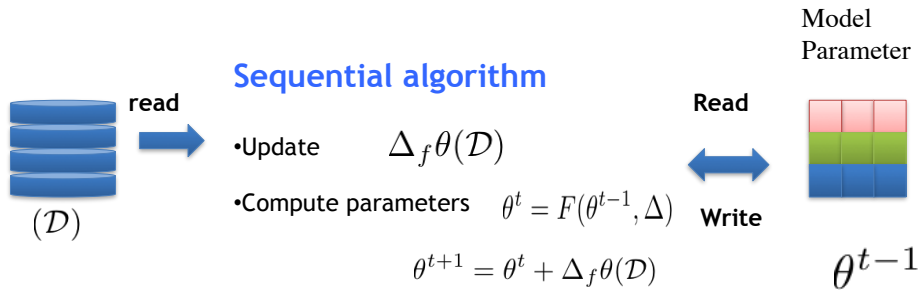
#### From a researcher's point of view

- Suppose an ideal system

# Representation of an iterative algorithm

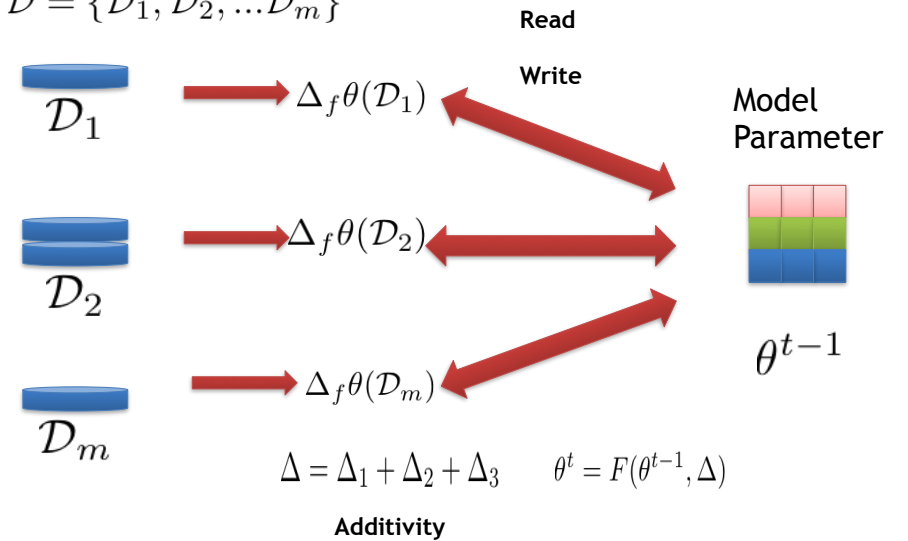
New model = old model + update (gradient, Data)

$$\theta^{t+1} = \theta^t + \Delta_f \theta(\mathcal{D})$$



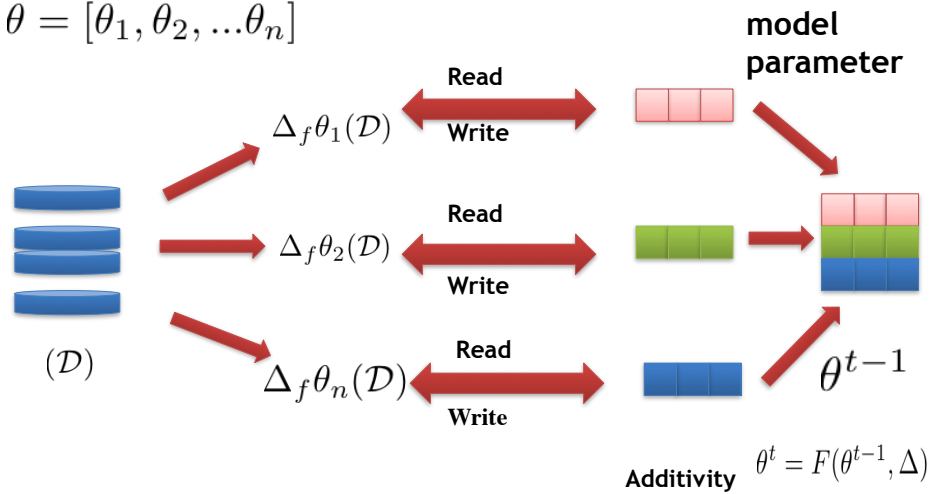
# Distributed Data (ensemble clustering)

$$\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_m\}$$

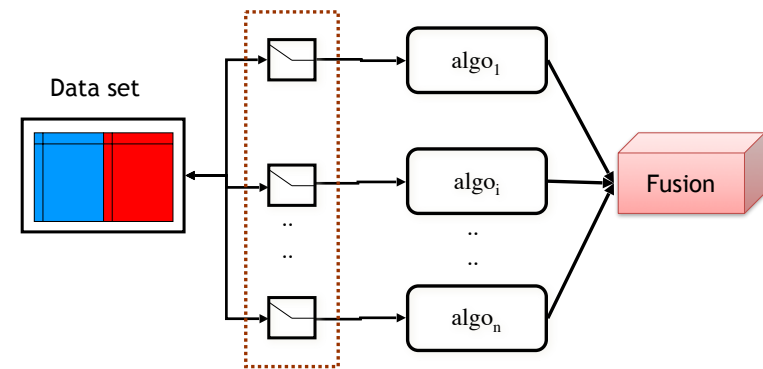


# Distributed Model

$$\theta = [\theta_1, \theta_2, \dots, \theta_n]$$



# Fusion (Clustering Ensembles)



# Question

- The distribution of data is often possible when the data is IID (independent and identically distributed (iid))
- The ideal is to work on an adaptable system (distributed, parallel, collaborative)

- What about tools / languages to manage data, models, parallelize / distribute, easily, quickly?



# Challenge - the same workspace

In Machine Learning



Workspace Matrix view

Database



Tables

Relational Database (Key, values)



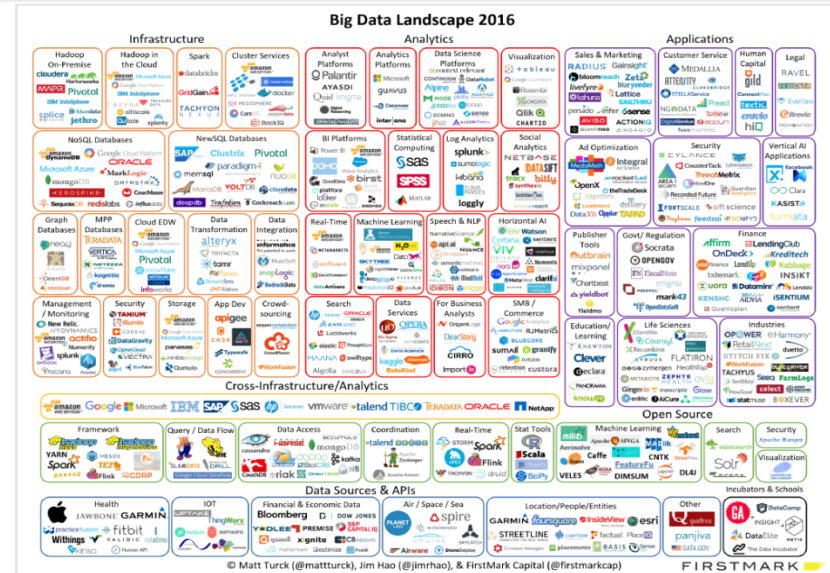
In distributed system



Unify the file management system with the analysis system taking into account the environment

We can not have two observations with the same key

# Outils - Languages



© Matt Turck (@mattturck), Jim Hao (@jimhao), & FirstMark Capital (@firstmarkcap)

FIRSTMARK

# The same workspace

- Manage Workspace as a single database



- Resilient Distributed Datasets (RDDs)

- Collection of objects stored in memory (or disk) through a cluster
- RDDs are distributed across workers
- Parallel processing (map, filter, ...)
- Rebuild automatically in case of failure



(Matei Zaharia)



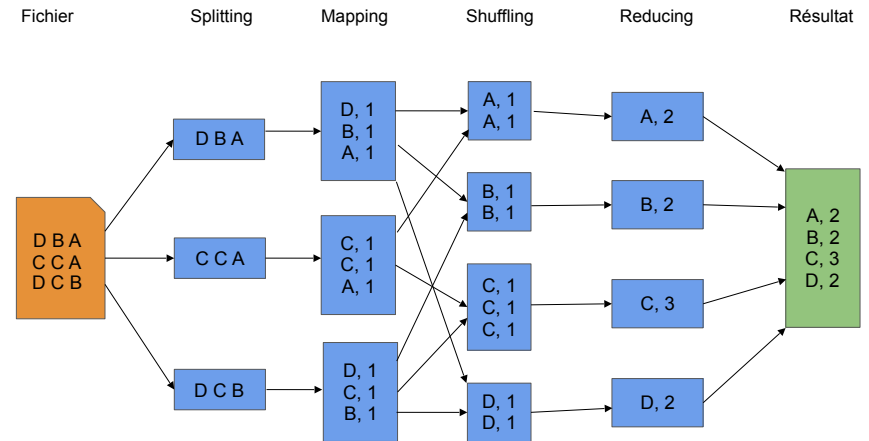


Spark/Hadoop = HDFS + MapReduce

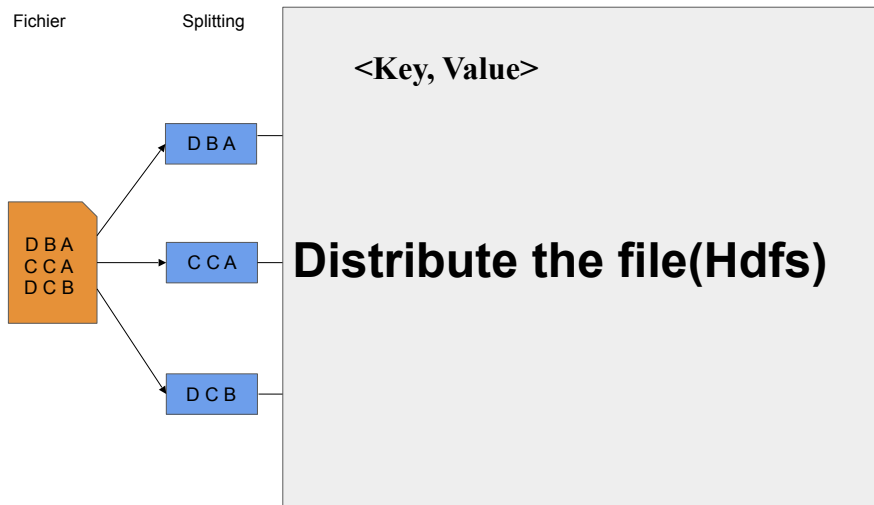
- HDFS : *Hadoop Distributed File System*
- MapReduce : is a programming model for processing a large data sets with a parallel, distributed algorithm on a cluster

- Architecture Scale-Out: Adding machines
- Security
- Data locality: Maximizes execution to the nearest data
- Scheduling: Automatic optimization of job scheduling
- Flexibility: Support for many languages and programming logic (Java, Scala, R, Python)
- Resiliency & High Availability: Unsuccessful jobs and tasks are restarted automatically without loss of the processes already performed

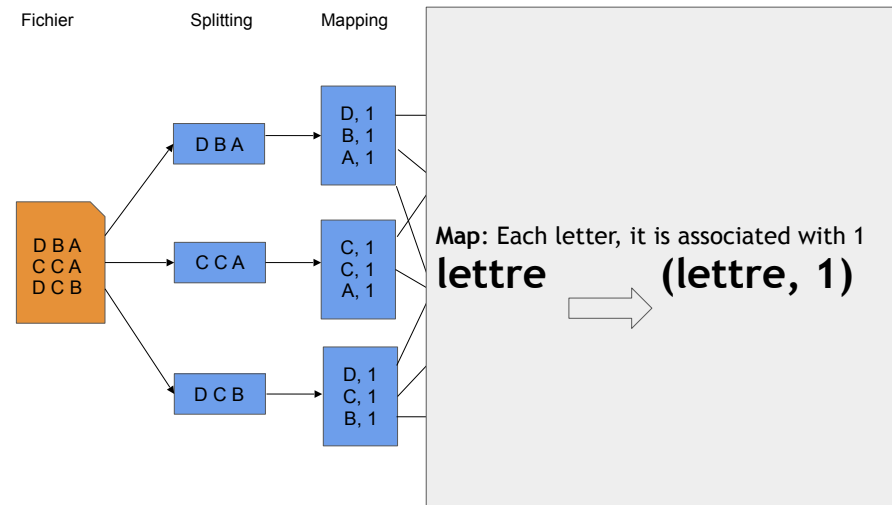
# MapReduce (wordcount)



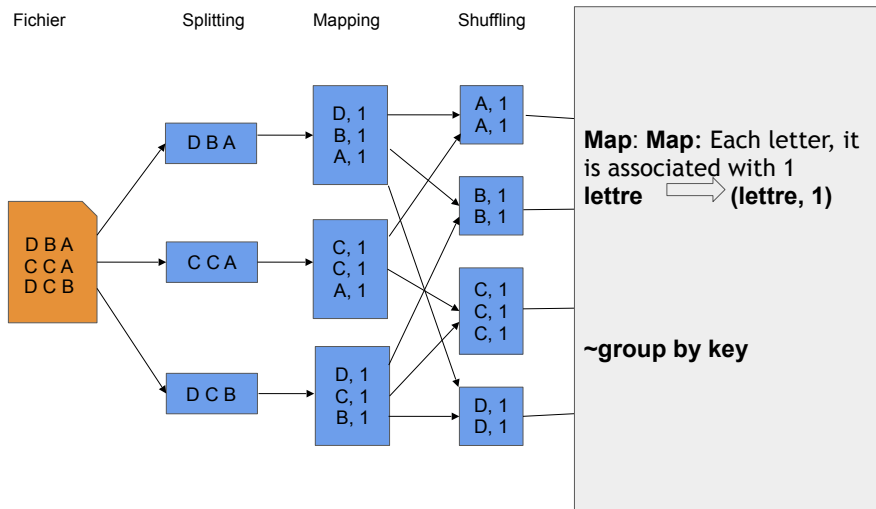
# MapReduce



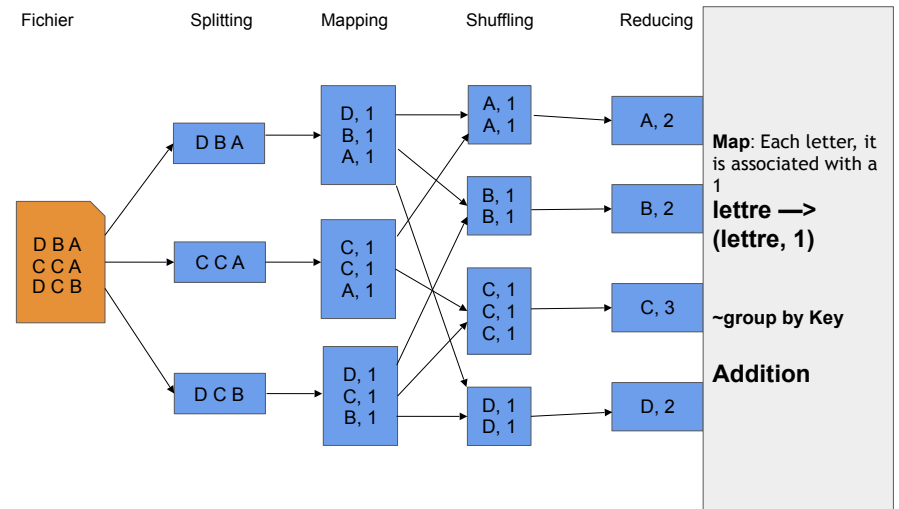
# MapReduce



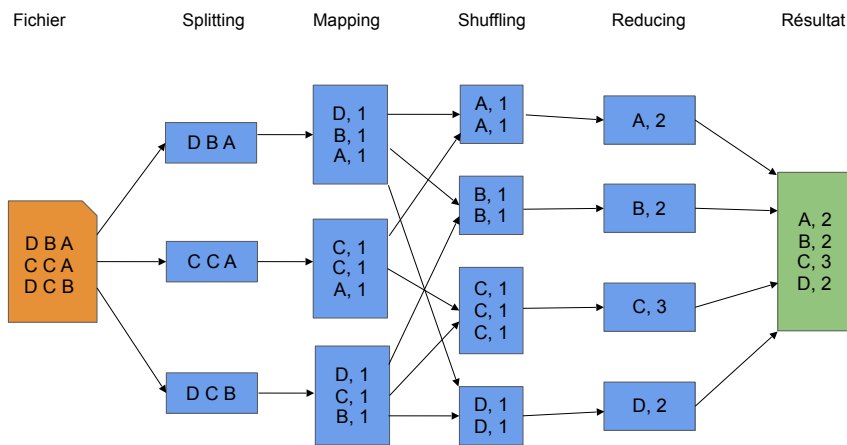
# MapReduce



# MapReduce



# MapReduce



# WordCount : Hadoop MR vs Spark



```

public class WordCount {
    public static class Map extends Mapper<LongWritable, Text, Text, IntWritable> {
        private final static IntWritable one = new IntWritable(1);
        private Text word = new Text();

        public void map(LongWritable key, Text value, Context context) {
            String line = value.toString();
            StringTokenizer tokenizer = new StringTokenizer(line);
            while (tokenizer.hasMoreTokens()) {
                word.set(tokenizer.nextToken());
                context.write(word, one);
            }
        }
    }

    public static class Reduce extends Reducer<Text, IntWritable, Text, IntWritable> {
        public void reduce(Text key, Iterable<IntWritable> values, Context context)
            throws IOException, InterruptedException {
            int sum = 0;
            for (IntWritable val : values) {
                sum += val.get();
            }
            context.write(key, new IntWritable(sum));
        }
    }
}
    
```

```

file = spark.textFile("hdfs://...")
file.flatMap(line => line.split(" "))
    .map(word => (word, 1))
    .reduceByKey(_ + _)
    
```

It's very close to the algorithm idea

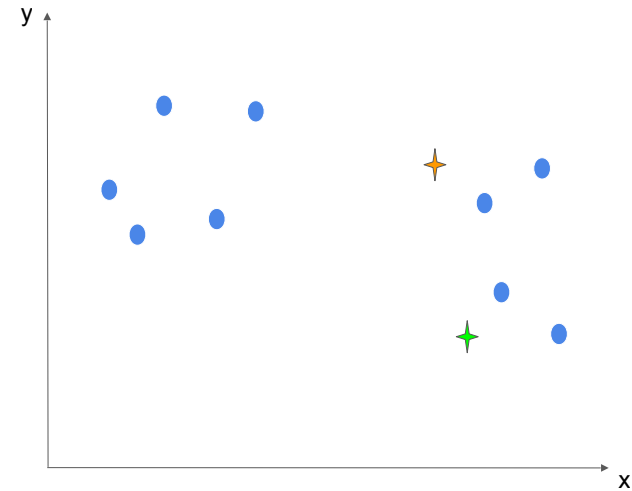
We discover a new language



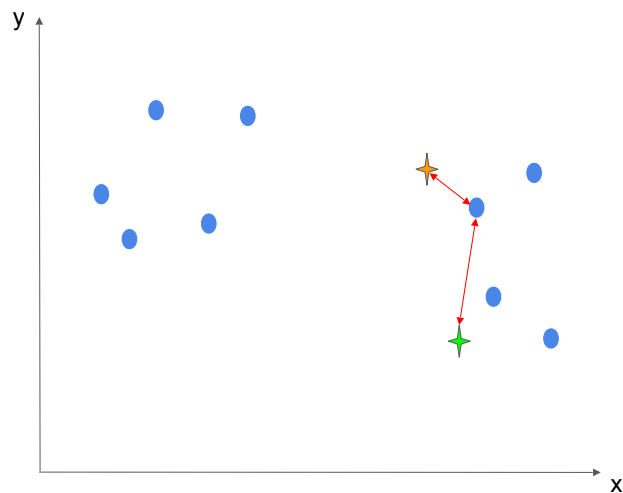
# K-means

Understand the K-means-MR from MLIB

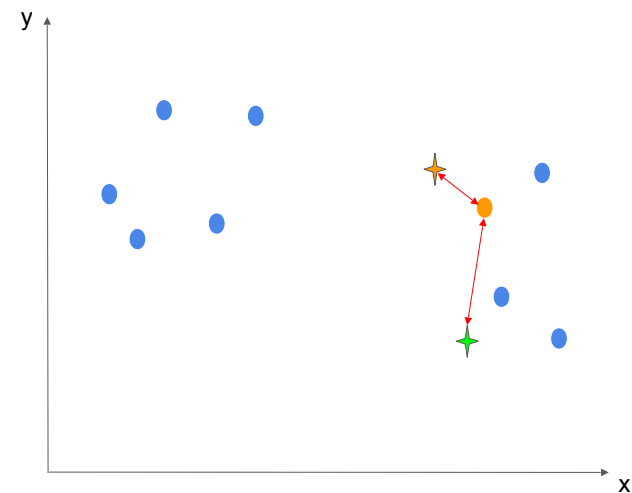
## K-means - initialisation



## K-means - compute distances with prototypes

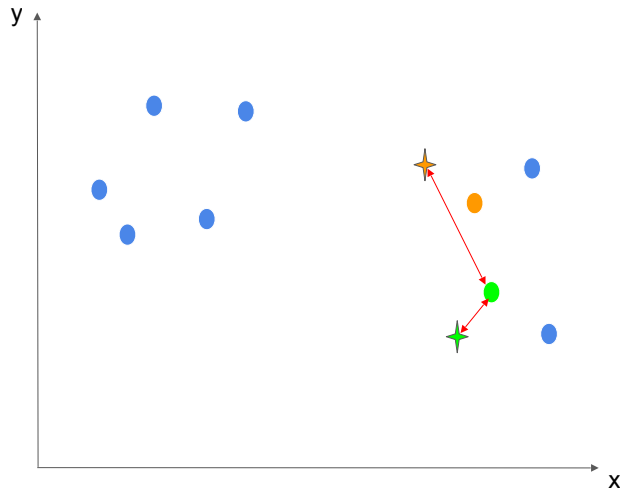


## K-means : Assignment

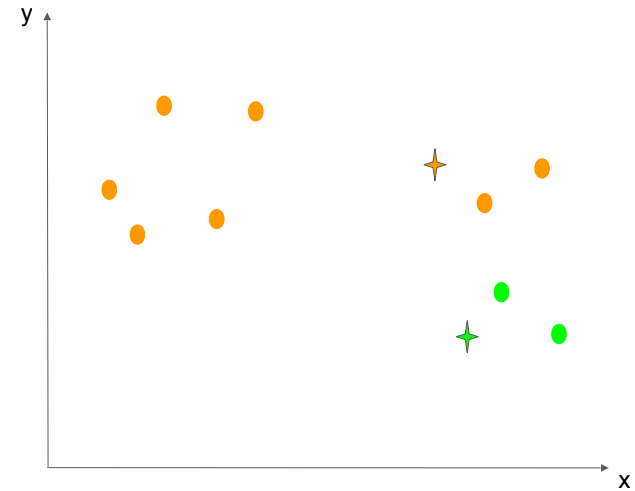




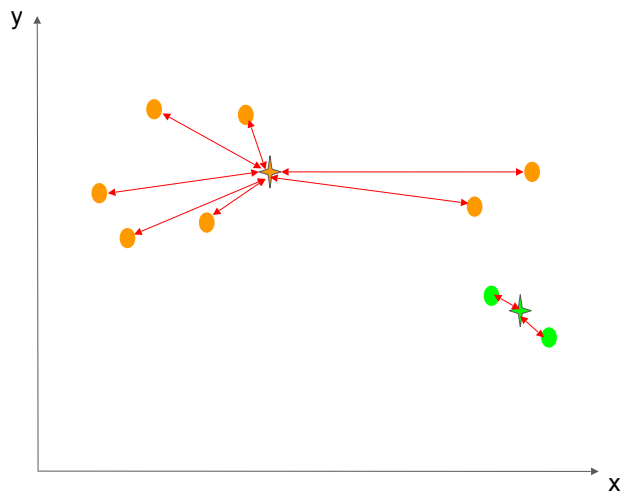
# K-means : Assignment



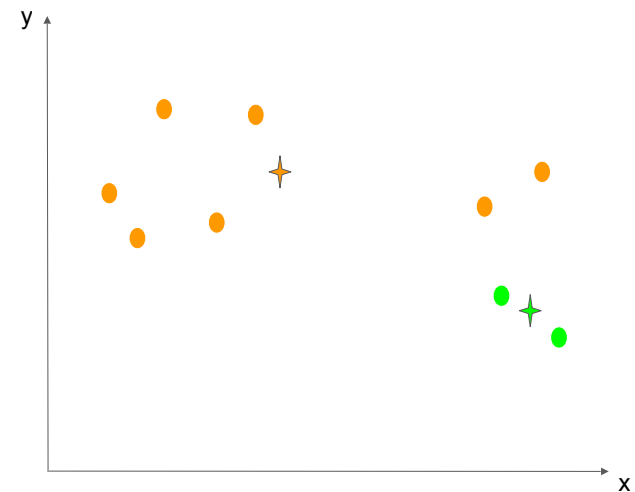
# K-means : Assignment



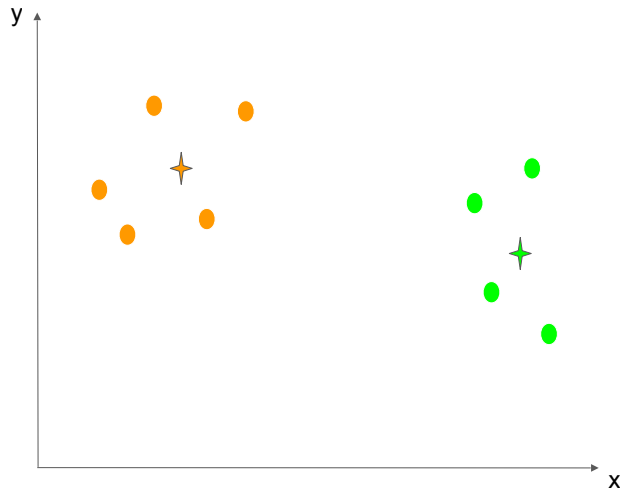
# K-means - Reduce (prototype update)



# K-means : iteration 1



# K-means : iteration 2



# K-means

• Iterative minimization which fixes the partition (c) alternately and then minimizes the inertia

### Assignment step:

For a set  $\mathbf{W}$  of fixed referents, the minimization of  $I$  with respect to  $\Phi$  is obtained by assigning each observation  $\mathbf{x}$  to the referent  $\mathbf{w}_c$  according to the new assignment function  $\Phi$

$$\phi(\mathbf{x}) = \arg \min_r \|\mathbf{x} - \mathbf{w}_r\|^2$$

### Quantisation step:

The partition  $\Phi$  is fixed. The function  $I(\mathbf{W}, \phi)$  is quadratic and convex with respect with  $\mathbf{W}$ . The minimum is

$$\frac{\partial I}{\partial \mathbf{W}} = \left[ \frac{\partial I}{\partial \mathbf{w}_1}, \frac{\partial I}{\partial \mathbf{w}_2}, \dots, \frac{\partial I}{\partial \mathbf{w}_p} \right]^p = 0 \quad \mathbf{w}_c = \frac{\sum_{\mathbf{x}_i \in P_c} \mathbf{x}_i}{|P_c|}$$

# K-means

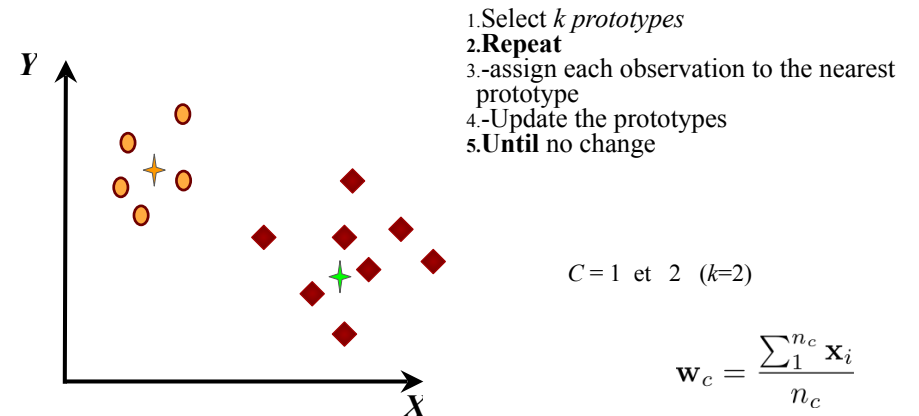
• Minimize the sum of local inertia with respect to  $\phi$  and  $\mathbf{W}$

$$I(\mathbf{W}, \phi) = \sum_{\mathbf{x}_i} \|\mathbf{x}_i - \mathbf{w}_{\phi(\mathbf{x}_i)}\|^2 = \sum_c \sum_{\mathbf{x}_i \in P_c} \|\mathbf{x}_i - \mathbf{w}_c\|^2$$

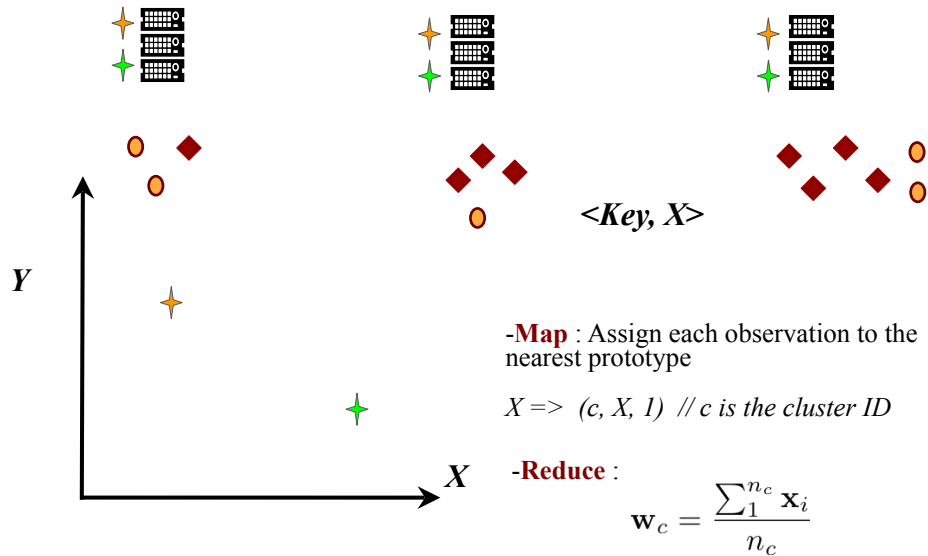


• The inertia  $I_c$  represents the quantisation error obtained if we replace each observation of  $P_c$  by its prototype  $\mathbf{w}_c$

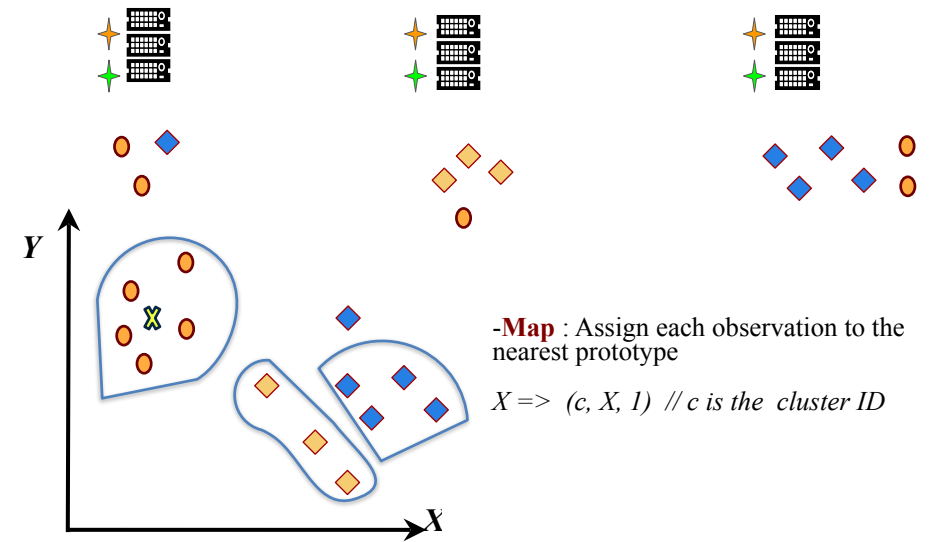
# Decomposition of K-Means



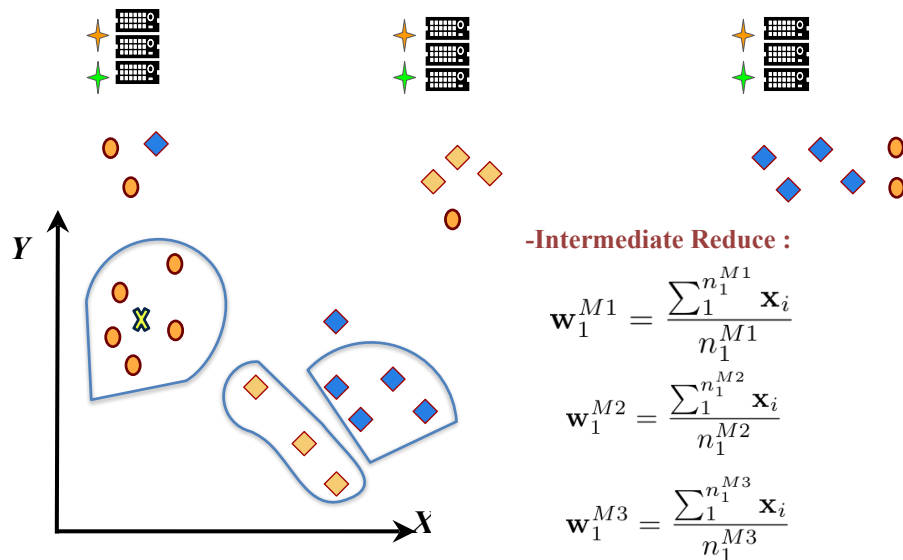
## Decomposition of K-Means (1st solution)



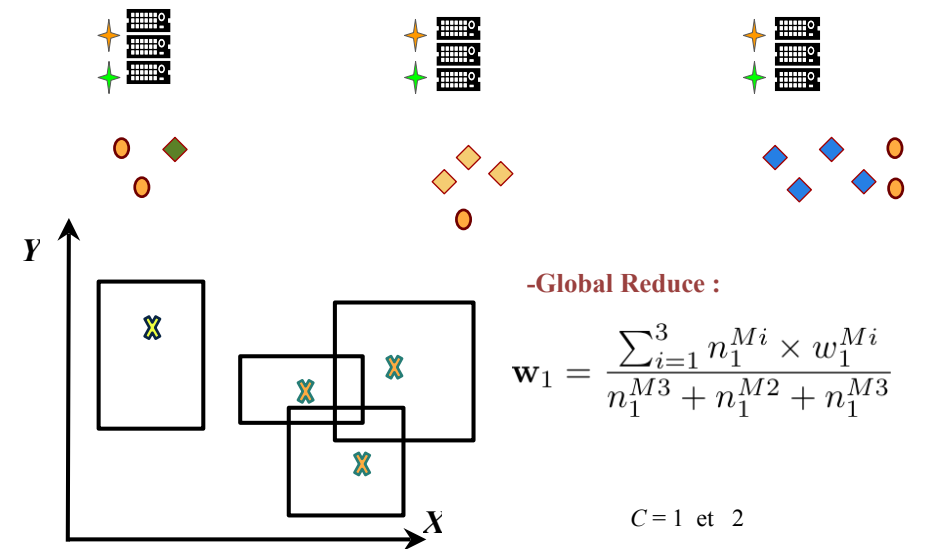
## Decomposition of K-Means (2sd solution)



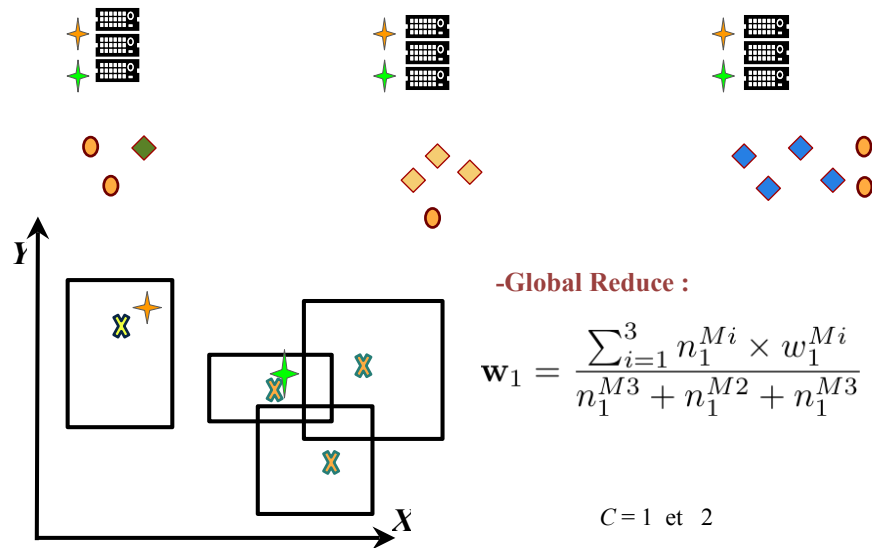
## Decomposition of K-Means (2sd solution)



## Decomposition of K-Means (2sd solution)



## Decomposition of K-Means (2sd solution)



Lebbah et al

45

## K-means, scala-Spark

```
for (i <- 1 until 10) { // 10 is the number of epoch

// Assignment step

val closest = data.map(x =>
  (closestCentroid(x, centroids), (x, 1))
)

// Quantisation step
// 1- addition of X and count the number of observation assigned
to each cluster

val pointStats=closest.reduceByKey{
  case ((x1, nb1), (x2, nb2)) => (x1 + x2, nb1 + nb2)
}
// 1- Compute the gravity centers (prototypes)
pointStats.foreach{case(id, value) =>
  centroids(id) = value._1 / value._2
}
}
```



Lebbah et al

47

## K-Means : The best solution

$$w_1 = \frac{\sum_{i=1}^3 n_1^{Mi} \times w_1^{Mi}}{n_1^{M3} + n_1^{M2} + n_1^{M3}}$$

$$w_c = \frac{\sum_1^{n_c} x_i}{n_c}$$

<Key, X>  
<N°cluster, X, I>

Do not forget we do not have  
two tuples with the same key

Initialisation : select  $K$  prototypes

**Repeat**

**-Map :** Assign each observation  $X$  to the nearest prototype  
 $X \Rightarrow (c, X, I)$  //  $c$  is the cluster ID

**-Reduce :** Addition of  $X$  belonging to the same cluster  $c$  and count  
the data

$(c, X1, I), (c, X2, I) \Rightarrow (c, X1+X2, I+1)$

// at the last reduce we obtain  $(c, \sum_{i=1}^{n_c} x_i, n_c)$

**-Collect :** Update the prototype  $w_c = \frac{\sum_1^{n_c} x_i}{n_c}$

**Until** no change

Lebbah et al

46

## K-means : MLlib - Spark

```
import org.apache.spark.mllib.clustering.{KMeans, KMeansModel}
import org.apache.spark.mllib.linalg.Vectors
```

// Load and parse the data

```
val data = sc.textFile("data/mllib/kmeans_data.txt")
val parsedData = data.map(s => Vectors.dense(s.split('
').map(_._toDouble))).cache()
```

// Cluster the data into two classes using KMeans

```
val numClusters = 2
val numIterations = 20
val clusters = KMeans.train(parsedData, numClusters, numIterations)
```

// Evaluate clustering by computing Within Set Sum of Squared Errors

```
val WSSSE = clusters.computeCost(parsedData)
println("Within Set Sum of Squared Errors = " + WSSSE)
```

// Save and load model

```
clusters.save(sc, "myModelPath")
val sameModel = KMeansModel.load(sc, "myModelPath")
```

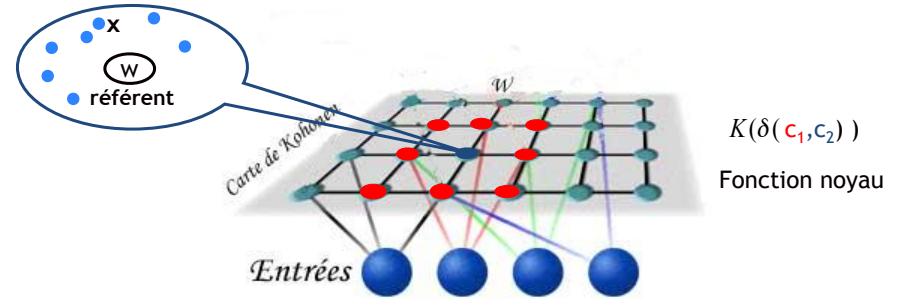
Lebbah et al

48

# First experience with MapReduce

Batch version  
SOM: Self-organizing Map

# Topological map :self-organizing Map



$\delta(c_1, c_2)$  *The short path on the graph*

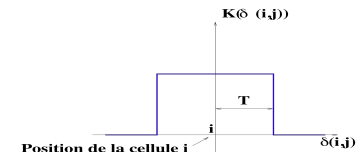
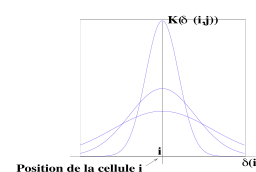
- It is a discrete space (C) of small dimension for visualization purposes (1-D, 2-D).
- C set of cells (nodes, neurons) connected by a non-oriented graph structure provided with a discrete distance  $\delta$  on C and a neighborhood structure

# SOM

# Topological Map

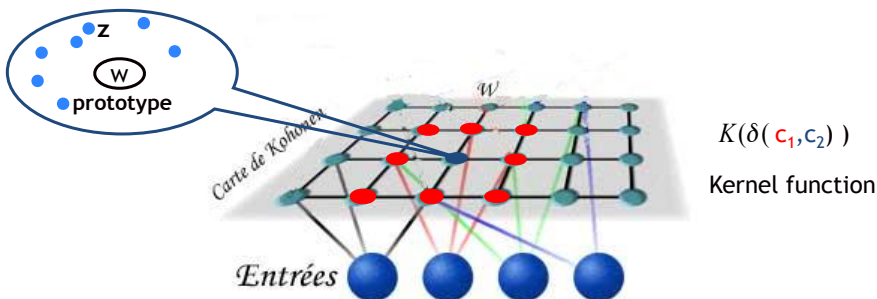
Cost Function

$$J_{som}^T(\mathcal{W}, \phi) = \sum_{\mathbf{x}_i \in \mathcal{A}} \sum_{c \in \mathcal{C}} K^T(\delta(c, \phi(\mathbf{x}_i))) \|\mathbf{x}_i - \mathbf{w}_c\|^2$$



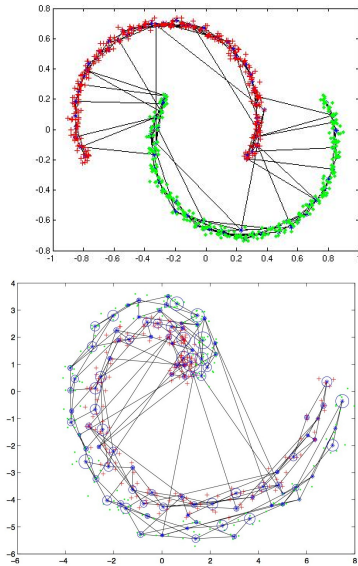
$$K^T(\delta(c, \phi(\mathbf{x}_i))) = \exp(-0.5 \frac{d}{T})$$

$V_c^T = \{r \in C / K^T(\delta(c, r)) > \alpha\}$ . The value of T determines the size of the neighborhood



$$J_{som}^T(\mathcal{W}, \phi) = \sum_{\mathbf{x}_i \in \mathcal{A}} \sum_{c \in \mathcal{C}} K^T(\delta(c, \phi(\mathbf{x}_i))) \|\mathbf{x}_i - \mathbf{w}_c\|^2$$

# Illustrations / visualisation



# Batch version

Iterative minimization  $J_{som}^T(\mathcal{W}, \phi)$  For a fixed parameter T:

## Assignment step

The set  $\mathcal{W}$  of the referents is fixed, the minimization is obtained by assigning each observation  $\mathbf{x}$  to the referent  $\mathbf{w}_c$  according to the new assignment function  $\Phi^T$

$$\phi^T(\mathbf{x}) = \arg \min_{r \in \mathcal{C}} \left( \sum_{c \in \mathcal{C}} K^T(\delta(c, r)) \|\mathbf{x}_i - \mathbf{w}_c\|^2 \right)$$

## Minimization step:

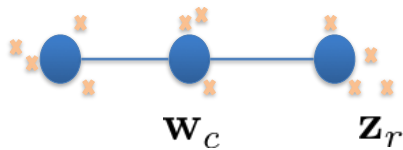
The partition is fixed  $\Phi^T$ .  $J_{som}^T(\mathcal{W}, \phi)$  is minimized with respect to all referents  $\mathcal{W}$ .

$$\mathbf{w}_c^T = \frac{\sum_{r \in \mathcal{C}} K(\delta(c, r)) \mathbf{Z}_r}{\sum_{r \in \mathcal{C}} K(\delta(c, r)) n_r}$$

$\mathbf{Z}_r$  represents the sum of all observations assigned to the cell  $r$

$n_r$ , number of observation

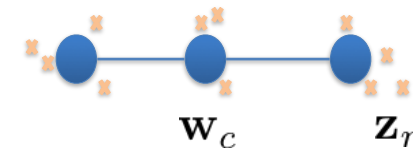
# SOM, MapReduce



$$\phi(\mathbf{x}_i)$$

$$\mathbf{w}_c = \frac{\sum_{r \in \mathcal{C}} \mathcal{K}(c, r) \sum_{\mathbf{x}_i \in \mathcal{C}_r} \mathbf{x}_i}{\sum_{r \in \mathcal{C}} \mathcal{K}(c, r) m_r}$$

# SOM, MapReduce



$$\phi(\mathbf{x}_i)$$

$$\mathbf{w}_c = \frac{\sum_{r \in \mathcal{C}} \mathcal{K}(c, r) \sum_{\mathbf{x}_i \in \mathcal{C}_r} \mathbf{x}_i}{\sum_{r \in \mathcal{C}} \mathcal{K}(c, r) m_r}$$

$$\mathbf{w}_c = \frac{\sum_{\mathbf{x}_i \in \mathcal{A}} \mathcal{K}(c, \phi(\mathbf{x}_i)) \mathbf{x}_i}{\sum_{\mathbf{x}_i \in \mathcal{A}} \mathcal{K}(c, \phi(\mathbf{x}_i))}$$

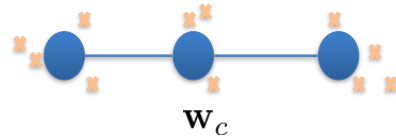






# SOM, MapReduce

$$w_c = \frac{\sum_{x_i \in \mathcal{A}} \mathcal{K}(c, \phi(x_i)) x_i}{\sum_{x_i \in \mathcal{A}} \mathcal{K}(c, \phi(x_i))}$$



Reduce Step :

**Numerator**

**Denominator**

$$\sum_{x_i \in \mathcal{A}} K(1, \phi(x_i)) x_i$$

$$\sum_{x_i \in \mathcal{A}} K(1, \phi(x_i))$$

; 1

$$\sum_{x_i \in \mathcal{A}} K(2, \phi(x_i)) x_i$$

$$\sum_{x_i \in \mathcal{A}} K(2, \phi(x_i))$$

; 2

$$\sum_{x_i \in \mathcal{A}} K(3, \phi(x_i)) x_i$$

$$\sum_{x_i \in \mathcal{A}} K(3, \phi(x_i))$$

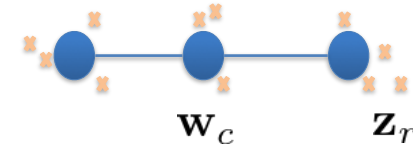
; 3

$$\sum_{x_i \in \mathcal{A}} K(C, \phi(x_i)) x_i$$

$$\sum_{x_i \in \mathcal{A}} K(C, \phi(x_i))$$

; C

# SOM, MapReduce

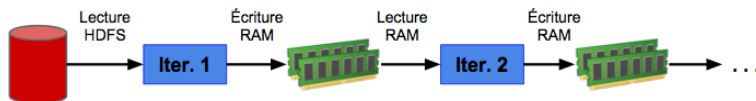


$$w_c = \frac{\sum_{r \in C} \mathcal{K}(c, r) \sum_{x_i \in C_r} x_i}{\sum_{r \in C} \mathcal{K}(c, r) m_r}$$

$$w_c = \frac{\sum_{x_i \in \mathcal{A}} \mathcal{K}(c, \phi(x_i)) x_i}{\sum_{x_i \in \mathcal{A}} \mathcal{K}(c, \phi(x_i))}$$

Numerator  
denominator

# MapReduce / Spark



SOM

Assignment

map<sub>1</sub>

Quantization

Reduce

BiTM

Assignment

map<sub>1</sub>

Quantization

Reduce

Assignment

map<sub>2</sub>

Reduce

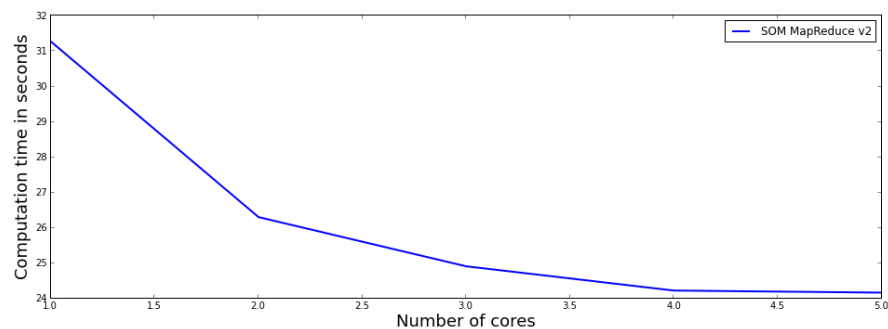
# SOM MR

```

prototypes = data.map{x =>
  val closestc = closestPrototype(prototypes, x)

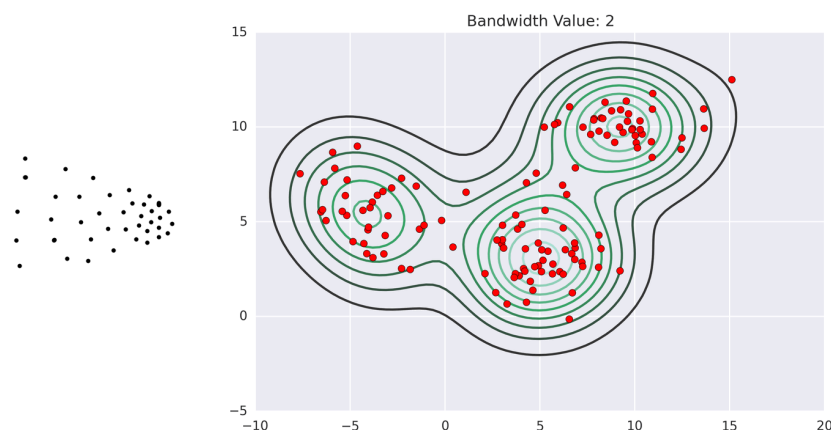
  prototypes.map{ c =>
    val t = K(c, closestc)
    (c, (x*t, t))
  }
}.reduce{ (a, b) =>
  a.zip(b).map{
    case ((aNum, aDenom), (bNum, bDenom)) =>
      (aNum+bNum, aDenom+bDenom)
  }
}.map(w => w._1 / w._2)
    
```

# Topological Map for Mixed Data



- 100 millions of observations
- SOM - Map : 10 x 10

# Illustration



# Second experience

## Mean-shift

# Modal Clustering

**Clusters** = basins of attractions of gradient ascent paths of probability density function (pdf) (Li et al, 2007, JMLR)

## Mean shift clustering

- Assume density  $f$  and gradient field  $Df$  are known
- Initial condition  $y_0$
- Update  $y_{j+1} = y_j + \mathbf{A} Df(y_j)/f(y_j)$ ,  $j = 1, 2, \dots$
- Sequence  $\{y_0, y_1, \dots\}$  converges to local mode in  $f$  (for any  $\mathbf{A}$ )

- Assume density  $f$  and gradient field  $Df$  are **unknown**
- If  $\hat{f}, D\hat{f}$  are  $k$ -nn estimators then update simplifies to

$$y_{j+1} = \frac{1}{k} \sum_{X_i \in k\text{-nn}(y_j)} X_i, \quad j = 1, 2, \dots$$

(Fukunaga & Hostetler, 1975)

- Update  $y_{j+1}$  = mean of  $k$  nearest neighbours of current iterate  $y_j$

# Modal Clustering

**Clusters** = basins of attractions of gradient ascent paths of probability density function (pdf) (Li et al, 2007, JMLR)

## Mean shift clustering

- Assume density  $f$  and gradient field  $Df$  are known
- Initial condition  $y_0$
- Update  $y_{j+1} = y_j + \mathbf{A} Df(y_j)/f(y_j)$ ,  $j = 1, 2, \dots$
- Sequence  $\{y_0, y_1, \dots\}$  converges to local mode in  $f$  (for any  $\mathbf{A}$ )
- Assume density  $f$  and gradient field  $Df$  are **unknown**
- If  $\hat{f}_j, D\hat{f}_j$  are  $k$ -nn estimators then update simplifies to

$$y_{j+1} = \frac{1}{k} \sum_{X_i \in k\text{-nn}(y_j)} X_i, \quad j = 1, 2, \dots$$

(Fukunaga & Hostetler, 1975)

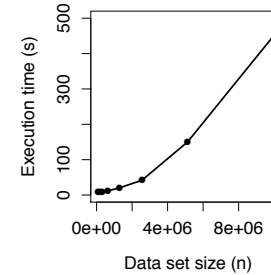
- Update  $y_{j+1}$  = mean of  $k$  nearest neighbours of current iterate  $y_j$

# Dissimilarity matrix

- Key to NNMS is  $n \times n$  dissimilarity matrix

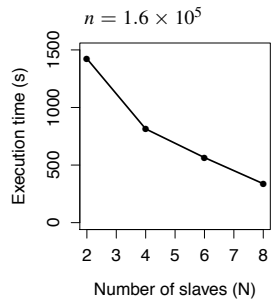
$$\begin{bmatrix} d(X_1, X_1) & d(X_1, X_2) & \dots & d(X_1, X_n) \\ d(X_2, X_1) & d(X_2, X_2) & \dots & d(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ d(X_n, X_1) & d(X_n, X_2) & \dots & d(X_n, X_n) \end{bmatrix}$$

- Dissimilarity matrix is  $O(n^2)$
- Within  $i$ -th column, sort distances
- Data point  $X_j, j \neq i$  with  $k$ -th smallest distances are  $k$ -nn to  $X_i$
- Quick sort =  $O(n \log n)$  for 1 column so  $k$ -nn is  $O(n^2 \log n)$



# Distributed way

- Each column in dissimilarity matrix can be computed independently
- Parallelise columns amongst  $N$  slave processes
- Time complexity reduced to  $O((n^2 \log n)/N)$  but still insufficient for large  $n$

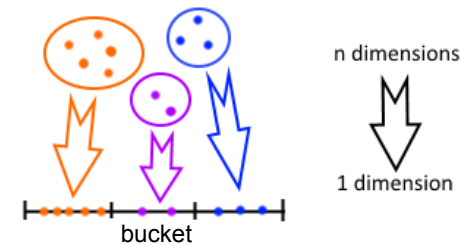


# LSH : Approximate the k-NN

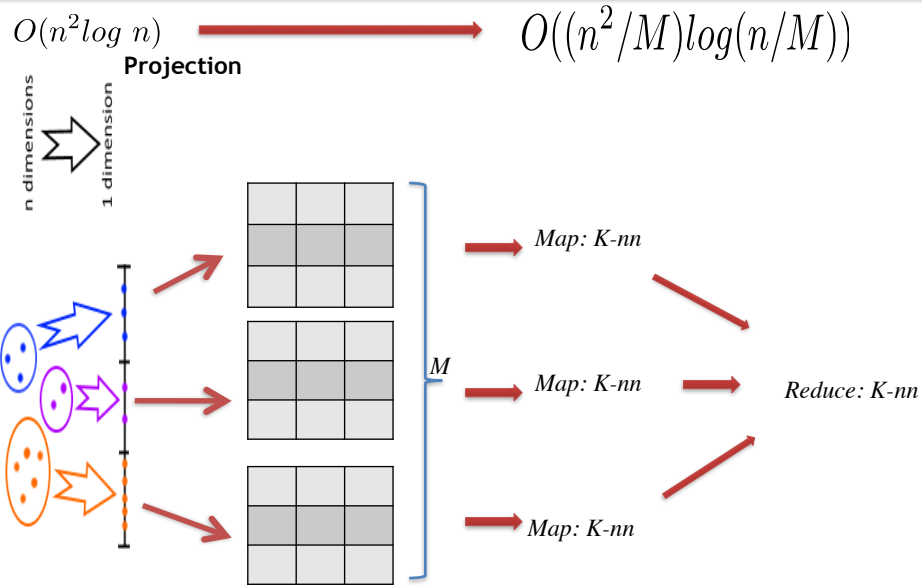
- Use approximate NN to reduce complexity from exact NN
- Locality Sensitive Hashing (LSH) based on random scalar projections of multivariate data

$$L(y) = Z^T y + U$$

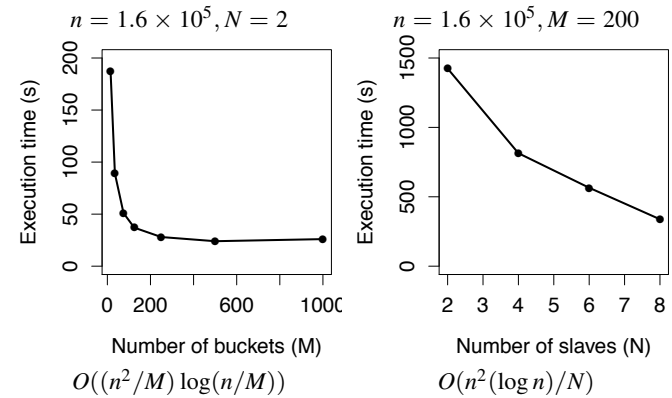
where  $Z$  is standard multivariate Gaussian,  $U$  is standard uniform (Indyk & Motwani, 1998, STOC '98)



# Illustration



- Spark Scala ecosystem
- Execution times for DNNMS-LSH-M



# Image Segmentation



Nearest neighbour mean-shift clustering (NNMS)  $O(n^2 \log n) \longrightarrow O((n^2/M) \log(n/M))$



**8h in R**   **5 min (4 cores)**  
**1min40 (24 cores)**

# Image Segmentation

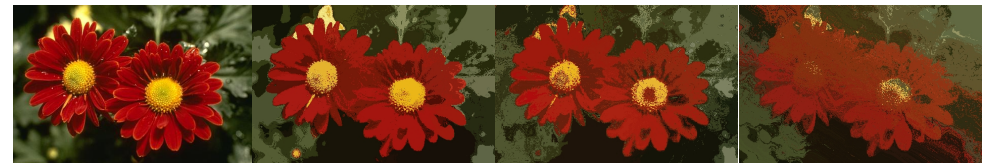


Image originale  
481x321

Image : 43 clusters  
200 buckets  
5 min (4 cores)  
1 min (20 cores)

Image : 40 clusters  
500 buckets  
8 min (4 cores)

Image : 34 clusters  
1000 buckets  
10 min (4 cores)

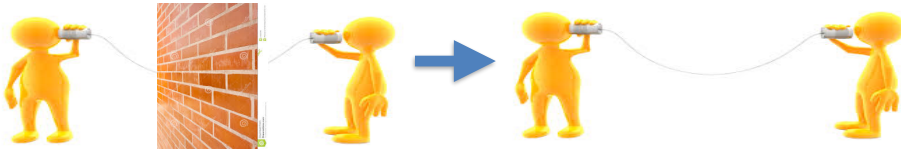


Mean Shift  
without LSH

A large number of buckets can impair quality.  
It must be adapted according to the size of the dataset.

## Conclusion : from a practical point of view

- It's time to use Scala in research and teaching. You will focus on your research problem



*With Scala and associated package, researchers, data scientists and developers can work in the same environment to build machine learning applications quickly!*

- The algorithms written with Scala are available on Spark-notebook  
<https://github.com/Spark-clustering-notebook/coliseum/wiki>



## Conclusion : from a research point of view

- The MapReduce paradigm is not the only solution! (but it helps)
- It is very important to model within the ecosystem : a new algorithm
- Your algorithm would be valid on a small scale and on a large scale

