# Università di Pisa

# Cloud Computing project:
# The $k$-means Clustering Algorithm in MapReduce

Candidati
**Alice Nannini**
**Fabio Malloggi**
**Marco Parola**
**Stefano Poleggi**

Relatore
**Dr. Nicola Tonellotto**

# Contents

# 1    Introduction

This project presents the implementation of the Kmeans algorithm based on a MapReduce version, using both the Hadoop framework both the Spark framework.

The two implementations of the kmeans algorithm developed must be performed with the following inputs:

- Name of the input file containing the dataset

- Number of centroids/clusters

- Output directory

- Number of total samples in the input dataset (the algorithm can be run assuming that you know this value)

The algorithm exit can occur due to two events:

- The maximum number of possible iteration has been reached

- The centroids calculated at i-th step and i+1-th step do not deviate beyond a certain threshold (Euclidean norm)

# 2    Dataset

The datasets for the final tests were generated with a python script, shown below and having the following format *'dataset_ numPoints_ kClusters_ dimPoints'*.

```python
import random

# inputs: n (records), k (clusters), d (dimensions)
numPoints = [1000,10000,100000]
kClusters = [7,13]
dimPoints = [3,7]


for n in numPoints:
    for k in kClusters:
        for d in dimPoints:
            # open a new file
            f = open("data/dataset_"+str(n)+"_"+str(k)+"_"+str(d)+".txt", "a")

            # compute the interval for creating the clusters
            interval = round(n/(2*k))
            count = 0
            print("dataset_"+str(n)+"_"+str(k)+"_"+str(d)+"; int: "+str(interval))

            # compute each point
            for i in range(n):
                if( (i%interval)==0 and i!=0):
                    count = count + 2

                x = ""
```

```
                for j in range(d):
                    x = x + str( interval*count + random.random()*interval )
                    x = x + " "
                x = x + "\n"
                # write the new point coordinates in the file
                f.write(x)

            f.close()
```

List of files generated from the previous code:

- dataset_100000_13_3.txt

- dataset_100000_13_7.txt

- dataset_100000_7_3.txt

- dataset_100000_7_7.txt

- dataset_10000_13_3.txt

- dataset_10000_13_7.txt

- dataset_10000_7_3.txt

- dataset_10000_7_7.txt

- dataset_1000_13_3.txt

- dataset_1000_13_7.txt

- dataset_1000_7_3.txt

- dataset_1000_7_7.txt

# 3    MapReduce pseudo-code

```
class MAPPER
  method MAP(sample_id id, sample_list l)
    for all sample s in sample_list l do
      dist = MAX_VALUE
      for all center c in cluster_centers cc do
        newDist <- computeDistance(s, c)
        if newDist < dist
          dist <- newDist
          clusterIndex = cc.index
      EMIT(index clusterIndex, sample s)



class REDUCER
  method REDUCE(index clusterIndex, samples [s1, s2,...])
    count <- 0
    center <- cluster_centers[clusterIndex]
    for all sample s in samples [s1, s2,...] do
      count <- count + 1
```

```
  for i in [0:size(s)] do
     newCenter[i] <- newCenter[i] + s[i]
for i in [0:size(newCenter)] do
  newCenter[i] <- newCenter[i] / count
EMIT(index clusterIndex, sample newCenter)
```

# 4  Hadoop Implementation

# 5  Spark Implementation

# 6  Tests and Results