

PC-2019/20 Implementazione dell'algoritmo K-means in Hadoop e OpenMP

Matteo Marulli

E-mail address

matteo.marulli@stud.unifi.it

Matteo Gemignani

E-mail address

matteo.gemignani@stud.unifi.it

Abstract

L'algoritmo di clustering K-means è uno dei più popolari algoritmi di cluster analysis usato per scoprire gruppi nei dati nel data mining e machine learning. Non è infrequente usare il Kmeans su grandi quantità di dati e in questo report proponiamo due versioni di K-means scritta una in OpenMP e l'altra in Hadoop. Per ogni versione descriveremo l'implementazione e l'analisi delle prestazioni.

Future Distribution Permission

The author(s) of this report give permission for this document to be distributed to Unifi-affiliated students taking future courses.

1. Introduzione

Il clustering è un area dell'unsupervised learning usato nel data mining come tecnica di analisi esplorativa dei dati molto utile quando i dati del dataset sono privi di una label.

Gli algoritmi che fanno parte di questa area ci consentono di organizzare i dati all'interno di gruppi senza avere alcuna precedente conoscenza delle appartenenze dei dati a tali gruppi.

I dati che finiscono in un gruppo sono molto simili tra di loro ma diversi dai dati che sono presenti in un altro gruppo.

Tra gli algoritmi di clustering il più famoso è senz'altro il K-means, la sua popolarità è dovuta alla sua semplicità e anche alla sua capacità di raggiungere buoni risultati nella pratica.

Non è raro usare il K-means su grandi

quantità di dati per questo partiremo dallo pseudo codice dell'algoritmo per analizzarlo e da questa analisi ricaveremo le implementazioni di OpenMP e Hadoop.

1.1. L'algoritmo

In questa sezione è riportato lo pseudo-codice del K-means.

Algorithm 1: K-means

Input: $\mathbb{D} \in \mathbb{R}^{n \times d}$, $k \in \mathbb{N} - \{0, 1\}$, $\epsilon \in \mathbb{R}$

$t = 0$

inizializza in modo casuale k centroidi

iniziali: $\mu_1^t, \mu_2^t, \dots, \mu_k^t \in \mathbb{R}^d$

repeat

$t \leftarrow t + 1$

$C_j \leftarrow \emptyset$ per ogni $j = 1, \dots, k$

 //Assegnamento ai
 cluster

foreach $x_j \in \mathbb{D}$ **do**

$j^* \leftarrow \operatorname{argmin}_i \{ \|x_j - \mu_j^t\|^2 \}$

 // Si assegna x_j al
 centroide più vicino

$C_{j^*} \leftarrow C_{j^*} \cup \{x_j\}$

end

 // Aggiornamento
 centroidi

foreach $i = 1, \dots, k$ **do**

$\mu_i^t \leftarrow \frac{1}{|C_i|} \sum_{x_j \in C_i} x_j$

end

until $\sum_{i=1}^k \|\mu_i^t - \mu_i^{t-1}\|^2 \leq \epsilon$;

Le parti più critiche sono la scelta dei centroidi iniziali, l'assegnamento di un punto ad

un cluster e l'aggiornamento dei centroidi.

Le tre operazioni saranno implementate in OpenMP e Hadoop.

1.2. Implementazione in OpenMP

Fra l'implementazione sequenziale e quella parallela non ci sono differenze tranne l'impostazione del numero dei Threads nel main che, nel caso dell'implementazione sequenziale viene impostata a 1.

1.2.1 Inizializzazione centroidi

Per OpenMp è stato scelto di creare i centroidi iniziali con questa strategia:

```
1 void KmeansOpenMP::initCentroids(double**  
    dataMatrix, int nRow, int nCol){  
    centroids = new double*[k];  
3  
    for (int i = 0; i < k; i++) {  
5        centroids[i] = new double[nCol];  
    }  
7    srand(seed);  
    for(int i = 0; i < k; i++){  
9        int row = int(rand() % nRow);  
        for(int j=0; j < nCol; j++){  
11            centroids[i][j] = dataMatrix[row  
                ][j];  
        }  
13    }  
}
```

initCentroids() sceglie casualmente k righe della matrice dei dati iniziale e queste saranno i centroidi di partenza.

Il costo di initCentroids() è $O(kd)$.

1.2.2 Assegnamento cluster

L'assegnamento di un punto ad un cluster è uno dei passi più costosi dell'algoritmo. Per ogni punto bisogna calcolare k distanze, cioè una per ogni centroide, quindi il costo di questa operazione diventa $O(nk)$.

```
int KmeansOpenMP::clusterAssignment(  
    double* point, double** centroids,  
    int nCol){  
2    int index = 0;  
    double distance = 0;
```

```
4    double oldDistance = DBL_MAX;  
    for(int i = 0; i < k; i++){  
6        double* centroid_i = centroids[i];  
        distance = getDistance(point,  
            centroid_i, nCol);  
8        if(distance < oldDistance){  
            index = i;  
            oldDistance = distance;  
10        }  
12    }  
    return index;  
14 }
```

Nel codice clusterAssignment() viene invocata la funzione che calcola la distanza tra due punti che ha costo $O(d)$, pertanto il costo dell'intero passaggio di assegnamento di un punto ad un cluster costa $O(knd)$

```
double KmeansOpenMP::getDistance(double *  
    point, double *centroid, int nCol) {  
2    double distance = 0;  
    for (int i = 0; i < nCol; i++) {  
4        distance += pow(point[i] - centroid  
            [i], 2);  
    }  
6    return sqrt(distance);  
}
```

Usando OpenMp la funzione clusterAssignment è stata parallelizzata nel seguente modo

```
1 #pragma omp parallel for schedule(static,  
    nRow/nThreads) default(none) shared(  
    labels) firstprivate(nThreads,  
    centroids, nCol, nRow, dataset)  
    for (int indexPoint = 0; indexPoint  
        < nRow; indexPoint++){  
3        labels[indexPoint] =  
            clusterAssignment(dataset[  
                indexPoint], centroids, nCol  
            );  
    }
```

Il codice è stato parallelizzato dividendo in modo equo il carico di lavoro tra i thread in modo statico.

I thread condividono l'array labels per riportare il risultato dell'assegnamento al cluster per ogni punto.

Inoltre, i thread, hanno una copia locale, inizializzata al valore del thread principale, delle variabili che identificano il numero dei

threads, i centroidi, il numero righe, il numero di colonne, e il dataset.

Per fare ciò è stata utilizzata la direttiva `firstprivate` così da evitare una copia locale inizializzata a 0 portando i thread ad non entrare nei cicli `for`.

1.2.3 Aggiornamento centroidi

L'aggiornamento dei centroidi è stato parallelizzato in base al numero di cluster:

```
void KmeansOpenMP::centroidUpdate(int *
    labels,int* labelsCount,double **
    centroids,double **dataset,int nCol,
    int nRow){
2   int pointForCluster;
    #pragma omp parallel default(none) shared
        (centroids,labelsCount) firstprivate
            (nCol,nRow,labels,dataset)
4   {
        #pragma omp for collapse(2)
6       for(int c=0; c<k;c++){
            for(int i=0;i<nCol;i++){
8                 centroids[c][i]=0; // azzerà il
                    centroide c-esimo
10            }
        #pragma omp for
12        for(int i=0;i<k;i++){
            labelsCount[i]=0;
14        }

16 #pragma omp for
        for(int c=0;c<k;c++){
18             for(int j=0;j<nRow;j++){
                if(c==labels[j]){
20                 labelsCount[c]++; // conta il
                    numero di elementi nel
                    cluster c-esimo
                for(int i=0;i<nCol;i++){
22                     centroids[c][i]+=dataset[j]
                        [i]; // aggiornamento
                            centroide c-esimo
                                componente per
                                    componente
24                 }
            }
26        }
    }
28 #pragma omp for collapse(2)
        for(int c=0; c<k;c++){
30             for(int i=0;i<nCol;i++){
```

```
centroids[c][i]=centroids[c][i]/
    labelsCount[c]; // dividi
    ogni componente per il
    numero di oggetti nel
    cluster
```

```
32     }
    }
34 }
}
```

Purtroppo non è stato possibile parallelizzare l'istruzione più costosa(Ciclo `for` Rig. 18 del codice precedente) causa di una dipendenza dovuta da `labelsCount[c]++`.

1.3. Implementazione in Hadoop

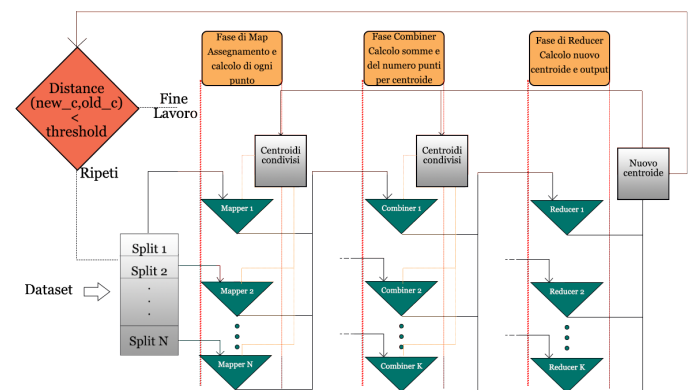


Figure 1. Flowchart dell'implementazione di K-means in Hadoop

Per Hadoop è stata scelta un'altra strategia di inizializzazione dei centroidi

```
private void initCentroids() {
2   for(int i=0;i<this.k;i++){
        double[] center= new double[this
            .nb_dimensions];
4       for(int j=0;j<this.nb_dimensions
            ;j++){
            center[j]=Math.floor(
                -1000+2000*random.
                    nextDouble())/100 ;
6         this.centroids.add(center);
8     }
}
```

Questa volta, invece di estrarre casualmente dei punti dal dataset, è stato deciso di generarli in modo aleatorio, evitando di leggere tutto il

dataset in modo sequenziale.

Per quanto riguarda le operazioni di assegnamento di un punto ad un cluster e l'aggiornamento dei centroidi, sono state rivisitate in modo da applicare il paradigma di programmazione di Hadoop: il "MapReduce", come mostrato in figura 1.

Il metodo di Map processa i punti del dataset in modo parallelo per coppie <key, value>, ugualmente i metodi Reduce e Combine lavorano in parallelo le coppie <key, List<value>> dove questa volta le chiavi saranno i centroidi e avranno associate delle liste di punti.

Per i tipi di key e value sono state usate le classi messe a disposizione da Hadoop, evitando la scrittura di una possibili nuove classi, Point e/o Centroids, che avrebbero dovuto implementare l'interfaccia WritableComparable e suoi relativi metodi.

1.3.1 Map

```
@Override
2   public void map(Object key, Text value
    , Context context) throws
      IOException, InterruptedException
    {
      String line = value.toString();
4     List<Double> coordinates = new
        ArrayList<Double>();
      StringTokenizer tokenizer = new
        StringTokenizer(line, ",");
6     while(tokenizer.hasMoreTokens())
        coordinates.add(Double.
            parseDouble(tokenizer.
                nextToken()));

8
      double distance = Double.MAX_VALUE;
10     double distanceTmp;
      long index = -1;
12     for(int i = 0; i<centroids.size();i
        ++){
          distanceTmp = getDistance(
              coordinates, centroids.get(i
                ));
14     if(distanceTmp<distance){
          distance = distanceTmp;
16     index= i;
```

```
    }
18   }
    if(index!=-1)
20     context.write(new LongWritable(
        index), value);
    else
22     logger.fatal("\n\nNessun cluster
        vicino trovato? min = "+
            distance+" coordinate = "+
            coordinates+"\n\n\n");
    }
}
```

Il metodo map esegue l'assegnamento di un punto ad cluster come la funzione clusterAssignment di OpenMP. Come possiamo vedere inizialmente il punto è salvato in una variabile Text e quindi per essere utilizzato, come un vettore di numeri reali, deve essere letto e trasformato in un Double.

Poi, esattamente come in C++, si calcola la distanza tra i centroidi e il punto estratto, scegliendo il centroide più vicino.

Infine, si passa il risultato al Combine scrivendo sul file contex l'indice del centroide più vicino e la lista delle coordinate del punto.

1.3.2 Combiner

La classe Combiner ha lo scopo di ridurre il trasferimento dati da Map a Reduce dal momento che i dati sono salvati sul disco locale dell'host.

Il Combiner prende una la lista di punti associati ad un cluster e ne calcola la somma parziale delle coordinate e la quantità.

```
1   public static int nb_dimension;
3
4   @Override
5   public void reduce(LongWritable key,
        Iterable<Text> values, Context
        context) throws IOException,
        InterruptedException {
6     double[] sum = new double[
            nb_dimension];
7     int numPoint = 0;
9     for(int i = 0; i<nb_dimension;i++)
```

```

11         sum[i] = 0;
12     for(Text record: values){
13         String line = record.toString();
14         List<Double> coordinates = new
15             ArrayList<Double>();
16         StringTokenizer tokenizer = new
17             StringTokenizer(line, ",");
18         while(tokenizer.hasMoreTokens())
19             coordinates.add(Double.
20                 parseDouble(tokenizer.
21                     nextToken()));
22
23         for(int i = 0; i<nb_dimension;i
24             ++){
25             sum[i] += coordinates.get(i);
26             numPoint++;
27         }
28         StringBuilder stringBuilder = new
29             StringBuilder();
30         for(int i = 0; i<nb_dimension;i++)
31             {
32                 stringBuilder.append(sum[i]);
33                 stringBuilder.append(",");
34             }
35         stringBuilder.append(numPoint);
36         context.write(key, new Text(
37             stringBuilder.toString()));
38     }

```

Infine passa al Reduce, sempre scrivendo sul file context, una coppia chiave valore, dove la chiave è ancora il centroide e il valore è un testo con la semisomma della coordinate e il numero di punti associati al centroide.

1.3.3 Reduce

Il Reduce calcola i nuovi centroidi sommando separatamente tutte le somme parziali delle coordinate e quelle dei numero di punti associati al centroide.

Fatto ciò, per ottenere il valore del nuovo centroide, il Reducer divide la somma ottenuta per il numero totale dei punti nel cluster.

Infine scrive sul file context la coppia indice del centroide e coordinate del centroide.

```

1 public static int nb_dimension;
2 public static List<double[]> centroids
3     = new ArrayList<>();

```

```

3 public static DecimalFormat dFormatter
4     = new DecimalFormat("#.###");
5
6 @Override
7 public void reduce(LongWritable key,
8     Iterable<Text> values, Context
9     context) throws IOException,
10     InterruptedException {
11     double[] center = new double[
12         nb_dimension];
13     int numPoint = 0;
14
15     for(int i = 0; i<nb_dimension;i++)
16         center[i] = 0;
17
18     for(Text record: values){
19         String line = record.toString();
20         List<Double> coordinates = new
21             ArrayList<Double>();
22         StringTokenizer tokenizer = new
23             StringTokenizer(line, ",");
24         for(int i = 0; i<nb_dimension;i
25             ++){
26             center[i] += Double.
27                 parseDouble(tokenizer.
28                     nextToken());
29             numPoint += Integer.parseInt(
30                 tokenizer.nextToken());
31         }
32
33         for(int i = 0; i<nb_dimension;i++)
34             center[i] /= numPoint;
35
36         StringBuilder center_sb = new
37             StringBuilder();
38         for(int i = 0; i< nb_dimension;i++)
39             {
40                 center_sb.append(dFormatter.
41                     format(center[i]));
42                 if(i< nb_dimension-1)
43                     center_sb.append(" ");
44             }
45         centroids.set((int)key.get(),
46             center);
47         context.write(key, new Text(
48             center_sb.toString()));
49     }

```

1.3.4 Driver

La classe driver contiene il metodo run che permette di eseguire il Kmeans.

Il metodo initJob contiene la configurazione del Job che determina le istruzioni che

permettono l'esecuzione del processo di Map-Reduce.

```
1  private void run() throws IOException,
    ClassNotFoundException,
    InterruptedException {
3      do{
        centroids.clear();
5        centroids.addAll(UpdateCentroids
            .centroids);
        ClusterAssignment.centroids.
            clear();
7        ClusterAssignment.centroids.
            addAll(centroids);
        UpdateCentroids.centroids.clear
            ();
9        UpdateCentroids.centroids.addAll
            (centroids);

11       initJob();
        deleteOutputDirIfExists();
13
        if(!job.waitForCompletion(true))
15             logger.info("job failed");
        n_iter++;
17
        logger.info("\n#iter = "+n_iter+
            "\n");
19        logger.info("\tOld_centers: "+
            centroidsToString(centroids)
            +"\n");
        logger.info("\tNew_centers: "+
            centroidsToString(
                UpdateCentroids.centroids)+"
            \n");
21
    }while(compare(centroids,
        UpdateCentroids.centroids));
23    logger.info("Finished in "+ n_iter
        + " iterations");
}
```

Infine, nel metodo run, viene effettuata una "pulitura" delle strutture dati static che contengono i centroidi nei metodi Map e Reduce; vengono svuotate ad ogni iterazione usando il metodo clear() per far posto ai nuovi centroidi che vengono aggiunti con l'istruzione addAll().

1.4. Avvio del programma in Hadoop

Per avviare il programma da terminale bisogna passare gli argomenti nel seguente ordine:

1. path dell'input
2. path dell'output
3. k che indica il numero di cluster desiderato
4. ϵ che indica la tolleranza

2. Analisi dei risultati

2.1. Hardware

Per condurre le analisi delle prestazioni è stata usata una macchina MSI con le seguenti specifiche:

- **Processore:** Intel® Core™ i7-10750H CPU @ 2.60GHz 6 core
- **Ram:** 16GB DDR4 3200Mhz
- **Hardisk:** SSD 1TB

2.2. Dataset di prova

Per testare il codice, è stato creato un dataset artificiale usando sklearn. Il metodo utilizzato è make_blob ed con i seguenti parametri:

- n_samples= 10.000.000;
- n_features= 10;
- centers= 10;

2.3. Parallelo Vs sequenziale

Viene riportata soltanto l'analisi delle prestazioni OpenMP in quanto per Hadoop servirebbero molti dati per poter osservare uno speedUP significativo.

Come è possibile vedere dalla Figura 2 la versione OpenMP si mantiene allineata con lo speedUP lineare fino a 3 thread per poi man mano distanziarsi leggermente fino a 6 thread.

Dopo 6 thread diverge completamente diventando sub-lineare, si può quindi concludere

che l'implementazione è al massimo 6 volte più veloce della versione sequenziale.

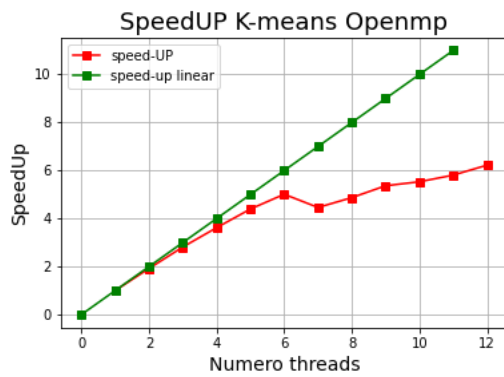


Figure 2. SpeedUP della versione OpenMP di Kmeans