

# Co-clustering with High Machine Scalability

April 29, 2015

## Abstract

*DisCo*, the clustering method using MapReduce, is introduced in 2008. It requires the main thread to collect the assignment of cluster to all rows and columns, and broadcast the assignment to all mappers before the each MapReduce task starts. The time complexity of this process is only related with number of rows and column. Therefore, although the matrix is sparse or not, gathering and broadcasting the cluster assignment takes  $O(m + n)$  times. In this paper, we introduce a improved version of DisCo, by deleting it's bottleneck processes.

## 1 Introduction

Clustering is one of the key issues of graph mining to find valuable knowledges from the raw data. As the graph size is getting bigger, the methods to manage these big dataset is also needed. The *DisCo* is one of the clustering methods managing big data by using Map-Reduce introduced in 2008. Each task of DisCo is consisted with broadcasting parameter  $r$  and  $c$ , MapReduce job, modify  $r(i)$ ,  $c(j)$  for all rows and columns. As  $m$  and  $n$  grows, the most burden part is broadcasting the  $r(i)$ ,  $c(j)$ . The biggest problem of DisCo is, broadcasting process has no machine scalability. Whether the matrix is sparse or not, amount of information that the main thread should broadcast is same as  $O(m + n)$ . Therefore the machine scalability is very poor as the row and column dimension getting bigger.

Our main contribution is giving DisCo to have higher machine scalability, i.e. removing broadcasting process. By this method, we removed the broadcasting, gathering process from the all MapReduce tasks, and having good machine scalability.

## 2 Background : DisCo

DisCo is consisted with two main phases. First phase is pre-processing, making row, column adjacency list as sequence file on HDFS as follows.

Key :  $i$   
Value :  $radj_i$

Second phase is Co-Clustering process. Before the DisCo Co-Clustering the row once, it broadcasts  $S$ ,  $c$  to all mapper. The space complexity of the parameters is  $O(k \times l + n)$ . Because the broadcasting process has no machine scalability, as the  $n$  getting bigger, the time to broadcasting parameters is higher whether the adjacency matrix is sparse or environment has many machines. This is main disadvantage of DisCo. Also, the CollectResult() is conducted by main thread which requires  $O(m)$  calculation.

Symbol	Definition
$A$	the $m$ by $n$ matrix
$m, n$	Number of rows and columns
$i, j$	Row and Column indices
$radj, cadj$	Row and column adjacency list
$radj_i, cadj_j$	Row adjacency list of row $i$ , column adjacency list of column $j$
$a_{i,j}$	$(i, j)$ element of $A$
$r, c$	Cluster assignment of row, column
$r(i)$	Cluster assignment of row $i$
$c(j)$	Cluster assignment of column $j$
$I(v)$	Set of row indices with $r(i) = v$
$J(w)$	Similar to $I(v)$ , but for columns
$R, C$	Size of row, column cluster
$R(v)$	Size of row cluster $v$
$C(w)$	Size of column cluster $w$
$S$	$k \times l$ statistic matrix
$k, l$	Size of row, column cluster, respectively
$v, w$	row and column cluster index
$S_{v,w}$	The $(v, w)$ element of $S$
$S_{v,:}$	$v$ th row vector of $S$
$S_{:,w}$	$w$ th column vector of $S$
$rst_i$	Row Statistic vector of size $l$ driven by $radj_i$ and $c$
$cst_j$	Similar to $rst_i$ , but for columns
$Rdata$	HDFS sequence file of row information
$Cdata$	HDFS sequence file of column information
$Rdata_i$	$(key, value)$ pair of sequence file $Rdata$ . $(key, value) = (i, (r(i), rst_i, radj_i))$
$Cdata_j$	$(key, value)$ pair of sequence file $Cdata$ . $(key, value) = (j, (c(j), cst_j, cadj_j))$

Table 1: Definition of Symbols

---

**Algorithm 1: CC( $A, k, l$ )**

---

```
1 Initialize  $r$  and  $c$ 
2 preprocess()
3 repeat
4   CC-MapReduce() begin
5     Broadcast  $c, S$  to all mapper.
6     CCRowMapper(key  $key$ , value  $value$ ) begin
7        $(key, value) = (i, rad_i)$ 
8        $rst_i = \text{row statistic}(rad_i, c)$  for  $v = 1..k$  do
9         if assigning  $i$  to  $v$  would lower cost then
10           $r(i) = v$ 
11        end
12      end
13       $Out(r(i), (rst_i, i))$ 
14    end
15    CCRowReducer(key  $key$ , values  $V$ )
16    begin
17       $v = key$ 
18       $S_v = \text{empty array of size } l$   $I_v = \emptyset$  for  $(rst_i, i) \in V$  do
19         $S_v = \text{combine statistics}(S_v, rst_i)$ 
20         $I_v = I_v \cup i$ 
21      end
22       $Out(v, (S_v, I_v))$ 
23    end
24  end
25  CollectResults() begin
26    /* Conducted by main thread */
27     $S = k$  by  $l$  empty matrix
28    for reducer output  $(v, (S_v, I_v))$  do
29       $S_{v,:} = \text{combine statistics}(S_{v,:}, S_v)$  for  $i \in I_v$  do
30         $r(i) = v$ 
31      end
32    end
33    Do the same for columns
34 until cost does not decrease
35 return  $r$  and  $c$ 
```

---

### 3 Co-Clustering without broadcasting

Our method is consisted with two phases. First is preprocessing the data. The difference between DisCo and our method in preprocessing data is that the initialized  $r, c$  is broadcasted to all mapper to set  $r(i)$  and calculate  $rst_i$  for all row index  $i$ .

Key :  $i$   
Value :  $r(i), rst_i, rad_i$

It is same for column. These two MapReduce tasks are the *only tasks* requires  $r$  and  $c$  need to be broadcasted.

The second phase is Co-Clustering. The difference between DisCo and our method is the parameter should be maintained, gathered and broadcasted is only  $S$ . The space complexity of the parameter is  $O(k \times l)$  while

the DisCo requires  $O(n)$  data to run the row coclustering task. The key idea to remove broadcasting  $c$  from all Co-Clustering the row is make the data already has the information of  $rst_i$ . In the line 6 of Algorithm1,  $rst_i$  is calculated by  $radj_i$  and  $c$ . If the data already has  $radj_i$ , then broadcasting  $c$  is unnecessary.

In summary, row co-clustering requires  $rst_i$  to change  $r(i)$  of all rows and  $cst_j$  of all columns. Our key idea is is let the data already has  $rst_i$ , and make the result change  $r(i)$  and  $cst_j$  within one row coclustering task. For acheiving this, our method contains one Map-only task, two MapReduce tasks.

1. CC-Mapper : Map-Only Tasks. Takes  $Rdata$ ,  $Cdata$  and return  $Rdata$  and following three types of data for oncoming MapReduce tasks
  - (a)  $S_{part}$  : set of  $(r(i), (1, rst_i))$  to update  $S$ ,  $R$
  - (b)  $Rdata_{part}$ : set of  $(r(i), radj_i)$  to calculate  $cst_j$  of all column  $j$
  - (c)  $Cdata_{part}$ : set of  $(j, (c(j), cadj_j))$  to reconstruct  $Cdata$
2. CC-Constructor : MapReduce Task. Takes  $Rdata_{part}$  and  $Cdata_{part}$  to construct  $Cdata$
3. CC-Collector : MapReduce Task. Takes  $S_{part}$  to update  $S$  and  $R$

---

**Algorithm 2: CC-new( $A, k, l$ )**

---

```

1 repeat
2   CC-Mapper()
3   CC-Constructor()
4   CC-Collector()
5   Do the same for columns
6 until cost does not decrease
7 return  $r$  and  $c$ 

```

---

**Algorithm 3: CC-Mapper**

---

```

1 CC-new-Mapper(key  $key$ , value  $value$ )/ * Map only task * /
   Input:  $Rdata_i, Cdata_j$ 
   Output:  $Rdata, S_{part}, Rdata_{part}, Cdata_{part}$ 
2 begin
3   if data form is  $Rdata$  then
4      $(i, (r(i), rst_i, radj_i)) = (key, value)$ 
5     for  $v = 1..k$  do
6       if assigning  $i$  to  $v$  would lower cost then
7          $r(i) = v$ 
8       end
9     end
10     $Rdata_i = (r(i), rst_i, radj_i)$ 
11    Out( $i, Rdata_i$ ) to  $Rdata$ 
12    Out( $r(i), (1, rst_i)$ ) to  $S_{part}$ 
13    Out( $r(i), (i, radj_i)$ ) to  $Rdata_{part}$ 
14  end
15  else
16     $(j, (c(j), cst_j, cadj_j)) = (key, value)$ 
17    Out( $j, (c(j), cadj_j)$ )
18  end
19 end

```

---

Note that the line 28 of Algorithm 4 calculate  $cst_j$  with  $Bag$ , which is hashmap having tuple as (cluster id, set of indices). DisCo calculate  $cst_j$  using  $cadj_j$  and  $r$ , but most of Co-Clustering algorithm such as *CrossAs-*

sociation needs only  $Bag$  to calculate  $cst_j$ .

---

**Algorithm 4: CC-Constructor**

---

**Input:**  $Rdata_{part}, Cdata_{part}$

**Output:**  $Cdata$

---

```

1 begin
2   CC-Const-Mapper(key  $key$ , value  $value$ ) begin
3     if data has form  $Rdata_i$  then
4        $(r(i), radj_i) = (key, value)$ 
5       for  $j \in radj_i$  do
6          $Out(j, (r(i), i))$ 
7       end
8     end
9     else
10       $Out(key, value)$ 
11    end
12  end
13  CC-const-Reducer(key  $key$ , values  $Values[]$ )
14  begin
15     $key = j$ 
16     $Bag$  = empty HashMap.
17     $cst_j$  =empty array of size  $k$ 
18    for  $value \in Values$  do
19      if  $value = (c(j), cadj_j)$  then
20        continue
21      end
22      else
23         $(v, i) = value$ 
24        add  $i$  to  $Bag.get(v)$ 
25      end
26    end
27     $Bag$  became set of tuple  $(v, \text{set of row index } i \text{ with } r(i) = v)$ 
28     $cst_j = column\ statistic(Bag)$ 
29     $Cdata_j = (c(j), cst_j, cadj_j)$ 
30     $Out(j, Cdata_j)$  to  $Cdata$ 
31  end
32 end

```

---

---

**Algorithm 5: CC-Collector**

---

**Input:**  $S_{part}$   
**Output:**  $S, R$

```
1 begin
2   CC-collect-Mapper(key  $key$ , value  $value$ ) begin
3      $(r(i), (1, rst_i)) = (key, value)$  for some row  $i$ 
4      $x =$  arbitrary value of 0 to 1
5      $key = (r(i) + x) \times NumOfReducer$ 
6     /* To make data spreads to all reducer equivalently */
7      $Out(key, value)$ 
8   end
9   CC-collect-Reducer(key  $key$ , values  $Values[]$ )
10  begin
11     $v = (key - (key \% NumOfReducer)) \div NumOfReducer$ 
12     $S_{v,:} =$  empty array of size  $l$ 
13     $R(v) = 0$ 
14    for  $value \in Values$  do
15       $value = (1, rst_i)$  for some row  $i$ 
16       $R(v) = R(v) + 1$ 
17       $S_{v,:} = combine\ statistics(S_{v,:}, rst_i)$ 
18    end
19     $Out(v, (R(v), S_{v,:}))$  to  $S$ 
20  end
21  Combine all reducer output to  $R, S$ 
22 end
23 Update  $R, S$ 
24 Update  $H$ 
```

---

## 4 Experiment

Our method require several intermediate data, but showed good machine scalability. DisCo also showed the good machine scalability in MapReduce task and it is faster than our methods if we compare only MapReduce task. But, DisCo takes a lot of time to broadcast the parameter as  $m, n$  gets bigger, even if the MapReduce task takes less time than ours. The number of machines, density of the data are the two factor of the Co-Clustering to decide which method to choose. The conclusion is, if the data is sparse and the environment has lots of machines, our method is better than DisCo.

## 5 Colclusion