

FunctionDescription of AdjClustBand_heap

Parismita Das

25 March 2017

General Description for adjClustBand_heap Implementation

The adjClustBand_heap function takes in the Cosine Similarity Matrix, diagonal elements along with its width and no of rows as input.

As constrained Hierarchical Clustering has a time complexity of $O(n^2)$ which is very costly. It took 22 minutes to calculate 32768 SNPs. Furthermore, it takes 4.5 hours (on a standard 2.2 Ghz single CPU) to analyze a whole genome of 500k simulated SNPs (for Affymetrix 500k arrays) genotyped on 100 individuals.

Hence the quadratic complexity of the adjacency-constrained hierarchical clustering remain unsatisfactory for many reasons.

To avoid this situation We introduce of a parameter h aims to control the maximum lag between items X_i and X_j for similarity calculations. Thus, the similarity measures are computed between these two items only if i and j differ by no more than h . To implement this we use a pencil trick which calculates the distance in complexity of $O(1)$

Furthermore we use MinHeap for finding the best fusion. Thus reduce the complexity from $O(p^2)$ to $O(p \log p)$.

In the adjClustBandheap.R file, the input is converted into a left and right matrix of $p \times h$ size. And Heap of $p-1$ size is build. The output is a hclust class object, Which gives similar results as rioja chclust(coniss linkage) and with better optimisation.

The Algorithm Used

Algorithm 8 The **optimized** algorithm of the Ward's constrained hierarchical clustering applied to a h -band similarity matrix

```

1: procedure CWARD( $\mathbf{X} \in \mathbb{R}^{n \times p}$ , Sim,  $h$ )
2:   Calculate the two  $p \times h$  arrays of pencils sums  $\triangleright \mathcal{O}(ph)$ 
3:   Initialize the chained array  $Tab$ 
4:    $heap \leftarrow \text{buildHeap}(1 : (p - 1), D)$   $\triangleright \mathcal{O}(p \log(p))$ 
5:    $jj \leftarrow p$ 
6:   for  $step = 1$  to  $p - 1$  do
7:     while ( $\neg Tab[valid, heap[1]]$ ) do
8:        $heap \leftarrow \text{deleteMin}(heap)$   $\triangleright \mathcal{O}(\log(p))$ 
9:     end while
10:     $posMin \leftarrow heap[1]$ 
11:     $i^* \leftarrow Tab[Cl1, posMin]$ 
12:     $heap \leftarrow \text{deleteMin}(heap)$   $\triangleright \mathcal{O}(\log(p))$ 
13:     $d_1 \leftarrow D(C_{i^*-1}, C_{i^*} \cup C_{i^*+1})$   $\triangleright \mathcal{O}(1)$ 
14:     $d_2 \leftarrow D(C_{i^*} \cup C_{i^*+1}, C_{i^*+2})$ 
15:    Add the distances  $d_1$  and  $d_2$  to  $Tab$ 
16:     $heap \leftarrow \text{insertHeap}(heap, jj, D)$   $\triangleright \mathcal{O}(\log(p))$ 
17:     $heap \leftarrow \text{insertHeap}(heap, jj + 1, D)$   $\triangleright \mathcal{O}(\log(p))$ 
18:    Update the neighbors of  $C_{i^*-1}$  and  $C_{i^*+2}$  in  $Tab$ 
19:    Set  $Tab[valid, posMin]$ ,  $Tab[valid, posL]$ 
      and  $Tab[valid, posR]$  to FALSE
20:     $jj \leftarrow jj + 2$ 
21:  end for

```

Function Description

- `percDown` : Rebalancing of elements after insertion, so that the Min-Heap property is maintained.
- `deleteMin_C` : Delete the minimum position after merging of the clusters.
- `insertHeap_C` : Insert the cluster after obtaining the minimum distance from `distance_C`
- `neighborCl_C` : Give the adjacent position of cluster.
- `neiNeighborPos_C` : Give the neighbour position of adjacent position of cluster.
- `pencil_C` : Calculate the pencil shaped area by using partial sum of similarity, used for distance calculation.
- `distance_C` : Calculate distance using Ward's criteria using Lance-William method and pencil implementation.
- `cWardHeaps` : The actual function that clusters according to the given similarity matrix using Min-heap for merging of clusters and ward's criteria for distance calculation required for merging of clusters.
- `.toMatL` : creates $p \times h$ matrix to calculate partial sum of similarity

- `.toMatR` : creates $p \times h$ matrix to calculate partial sum of similarity

The output has attributes such as `traceW`, `gain`, `merge`, `height`, `label`, `order` which is calculated in `cWardHeaps` function. That are used to create a 'hclust' object in the R interface of the implementation.

The Functions Dependency Graph

