

Robust Single Linkage Algorithm and Extract Flat Clustering

Taoran Xue

George Washington University

April 25, 2017

- There are many sources of almost unlimited data:
 - ① Images from the web.
 - ② Speech recorded by a microphone.
 - ③ Records of credit card or other transactions.
- Noise points occasionally draws between two cluster

Minimum spanning tree

For each i , set $r(x_i)$ to the distance from x_i to its k th nearest neighbor.

As r grows from 0 to ∞ :

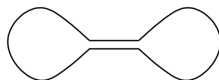
- 1 Construct a graph G_r with nodes $\{x_i : r(x_i) \leq r\}$. Include edge (x_i, x_j) if $\|x_i - x_j\| \leq \alpha r$.
- 2 Let $\mathbb{C}_n(r)$ be the connected components of G_r .

Definition

Set $r_k(x_i)$ to the distance to k th nearest neighbor. For any $r = \max\{r_k(x_i)\}$, connect points x_i and x_j , if $\|x_i - x_j\| \leq \alpha r$.

New distance function

Effect 1: thin bridges



For any set Z , let Z_σ be all points within distance σ of it.

Figure 1: “Thin bridge” effect.

Definition

Set $core_k(x_i)$ to the distance to k th nearest neighbor.

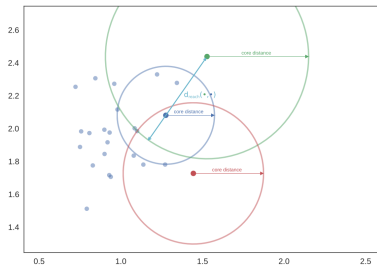
$$d_{\text{mrd}_k}(x_i, x_j) = \max\{core_k(x_i), core_k(x_j), ||x_i - x_j||\}$$

Mutual reachability distance

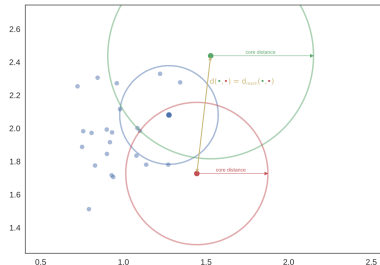
Definition

Set $core_k(x_i)$ to the distance to k th nearest neighbor.

$$d_{\text{mrd}_k}(x_i, x_j) = \max\{core_k(x_i), core_k(x_j), ||x_i - x_j||\}$$



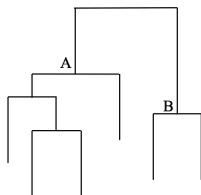
(a)



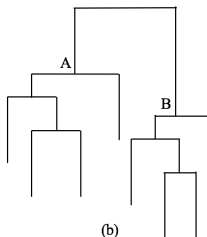
(b)

Condense tree

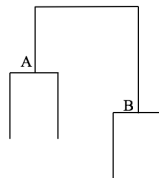
- If left child cluster point number is greater than minimum cluster size, but right side is not, Figure(a), keep the left branch and ignore right cluster;
- If left and right child clusters are both greater than minimum cluster size, Figure(b), we consider that a cluster split and let the split persist the whole tree;
- If left and right child clusters are both fewer than minimum cluster size, Figure(c), we ignore the two cluster.



(a)



(b)



(c)

Condense tree algorithm

```
Input:  $H[m] \leftarrow$  hierarchy tree  
Output:  $T[n] \leftarrow$  condense tree  
 $nodeList \leftarrow \text{BFS}(H);$   
for  $r \leftarrow 0$  to  $m$  do  
    if  $nodeList[r]$  is ignored then  
        Pass;  
    end  
     $left \leftarrow nodeList[r].child1;$   
     $leftCount \leftarrow H[left].childrenSize;$   
     $right \leftarrow nodeList[r].child2;$   
     $rightCount \leftarrow H[right].childrenSize;$   
    if  $leftCount \geq \text{minClusterSize}$  and  $rightCount \geq \text{minClusterSize}$  then  
         $T[p++] \leftarrow (\text{nextLabel}, ++\text{nextLabel}, 1/\text{distance}, \text{leftCount});$   
         $T[p++] \leftarrow (\text{nextLabel}, ++\text{nextLabel}, 1/\text{distance}, \text{rightCount});$   
    end  
    if  $leftCount < \text{minClusterSize}$  then  
        for  $tmpNode$  in  $\text{BFS}(left)$  do  
            if  $tmpNode$  is leaf then  
                 $T[p++] \leftarrow (\text{nextLabel}, tmpNode, 1/\text{distance}, 1);$   
            end  
            Ignore  $tmpNode;$   
        end  
    end  
    if  $rightCount < \text{minClusterSize}$  then  
        for  $tmpNode$  in  $\text{BFS}(right)$  do  
            if  $tmpNode$  is leaf then  
                 $T[p++] \leftarrow (\text{nextLabel}, tmpNode, 1/\text{distance}, 1);$   
            end  
            Ignore  $tmpNode;$   
        end  
    end  
end
```

Definition

Let $\lambda = \frac{1}{d_{\text{mrd}_k}}$. For each cluster we give λ_{birth} and λ_p to be the lambda value when the cluster split off then became it's own cluster, and the lambda value (if any) when the cluster split into smaller clusters respectively.

$$S(C) = \sum_{p \in C} (\lambda_p - \lambda_{\text{birth}})$$

Cluster stability algorithm

Input: $T[n] \leftarrow$ condense tree in reverse topological order contains a tuple of (parent,child, λ ,childrenSize)

Output: $S[n] \leftarrow$ stability of every node in condense tree

for $r \leftarrow 0$ **to** n **do**

 currChild $\leftarrow T[r].\text{child};$

 curr $\lambda \leftarrow T[r].\lambda;$

if currChild = prevChild **then**

 min $\lambda \leftarrow \text{Min}(\text{min}\lambda, \text{curr}\lambda);$

else

 birth $\lambda[\text{currChild}] \leftarrow \text{min}\lambda;$

 prevChild $\leftarrow \text{currChild};$

 min $\lambda \leftarrow \text{curr}\lambda;$

end

end

for $r \leftarrow 0$ **to** n **do**

$S[r] \leftarrow S[r] + T[r].\lambda - \text{birth}\lambda[T[r].\text{parent}] \times T[r].\text{childrenSize};$

end

Definition

Set $SC(C)$ is the sum of the stabilities of the child cluster of cluster C .

$$SC(C) = \sum_{c \in C} S(c)$$

Flat clustering algorithm

Input: $S[n] \leftarrow$ stability condense tree sorted in reverse topological order

Output: $L[n] \leftarrow$ **True** cluster is selected; **False** otherwise

$L \leftarrow \{\mathbf{True}\};$

for $r \leftarrow 0$ **to** n **do**

$\text{childList} \leftarrow \{\text{list of node whose parent is } r\};$

$\text{subtreeStabilities} \leftarrow \sum_{c \in \text{childList}} S[c];$

if $\text{subtreeStabilities} > S[r]$ **then**

$L[r] \leftarrow \mathbf{False};$

$S[r] \leftarrow \text{subtreeStabilities};$

else

for tmpNode **in** $\text{BFS}(r)$ **do**

if $\text{tmpNode} \neq r$ **then**

$L[\text{tmpNode}] \leftarrow \mathbf{False}$

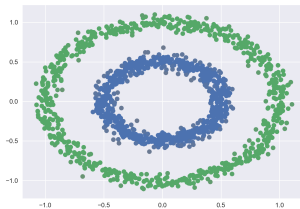
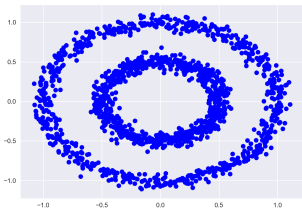
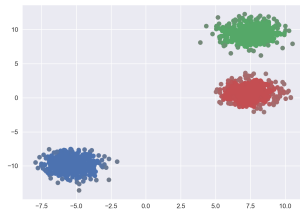
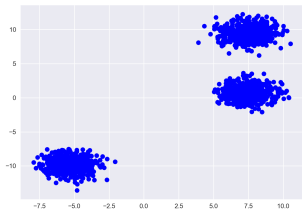
end

end

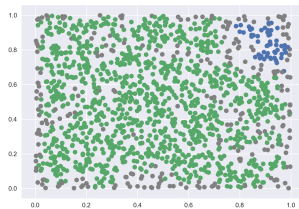
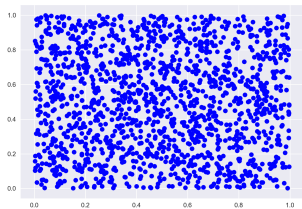
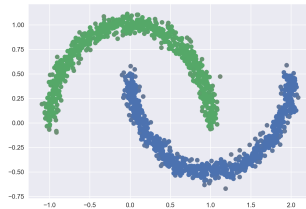
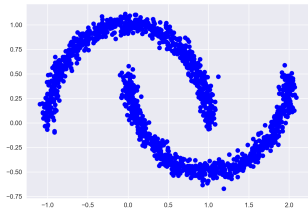
end

end

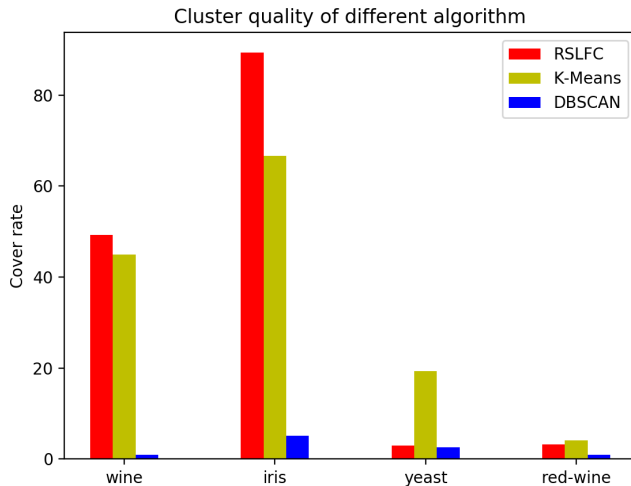
Experiment



Experiment (cont.)



Experiment (cont.)



- Campello, R. J., Moulavi, D., & Sander, J. (2013). Density-based clustering based on hierarchical density estimates. In *Pacific-asia conference on knowledge discovery and data mining* (pp. 160–172).
- Chaudhuri, K., & Dasgupta, S. (2010). Rates of convergence for the cluster tree. In *Advances in neural information processing systems* (pp. 343–351).
- McInnes, L., Healy, J., & Astels, S. (2017, mar). hdbscan: Hierarchical density based clustering. *The Journal of Open Source Software*, 2(11).