

Dive a Bit Deeper into Hierarchical Clustering

Presenter | Sijie Li

2025.08.25



Fast optimal leaf ordering for hierarchical clustering

Ziv Bar-Joseph¹, David K. Gifford^{1,2} and Tommi S. Jaakkola²

¹Laboratory for Computer Science, MIT, 545 Technology Square, Cambridge, MA, 02139, USA and ²Artificial Intelligence Laboratory, MIT, 545 Technology Square, Cambridge, MA, 02139, USA

Received on February 5, 2001; revised and accepted on March 26, 2001

BIOINFORMATICS APPLICATIONS NOTE

Vol. 24 no. 5 2008, pages 719–720
doi:10.1093/bioinformatics/btm563

Gene expression

Defining clusters from a hierarchical cluster tree: the Dynamic Tree Cut package for R

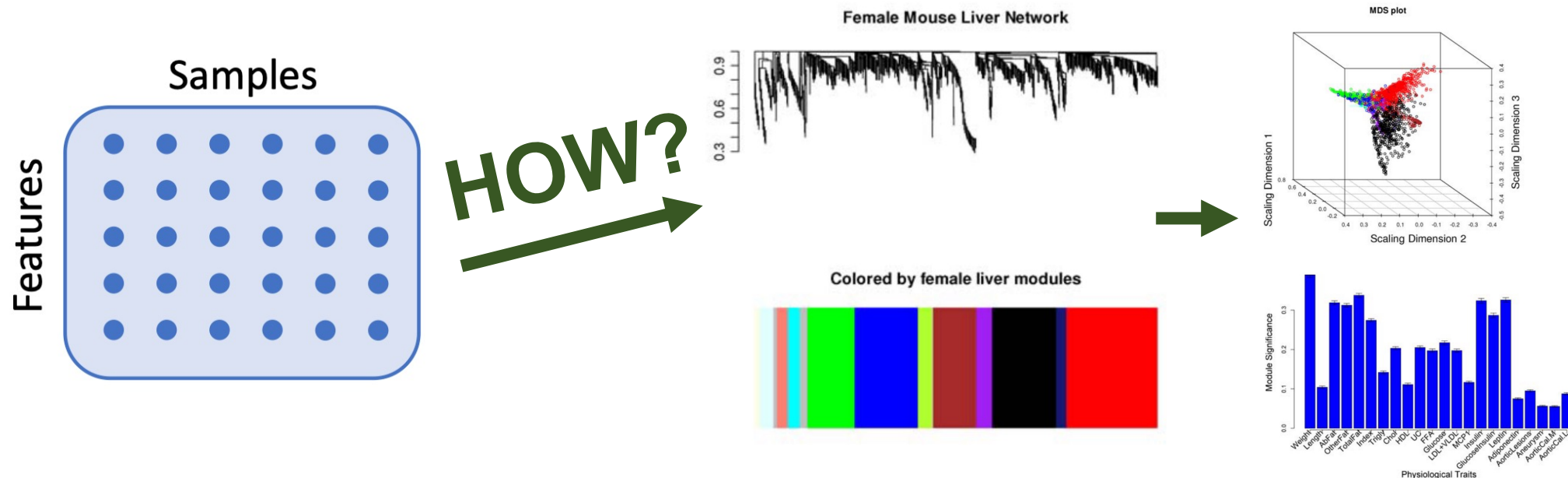
Peter Langfelder^{1,†}, Bin Zhang^{2,†} and Steve Horvath^{1,*}

Bioinformatics

Why should we care?

I want to know deeper out of curiosity.

- What's behind the concise codes?
- Widely used in biology → may be useful to know more.
- Utilize to develop useful bioinformatic tools...



* Also, taste papers form unfamiliar realms

(Hovarth et al., 2006)

Why should we care?

Wisdom from the classical machine learning algorithms may sometimes outweighs the fancy deep-learning.

- Absorb the wisdom from old papers to get new ideas.
- Make the most of them!

nature methods



Brief Communication

<https://doi.org/10.1038/s41592-025-02772-6>

Deep-learning-based gene perturbation effect prediction does not yet outperform simple linear baselines

(Anders et al., 2025)

Distance matrix records the dissimilarity.

```
from scipy.spatial.distance import pdist  
distance = pdist(rna_bulk, "correlation")
```

$$\rho = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E[(x_i - \bar{x})(y_i - \bar{y})]}{E[(x_i - \bar{x})^2]^{1/2} \cdot E[(y_i - \bar{y})^2]^{1/2}}$$

$$\begin{array}{c} \text{sample 1} \quad \text{sample 2} \quad \text{sample 3} \\ \begin{matrix} \text{gene 1} \\ \text{gene 2} \\ \text{gene 3} \end{matrix} \begin{bmatrix} & x_1^T & \\ & x_2^T & \\ & x_3^T & \end{bmatrix} \end{array} \xrightarrow{\text{pdist}} \begin{array}{c} \text{gene 1} \quad \text{gene 2} \quad \text{gene 3} \\ \begin{matrix} \text{gene 1} \\ \text{gene 2} \\ \text{gene 3} \end{matrix} \begin{bmatrix} 1 & 1-\rho_{12} & 1-\rho_{13} \\ & 1 & 1-\rho_{23} \\ & & 1 \end{bmatrix} \end{array}$$

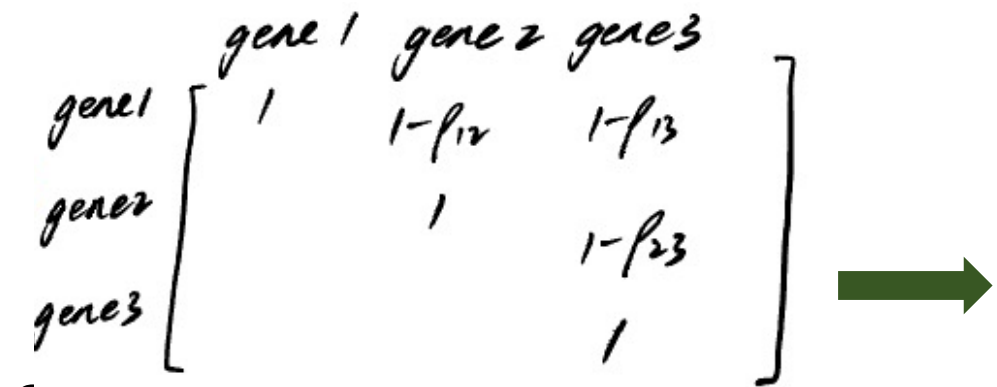
* Supplement

- There are many ways to define “distance”
- Also, other dissimilarity metrics are introduced

e.g. topological overlap dissimilarity measure (dTOM) in WGCNA package. It even considers the scale-free property of the gene network!

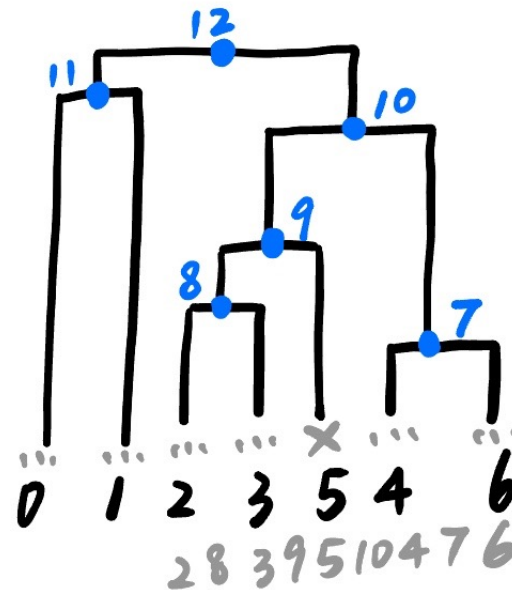
Linkage records the hierarchically merging process.

```
link = linkage(distance, "average")
```



Steps

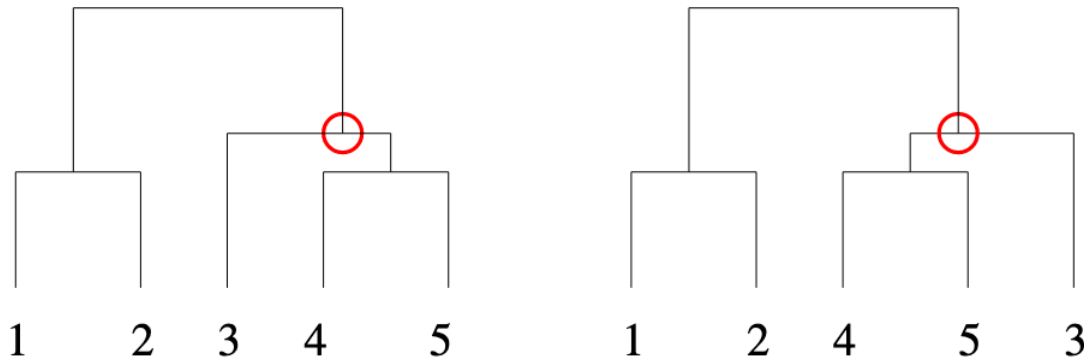
1. Iterate to find minimum
2. merge and replace the distance by the avg
3. [repeat] Iterate to find minimum...



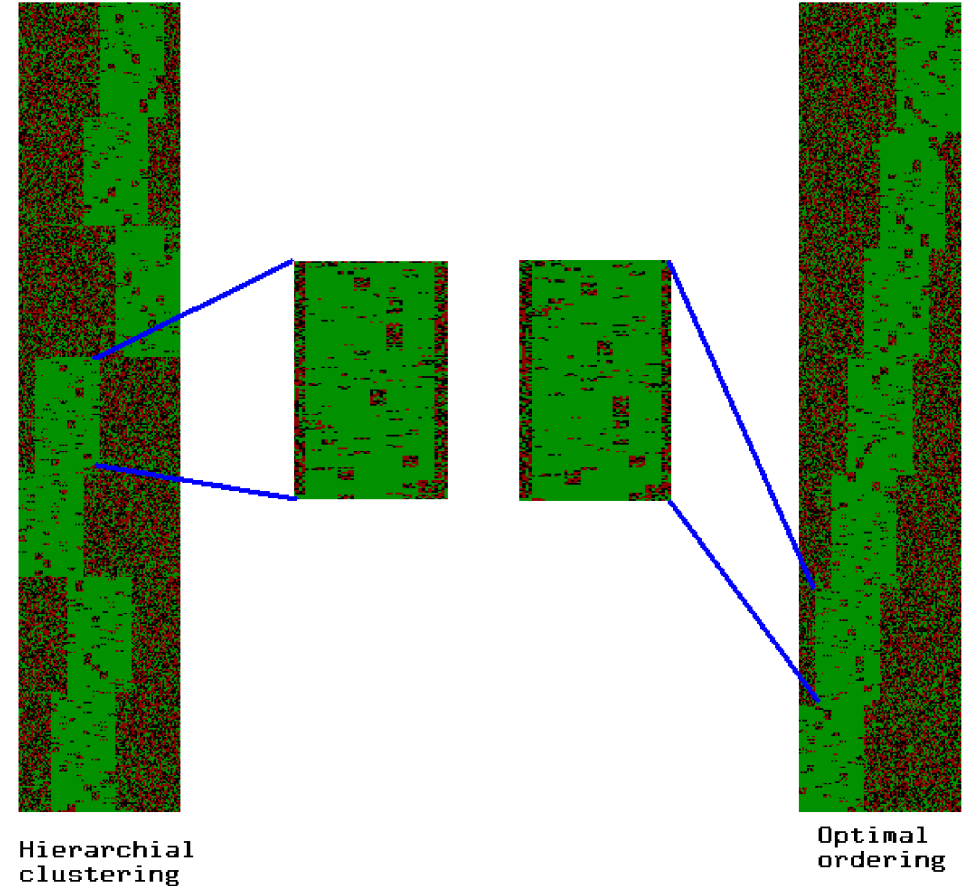
node1	node2	merging height	# object
4	6	0.1	2
2	3	0.14	2
8	5	0.2	3
9	7	0.3	5
...			

link

There are 2^{n-1} ways to arrange a dendrogram.

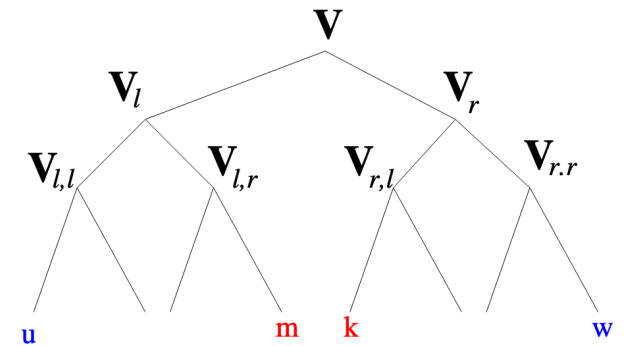


- What is the best ordering?
 - ✓ Max similarity between neighbors
- How to get the best ordering?



Dendrogram revealing biological structures is obtained in a bottom up manner.

```
optOrdering( $v, S$ ) {  
  If  $|v| = 1$  {                                //  $v$  has only one leaf  
     $M(v, u, u) = 0$                             //  $u$  is the only leaf in  $v$   
    return  $M(v, u, u)$   
  Else                                          //  $v$  has more than one leaf  
     $M(v_l, L, R) = \text{optOrdering}(v_l, S)$       //  $v_l$  is the left subtree of  $v$   
     $M(v_r, L, R) = \text{optOrdering}(v_r, S)$       //  $v_r$  is the right subtree of  $v$   
    For all leaves  $u \in v_l$  {  
      For all leaves  $w \in v_r$  {  
         $M(v, u, w) = \max_{m \in v_{l,r}, k \in v_{r,l}} M(v_l, u, m) + M(v_r, w, k) + S(m, k)$   
         $M(v, w, u) = M(v, u, w)$   
      }  
    }  
    return  $M(v, L, R)$                         //  $L, R$  stands for all possible pairs of leaves from  $v_l$  and  $v_r$   
}
```



V: subtree/internal node

(Bar-Joseph et al., 2001)

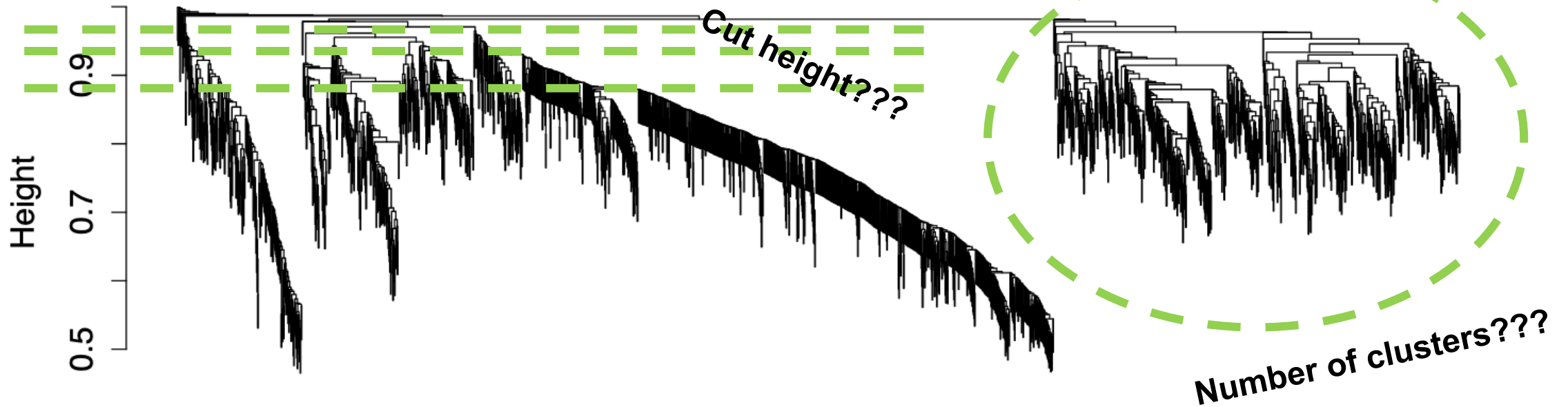
Stopping earlier can shorten the time to find the optimal ordering.

```
curMax =  $-\infty$ 
For all leaves m according to the order of  $M(v_l, u, R)$  {
  if  $M(v_l, u, m) + M(v_r, w, k_0) + C(v_l, v_r) \leq \textit{curMax}$       //  $k_0$  is first in the order of  $M(v_r, w, L)$ 
     $M(v, u, w) = \textit{curMax}$ , terminate the search
  For all leaves k according to the order of  $M(v_r, w, L)$ 
    if  $M(v_l, u, m) + M(v_r, w, k) + C(v_l, v_r) \leq \textit{curMax}$ 
      break // terminate the inner loop
    if  $\textit{curMax} < M(v_l, u, m) + M(v_r, w, k) + S(m, k)$ 
       $\textit{curMax} = M(v_l, u, m) + M(v_r, w, k) + S(m, k)$ 
  } // end of the inner loop
} // end of the outer loop
 $M(v, u, w) = \textit{curMax}$ 
```

(Bar-Joseph et al., 2001)

To get clusters!

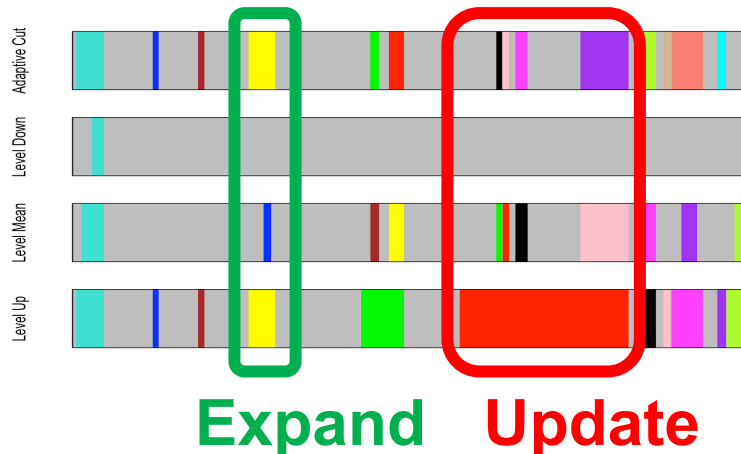
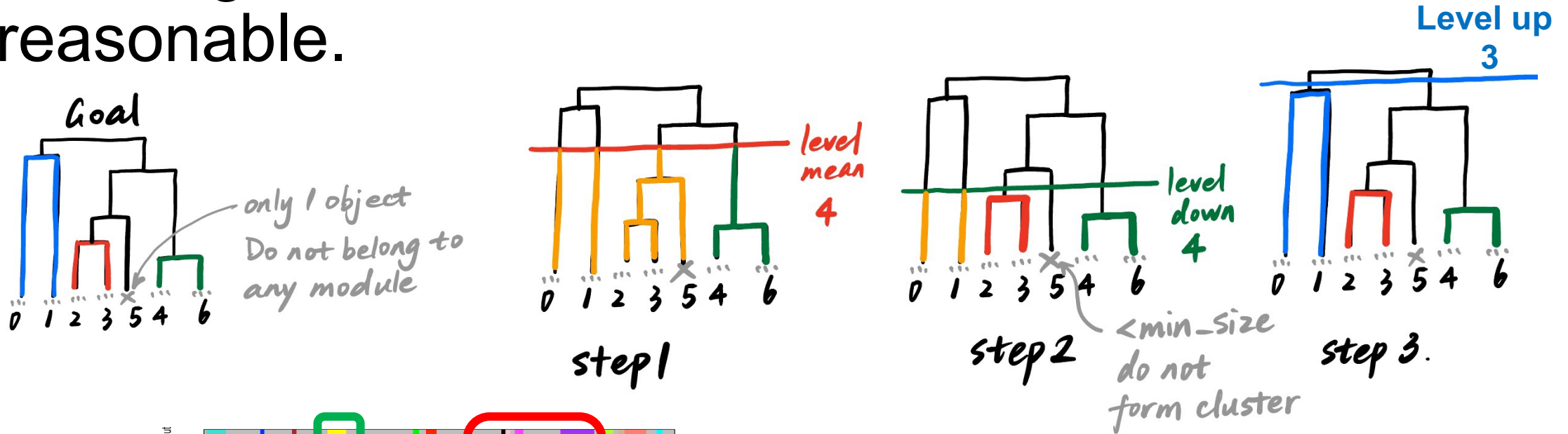
Specifying cut height or cluster number manually may miss biological insights.



- Manually specifying parameters places the burden of choosing good settings on the user! E.g. adjusting the resolution when doing single cell clustering 😞
- Makes the whole process hard to automate.

(Horvath et al., 2008)₈

Cut height can be softened to make clusters more reasonable.

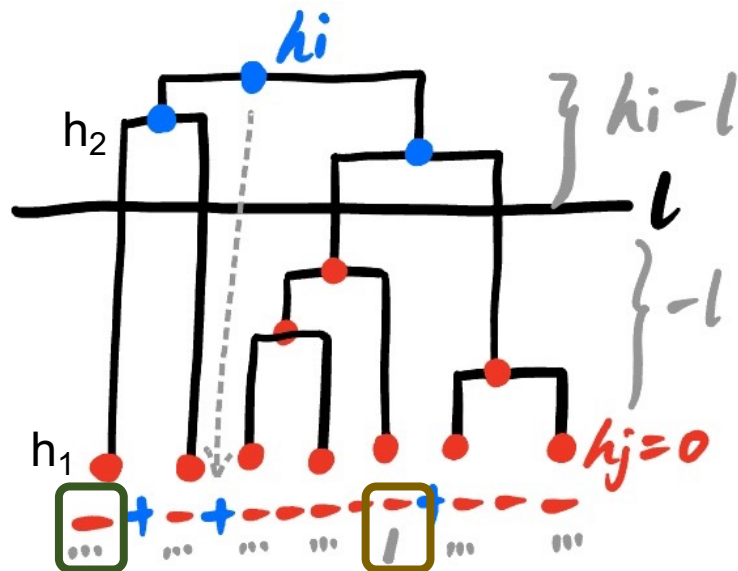


$$l_m = \frac{1}{n} \sum_{i=1}^n h_i,$$

$$l_u = \frac{1}{2} (l_m + \max\{h_1, h_2, \dots, h_n\}),$$

$$l_d = \frac{1}{2} (l_m + \min\{h_1, h_2, \dots, h_n\}).$$

```
def TreeCut(H,l,min_size):
```



- Nodes here including merging nodes
- # Nodes in the cluster = # **continuous** minus-
- If # > min size → form a cluster

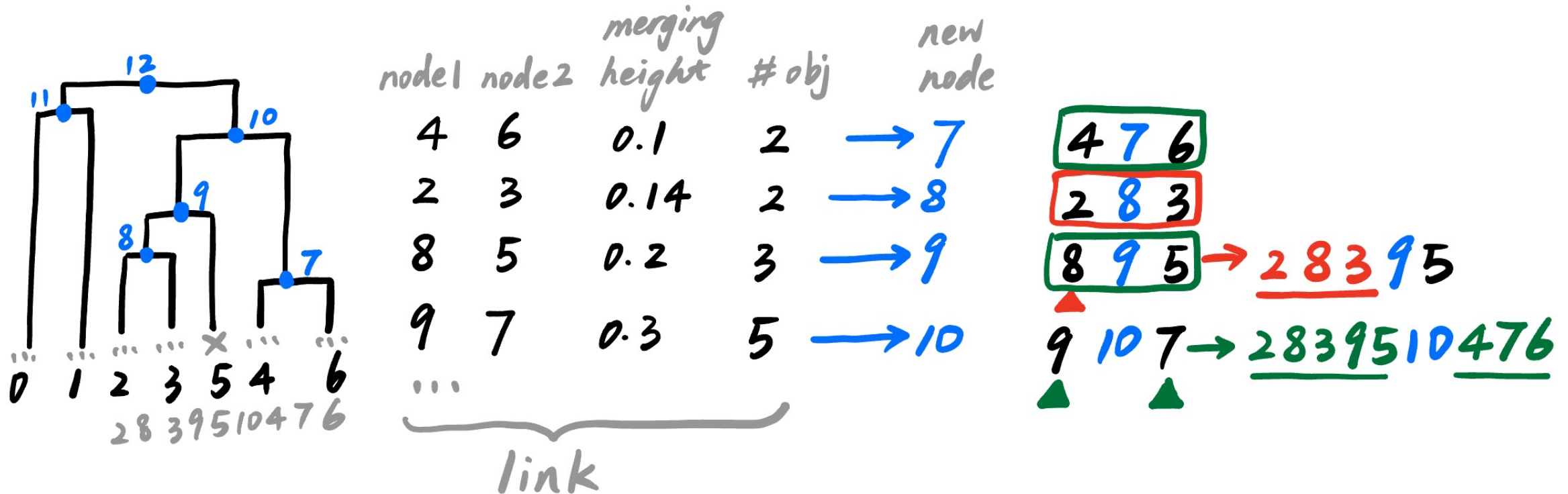
- `hclust()` and `linkage()` only do return the merging process
- Order is the **same as dendrogram!**

Arduous!
Confusions in the paper!

10

Reproduce DynamicTreeCut

def node_order(link):

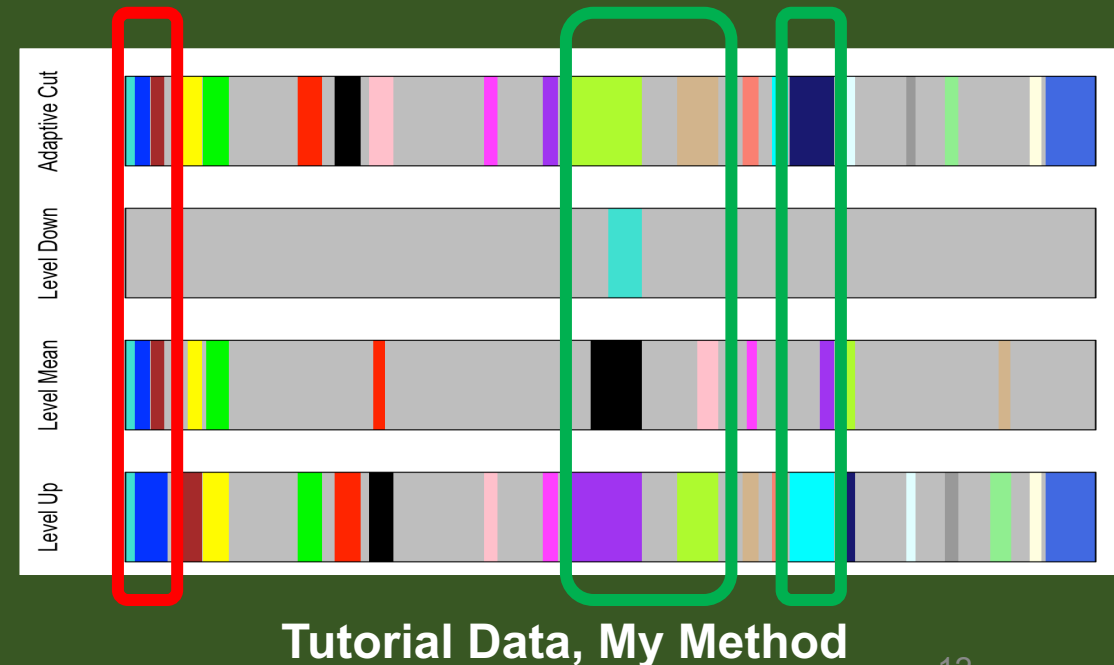
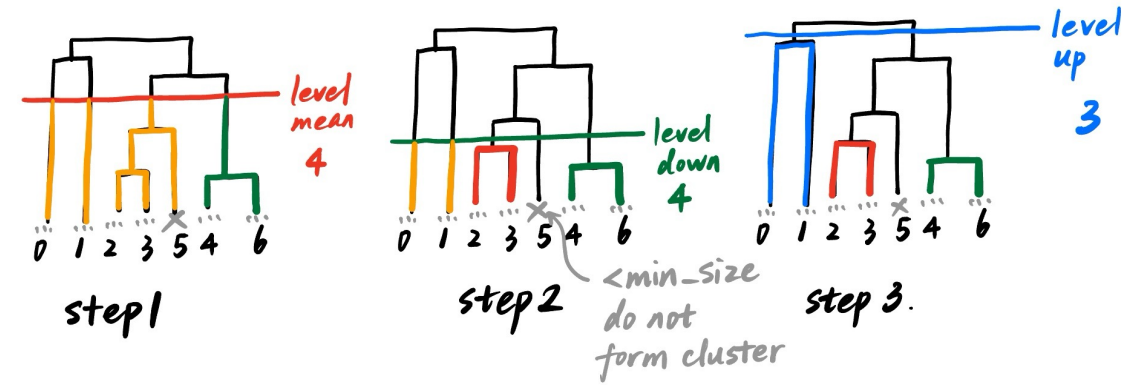
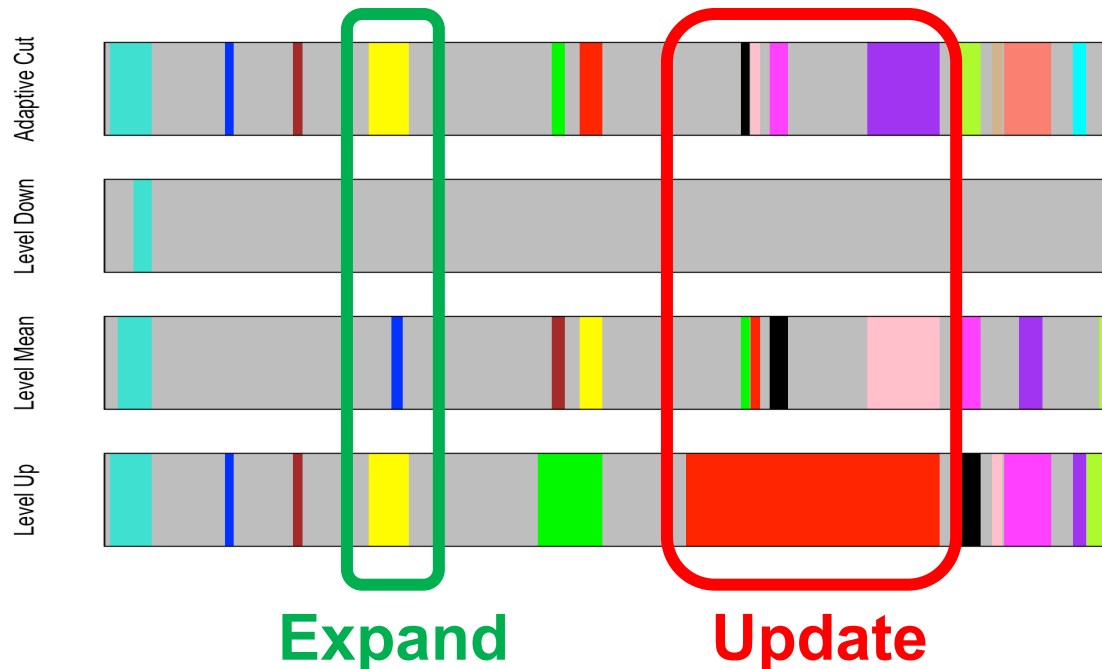


Then **map the height** to the node with link get height_order

Reproduce DynamicTreeCut

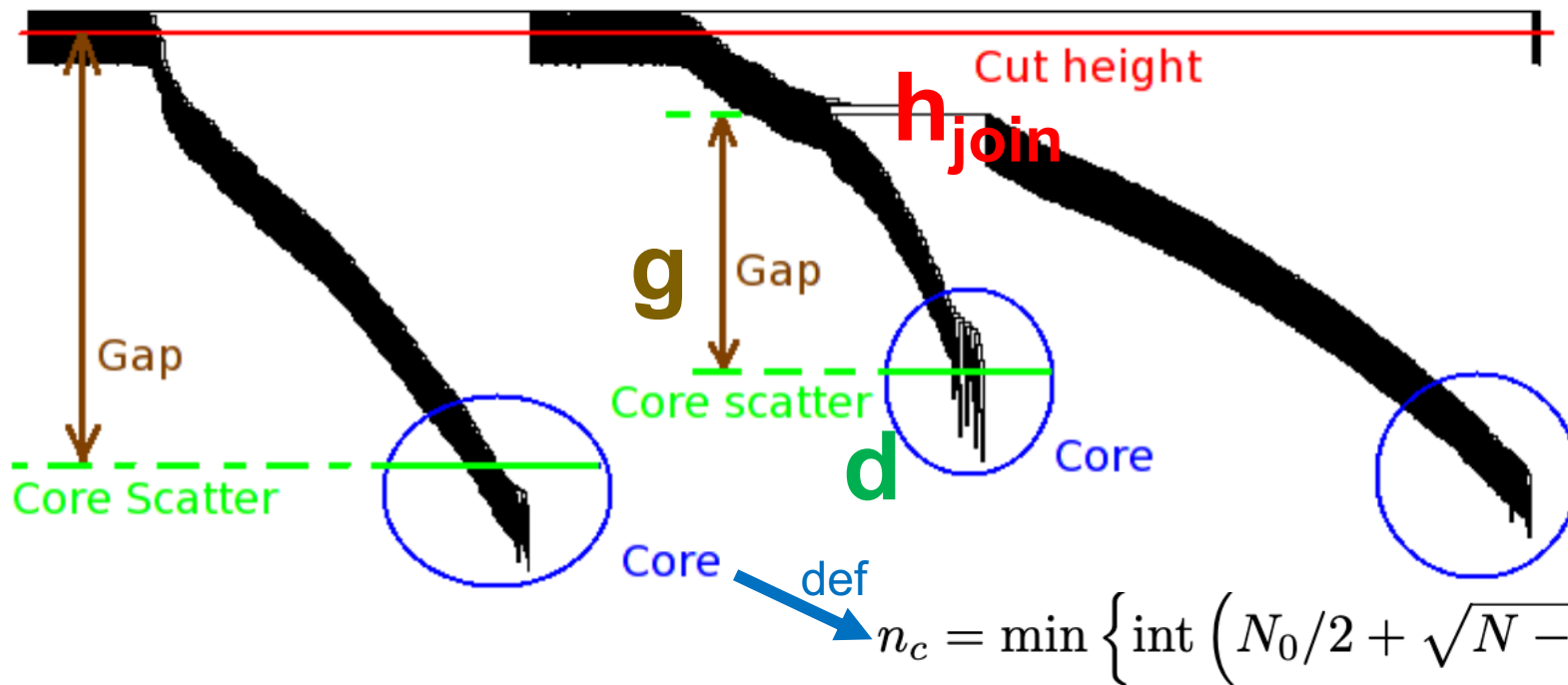
Def AdaptiveTreeCut(H):

- If ≥ 2 sub-clusters in the lower cut level result, **update**
- Or else, **expand** the cluster range



Clusters are defined based on the structure of the dendrogram.

- We can detect clusters easily with our bare eye! Shape matters!



Criteria to be a cluster

- ✓ # leaves $\geq N_0$
- ✓ $h_{\text{join}} \leq h_{\text{max}}$
- ✓ $g \geq g_{\text{min}}$
- ✓ $d \geq d_{\text{max}}$

Our eyes judge like this!

(Horvath et al., 2008)

Flat clustering methods can be included for optimization.

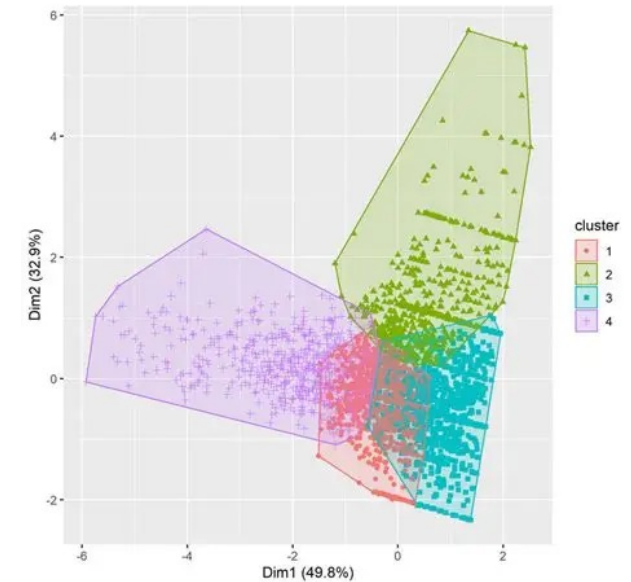
Apply to assign

- unlabeled leaves
- unlabeled tiny cluster ($\# \text{ leaves} < N_0$)

PAM (partition around medoid)-like method

- Method 1: assign according to avg(leaves-leaves) distance
- Method 2: assign according to leaves-medoid/medoid-medoid distance

PAM

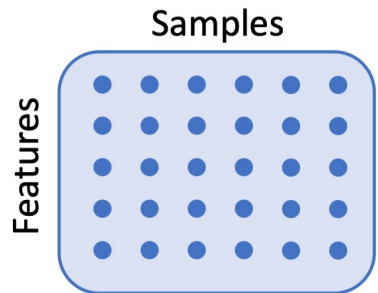


* **Flat clustering methods** always have resolution problem, but **hierarchical clustering** do not have.
Eg. Specify radius for DBSCAN; select number of clusters for k-means

In summary

It's so satisfying to finally figure all of this out. Hope you feel the same! 😊

Now, what's happen behind these codes and the principles to get the results are clear!



```
from scipy.cluster import hierarchy
from scipy.spatial.distance import pdist
distance = pdist(rna_bulk, "correlation")
link = hierarchy .linkage(distance, "average")
dn = hierarchy .dendrogram(link)
```

```
from dynamicTreeCut import cutreeHybrid
clusters = cutreeHybrid(link, distance)
```

