

Python package for proteomics data clustering

☐ Introduction

☐ Preprocessing

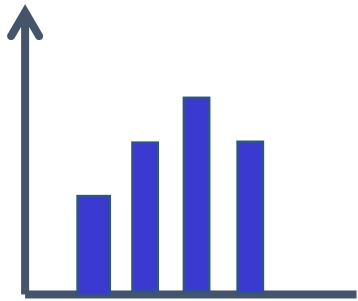
☐ Network Enhancement

☐ Clustering

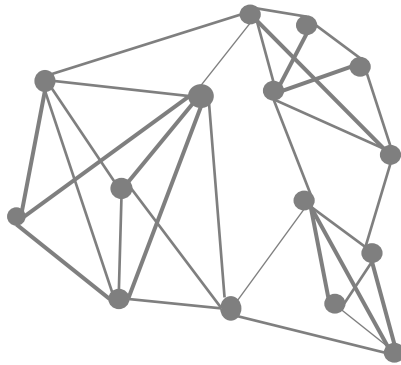
☐ Usage

☐ Reference

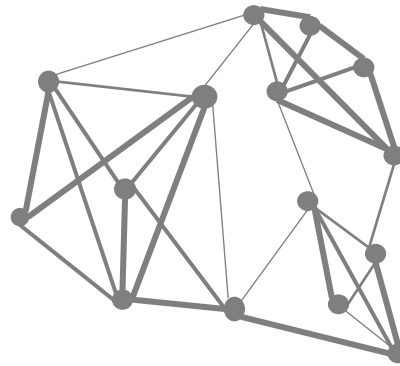
Overview



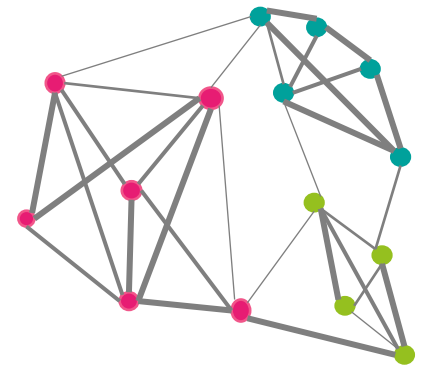
TMT Data



Graph



Denoised
Graph



Modules

Annotation of a graph

- Graph G is consisted of 3 components :
Nodes V , Edges E , and Weights W

$$G = (V, E, W)$$

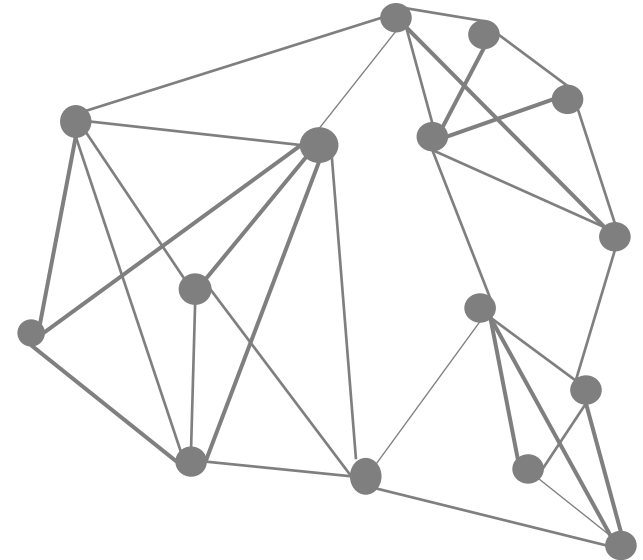
- Define the number of nodes as n

$$|V| = n$$

- Edges are defined as pairs of nodes

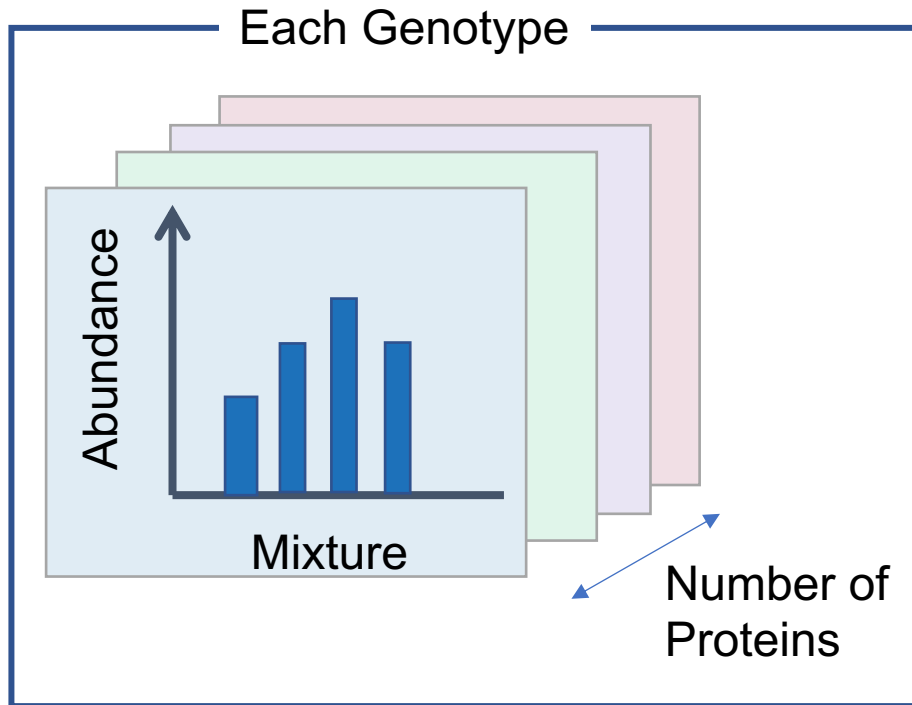
$$e_{i,j} = (v_i, v_j) \quad e_{i,j} \in E, v_i, v_j \in V$$

- If each edges have weight $w_{i,j}$ as a scalar value, weights are defined as a nxn matrix $W \in \mathbb{R}^{n \times n}$



Preprocessing

Data structure of TMT



Annotation

- Proteins
 $p \in \mathbf{P}, |\mathbf{P}| = n$
- Mixture/ Fraction
 $m \in \mathbf{M}$
- Genotype
 $g \in \mathbf{G}$
- Abundance vector of
 i-th protein in j-th genotype
 $\mathbf{a}_{p_i, g_j} = (a_{p_i, g_j, m} \mid m \in \mathbf{M})$

Correlation matrix

Think about single genotype g .

The correlation coefficient between i-th and j-th protein can be defined as :

$$\text{corr}(\mathbf{a}_{p_i}, \mathbf{a}_{p_j}) = \frac{\sum_{m \in M} (a_{m,p_i} - \widehat{\mathbf{a}}_{p_i})(a_{m,p_j} - \widehat{\mathbf{a}}_{p_j})}{\sqrt{\sum_{m \in M} (a_{m,p_i} - \widehat{\mathbf{a}}_{p_i})^2 \sum_{m \in M} (a_{m,p_j} - \widehat{\mathbf{a}}_{p_j})^2}}$$

Correlation matrix $W \in \mathbb{R}^{n \times n}$ is defined as:

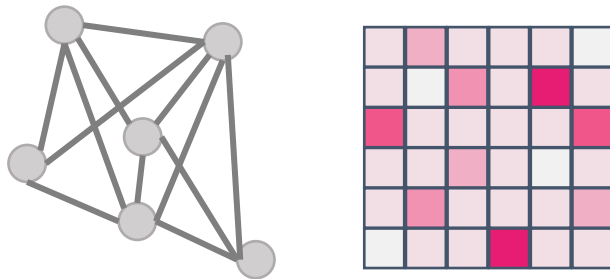
$$W_{i,j} = \text{corr}(\mathbf{a}_{p_i}, \mathbf{a}_{p_j})$$

By regarding correlation coefficients as weights and proteins as nodes, the sets (P, E, W) is regarded as a graph.

Network Enhancement

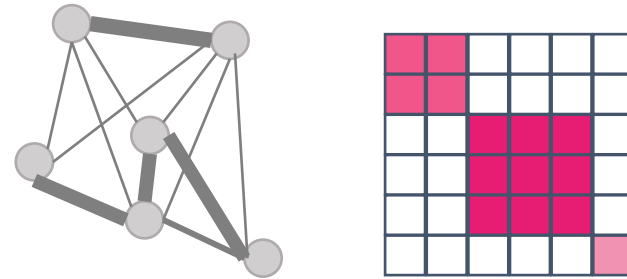
Goal : Getting new wights matrix W^*

$$G = (P, E, W)$$



Initial graph

$$G^* = (P, E, W^*)$$



Weighted graph

- Weights should be positive : $W_{i,j} \geq 0$
- "Heavy" edge should be heavier, "right" edge should be lighter
Increasing the eigengap of W while preserving the eigenvector

3 steps

1. Define probability transition matrix P as:

$$P_{i,j} = \frac{W_{i,j}}{\sum_{k \in \mathcal{N}_i} W_{i,k}} * I\{j \in \mathcal{N}_i\}$$

\mathcal{N}_i : k-th neighborhoods of i-th vertex

$$I\{j \in \mathcal{N}_i\} = \begin{cases} 1 & \text{if } j \in \mathcal{N}_i \\ 0 & \text{else} \end{cases}$$

2. Yield doubly stochastic matrix \mathcal{T} as:

$$\mathcal{T}_{i,j} \leftarrow \sum_{k=1}^n \frac{P_{i,k} P_{j,k}}{\sum_{v=1}^n P_{v,k}}$$

3. Update weight

$$W_{t+1} = \alpha \mathcal{T} \times W_t \times \mathcal{T} + (1 - \alpha) \mathcal{T}$$

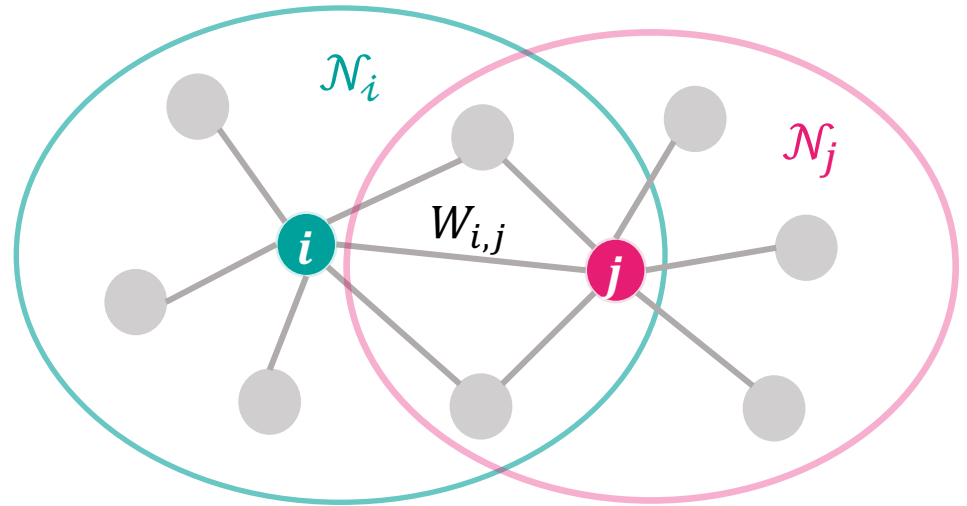
α : Regularization parameter
 t : Iteration number

Probability transition matrix

$$P_{i,j} = \frac{W_{i,j}}{\sum_{k \in \mathcal{N}_i} W_{i,k}} * I\{j \in \mathcal{N}_i\}$$

j is in $\mathcal{N}_i : P_{i,j} > 0$

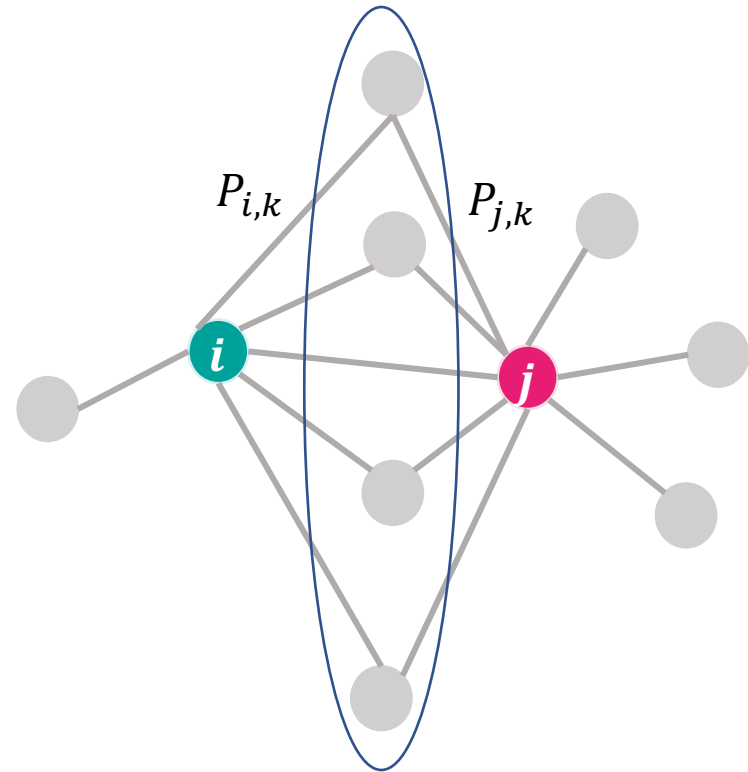
i isn't in $\mathcal{N}_j : P_{j,i} = 0$



- P isn't necessarily symmetric
- P is sparse

Doubly stochastic matrix

$$\begin{aligned} \mathcal{T}_{i,j} &= \sum_{k=1}^n \frac{P_{i,k} P_{j,k}}{\sum_{v=1}^n P_{v,k}} \\ &= \frac{\sum_{k=1}^n P_{i,k} P_{j,k}}{\sum_{k=1}^n \sum_{v=1}^n P_{v,k}} \end{aligned}$$



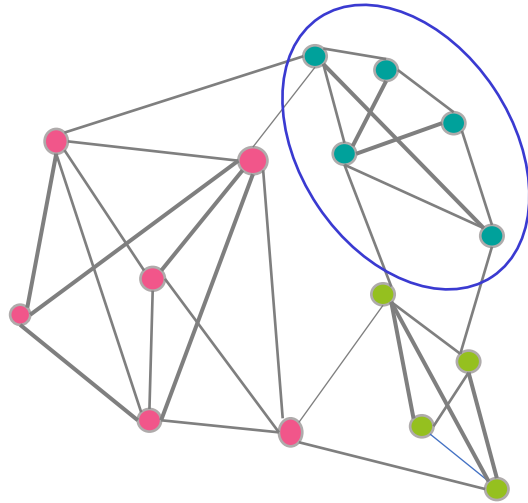
- T is symmetric
- Summation of all rows and columns of T is 1

Update of weights

$$\begin{aligned}
 W_t &= \alpha \mathcal{T} \times W_{t-1} \times \mathcal{T} + (1 - \alpha) \mathcal{T} \\
 &= \alpha \mathcal{T} \times (\alpha \mathcal{T} \times W_{t-2} \times \mathcal{T} + (1 - \alpha) \mathcal{T}) \times \mathcal{T} + (1 - \alpha) \mathcal{T} \\
 &= \alpha^2 \mathcal{T}^2 W_{t-2} \mathcal{T}^2 + \alpha \mathcal{T} (1 - \alpha) \mathcal{T} \times \mathcal{T} + (1 - \alpha) \mathcal{T} \\
 &= \alpha^2 \mathcal{T}^2 W_{t-2} \mathcal{T}^2 + (1 - \alpha) \mathcal{T} (\alpha \mathcal{T}^2 + 1) \\
 &\dots \\
 &= \alpha^t \mathcal{T}^t W_0 \mathcal{T}^t + (1 - \alpha) \mathcal{T} \sum_{k=0}^{t-1} (\alpha \mathcal{T}^2)^k
 \end{aligned}$$

$$W_{t \rightarrow \infty} = (1 - \alpha) \mathcal{T} (\mathcal{I} - \alpha \mathcal{T}^2)^{-1}$$

Overview and basic annotation



Community $C_i \subseteq V$

s.t.

$V = \cup_i C_i, C_i \cap C_j = \phi$ when $i \neq j$

Partition $P = \{C_1, C_2 \dots C_r\}$

Graph $G = (V, E, W)$

Node movement $P(v \rightarrow C)$

Goal : To obtain good partition \rightarrow Optimization of quality function $H(G, P)$

$$H(G, P) = \sum_{C \in P} [E(C, C) - \gamma \binom{|C|}{2}]$$

$$H(G, P) = \sum_{C \in P} [E(C, C) - \frac{\gamma}{2m} \binom{|C|}{2}]$$

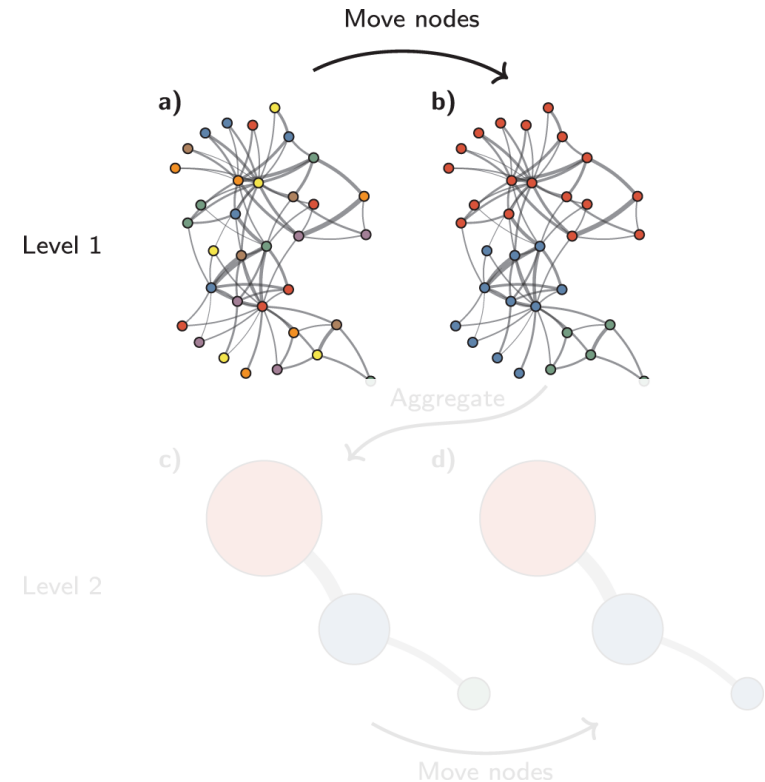
Louvain algorithm

```

function Louvain(Graph G, Partition P)
  do
    P  $\leftarrow$  MoveNodes(G, P)
    done  $\leftarrow$  |P| = |V (G)|
    if not done then
      G  $\leftarrow$  AggregateGraph(G, P)
      P  $\leftarrow$  SingletonPartition(G)
    end if
  while not done
  return flat* (P)
end function
  
```

```

function MoveNodes(Graph G, Partition P)
  do
     $H_{old} = H(P)$ 
    for  $v \in V (G)$  do
       $C' \leftarrow \operatorname{argmax}_{C \in P \cup \emptyset} \Delta H_P(v \rightarrow C)$  .
      if  $\Delta H_P(v \rightarrow C') > 0$  then
         $v \rightarrow C'$ 
      end if
    end for
  while  $H(P) > H_{old}$ 
  return P
end function
  
```



Assign nodes to communities which maximize quality function

Louvain algorithm

```

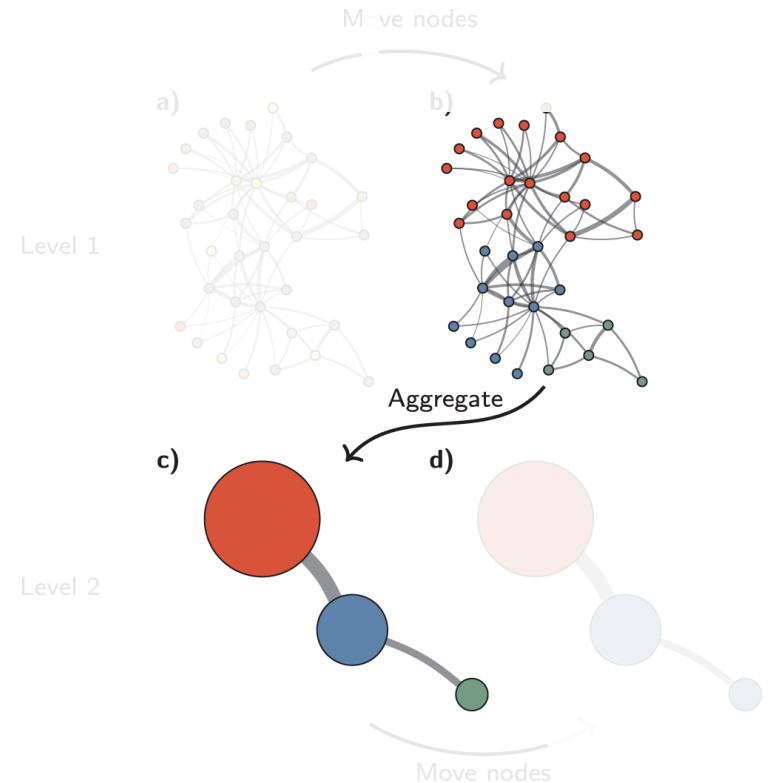
function Louvain(Graph G, Partition P)
  do
     $P \leftarrow \text{MoveNodes}(G, P)$ 
     $\text{done} \leftarrow |P| = |V(G)|$ 
    if not done then
       $G \leftarrow \text{AggregateGraph}(G, P)$ 
       $P \leftarrow \text{SingletonPartition}(G)$ 
    end if
  while not done
  return flat* (P)
end function
  
```

```

function AggregateGraph(Graph G, Partition P)
   $V \leftarrow P$ 
   $E \leftarrow \{(C, D) \mid (u, v) \in E(G), u \in C \in P, v \in D \in P\}$ 
  return Graph(V, E)
end function
  
```

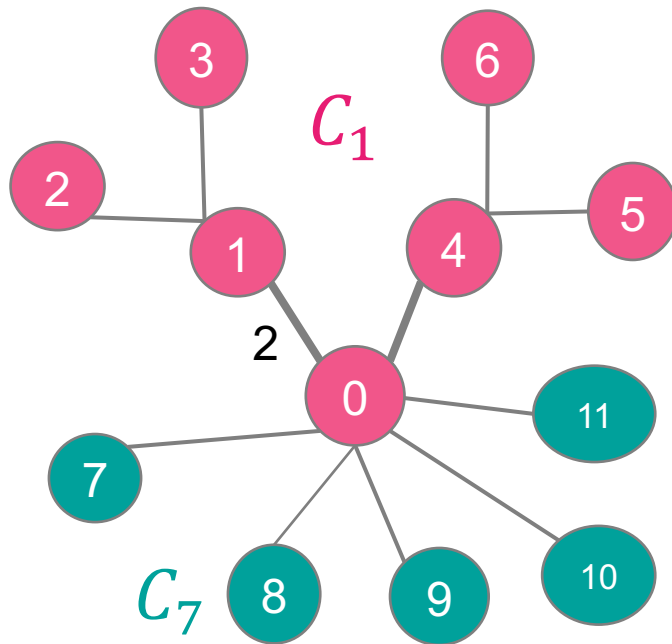
```

function SingletonPartition(Graph G)
  return  $\{\{v\} \mid v \in V(G)\}$ 
end function
  
```

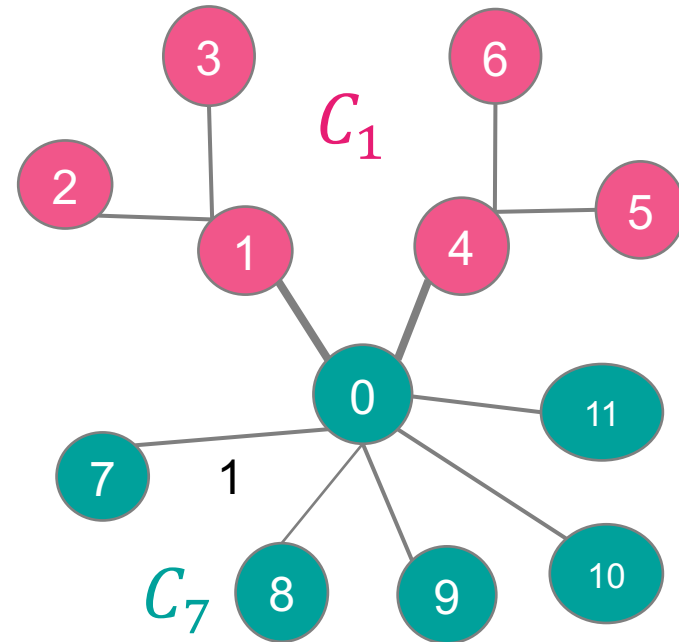


◀ After all nodes are assigned to communities, each communities are assumed as a node

Disconnection problem of Louvain algorithm



$$\Delta H(0 \rightarrow C_1) = 2 \times 2 - 6\gamma$$



$$H(0 \rightarrow C_7) = 5 - 5\gamma$$

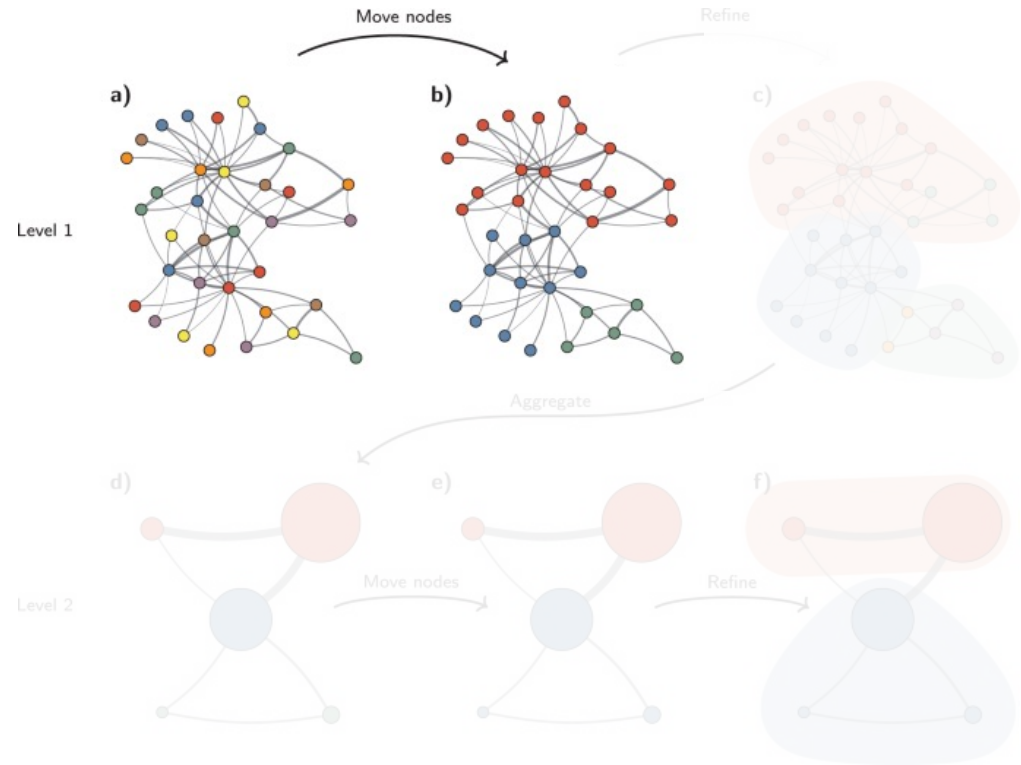
Leiden algorithm

```

function Leiden(Graph G, Partition P)
  do
     $P \leftarrow \text{MoveNodesFast}(G, P)$ 
     $\text{done} \leftarrow |P| = |V(G)|$ 
    if not done then
       $P_{\text{refined}} \leftarrow \text{RefinePartition}(G, P)$ 
       $G \leftarrow \text{AggregateGraph}(G, P_{\text{refined}})$ 
       $P \leftarrow \{ \{v \mid v \subseteq C, v \in V(G)\} \mid C \in P \}$ 
    end if
  while not done
  return flat* (P)
end function
  
```

```

function MoveNodesFast(Graph G, Partition P)
   $Q \leftarrow \text{Queue}(V(G))$ 
  do
     $v \leftarrow Q.\text{remove}()$ 
     $C' \leftarrow \text{argmax}_{C \in P \cup \emptyset} \Delta H_P(v \rightarrow C)$ 
    if  $\Delta H_P(v \rightarrow C') > 0$  then
       $v \rightarrow C'$ 
       $N \leftarrow \{u \mid (u, v) \in E(G), u \notin C'\}$ 
       $Q.\text{add}(N - Q)$ 
    end if
  while  $Q \neq \emptyset$ 
  return P
end function
  
```



◀ Make sure to visit all nodes by using queue

Leiden algorithm

```

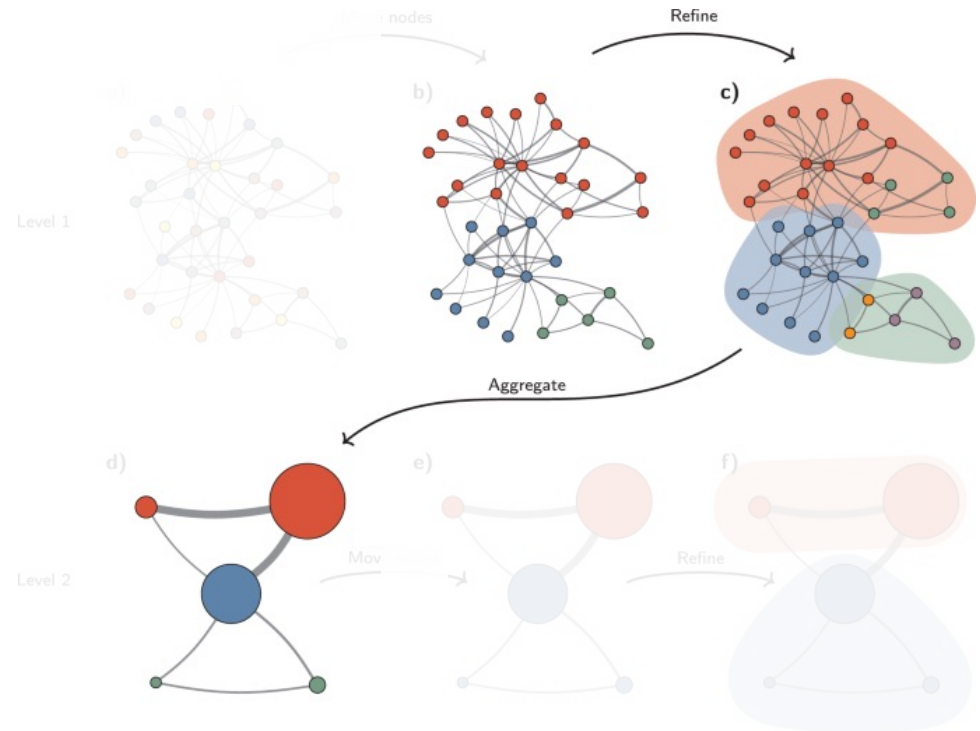
function Leiden(Graph G, Partition P)
  do
     $P \leftarrow \text{MoveNodesFast}(G, P)$ 
     $\text{done} \leftarrow |P| = |V(G)|$ 
    if not done then
       $P_{\text{refined}} \leftarrow \text{RefinePartition}(G, P)$ 
       $G \leftarrow \text{AggregateGraph}(G, P_{\text{refined}})$ 
       $P \leftarrow \{\{v \mid v \subseteq C, v \in V(G)\} \mid C \in P\}$ 
    end if
  while not done
  return flat*(P)
end function
  
```

```

function RefinePartition(Graph G, Partition P)
   $P_{\text{refined}} \leftarrow \text{SingletonPartition}(G)$ 
  for  $C \in P$  do
     $P_{\text{refined}} \leftarrow \text{MergeNodesSubset}(G, P_{\text{refined}}, C)$ 
  end for
  return  $P_{\text{refined}}$ 
end function
  
```

```

function AggregateGraph(Graph G, Partition P)
   $V \leftarrow P$ 
   $E \leftarrow \{(C, D) \mid (u, v) \in E(G), u \in C, v \in D \in P\}$ 
  return Graph(V, E)
end function
  
```



◀ Divide partitions into subsets (Refinement)

◀ Aggregation is based on P_{refined} , but these subsets will be assigned to the same community

Code availability

<https://github.com/lo-Saito/ProteinClustering> MIT License

Requirement

Python \geq 3.7

TODO

Interface for .rda

References

- [1] [Genetic Disruption of WASHC4 Drives Endo-lysosomal Dysfunction and Cognitive-Movement Impairments in Mice and Humans.](#) Courtland J.L., Bradshaw T.W.A., Waitt G., Soderblom E., Ho T., Rajab A., Vancini R., Kim I.H., Soderling S.H. (2021). eLife; 10:e61590 doi: 10.7554/eLife.61590

- [2] [From Louvain to Leiden: guaranteeing well-connected communities.](#) Traag, V.A., Waltman. L., Van Eck, N.-J. (2018). Scientific reports, 9(1), 5233. 10.1038/s41598-019-41695-z

- [3] [Network Enhancement as a general method to denoise weighted biological networks.](#) Wang B., Pourshafeie A., Zitnik M., Zhu J., Bustamante C.D., Batzoglou S., Leskovec J. (2018). Nature Communications, 9, 3108. 10.1038/s41467-018-05469-x

Appendix

Refine partition

```

function MergeNodesSubset(Graph G, Partition P, Subset S)
  R = {v | v ∈ S, E(v, S - v) ≥ γ||v|| · (||S|| - ||v||)}
  for v ∈ R do
    if v in singleton community then
      T ← {C | C ∈ P, C ⊆ S, E(C, S - C) ≥ γ||C|| · (||S|| - ||C||)}

      
$$\Pr(C' = C) \sim \begin{cases} \exp\left(\frac{1}{\theta} \Delta H_P(v \rightarrow C)\right) & \text{if } \Delta H_P(v \rightarrow C) \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad \text{for } C \in T$$


      v → C'
    end if
  end for
  return P
end function

```