

Clusters and Communities

Internet Analytics (COM-308)

Prof. Matthias Grossglauser
School of Computer and Communication
Sciences



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Overview

- Clustering:
 - Given: set of points with a distance metric
 - Find sets of points that are close to each other, but far from other points
- Community detection:
 - Given: network
 - Find sets of nodes that are highly interconnected, but poorly connected to other nodes
- Many important applications:
 - Data analysis
 - Problem decomposition
 - Resource allocation
 - ...

Clustering: goal and definition

- Find organization in data:
 - Image compression: each pixel has a color → find small set of colors so that each pixel is close to one color
 - Mobility: point = GPS trace with noise → find representative set of trajectories
- Definition:
 - Given: set of points (vectors) with a distance function (metric space)
 - Find: partition (hard or soft) of points into clusters; plus potentially more information (characterization of clusters)

K-means clustering algorithm

- Input:
 - N data points x_1, \dots, x_N
 - K : number of clusters
- Output:
 - K cluster centers μ_k
 - r_{nk} : point-cluster assignment indicator
 - $r_{nk} = 1$ means point x_n is in cluster k
- Cost function:
 - $J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$
- Optimal K-means: NP-hard
- Solution: iterative heuristic to approximate solution

K-means: iterative approximation

- Initialize μ
- Until convergence:
 - Minimize J w.r.t. $\{r_{nk}\}$:
 - $r_{nk} = 1$ only for
 $k = \operatorname{argmin} \|x_n - \mu_k\|$
 - Minimize J w.r.t. μ :
 - Set derivative of J w.r.t. μ_k to zero
 - $2 \sum_n r_{nk} (x_n - \mu_k) = 0$
 - Solve for μ : $\mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$

E-step:
Attribute each x_n
to closest center

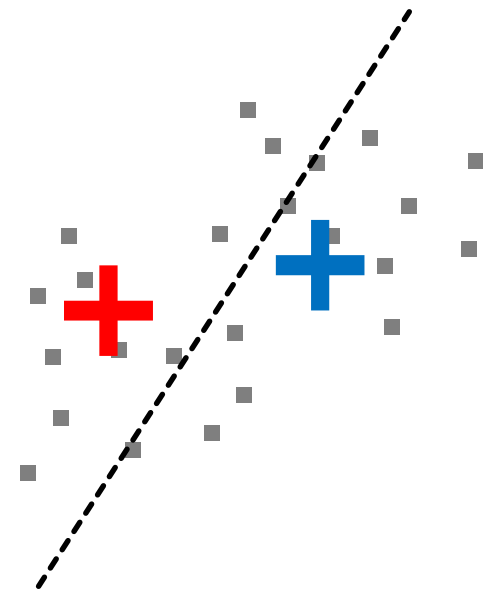
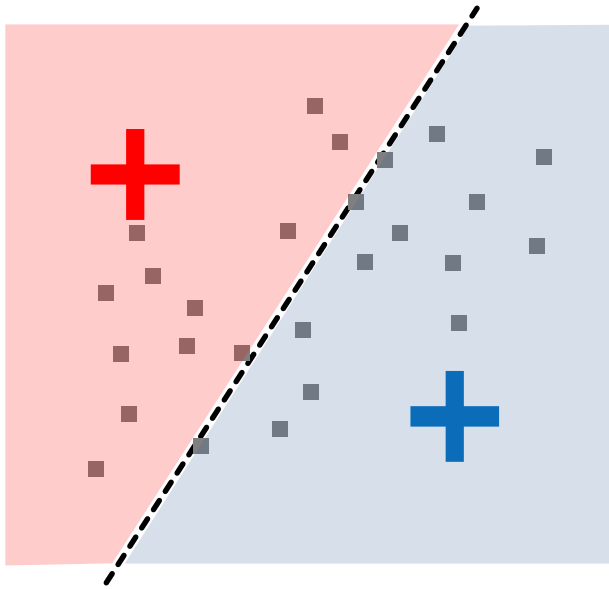
M-step:
New cluster center μ_k
= center of mass of
points of cluster k

K-means: example

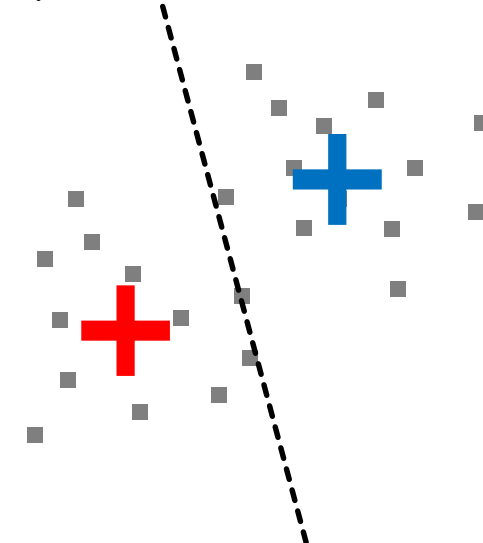
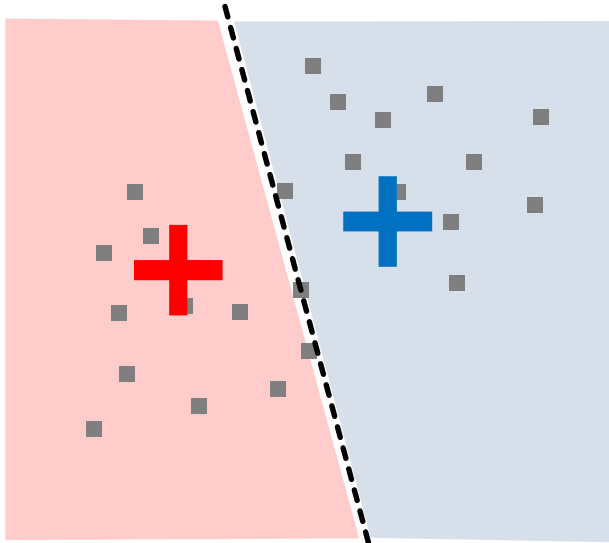
E-step

M-step

Iter 1



Iter 2



From K-means to Mixtures of Gaussians

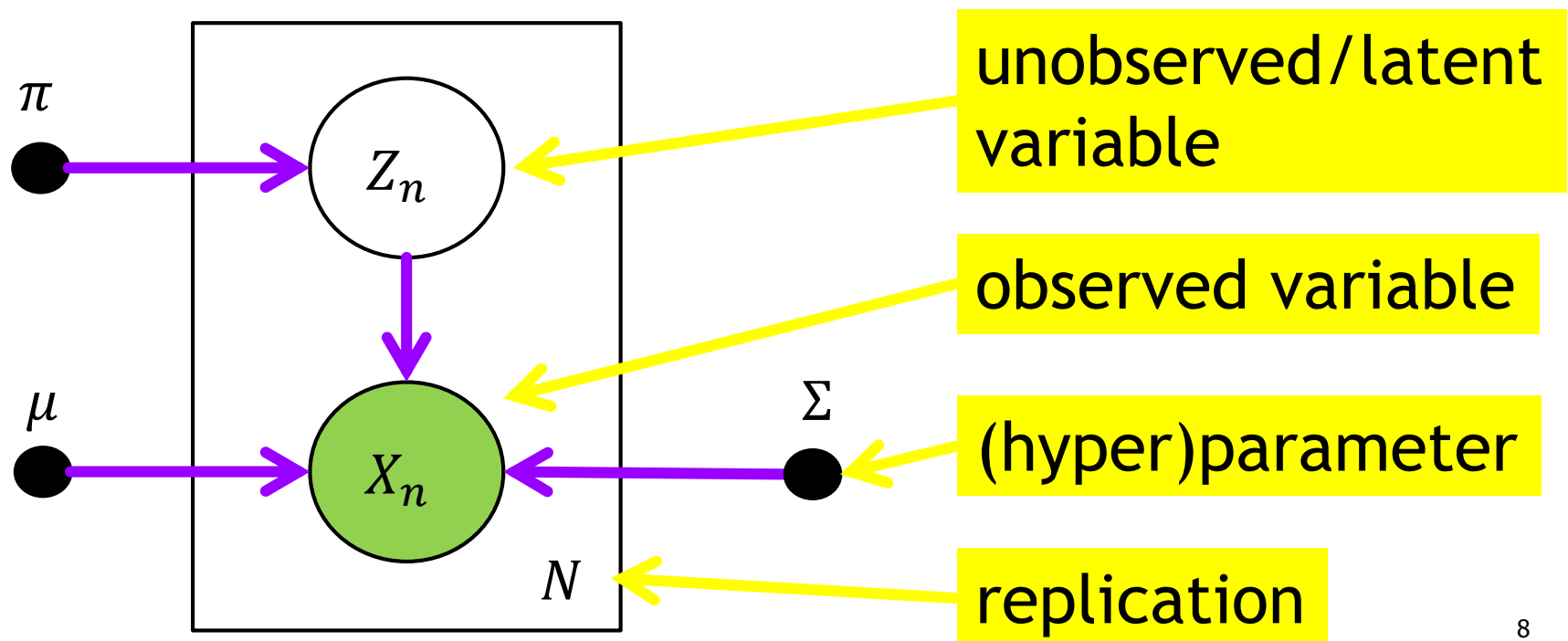
- Convergence:
 - J non-increasing \rightarrow finds local min
 - But optimal clustering not guaranteed
- K-Means: Can be generalized to non-Euclidean distance functions
- Features:
 - Each point attributed to exactly one cluster
 - Not a generative model: cannot “simulate” data based on learned $\{r_{nk}\}, \{\mu_k\}$ (no distribution for new values)
- Improvement:
 - Soft attribution: a point can belong to several clusters
 - Generative model: distribution over points

Gaussian mixture model (GMM)

- GMM: distribution of single data point

$$p(x) = \sum_{k=1}^K \pi_k N(x|\mu_k, \Sigma_k)$$

- Random variable Z_k : point belongs to cluster k
 - $p(Z_k = 1) = \pi_k$: mixing coefficients



GMM

- Latent variable $Z = (Z_1, Z_2, \dots, Z_k, \dots, Z_K)$:
 - $p(Z) = \prod_k \pi_k^{Z_k}$
 - “One-hot”: $Z = (0, 0, 0, \dots, 1, \dots, 0, 0)$
- Conditional distribution of data point:
 - $p(X|Z_k = 1) = N(x|\mu_k, \Sigma_k)$, i.e.,
 - $p(X|Z) = \prod_k N(x|\mu_k, \Sigma_k)^{Z_k}$
- Data distribution:
 - $p(X) = \sum_z p(Z)p(X|Z) = \sum_k \pi_k N(x|\mu_k, \Sigma_k)$
- Conclusion:
 - Gaussian mixture can be viewed as follows: choose a cluster k with distribution $\{\pi\}$; then generate a point according to Gaussian $N(X|\mu_k, \Sigma_k)$ of the chosen cluster

GMM: posterior

- Finding clusters = computing posterior, ie, distribution of Z given data X
- Def: $\gamma_k = p(Z_k = 1|X)$
- Bayes' theorem:
 - $$\gamma_k = \frac{p(Z_k=1)p(X|Z_k=1)}{\sum_j p(Z_j=1)p(X|Z_j=1)} = \frac{\pi_k N(X|\mu_k, \Sigma_k)}{\sum_j \pi_j N(X|\mu_j, \Sigma_j)}$$
- Interpretation:
 - For fixed $\{\mu_k, \Sigma_k\}$,
 π_k is the prior for the cluster Z of point X , and
 γ_k is the posterior

GMM: ML estimator for μ

- Log-likelihood function:

- $L = \log p(X_1, \dots, X_n | \pi, \mu, \Sigma) = \sum_n \log \sum_k \pi_k N(x_n | \mu_k, \Sigma_k)$

- Maximizing w.r.t. μ :

- $\frac{\partial L}{\partial \mu_k} = 0 \Rightarrow$

$$\sum_n \underbrace{\frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_j \pi_j N(x_n | \mu_j, \Sigma_j)}}_{\gamma_{nk}} \Sigma_k^{-1} (x_n - \mu_k) = 0$$

- Solution: $= \gamma_{nk}$: posterior of Z given $X = X_n$
- $\mu_k = \frac{1}{N_k} \sum_n \gamma_{nk} x_n$ (roughly, weighted center of mass)
 - with $N_k = \sum_n \gamma_{nk}$ (roughly, # of points in class k)

GMM: ML estimator for $\{\Sigma\}$ and π

- Maximizing w.r.t. Σ :
 - $\Sigma_k = \frac{1}{N_k} \sum_n \gamma_{nk} (x_n - \mu_k)(x_n - \mu_k)^T$
 - (roughly, weighted empirical covariance matrix within class k)
- Maximizing w.r.t. π :
 - $\pi_k = \frac{N_k}{N}$
 - (roughly, number of points attributed to cluster k)

EM algorithm for GMM

E-step:

Compute posterior of latent variables Z given parameters from

M-step:

$$\gamma_{nk} = \frac{\pi_k N(X_n | \mu_k, \Sigma_k)}{\sum_j \pi_j N(X_n | \mu_j, \Sigma_j)}$$

M-step:

Compute new parameters using distribution of latent variables from E-step:

$$\mu_k = \frac{1}{N_k} \sum_n \gamma_{nk} x_n$$

$$\Sigma_k = \frac{1}{N_k} \sum_n \gamma_{nk} (x_n - \mu_k)(x_n - \mu_k)^T$$

$$\pi_k = \frac{N_k}{N}$$

E-step 2

E-step 3

M-step 2

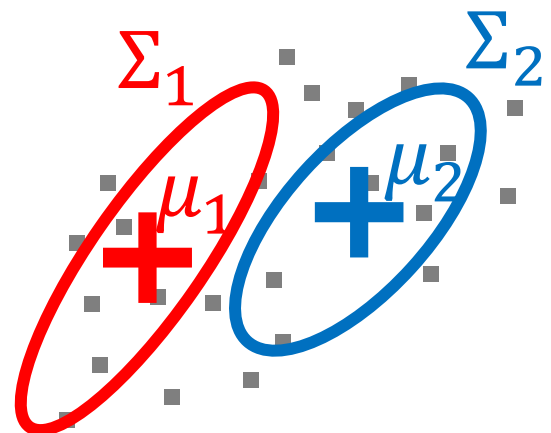
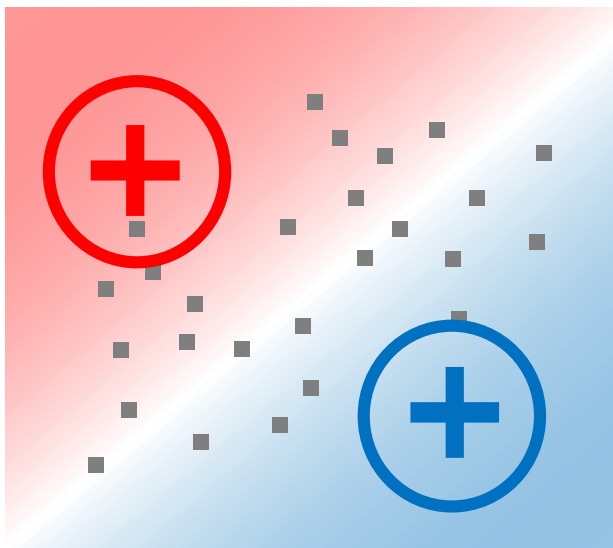
Until convergence of likelihood
 $L = \log p(X_1, \dots, X_n | \pi, \mu, \Sigma)$

EM for GMM: example

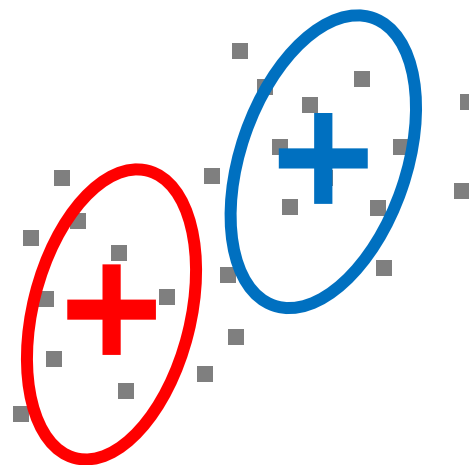
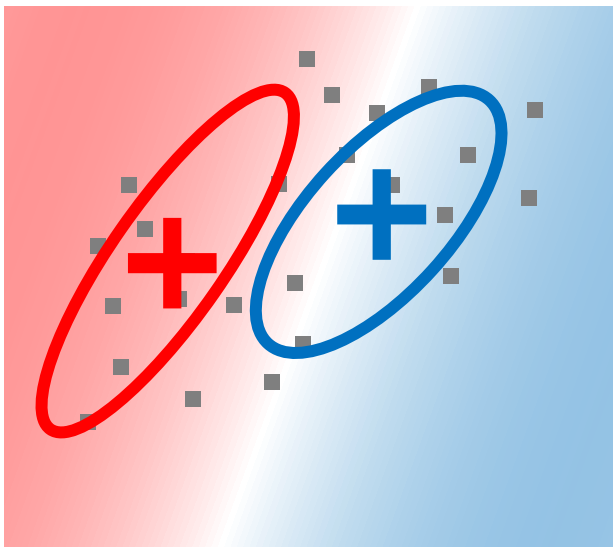
E-step

M-step

Iter 1



Iter 2



K-means vs GMM

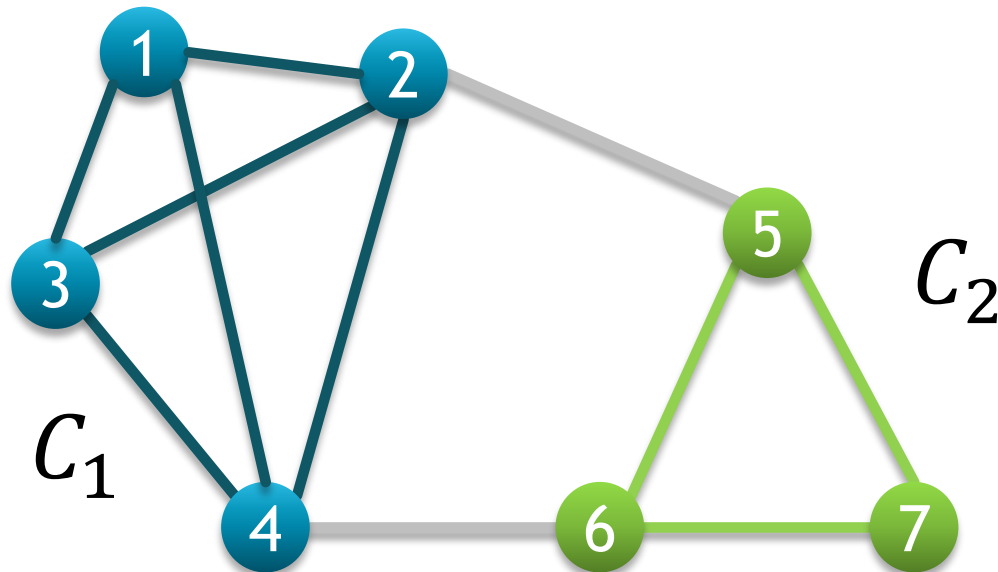
	K-Means	GMM
Membership	Hard	Soft
Generative	No	Yes
E-step updates	r_{nk} = closest center	$\gamma_{nk} = P(Z_k X_n)$
M-step updates	μ_k	μ_k, Σ_k, π_k
Convergence	Guaranteed	Guaranteed
Optimal	No	No
Characterization	Centers, membership	Centers, weights, shapes

Community detection: goal and def

- Find organization in graphs:
 - Email or phone graph → find organizational units
 - Citation networks → scientific topics and their relationships
 - Social networks → groups with shared interests (language, etc.)
- Definition:
 - Given: a network $G(V, E)$
 - Find: a partition (hard or soft), or a hierarchy, such that node in same community are more “meshed” than other nodes

Modularity: strength of communities

- Def:
 - Partitioning of nodes into communities: $\{C_i\}$
 - $Q = \frac{1}{2m} \sum_{C_i \in \mathcal{C}} \sum_{u,v \in C_i} \left(\mathbb{1}_{uv} - \frac{d_u d_v}{2m} \right)$



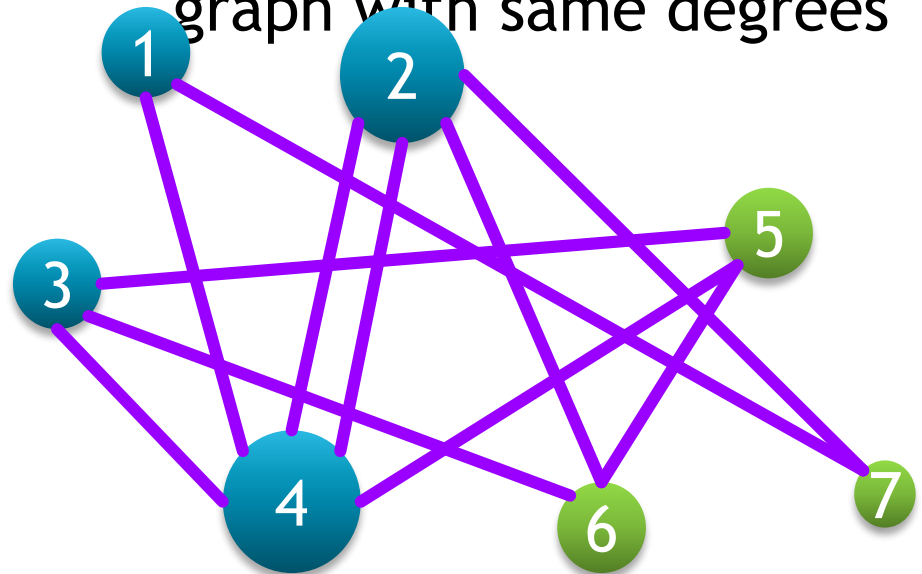
Note: inner sum is over all *ordered* (u, v) , and includes (u, u)

Modularity: interpretation

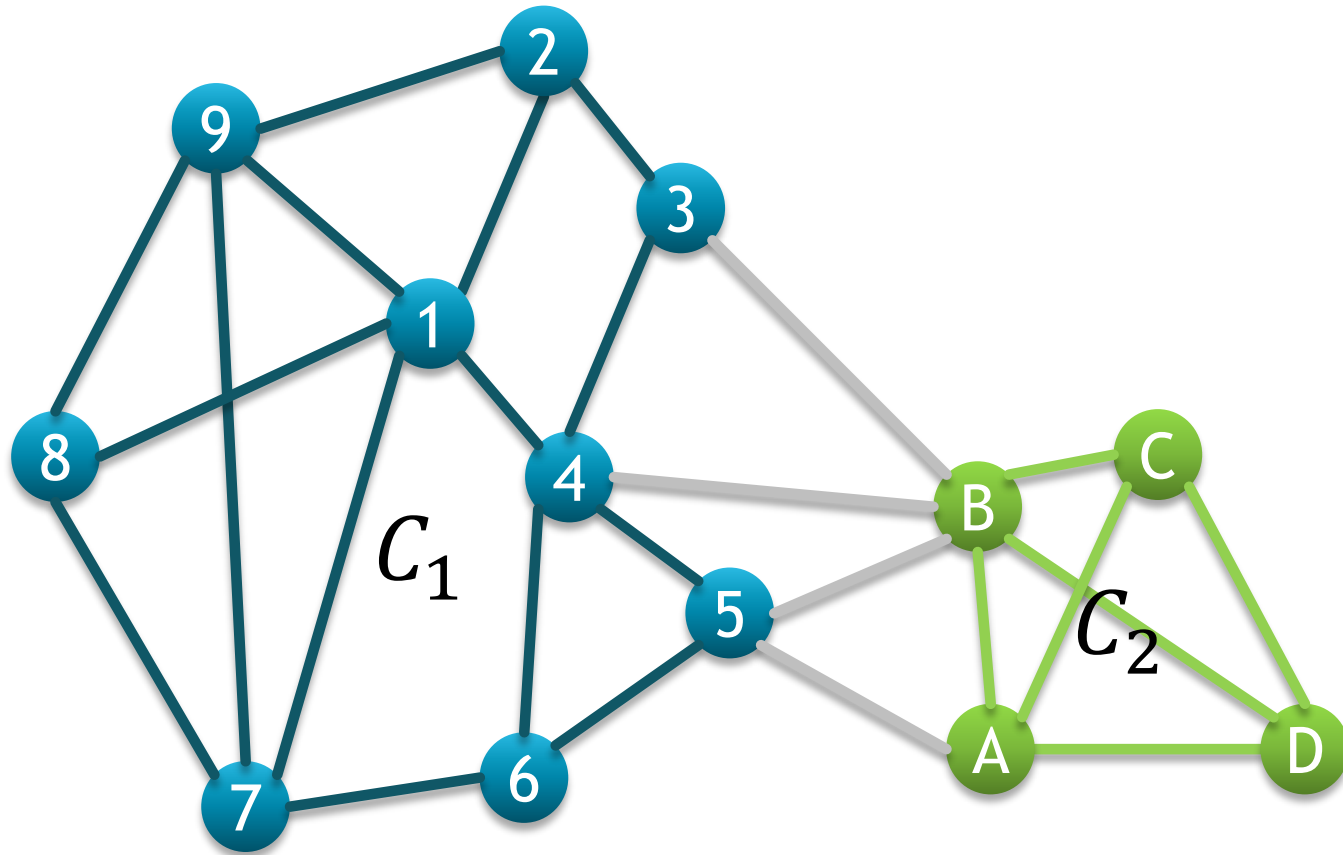
- Number of stubs: $2m$
- Expected # edges (u, v) in C_i :
$$\frac{1}{2} \sum_{u,v \in C_i} \frac{d_u d_v}{2m} = e_i$$
- Actual # edges in community C_i :
$$\frac{1}{2} \sum_{u,v \in C_i} 1_{uv} = m_i$$
- Modularity: compares actual graph with unstructured graph with same node

weights

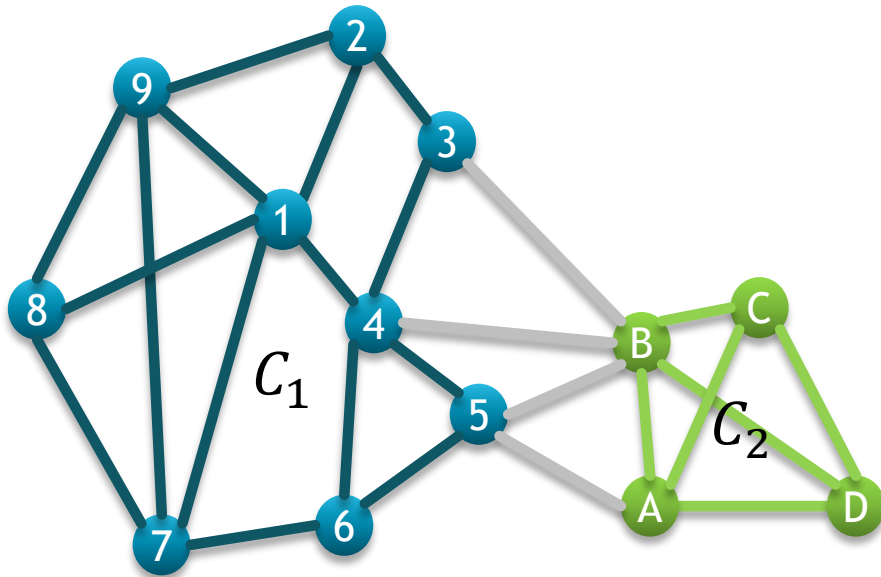
- $Q = \frac{\sum_{C_i} (m_i - e_i)}{m}$
 - m_i : # edges in community C_i
 - e_i : expected # edges in community C_i in random graph with same degrees



Modularity: example



Modularity: example



C_1 :

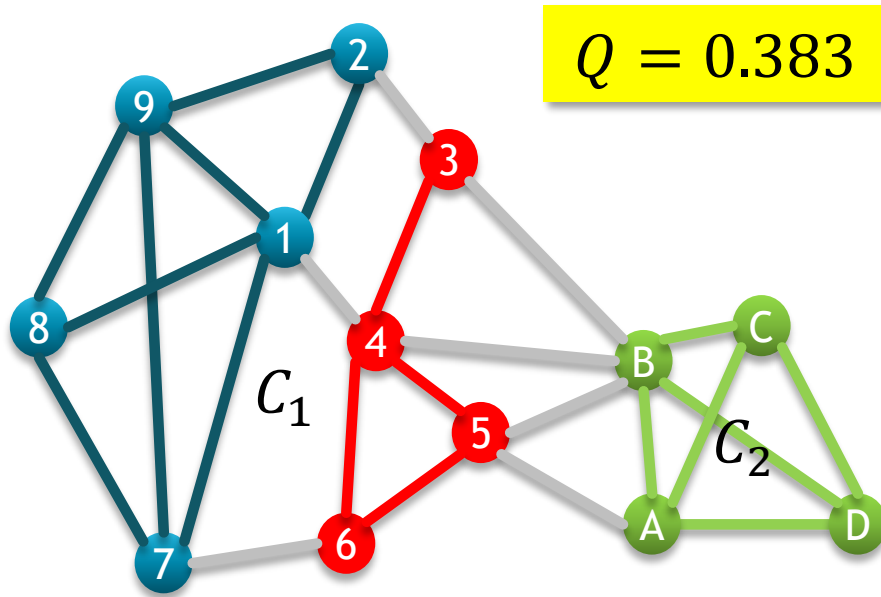
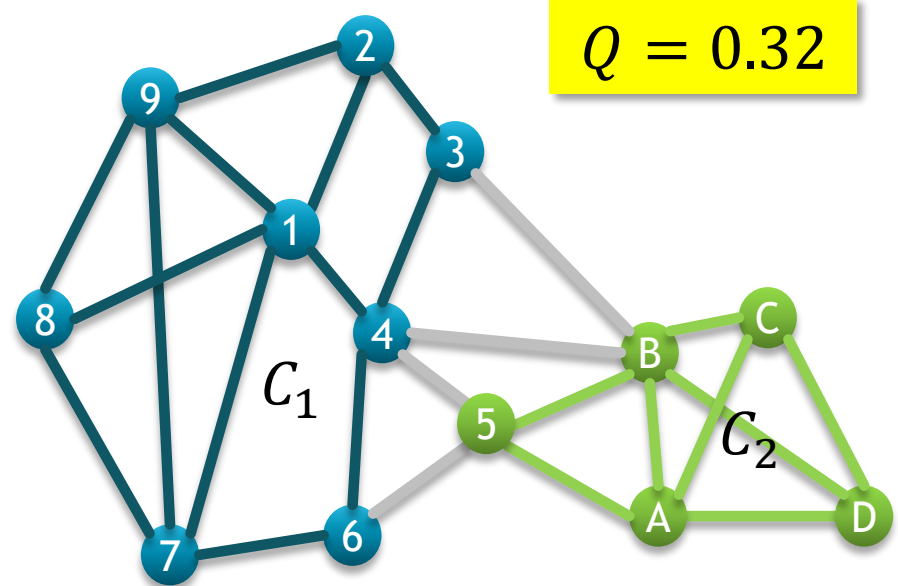
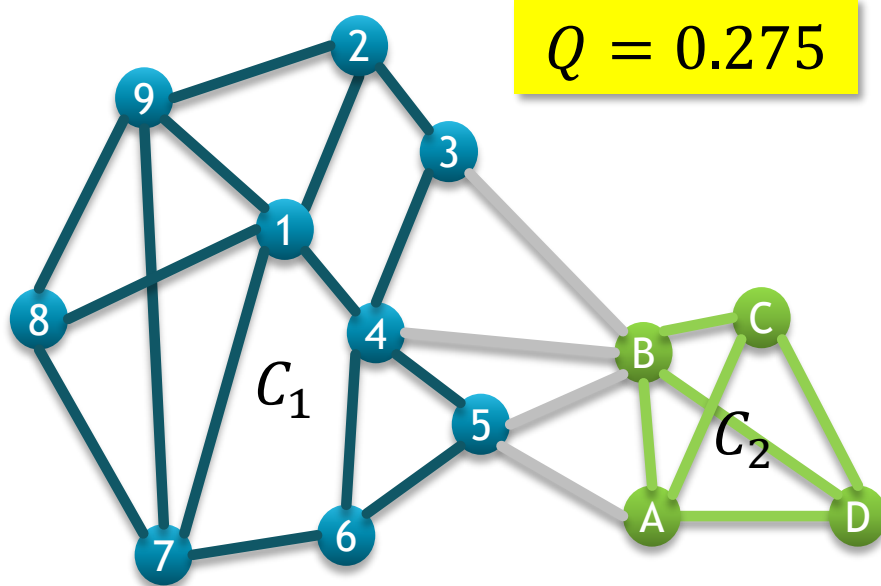
- $m_1 = 15$
- $e_1 = 11.56$

C_2 :

- $m_2 = 6$
- $e_2 = 2.56$

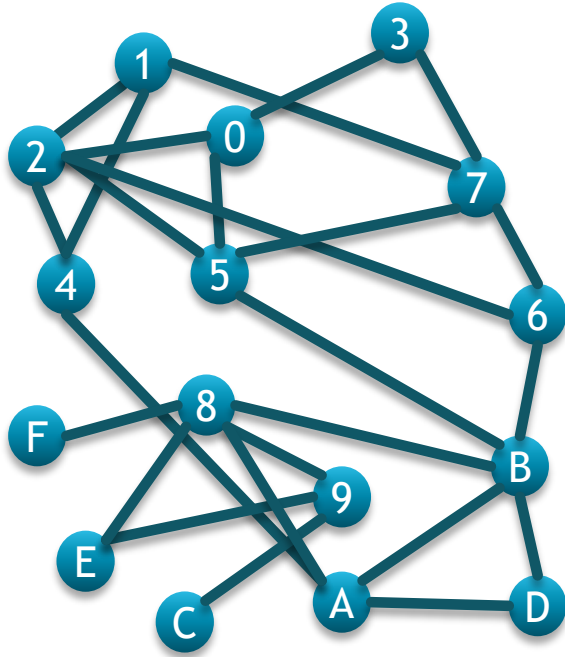
- $m = 25$
- Edge (4,8):
 - $d_4 = 5; d_8 = 3$
 - Expected $\# = \frac{15}{50} \sim 0.3$
- $Q = \frac{(15+6)-(11.56+2.56)}{25} = 0.275$

Modularity: example



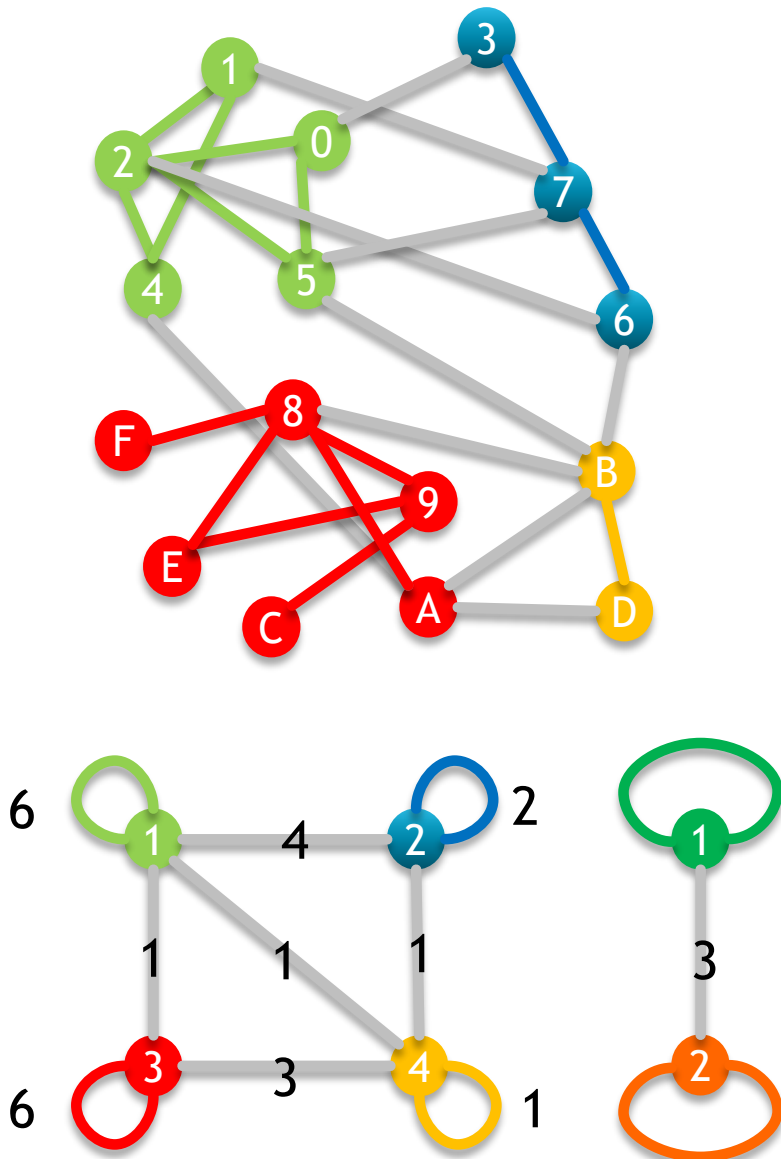
- Max-modularity is NP-hard
- Need efficient heuristics

Louvain method



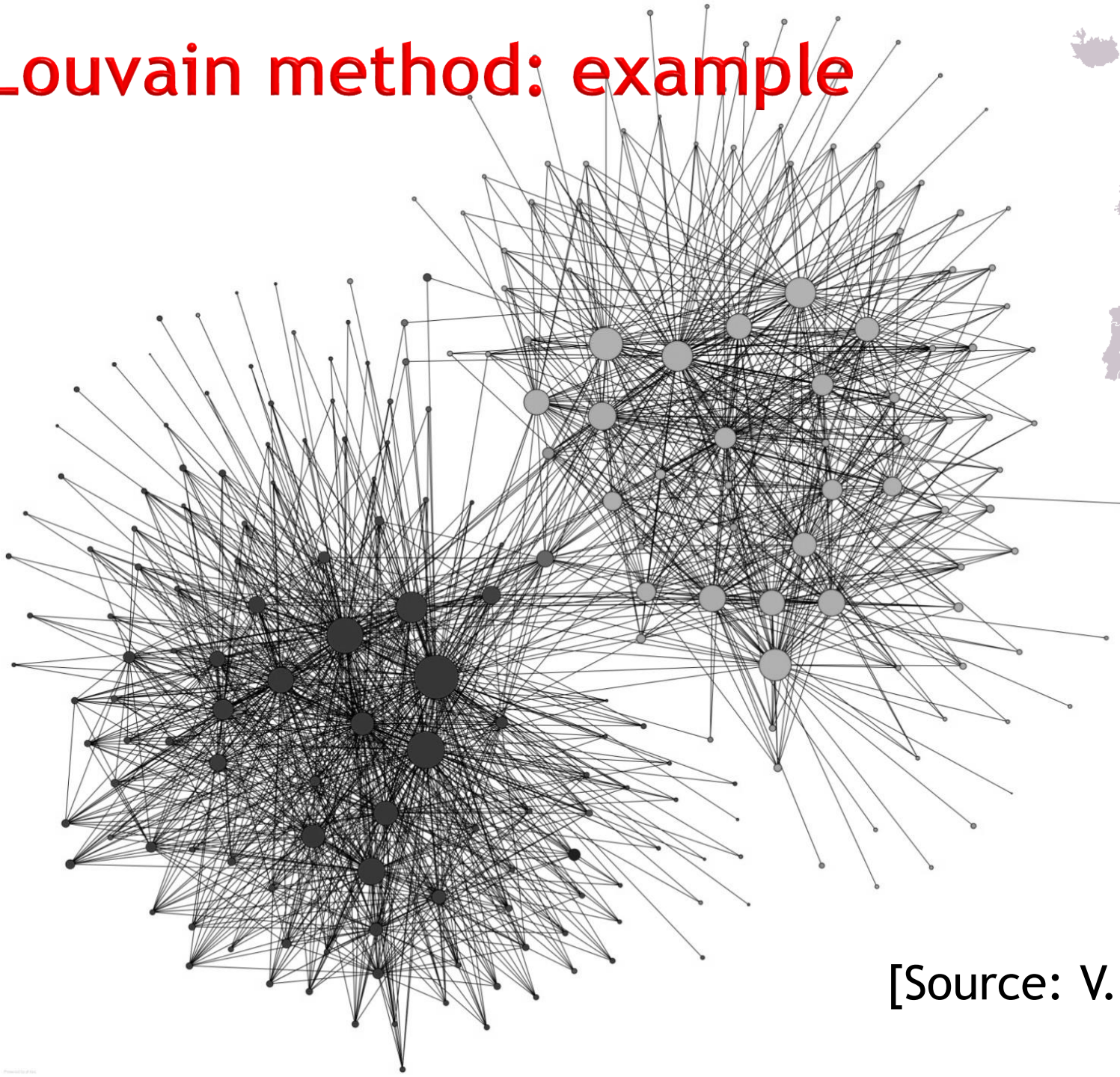
- Idea:
 - Building hierarchy of communities
 - Bottom-up: start with each node a separate community, then coalesce communities
- For every node u :
 - Compare modularity if u is added to the community of a neighbor v
 - Choose neighbor v that increases modularity most; if none, leave u in current community

Louvain method



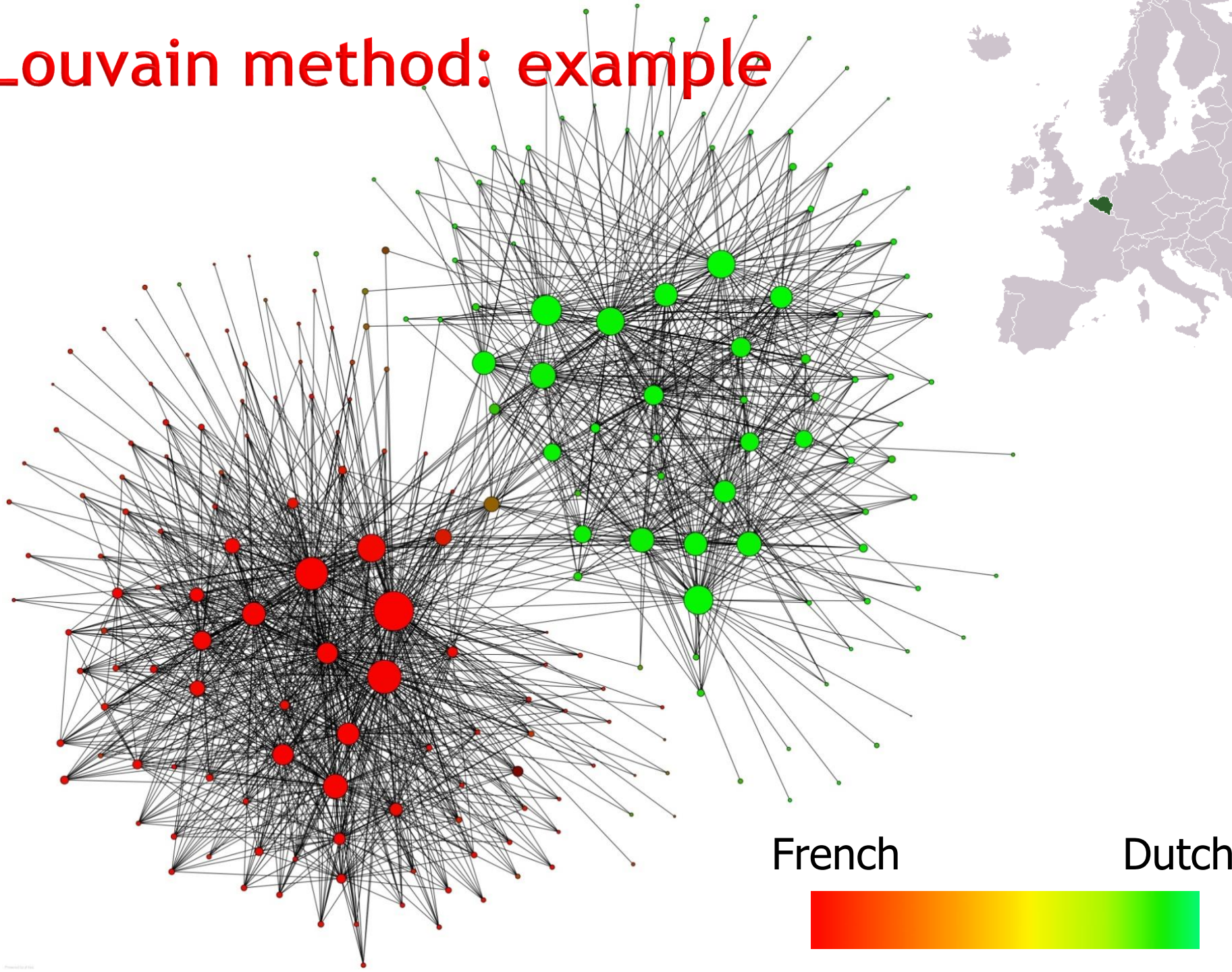
- Iterate through all nodes u (possibly several times) until no more modularity increases possible
 - Local maximum
- Form a new graph capturing the network of communities
 - Link weights = # edges between communities
 - Self-loops: internal edges
- Repeat the procedure until convergence

Louvain method: example



[Source: V. Blondel]

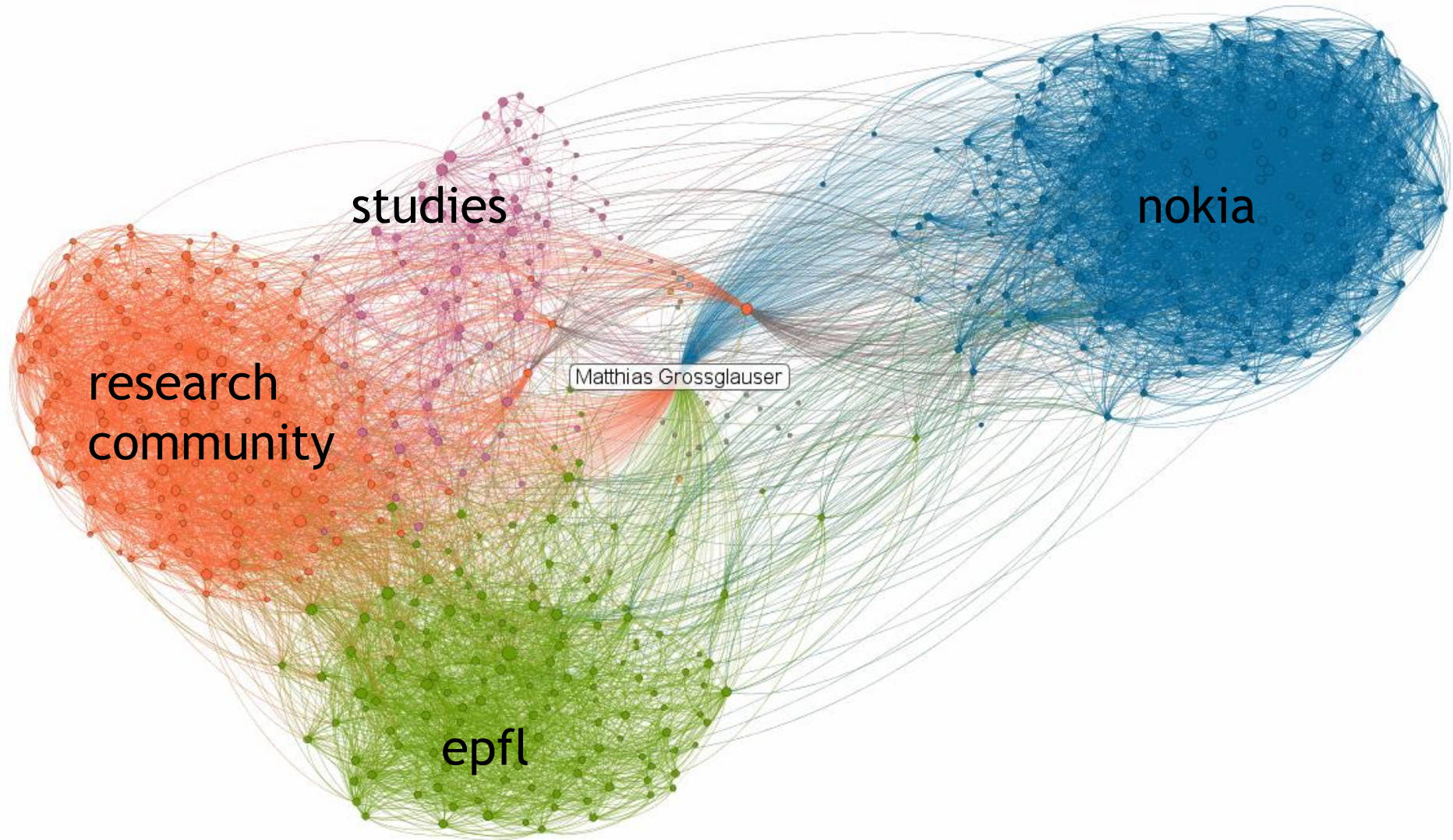
Louvain method: example



Louvain method in LinkedIn Labs

LinkedIn Maps

Matthias Grossglauser's Professional Network
as of February 16, 2013



Summary

- Unsupervised techniques for grouping data
- Clusters: set of points close in distance, far from other clusters
 - Criterion: distances to center, likelihood
 - GMM: Gaussian parameters characterize cluster
- Community: set of nodes with high edge density, low edge density to other communities
 - Criterion = modularity
- In general, no optimal solutions
 - Exponential computational cost
- Heuristics:
 - Expectation Maximization for mixture models
 - Louvain method: build hierarchy bottom-up

References

- [Ch. M. Bishop, Pattern Recognition and Machine Learning, Springer, 2006 (chapter 9)]
- [V. Blondel, lecture notes on community detection, 2013]
- [M. Newman, Networks, Oxford UP, 2010]