

Variational EM and K-Means

Machine Learning
CSx824/ECEx242
Bert Huang
Virginia Tech

Outline

- Bounding the log marginal likelihood
- Families of variational distributions
 - Fully factorized distributions
 - Point distributions
- K-means as “hard”-EM

Marginal Likelihood

$$\begin{aligned} p(X|\theta) &= \int_Z p(X, Z|\theta) dZ \\ &= \sum_Z p(X, Z|\theta) \end{aligned}$$

e.g., $X = \{x_1, \dots, x_n\}$

$$\theta = \{\mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K, \dots p(c)\}$$

$$Z = \{z_1, \dots, z_n\} \quad (\text{cluster memberships})$$

$$p(X, Z|\theta) = \prod_{i=1}^n p(z_i) \mathcal{N}(x_i | \mu_{z_i}, \Sigma_{z_i})$$

$$\log p(X|\theta) = \log \sum_Z p(X, Z|\theta)$$

log marginal likelihood

$$\operatorname{argmax}_{\theta} \log \sum_Z p(X, Z|\theta)$$

learning objective

Jensen's Inequality

For any convex function φ ,

$$\varphi(\mathbb{E}[X]) \leq \mathbb{E}[\varphi(X)]$$

For any concave function ϕ ,

$$\phi(\mathbb{E}[X]) \geq \mathbb{E}[\phi(X)]$$

$$\log(\mathbb{E}[X]) \geq \mathbb{E}[\log X]$$

Variational Bound

$$\log(\mathbb{E}[X]) \geq \mathbb{E}[\log X]$$

$$\begin{aligned} \log \sum_Z p(X, Z|\theta) &= \log \sum_Z \frac{q(Z)}{q(Z)} p(X, Z|\theta) & \sum_Z q(Z) &= 1 \\ &= \log \sum_Z q(Z) \frac{p(X, Z|\theta)}{q(Z)} \\ &\geq \sum_Z q(Z) \log \frac{p(X, Z|\theta)}{q(Z)} \\ &= \sum_Z q(Z) \log p(X, Z|\theta) - \sum_Z q(Z) \log q(Z) \end{aligned}$$

Variational Bound

expectation

entropy

$$\log \sum_Z p(X, Z|\theta) \geq \sum_Z q(Z) \log p(X, Z|\theta) - \sum_Z q(Z) \log q(Z)$$

We can pick any q distribution and the bound holds

$$\operatorname{argmax}_{\theta, q \in Q} \sum_Z q(Z) \log p(X, Z|\theta) - \sum_Z q(Z) \log q(Z)$$

$$q(Z) = \prod_{i=1}^n q(z_i) \quad \sum_{z_i} q(z_i) = 1$$

Fully Factorized Variational Family

$$\operatorname{argmax}_{\theta, q \in Q} \sum_Z q(Z) \log p(X, Z | \theta) - \sum_Z q(Z) \log q(Z)$$

$$q(Z) = \prod_{i=1}^n q(z_i) \quad \sum_{z_i} q(z_i) = 1$$

$$\operatorname{argmax}_{\theta, q \in Q} \sum_{i=1}^n \sum_{z_i} q(z_i) \log p(x_i, z_i | \theta) - q(z_i) \log q(z_i)$$

Point Distributions

$$\operatorname{argmax}_{\theta, q \in Q} \sum_Z q(Z) \log p(X, Z | \theta) - \sum_Z q(Z) \log q(Z)$$

$$q(Z) = \prod_{i=1}^n q(z_i) \quad q(z_i) = \begin{cases} 1 & \text{if } z_i = \hat{z}_i \\ 0 & \text{otherwise} \end{cases}$$

$$\operatorname{argmax}_{\theta, q \in Q} \sum_{i=1}^n \sum_{z_i} q(z_i) \log p(x_i, z_i | \theta) - q(z_i) \log q(z_i)$$

$$\operatorname{argmax}_{\theta, q \in Q, \hat{Z}} \sum_{i=1}^n \log p(x_i, \hat{z}_i | \theta)$$

point distributions are often easier to compute, but less robust

Point Distributions for GMMs

$$\operatorname{argmax}_{\theta, q \in Q, \hat{Z}} \sum_{i=1}^n \log p(x_i, \hat{z}_i | \theta)$$
$$\sum_{i=1}^n \log \mathcal{N}(x_i | \mu_{\hat{z}_i}, \Sigma_{\hat{z}_i})$$

$$\hat{z}_i \leftarrow \operatorname{argmax}_z \log \mathcal{N}(x_i | \mu_z, \Sigma_z)$$

$$\mu_z \leftarrow \frac{\sum_{i; \hat{z}_i=z} x_i}{\sum_{i; \hat{z}_i=z} 1} \quad \Sigma_z \leftarrow \frac{\sum_{i; \hat{z}_i=z} (x_i - \mu_i)(x_i - \mu_i)^\top}{\sum_{i; \hat{z}_i=z} 1}$$

K-means

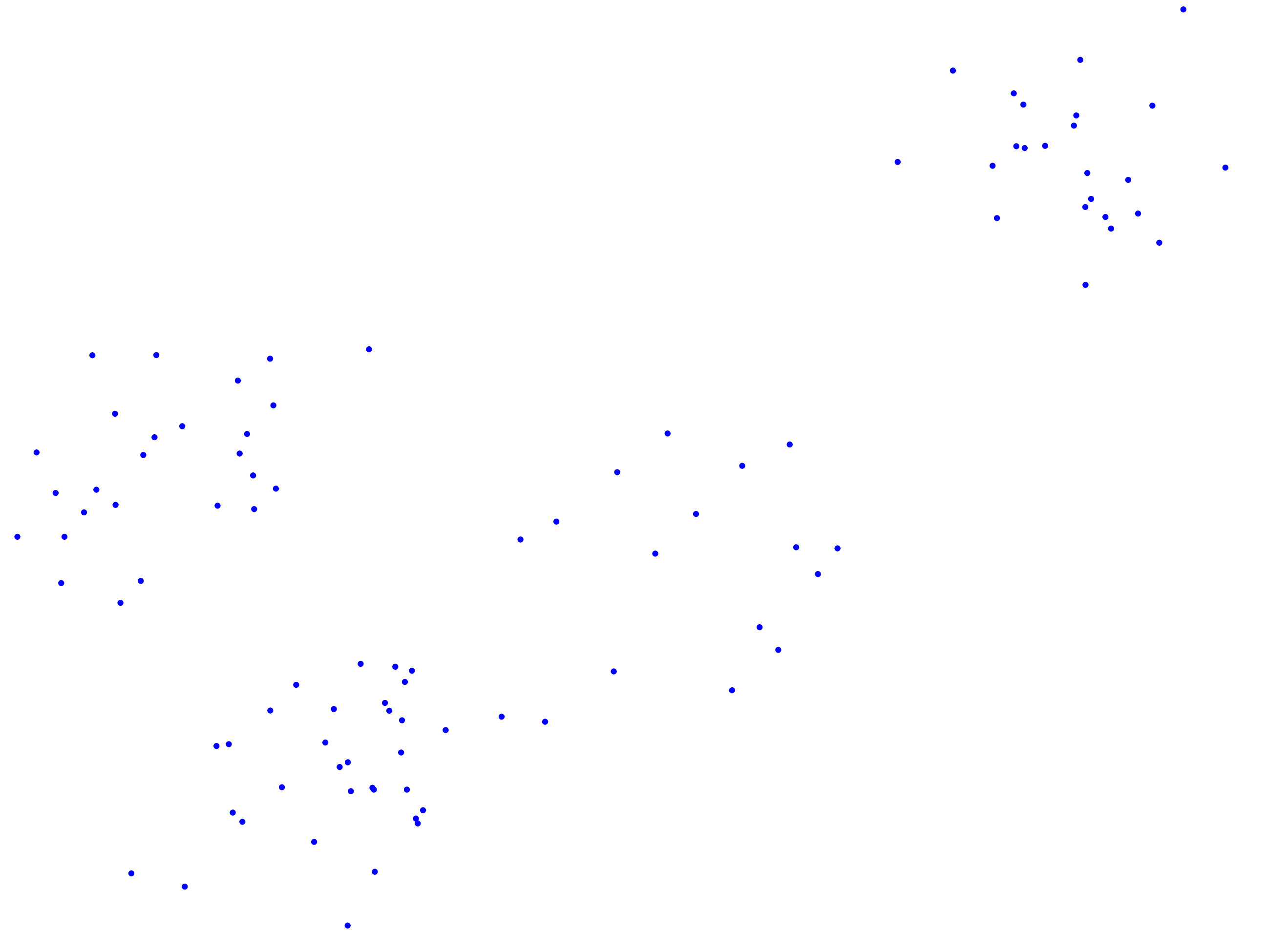
$$\hat{z}_i \leftarrow \operatorname{argmin}_z \|x_i - \mu_z\|$$

assign points to
closest mean

$$\mu_z \leftarrow \frac{\sum_{i; \hat{z}_i=z} x_i}{\sum_{i; \hat{z}_i=z} 1}$$

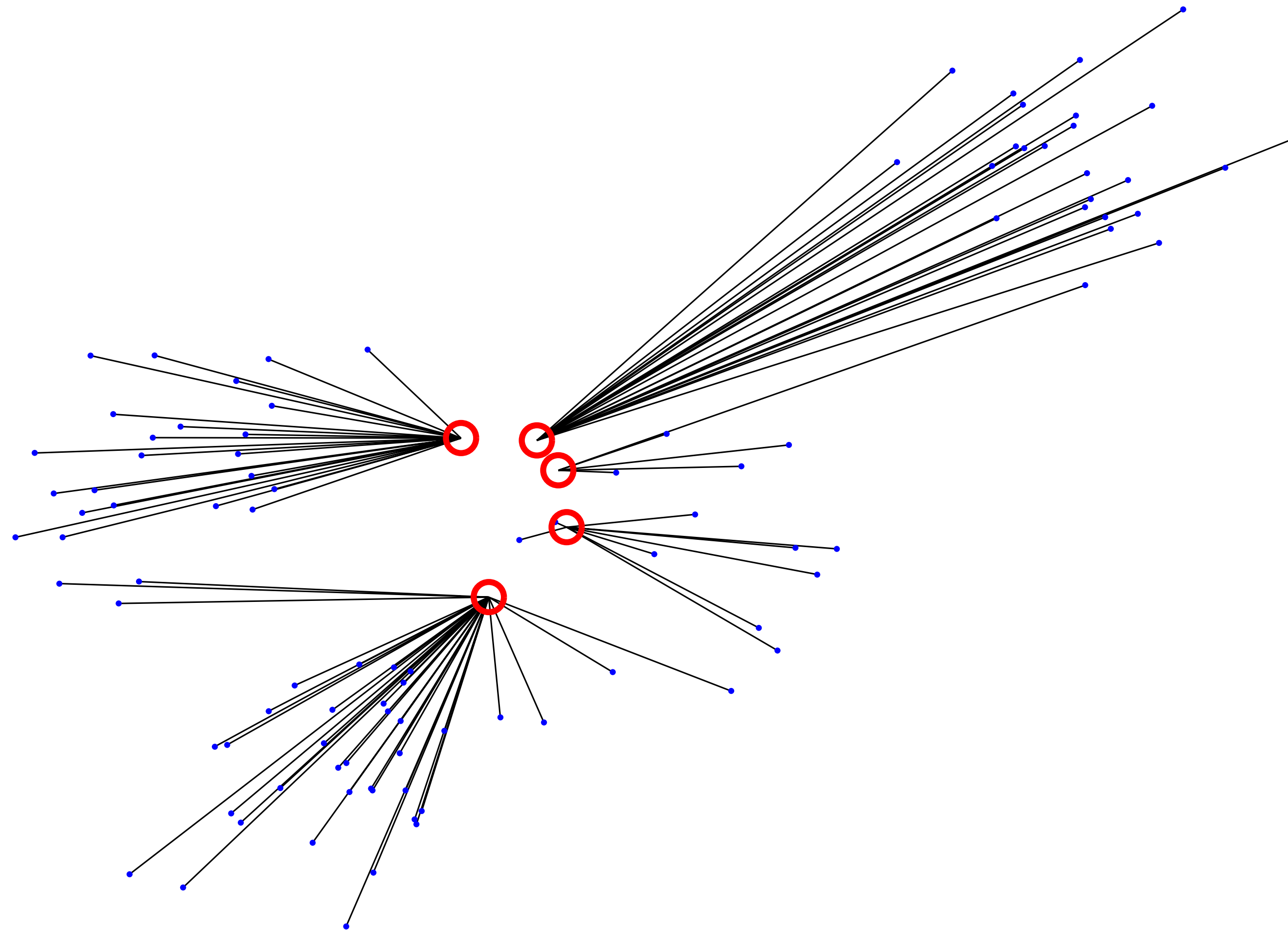
set means to average
of points in cluster

Example



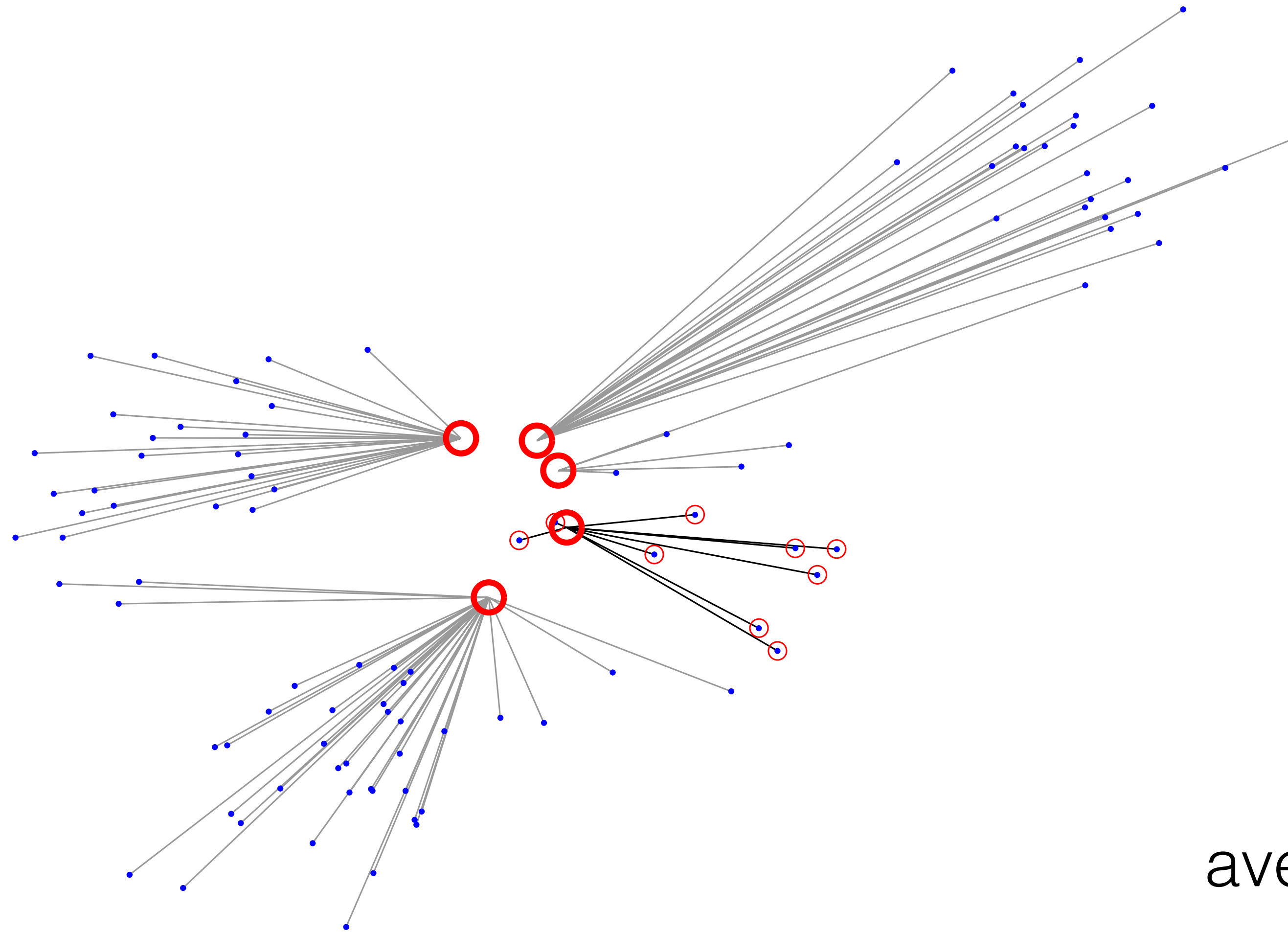
input data

Example



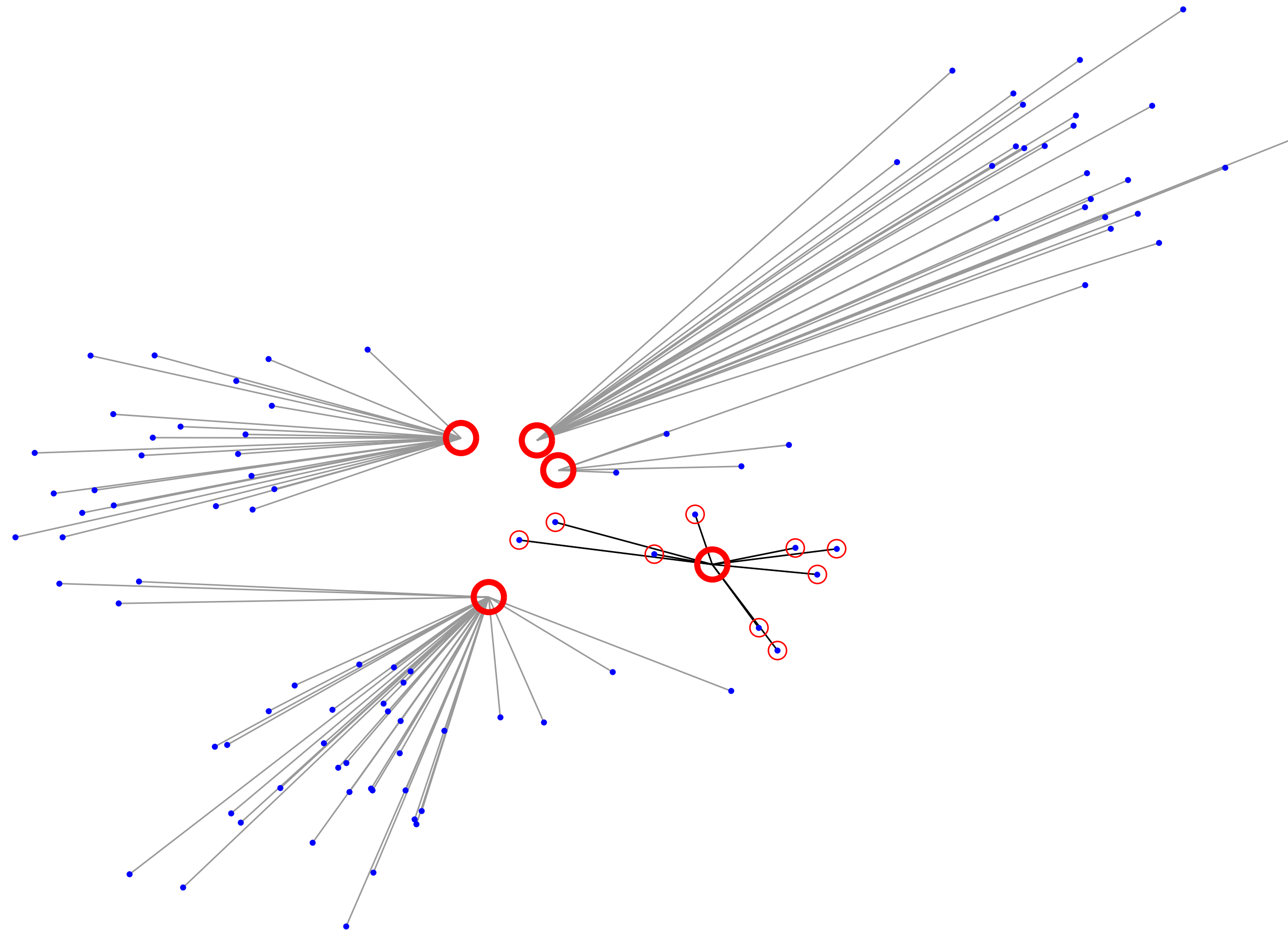
assign points to
initialized means

Example



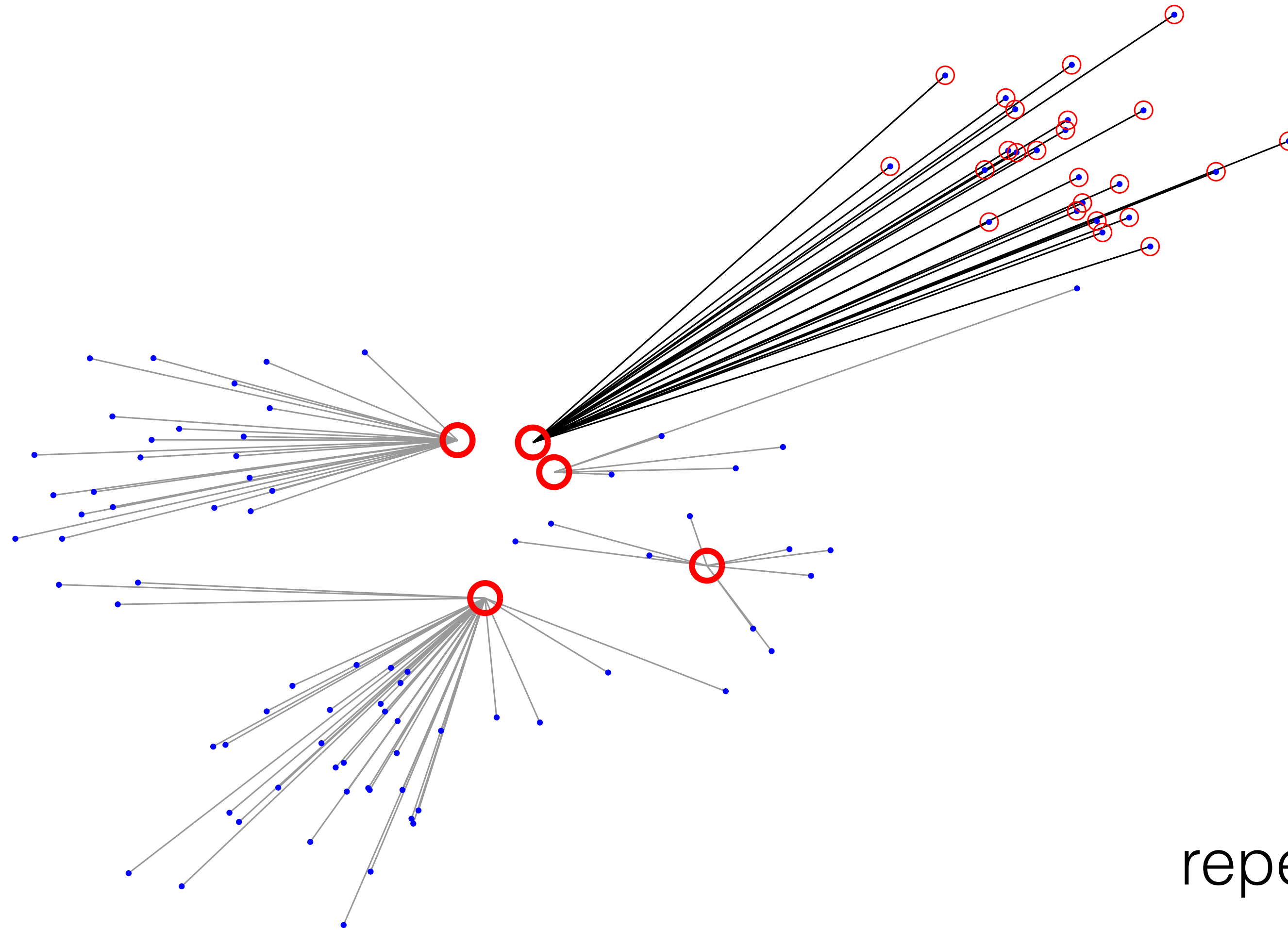
average each cluster

Example



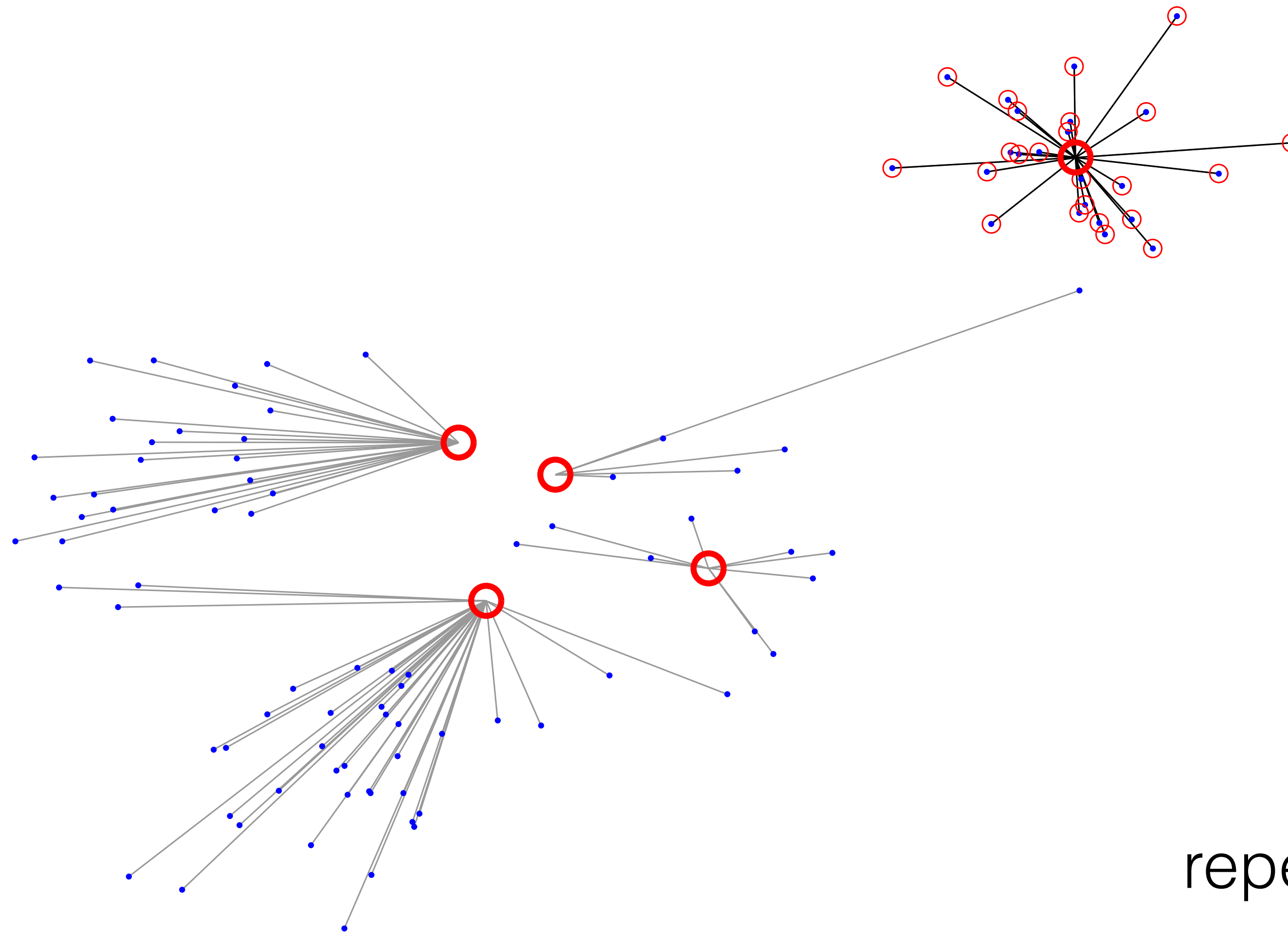
update mean

Example



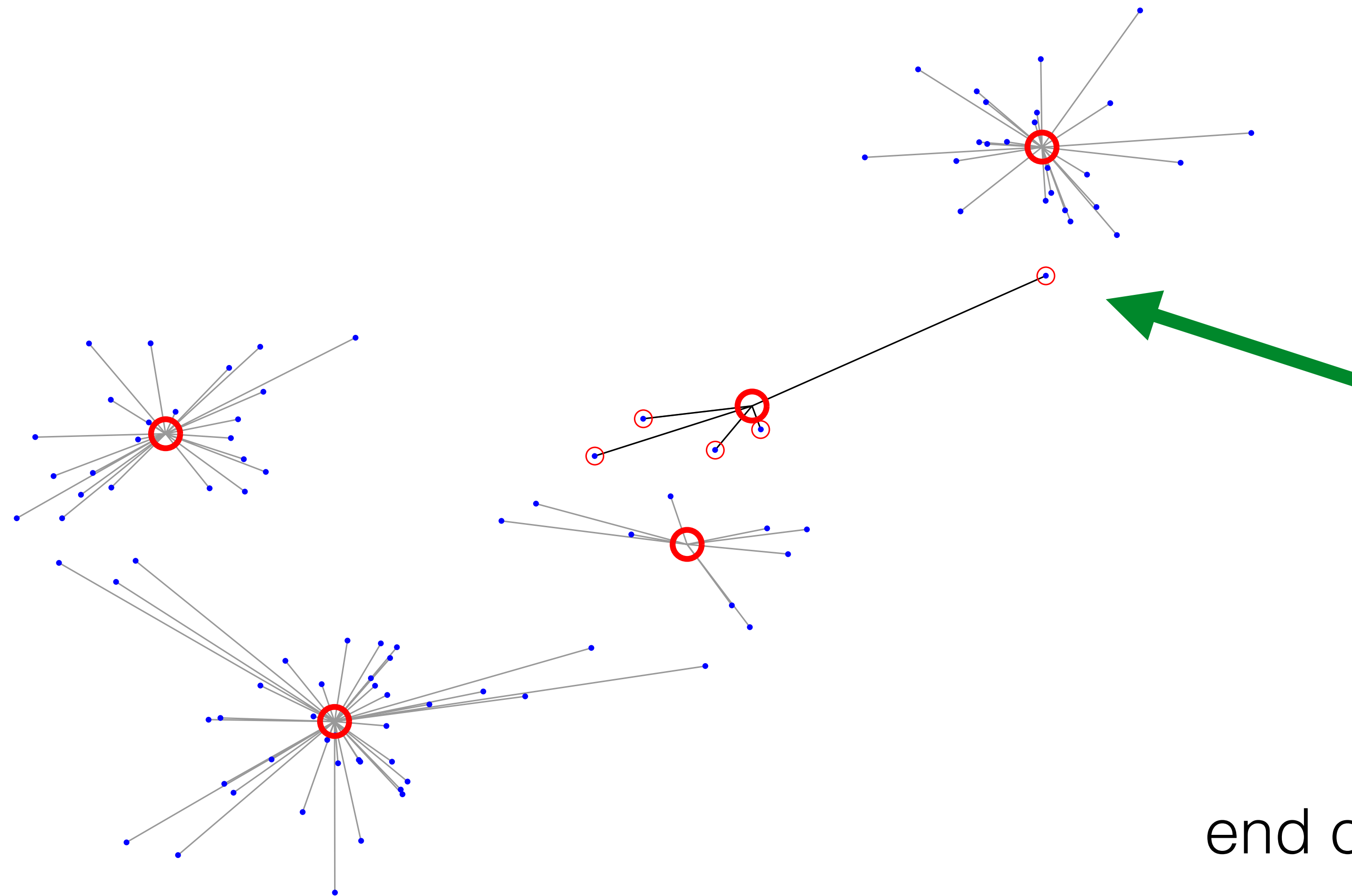
repeat for all clusters

Example

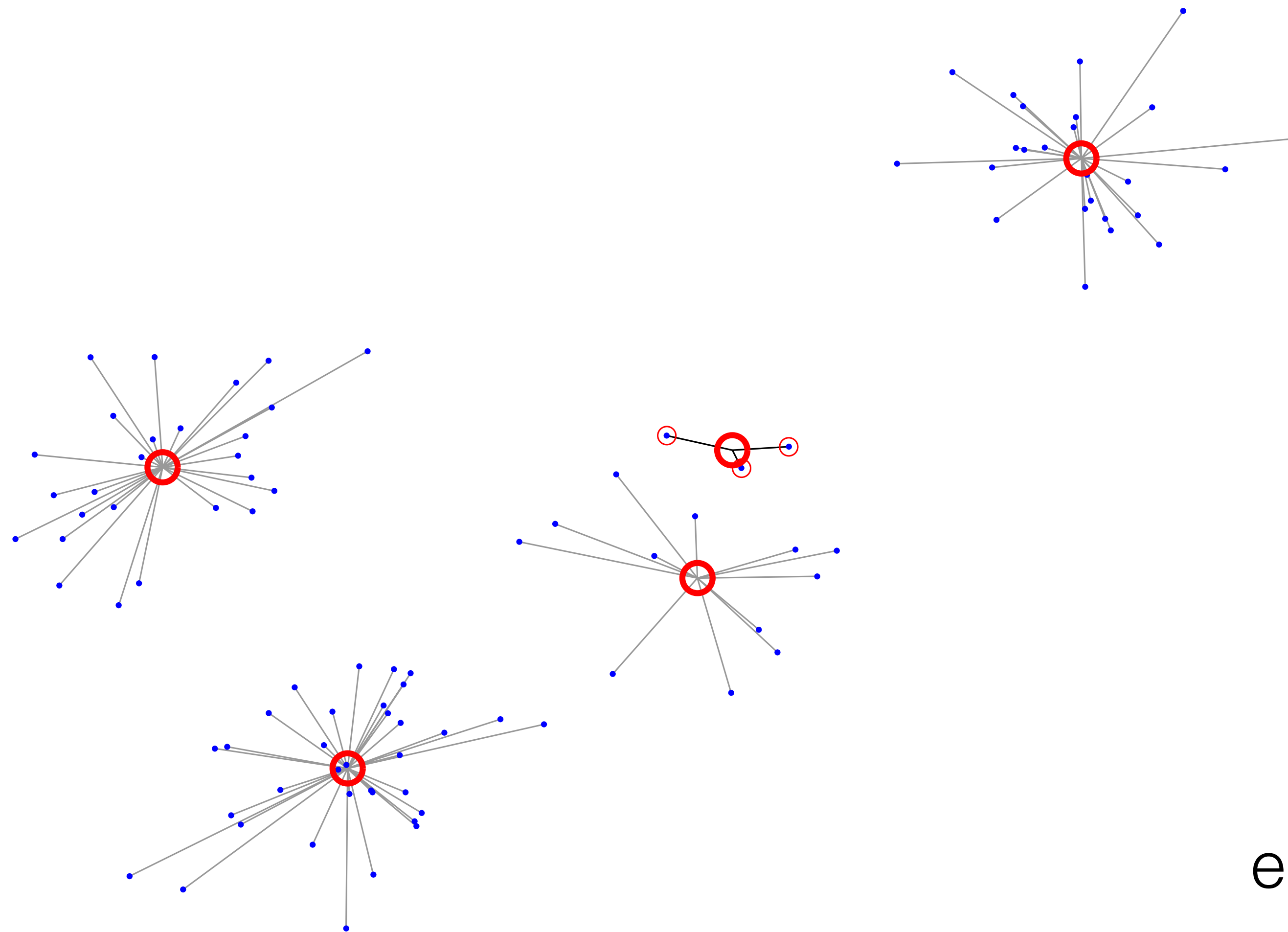


repeat for all clusters

Example

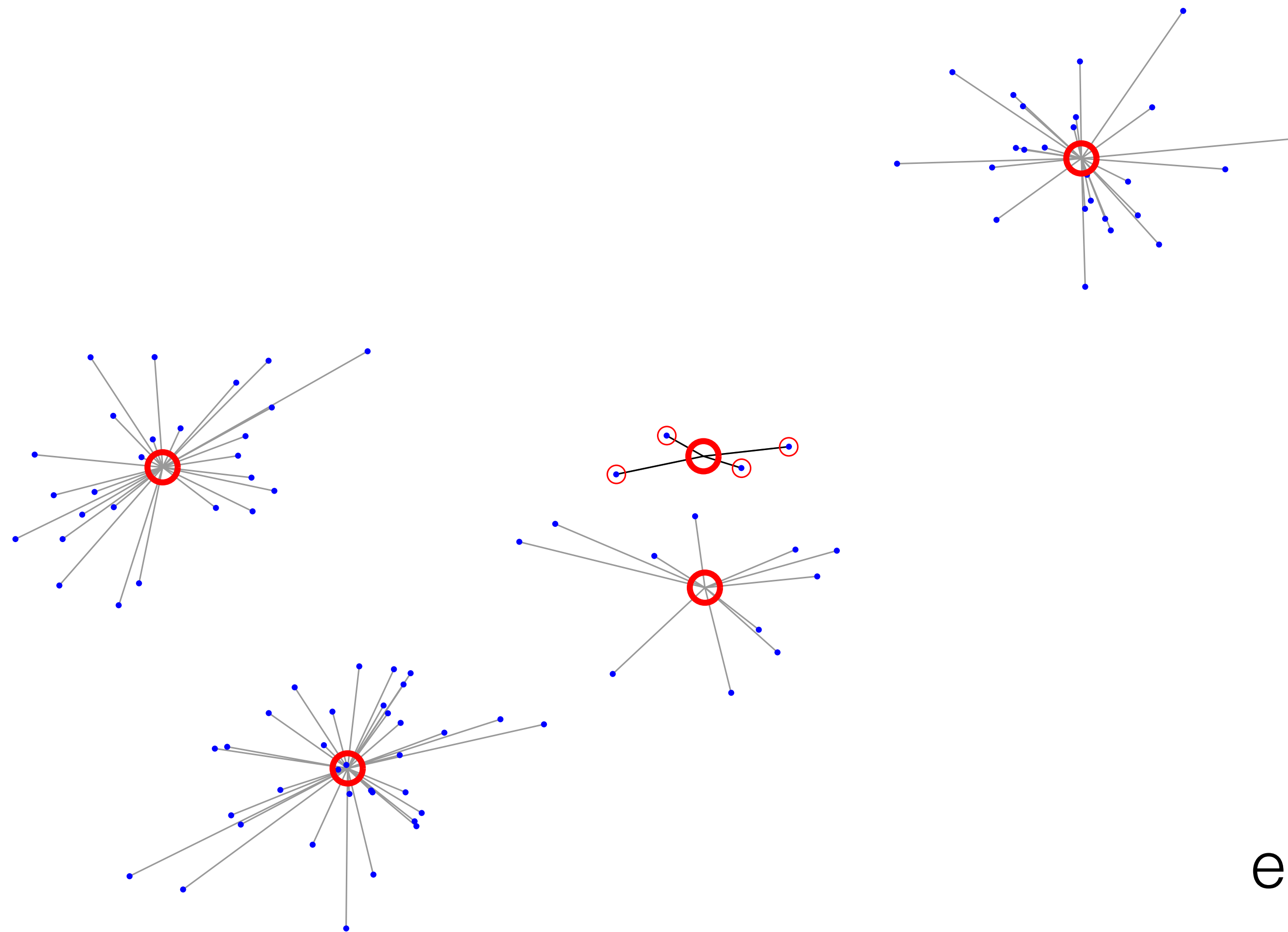


Example



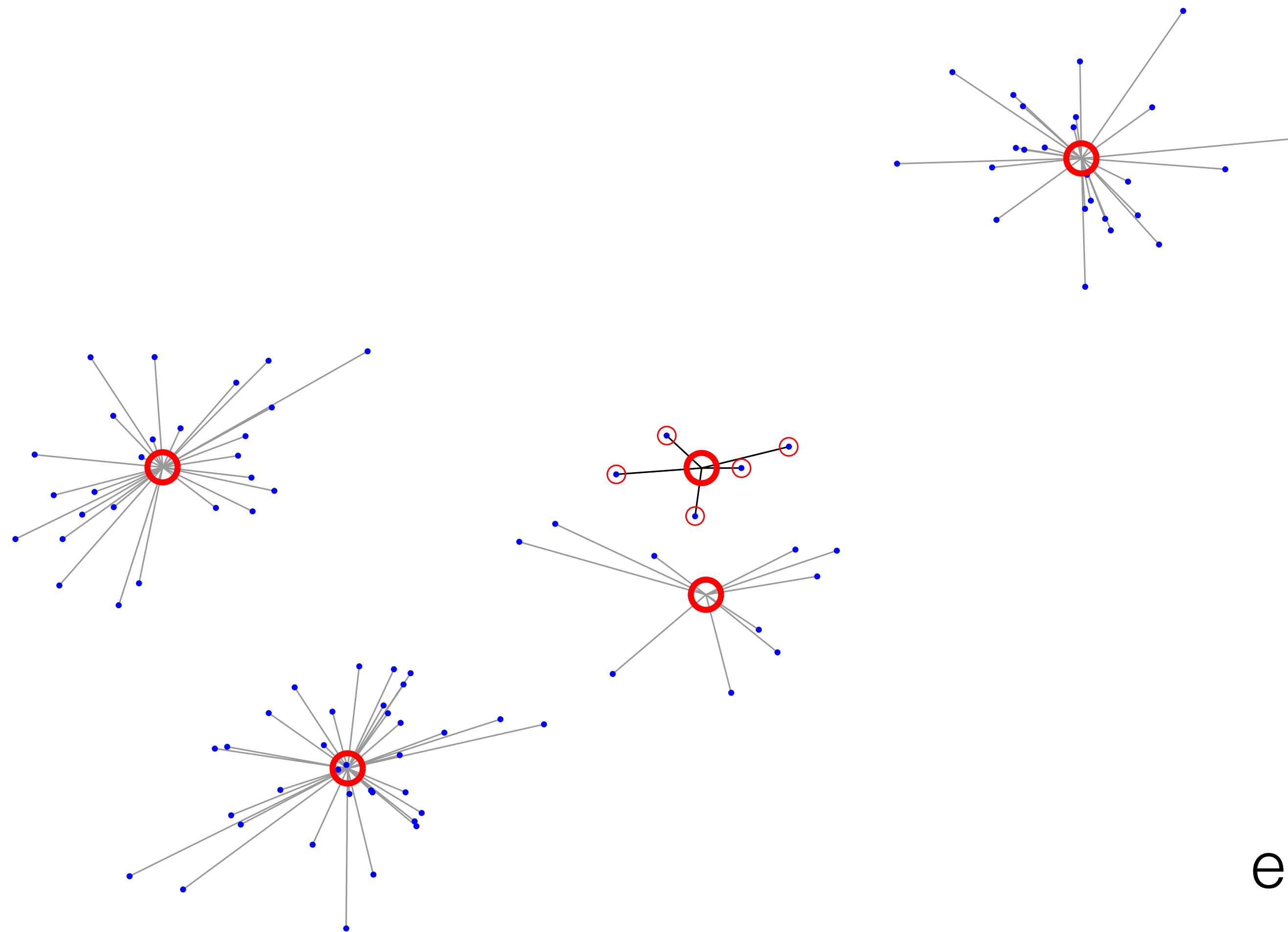
end of iteration 2

Example



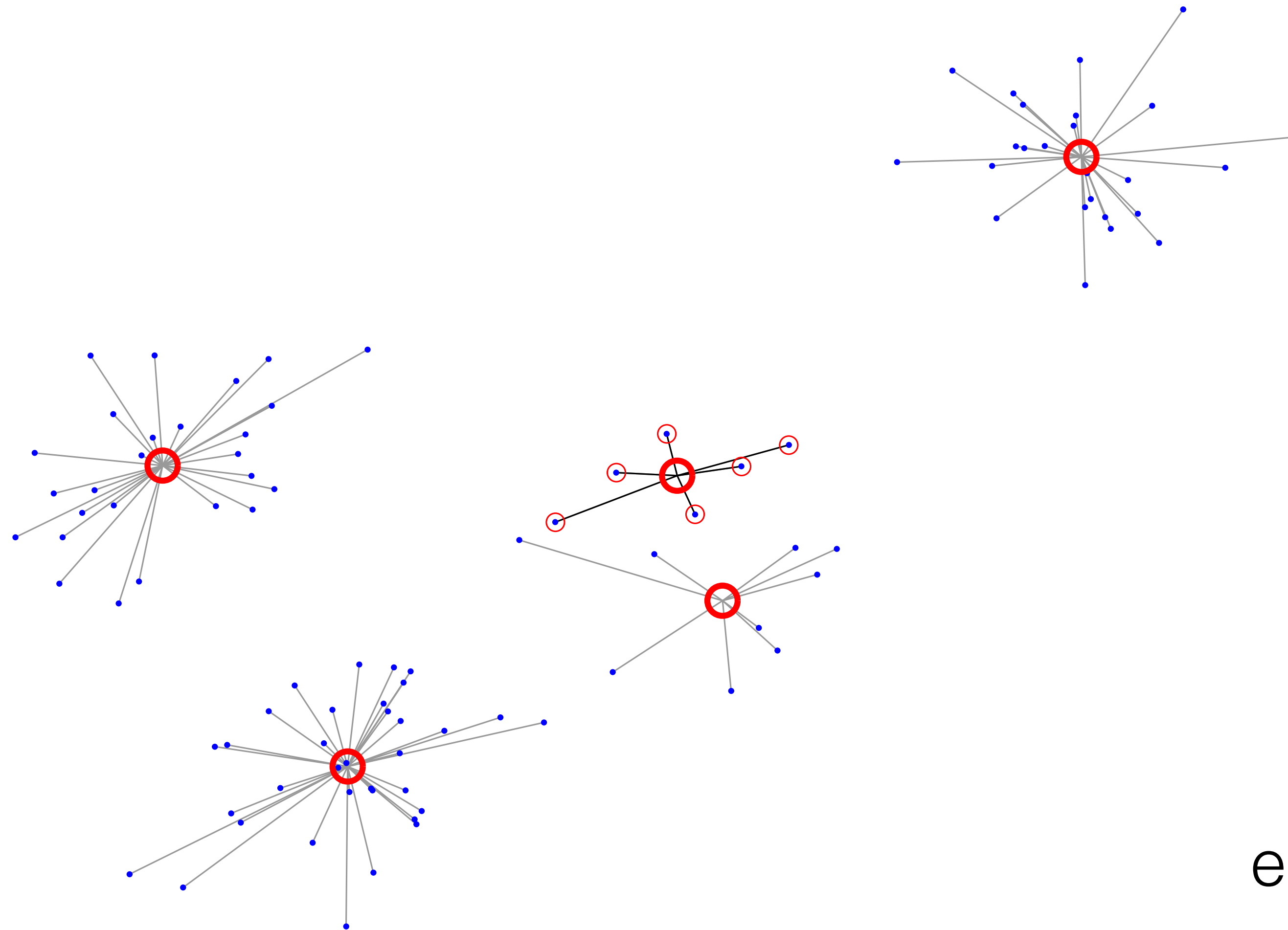
end of iteration 3

Example



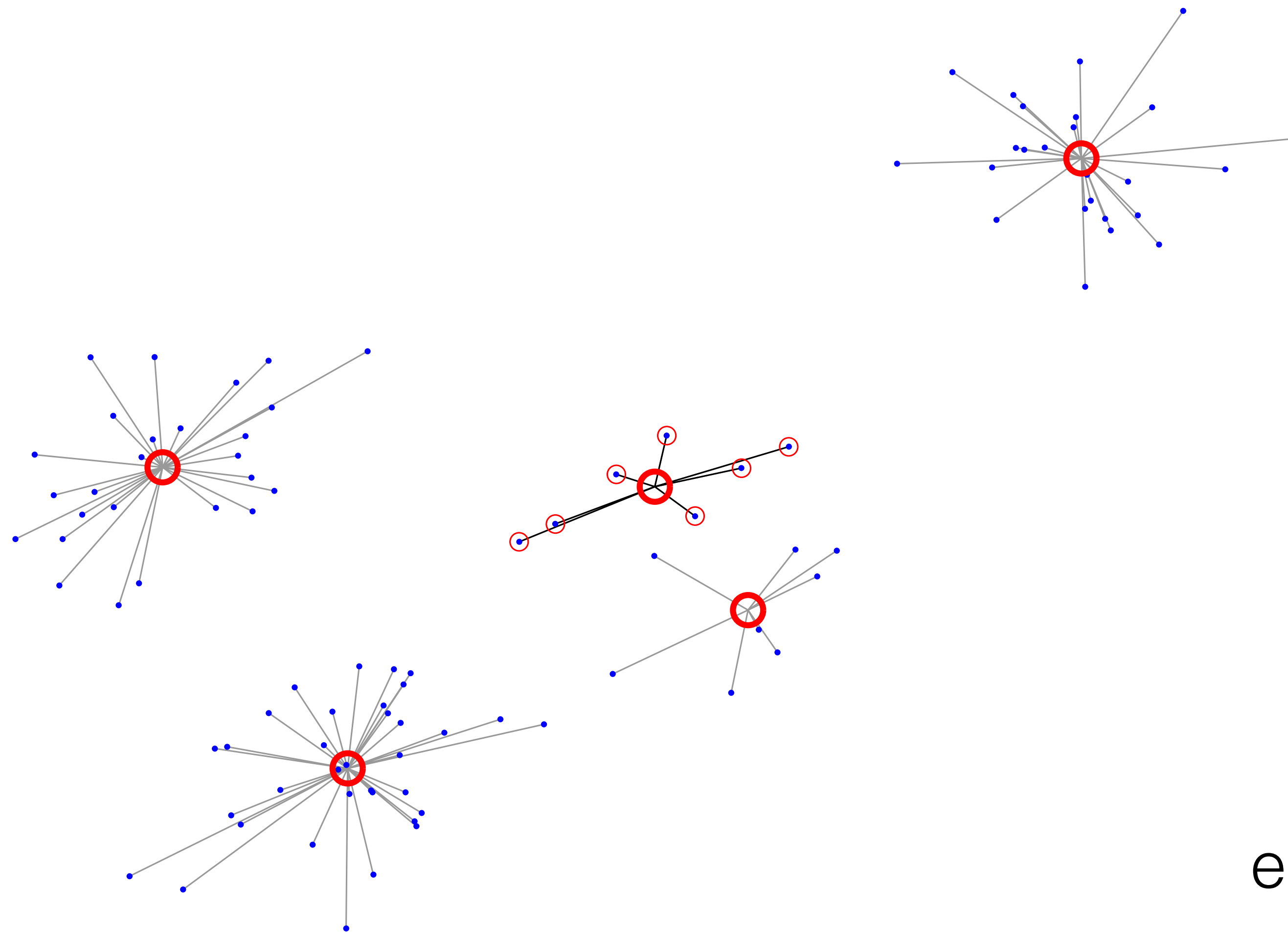
end of iteration 4

Example



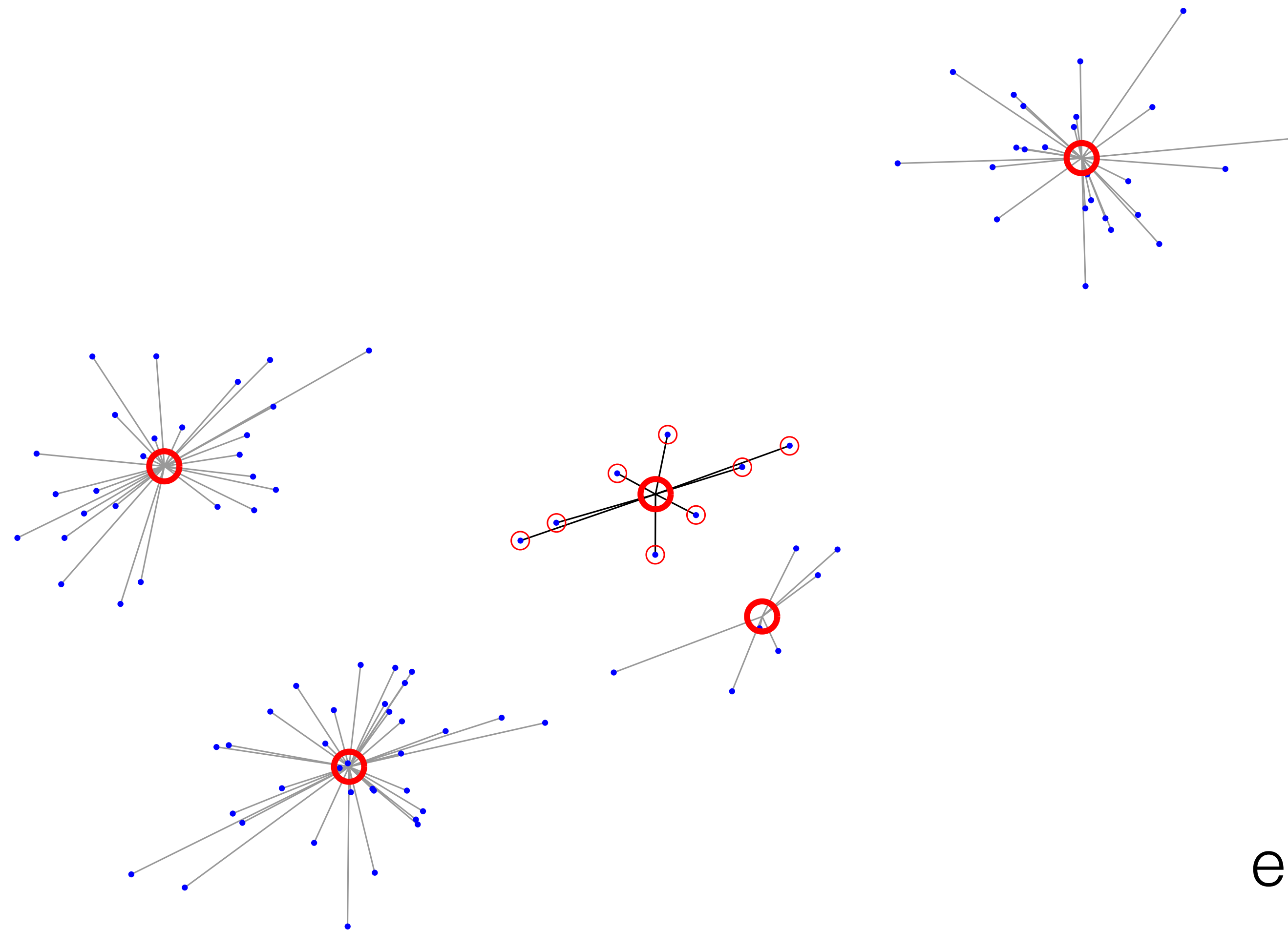
end of iteration 5

Example



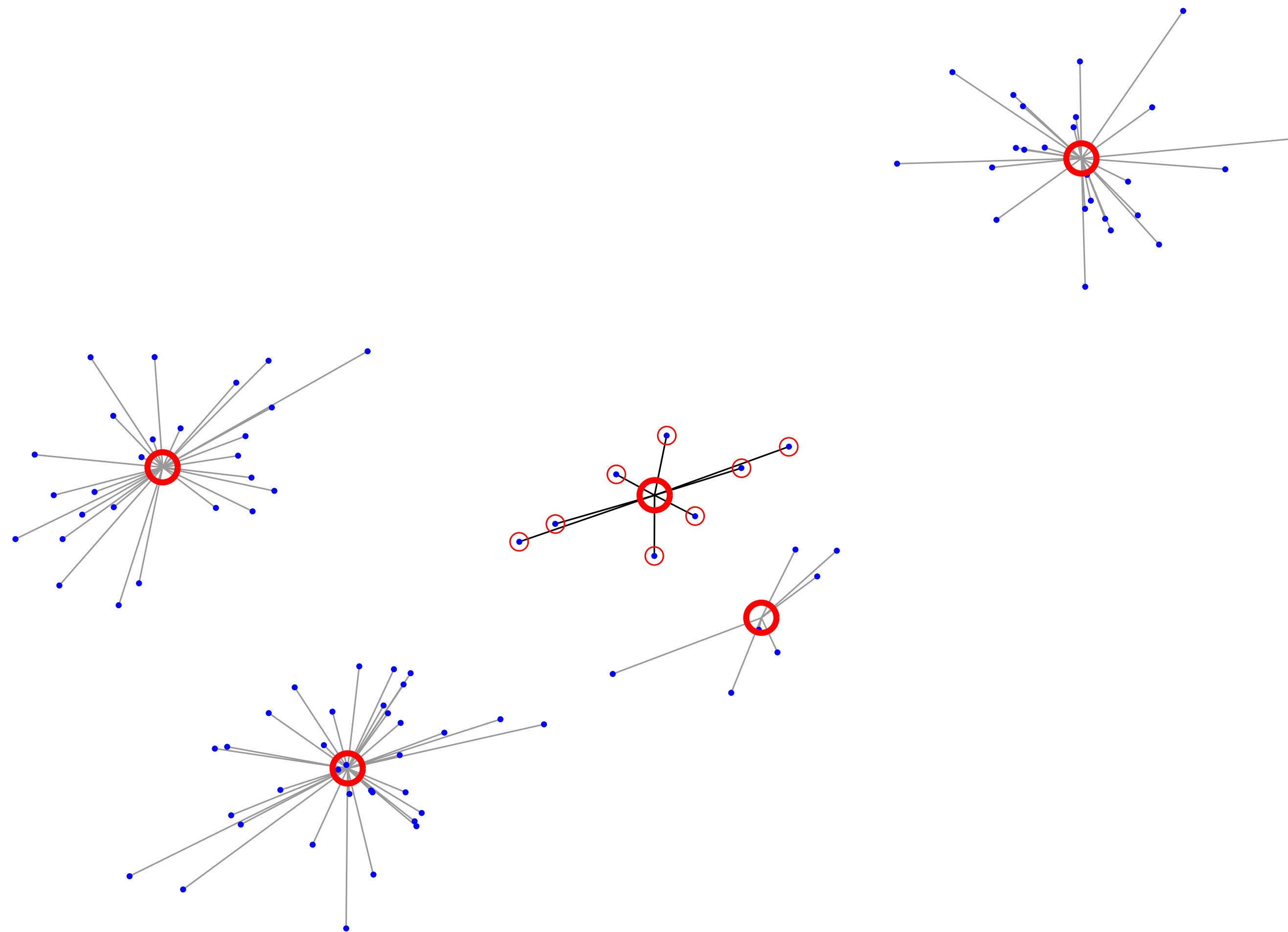
end of iteration 6

Example



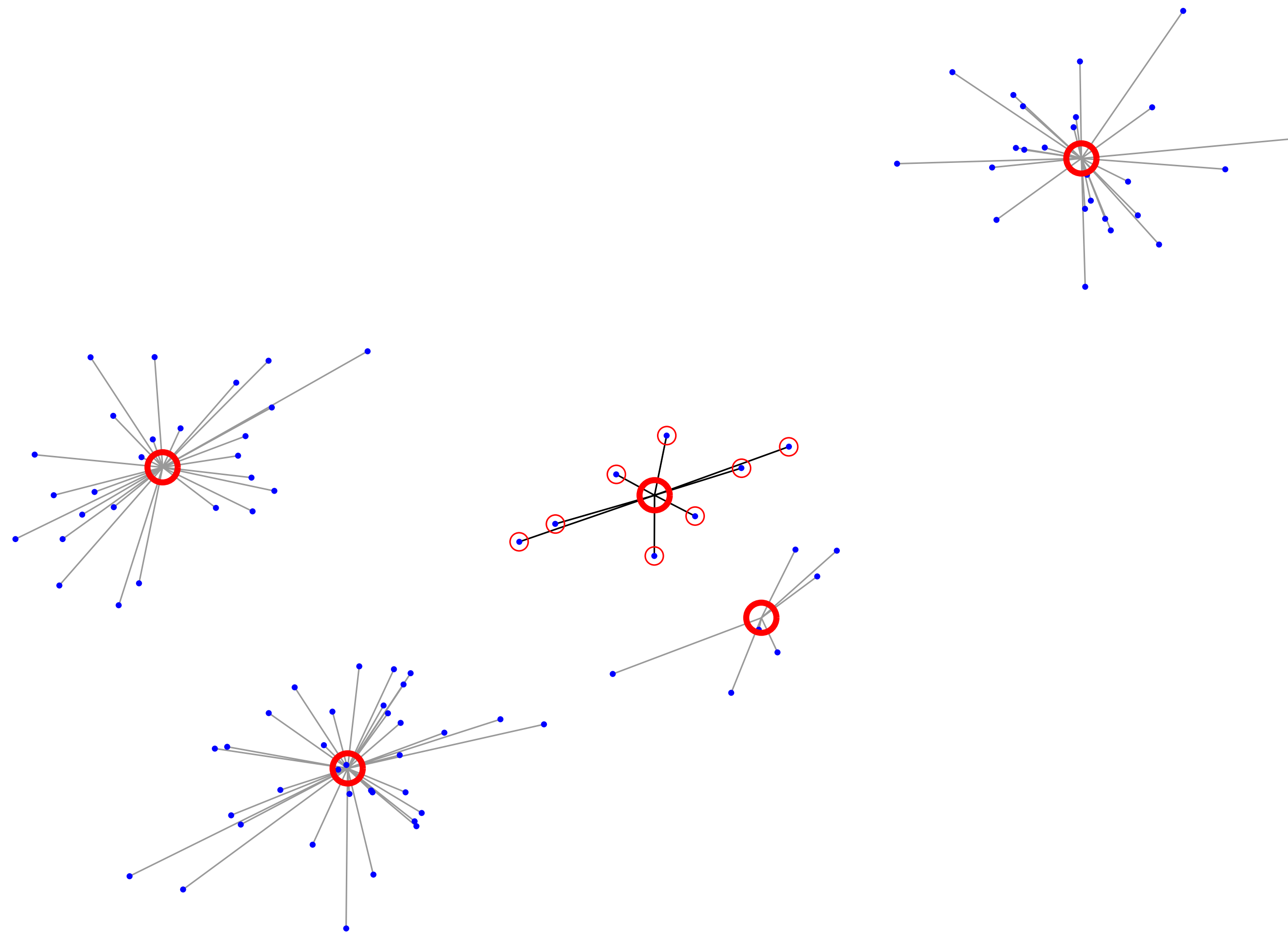
end of iteration 7

Example



converged at
iteration 8

Example



Summary

- Used Jensen's inequality to derive lower bound on log marginal likelihood
- Bound uses variational distribution q . We get to choose what family of q distributions to consider
- Using fully-factorized multinomial distributions for q gets EM
- Fully-factorized point distributions gets "hard"-EM, and using fixed, spherical covariance gets K-means