

# An introduction to graph analysis and modeling

## The Stochastic block model

MSc in Statistics for Smart Data – ENSAI

Autumn semester, 2018

<https://github.com/jchiquet/CourseStatNetwork>



# Motivations

Last time: find an underlying organization in a observed network

Spectral or hierachical clustering for network data

~> Not model-based, thus no statistical inference possible

Today: clustering of network based on a probabilistic model of the graph

Become familiar with

- the stochastic block model, a random graph model tailored for clustering vertices,
- the variational EM algorithm used to infer SBM from network data.

hierarchical clustering  $\leftrightarrow$  Gaussian mixture models



hierarchical/spectral clustering for network  $\leftrightarrow$  Stochastic block model

# Motivations

Last time: find an underlying organization in a observed network

Spectral or hierachical clustering for network data

~> Not model-based, thus no statistical inference possible

Today: clustering of network based on a probabilistic model of the graph

Become familiar with

- the stochastic block model, a random graph model tailored for clustering vertices,
- the variational EM algorithm used to infer SBM from network data.

hierarchical clustering  $\leftrightarrow$  Gaussian mixture models



hierarchical/spectral clustering for network  $\leftrightarrow$  Stochastic block model

# Motivations

Last time: find an underlying organization in a observed network

Spectral or hierachical clustering for network data

~> Not model-based, thus no statistical inference possible

Today: clustering of network based on a probabilistic model of the graph

Become familiar with

- the stochastic block model, a random graph model tailored for clustering vertices,
- the variational EM algorithm used to infer SBM from network data.

hierarchical clustering  $\leftrightarrow$  Gaussian mixture models



hierarchical/spectral clustering for network  $\leftrightarrow$  Stochastic block model

# Outline

## ① Background: mixture models and EM

- Mixture models

- Expectation-Maximization algorithm

- Example: mixture of Gaussians

## ② The Stochastic Block Model (SBM)

- Some Graphs Models and their limitations

- Mixture of Erdős-Rényi and the SBM

- Inference in SBM with variational EM

# Outline

- 1 Background: mixture models and EM
  - Mixture models
  - Expectation-Maximization algorithm
  - Example: mixture of Gaussians
- 2 The Stochastic Block Model (SBM)

# References



Pattern recognition and machine learning,  
Christopher Bishop

Chapter 9: Mixture Models and EM

<http://users.isr.ist.utl.pt/~wurmd/Livros/school/>



Models with Hidden Structure with Applications in Biology and  
Genomics,

Stéphane Robin

Master MathSV Course

[https:](https://www6.inra.fr/mia-paris/content/download/4587/42934/version/1/file/ModelsHiddenStruct-Biology.pdf)

[//www6.inra.fr/mia-paris/content/download/4587/42934/version/1/file/ModelsHiddenStruct-Biology.pdf](https://www6.inra.fr/mia-paris/content/download/4587/42934/version/1/file/ModelsHiddenStruct-Biology.pdf)



Classification non-supervisées,

É. Lebarbier, T. Mary-Huard

Chapitre 3 - méthode probabiliste: le modèle de mélange

<https://www.agroparistech.fr/IMG/pdf/ClassificationNonSupervisee-AgroParisTech.pdf>

# Outline

- ① Background: mixture models and EM
  - Mixture models**
  - Expectation-Maximization algorithm
  - Example: mixture of Gaussians
- ② The Stochastic Block Model (SBM)



# Latent variables models

## Definition

A **latent variable model** is a statistical model that relates, for  $i = 1, \dots, n$  individuals,

- a set of **manifest** (observed) variables  $\mathbf{X} = (X_i, i = 1, \dots, n)$  to
- a set of **latent** (unobserved) variables  $\mathbf{Z} = (Z_i, i = 1, \dots, n)$ .

Common assumption: conditional independence

$$\mathbb{P}((X_1, \dots, X_n) | (Z_1, \dots, Z_n)) = \prod_{i=1}^n \mathbb{P}(X_i | Z_i).$$

Famous examples

- $(Z_i, i \geq 1)$  is Markov chain: **Markov models**
- $Z_i$  categorical and independent: **mixture models**

# Latent variables models

## Definition

A **latent variable model** is a statistical model that relates, for  $i = 1, \dots, n$  individuals,

- a set of **manifest** (observed) variables  $\mathbf{X} = (X_i, i = 1, \dots, n)$  to
- a set of **latent** (unobserved) variables  $\mathbf{Z} = (Z_i, i = 1, \dots, n)$ .

Common assumption: conditional independence

$$\mathbb{P}((X_1, \dots, X_n) | (Z_1, \dots, Z_n)) = \prod_{i=1}^n \mathbb{P}(X_i | Z_i).$$

## Famous examples

- $(Z_i, i \geq 1)$  is Markov chain: **Markov models**
- $Z_i$  categorical and independent: **mixture models**
- what if  $X_i = X_{i,j}$  is a collection of edges in a graph?

# Latent variables models

## Definition

A **latent variable model** is a statistical model that relates, for  $i = 1, \dots, n$  individuals,

- a set of **manifest** (observed) variables  $\mathbf{X} = (X_i, i = 1, \dots, n)$  to
- a set of **latent** (unobserved) variables  $\mathbf{Z} = (Z_i, i = 1, \dots, n)$ .

Common assumption: conditional independence

$$\mathbb{P}((X_1, \dots, X_n) | (Z_1, \dots, Z_n)) = \prod_{i=1}^n \mathbb{P}(X_i | Z_i).$$

## Famous examples

- $(Z_i, i \geq 1)$  is Markov chain: **Markov models**
- $Z_i$  categorical and independent: **mixture models**
- what if  $X_i = X_{i'j'}$  is a collection of edges in a graph?

## Mixture models: the latent variables

When  $(Z_1, \dots, Z_n)$  are independent categorical variables, they give a **natural (latent) classification of the observations**  $(X_1, \dots, X_n)$  – or labels.

### Notations

Let  $(Z_1, \dots, Z_n)$  be *iid* categorical variables with distribution

$$\mathbb{P}(i \in q) = \mathbb{P}(Z_i = q) = \alpha_q, \quad \text{s.t.} \sum_{q=1}^Q \alpha_q = 1.$$

### Alternative (equivalent) notation

Let  $Z_i = (Z_{i1}, \dots, Z_{iQ})$  be an indicator vector of label for  $i$ :

$$\mathbb{P}(i \in q) = \mathbb{P}(Z_{iq} = 1) = \alpha_q, \quad \text{s.t.} \sum_{q=1}^Q \alpha_q = 1.$$

By definition,  $Z_i \sim \mathcal{M}(1, \alpha)$ , with  $\alpha = (\alpha_1, \dots, \alpha_Q)$ .

# Mixture models: the latent variables

When  $(Z_1, \dots, Z_n)$  are independent categorical variables, they give a **natural (latent) classification of the observations**  $(X_1, \dots, X_n)$  – or **labels**.

## Notations

Let  $(Z_1, \dots, Z_n)$  be *iid* categorical variables with distribution

$$\mathbb{P}(i \in q) = \mathbb{P}(Z_i = q) = \alpha_q, \quad \text{s.t.} \sum_{q=1}^Q \alpha_q = 1.$$

Alternative (equivalent) notation

Let  $Z_i = (Z_{i1}, \dots, Z_{iQ})$  be an indicator vector of label for  $i$ :

$$\mathbb{P}(i \in q) = \mathbb{P}(Z_{iq} = 1) = \alpha_q, \quad \text{s.t.} \sum_{q=1}^Q \alpha_q = 1.$$

By definition,  $Z_i \sim \mathcal{M}(1, \alpha)$ , with  $\alpha = (\alpha_1, \dots, \alpha_Q)$ .

# Mixture models: the latent variables

When  $(Z_1, \dots, Z_n)$  are independent categorical variables, they give a **natural (latent) classification of the observations**  $(X_1, \dots, X_n)$  – or **labels**.

## Notations

Let  $(Z_1, \dots, Z_n)$  be *iid* categorical variables with distribution

$$\mathbb{P}(i \in q) = \mathbb{P}(Z_i = q) = \alpha_q, \quad \text{s.t.} \sum_{q=1}^Q \alpha_q = 1.$$

## Alternative (equivalent) notation

Let  $Z_i = (Z_{i1}, \dots, Z_{iQ})$  be an indicator vector of label for  $i$ :

$$\mathbb{P}(i \in q) = \mathbb{P}(Z_{iq} = 1) = \alpha_q, \quad \text{s.t.} \sum_{q=1}^Q \alpha_q = 1.$$

By definition,  $Z_i \sim \mathcal{M}(1, \alpha)$ , with  $\alpha = (\alpha_1, \dots, \alpha_Q)$ .

# Mixture models: the manifest variables

A mixture model represents the **presence of subpopulations** within an overall population as follows:

$$\mathbb{P}(X_i) = \sum_{z_i \in \mathcal{Z}_i} \mathbb{P}(X_i, Z_i) = \sum_{Z_i \in \mathcal{Z}_i} \mathbb{P}(X_i | Z_i) \mathbb{P}(Z_i).$$

## Conditional distribution of the manifest variables

We assume a **parametric distribution** of  $X$  in each subpopulation

$$X_i | \{Z_i = q\} \sim \mathbb{P}_{\theta_q} \quad \left( \Leftrightarrow X_i | \{Z_{iq}\} = 1 \sim \mathbb{P}_{\theta_q} \right)$$

The specificity of each class is handled by  $\{\theta_q\}_{q=1}^Q$ .

# Mixture models: likelihoods

## The complete-data likelihood

It is the join distribution of  $(X_i, Z_i)$ :

$$\mathbb{P}(X_i, Z_i) = \alpha_{Z_i} \mathbb{P}_{\theta_q}(X_{Z_i})$$

## The incomplete-data likelihood

It is the marginal distribution of  $X_i$  once  $Z_i$  integrated:

$$\mathbb{P}(X_i) = \sum_{q=1}^Q \mathbb{P}(X_i, Z_i = q) = \sum_{q=1}^Q \alpha_q \mathbb{P}_{\theta_q}(X_i)$$

↪ A mixture model is a sum of distributions weighed by the proportion of each subpopulation.



# Mixture models: likelihoods

## The complete-data likelihood

It is the joint distribution of  $(X_i, Z_i)$ :

$$\mathbb{P}(X_i, Z_i) = \alpha_{Z_i} \mathbb{P}_{\boldsymbol{\theta}_q}(X_{Z_i})$$

## The incomplete-data likelihood

It is the marginal distribution of  $X_i$  once  $Z_i$  integrated:

$$\mathbb{P}(X_i) = \sum_{q=1}^Q \mathbb{P}(X_i, Z_i = q) = \sum_{q=1}^Q \alpha_q \mathbb{P}_{\boldsymbol{\theta}_q}(X_i)$$

↪ A **mixture model** is a sum of distributions weighed by the proportion of each subpopulation.

# Outline

- 1 Background: mixture models and EM
  - Mixture models
  - Expectation-Maximization algorithm**
  - Example: mixture of Gaussians
- 2 The Stochastic Block Model (SBM)

# Intractability of the Likelihood

## Maximum Likelihood Estimator

The MLE aims to maximize the (marginal) likelihood of the observations:

$$L(\boldsymbol{\theta}; \mathbf{X}) = \mathbb{P}_{\boldsymbol{\theta}}((X_1, \dots, X_n)) = \int_{\mathbf{Z} \in \mathcal{Z}} \mathbb{P}_{\boldsymbol{\theta}}(\mathbf{X}, \mathbf{Z}) d\mathbf{Z}$$

Integrations are summation over  $\{1, \dots, Q\}$ : we have  $Q^n$  terms !

## Intractable summation

With mixture models, for  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_Q)$  we have

$$\log L(\boldsymbol{\theta}; \mathbf{X}) = \sum_{i=1}^n \log \left\{ \sum_{q=1}^Q \alpha_q \mathbb{P}_{\boldsymbol{\theta}_q}(X_i) \right\}.$$

↪ Direct maximization of the likelihood is impossible in practice

# Intractability of the Likelihood

## Maximum Likelihood Estimator

The MLE aims to maximize the (marginal) likelihood of the observations:

$$L(\boldsymbol{\theta}; \mathbf{X}) = \mathbb{P}_{\boldsymbol{\theta}}((X_1, \dots, X_n)) = \int_{\mathbf{Z} \in \mathcal{Z}} \mathbb{P}_{\boldsymbol{\theta}}(\mathbf{X}, \mathbf{Z}) d\mathbf{Z}$$

Integrations are summation over  $\{1, \dots, Q\}$ : we have  $Q^n$  terms !

## Intractable summation

With mixture models, for  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_Q)$  we have

$$\log L(\boldsymbol{\theta}; \mathbf{X}) = \sum_{i=1}^n \log \left\{ \sum_{q=1}^Q \alpha_q \mathbb{P}_{\boldsymbol{\theta}_q}(X_i) \right\}.$$

⇒ Direct maximization of the likelihood is impossible in practice

# Bayes decision rule / Maximum *a posteriori*

## Principle

Affect an individual  $i$  to the subpopulation which is the most likely according to the data:

$$\tau_{iq} = \mathbb{P}(Z_{iq} = 1 | X_i = x_i)$$

This is the **posterior probability** for  $i \in q$ .

## Application of the Bayes Theorem

It is straightforward to show that

$$\tau_{iq} = \frac{\alpha_q \mathbb{P}_{\theta_q}(x_i)}{\sum_{q=1}^Q \alpha_q \mathbb{P}_{\theta_q}(x_i)}$$

# Principle of the EM algorithm

If  $\theta$  were known

...estimating the **posterior probability**  $\mathbb{P}(Z_i|\mathbf{X})$  of  $\mathbf{Z}$  should be easy

*By means of the Bayes decision rule*

If  $\mathbf{Z}$  were known...

...estimating the **best set of parameter**  $\theta$  should be easy

*This is close to usual maximum likelihood estimation*

EM principle

Maximize the marginal likelihood iteratively:

- ① Initialize  $\theta$
- ② Compute the probability of  $\mathbf{Z}$  given  $\theta$
- ③ Get a better  $\theta$  with the new  $\mathbf{Z}$
- ④ Iterate until convergence

# Principle of the EM algorithm

If  $\theta$  were known

...estimating the **posterior probability**  $\mathbb{P}(Z_i|\mathbf{X})$  of  $\mathbf{Z}$  should be easy

*By means of the Bayes decision rule*

If  $\mathbf{Z}$  were known...

...estimating the **best set of parameter**  $\theta$  should be easy

*This is close to usual maximum likelihood estimation*

## EM principle

Maximize the marginal likelihood iteratively:

- 1 Initialize  $\theta$
- 2 Compute the probability of  $\mathbf{Z}$  given  $\theta$
- 3 Get a better  $\theta$  with the new  $\mathbf{Z}$
- 4 Iterate until convergence

# Formal algorithm

**Initialization:** start from a good guess either of  $\mathbf{Z}$  or  $\boldsymbol{\theta}$ , then iterate 1-2

## 1. Expectation step

Calculate the expected value of the loglikelihood under the current  $\boldsymbol{\theta}$

$$Q\left(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}\right)=\mathbb{E}_{\mathbf{Z}|\mathbf{X};\boldsymbol{\theta}^{(t)}}\left[\log L(\boldsymbol{\theta};\mathbf{X},\mathbf{Z})\right] \quad (\text{needs } \mathbb{P}_{\boldsymbol{\theta}^{(t)}}(\mathbf{Z}|\mathbf{X}))$$

## 2. Maximization step

Find the parameters that maximize this quantity

$$\boldsymbol{\theta}^{(t+1)}=\arg \max _{\boldsymbol{\theta}} Q\left(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}\right)$$

Stop when  $\left\|\boldsymbol{\theta}^{(t+1)}-\boldsymbol{\theta}^{(t)}\right\|<\varepsilon$  or  $\left\|Q^{(t+1)}-Q^{(t)}\right\|<\varepsilon$



# (Basic) Convergence analysis

## Theorem

*At each step of the EM algorithm, the loglikelihood increases. EM thus reaches a local optimum.*

## Proof.

On board.



# Choosing the number of component

## Reminder: Bayesian Information Criterion

The BIC is a model selection criterion which penalizes the adjustment to the data by the number of parameter in model  $\mathcal{M}$  as follows:

$$\text{BIC}(\mathcal{M}) = \log L(\hat{\boldsymbol{\theta}}; \mathbf{X}) - \frac{1}{2} \log(n) \text{df}(\mathcal{M}).$$

## Integrated Classification Criterion

It is an adaptation working with the complete-data likelihood:

$$\begin{aligned} \text{ICL}(\mathcal{M}) &= \log L(\hat{\boldsymbol{\theta}}; \mathbf{X}, \hat{\mathbf{Z}}) + \frac{1}{2} \log(n) \text{df}(\mathcal{M}) \\ &= \text{BIC} - \mathcal{H}(\mathbb{P}(\hat{\mathbf{Z}}|\mathbf{X})), \end{aligned}$$

where the entropy  $\mathcal{H}$  measures the separability of the subpopulations.

⇒ We choose  $\mathcal{M}(Q)$  that maximizes either BIC or ICL

# Choosing the number of component

## Reminder: Bayesian Information Criterion

The BIC is a model selection criterion which penalizes the adjustment to the data by the number of parameter in model  $\mathcal{M}$  as follows:

$$\text{BIC}(\mathcal{M}) = \log L(\hat{\boldsymbol{\theta}}; \mathbf{X}) - \frac{1}{2} \log(n) \text{df}(\mathcal{M}).$$

## Integrated Classification Criterion

It is an adaptation working with the complete-data likelihood:

$$\begin{aligned} \text{ICL}(\mathcal{M}) &= \log L(\hat{\boldsymbol{\theta}}; \mathbf{X}, \hat{\mathbf{Z}}) + \frac{1}{2} \log(n) \text{df}(\mathcal{M}) \\ &= \text{BIC} - \mathcal{H}(\mathbb{P}(\hat{\mathbf{Z}}|\mathbf{X})), \end{aligned}$$

where the entropy  $\mathcal{H}$  measures the separability of the subpopulations.

⇒ We choose  $\mathcal{M}(Q)$  that maximizes either BIC or ICL

# Choosing the number of component

## Reminder: Bayesian Information Criterion

The BIC is a model selection criterion which penalizes the adjustment to the data by the number of parameter in model  $\mathcal{M}$  as follows:

$$\text{BIC}(\mathcal{M}) = \log L(\hat{\boldsymbol{\theta}}; \mathbf{X}) - \frac{1}{2} \log(n) \text{df}(\mathcal{M}).$$

## Integrated Classification Criterion

It is an adaptation working with the complete-data likelihood:

$$\begin{aligned} \text{ICL}(\mathcal{M}) &= \log L(\hat{\boldsymbol{\theta}}; \mathbf{X}, \hat{\mathbf{Z}}) + \frac{1}{2} \log(n) \text{df}(\mathcal{M}) \\ &= \text{BIC} - \mathcal{H}(\mathbb{P}(\hat{\mathbf{Z}}|\mathbf{X})), \end{aligned}$$

where the entropy  $\mathcal{H}$  measures the separability of the subpopulations.

⇒ We choose  $\mathcal{M}(Q)$  that maximizes either BIC or ICL

# Outline

- ① Background: mixture models and EM
  - Mixture models
  - Expectation-Maximization algorithm
  - Example: mixture of Gaussians
- ② The Stochastic Block Model (SBM)

# Mixture of Gaussians

Calculs in the univariate case: complete likelihood

The distribution of  $X_i$  conditional on the label of  $i$  is assumed to be a univariate Gaussian distribution with unknown parameters:

$$X_i | Z_{iq} = 1 \sim \mathcal{N}(\mu_q, \sigma_q^2)$$

complete Likelihood ( $\mathbf{X}, \mathbf{Z}$ )

The model complete loglikelihood is

$$\log L(\boldsymbol{\mu}, \boldsymbol{\sigma}^2; \mathbf{X}, \mathbf{Z}) = \sum_{i=1}^n \sum_{q=1}^Q Z_{iq} \left( \log \alpha_q - \log \sigma_q - \log(\sqrt{2\pi}) - \frac{1}{2\sigma_q^2} (x_i - \mu_q)^2 \right)$$

# Mixture of Gaussians

Calculs in the univariate case: E-step

## E-step

For fixed values of  $\mu_q, \sigma_q^2$  and  $\alpha_q$ , the estimates of the posterior probabilities  $\hat{\tau}_{iq} = \mathbb{P}(Z_{iq} = 1 | X_i)$  are

$$\hat{\tau}_{iq} = \frac{\alpha_q \mathcal{N}(x_i; \mu_q, \sigma_q^2)}{\sum_{q=1}^Q \alpha_q \mathcal{N}(x_i; \mu_q, \sigma_q^2)},$$

where  $\mathcal{N}$  is the density of the normal distribution.

# Mixture of Gaussians

Calculs in the univariate case: M-step

## M-step

For fixed values of  $\tau_{iq}$ , the estimates of the model parameters are

$$\hat{\alpha}_q = \frac{\sum_{i=1}^n \tau_{iq}}{\sum_{i=1}^n \sum_{q=1}^Q \tau_{iq}} \quad \hat{\mu}_q = \frac{\sum_i \tau_{iq} x_i}{\sum_i \tau_{iq}} \quad \hat{\sigma}_q^2 = \frac{\sum_{i=1}^n \tau_{iq} (x_i - \mu_q)^2}{\sum_{i=1}^n \tau_{iq}}$$



## R code: auxiliary functions

We start by defining functions to compute the complete model loglikelihood, perform the E step and the M step.

```
get.cloglik <- function(X, Z, theta) {  
  alpha <- theta$alpha; mu <- theta$mu; sigma <- theta$sigma  
  xs <- scale(matrix(X,length(x),length(alpha)),mu,sigma)  
  return(sum(Z*(log(alpha)-log(sigma)-.5*(log(2*pi)+xs^2))))  
}  
  
M.step <- function(X, tau) {  
  n <- length(X); Q <- ncol(tau)  
  alpha <- colMeans(tau)  
  mu <- colMeans(tau * matrix(X,n,Q)) / alpha  
  sigma <- sqrt(colMeans(tau*sweep(matrix(X,n,Q),2,mu,"-")^2)/alpha)  
  return(list(alpha=alpha, mu=mu, sigma=sigma))  
}  
  
E.step <- function(X, theta) {  
  tau <- mapply(function(alpha, mu, sigma) {  
    alpha*dnorm(X,mu,sigma)  
  }, theta$alpha, theta$mu, theta$sigma)  
  return(tau / rowSums(tau))  
}
```

## R code: EM for univariate mixture

```
EM.mixture <- function(X, Q,
                       init.cl=sample(1:Q,n,rep=TRUE), max.iter=100, eps=1e-5) {
  n <- length(X); tau <- matrix(0,n,Q); tau[cbind(1:n,init.cl)] <- 1
  Eloglik <- vector("numeric", max.iter)
  iter <- 0; cond <- FALSE

  while (!cond) {
    iter <- iter + 1
    ## M step
    theta <- M.step(X, tau)
    ## E step
    tau <- E.step(X, theta)
    ## check consistency
    Eloglik[iter] <- get.cloglik(X, tau, theta)
    if (iter > 1)
      cond <- (iter>=max.iter) | Eloglik[iter]-Eloglik[iter-1] < eps
  }

  return(list(alpha = theta$alpha, mu = theta$mu, sigma = theta$sigma,
             tau = tau, cl = apply(tau, 1, which.max),
             Eloglik = Eloglik[1:iter]))
}
```

## Example: data generation

We first generate data with 4 components:

```
mu1 <- 5    ; sigma1 <- 1; n1 <- 100
mu2 <- 10   ; sigma2 <- 1; n2 <- 200
mu3 <- 15   ; sigma3 <- 2; n3 <- 50
mu4 <- 20   ; sigma4 <- 3; n4 <- 100
cl <- rep(1:4,c(n1,n2,n3,n4))
x <- c(rnorm(n1,mu1,sigma1),rnorm(n2,mu2,sigma2),
      rnorm(n3,mu3,sigma3),rnorm(n4,mu4,sigma4))
n <- length(x)

## we randomize the class ordering
rnd <- sample(1:n)
cl <- cl[rnd]
x <- x[rnd]

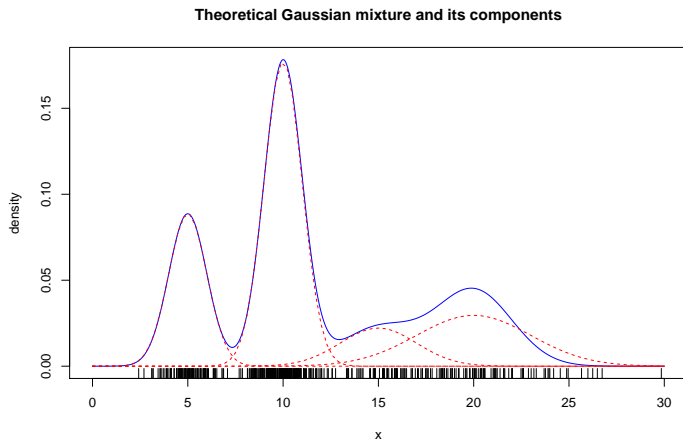
alpha <- c(n1,n2,n3,n4)/n
```

## Example: data generation - plot I

Let us plot the data and the theoretical mixture.

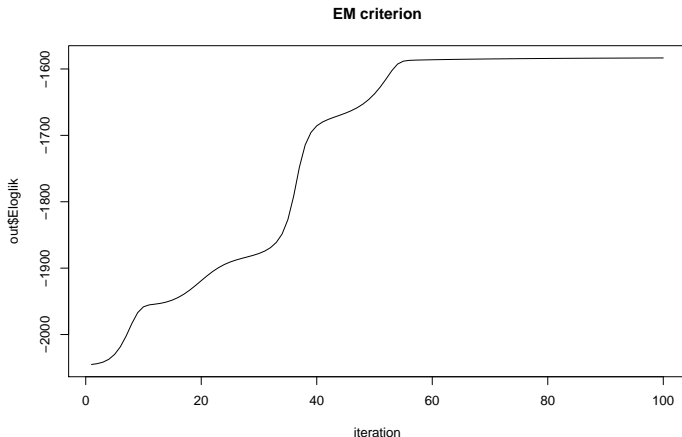
```
curve(alpha[1]*dnorm(x,mu1,sigma1) +  
      alpha[2]*dnorm(x,mu2,sigma2) +  
      alpha[3]*dnorm(x,mu3,sigma3) +  
      alpha[4]*dnorm(x,mu4,sigma3),  
      col="blue", lty=1, from=0,to=30, n=1000,  
      main="Theoretical Gaussian mixture and its components",  
      xlab="x", ylab="density")  
curve(alpha[1]*dnorm(x,mu1,sigma1), col="red", add=TRUE, lty=2)  
curve(alpha[2]*dnorm(x,mu2,sigma2), col="red", add=TRUE, lty=2)  
curve(alpha[3]*dnorm(x,mu3,sigma3), col="red", add=TRUE, lty=2)  
curve(alpha[4]*dnorm(x,mu4,sigma4), col="red", add=TRUE, lty=2)  
rug(x)
```

## Example: data generation - plot II



## Example: adjustment

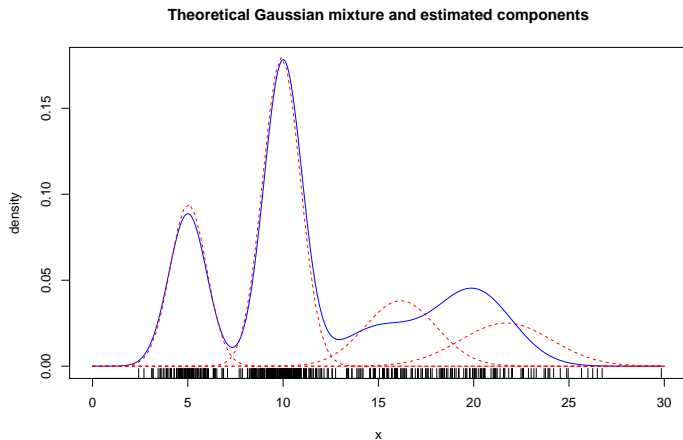
```
out <- EM.mixture(x, Q=4, init.cl=sample(1:4,n,rep=TRUE))  
plot(out$Eloglik, main="EM criterion", type="l", xlab="iteration")
```



## Example: adjustment - plot I

```
out <- EM.mixture(x, Q=4, init.cl=kmeans(x,4)$cl)
curve(alpha[1]*dnorm(x,mu1,sigma1) +
      alpha[2]*dnorm(x,mu2,sigma2) +
      alpha[3]*dnorm(x,mu3,sigma3) +
      alpha[4]*dnorm(x,mu4,sigma3), col="blue",
      lty=1, from=0,to=30, n=1000,
      main="Theoretical Gaussian mixture and estimated components",
      xlab="x", ylab="density")
curve(out$alpha[1]*dnorm(x,out$mu[1],out$sigma[1]), col="red", add=TRUE, lty=2)
curve(out$alpha[2]*dnorm(x,out$mu[2],out$sigma[2]), col="red", add=TRUE, lty=2)
curve(out$alpha[3]*dnorm(x,out$mu[3],out$sigma[3]), col="red", add=TRUE, lty=2)
curve(out$alpha[4]*dnorm(x,out$mu[4],out$sigma[4]), col="red", add=TRUE, lty=2)
rug(x)
```

# Example: adjustment - plot II





## Example: adjustment - classification I

```
table(cl, out$cl)
```

```
##
```

```
## cl      1      2      3      4
```

```
##  1      0      0 100      0
```

```
##  2 197      1      2      0
```

```
##  3      2  45      0      3
```

```
##  4      0  36      0  64
```

```
aricode::ARI(cl, out$cl)
```

```
## [1] 0.8686458
```

# Outline

① Background: mixture models and EM

② The Stochastic Block Model (SBM)

Some Graphs Models and their limitations

Mixture of Erdős-Rényi and the SBM

Inference in SBM with variational EM

# References



Statistical Analysis of Network Data: Methods and Models

Eric Kolaczyk

Chapters 5 and 6



Mixture model for random graphs, Statistics and Computing

Daudin, Robin, Picard

[pbil.univ-lyon1.fr/members/fpicard/franckpicard\\_fichiers/pdf/DPR08.pdf](http://pbil.univ-lyon1.fr/members/fpicard/franckpicard_fichiers/pdf/DPR08.pdf)



Analyse statistique de graphes,

Catherine Matias

Chapitre 4, Section 4

# Outline

① Background: mixture models and EM

② The Stochastic Block Model (SBM)

Some Graphs Models and their limitations

Mixture of Erdős-Rényi and the SBM

Inference in SBM with variational EM

# A mathematical model: Erdős-Rényi graph

## Definition

Let  $\mathcal{V} = 1, \dots, n$  be a set of fixed vertices. The (simple) Erdős-Rényi model  $\mathcal{G}(n, \pi)$  assumes random edges between pairs of nodes with probability  $\pi$ . In other word, the (random) adjacency matrix  $\mathbf{X}$  is such that

$$X_{ij} \sim \mathcal{B}(\pi)$$

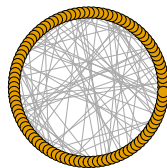
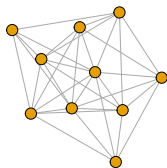
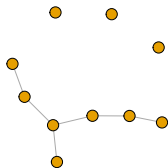
## Proposition (degree distribution)

*The (random) degree  $D_i$  of vertex  $i$  follows a binomial distribution:*

$$D_i \sim b(n - 1, \pi).$$

# Erdős-Rényi - example

```
G1 <- igraph::sample_gnp(10, 0.1)
G2 <- igraph::sample_gnp(10, 0.9)
G3 <- igraph::sample_gnp(100, .02)
par(mfrow=c(1,3))
plot(G1, vertex.label=NA) ; plot(G2, vertex.label=NA)
plot(G3, vertex.label=NA, layout=layout.circle)
```



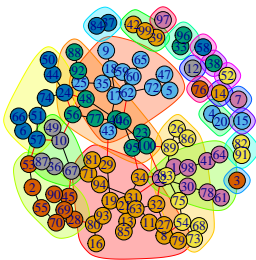
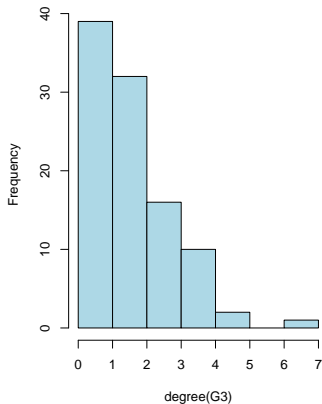
# Erdős-Rényy - limitations: very homogeneous

```
average.path.length(G3); diameter(G3)
```

```
## [1] 5.488342
```

```
## [1] 12
```

Histogram of degree(G3)



# Mechanism-based model: preferential attachment

The graph is defined dynamically as follows

## Definition

Start from a initial graph  $\mathcal{G}_0 = (\mathcal{V}_0, \mathcal{E}_0)$ , then for each time step,

- ① At  $t$  a new node  $V_t$  is added
- ②  $V_t$  is connected to  $i \in V_{t-1}$  with probability

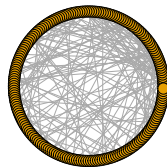
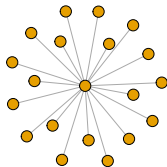
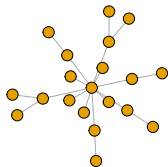
$$D_i^\alpha + \text{cst.}$$

$\rightsquigarrow$  Nodes with high degree get more connections thus **richers get richers**



# Preferential attachment - example

```
G1 <- igraph::sample_pa(20, 1, directed=FALSE)
G2 <- igraph::sample_pa(20, 5, directed=FALSE)
G3 <- igraph::sample_pa(200, directed=FALSE)
par(mfrow=c(1,3))
plot(G1, vertex.label=NA) ; plot(G2, vertex.label=NA)
plot(G3, vertex.label=NA, layout=layout.circle)
```



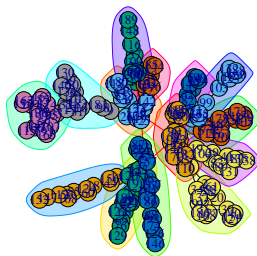
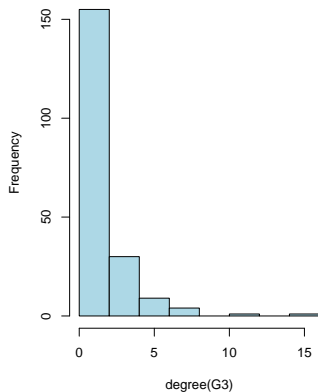
# Preferential attachment - limitations

```
average.path.length(G3); diameter(G3)
```

```
## [1] 6.098844
```

```
## [1] 14
```

Histogram of degree(G3)



# Limitations

- Erdős-Rényi

The ER model does not fit well real world network

- As can be seen from its degree distribution
- ER is generally too homogeneous

- Preferential attachment

- Is defined through an algorithm so performing statistics is complicated
- Is stucked to the power-law distribution of degrees

## The Stochastic Block Model

The SBM<sup>1</sup> generalizes ER in a mixture framework. It provides

- a statistical framework to adjust and interpret the parameters
- a flexible yet simple specification that fits many existing network data

---

<sup>1</sup>Other models exist (e.g. exponential model for random graphs) but less popular.

# Outline

① Background: mixture models and EM

② The Stochastic Block Model (SBM)

Some Graphs Models and their limitations

Mixture of Erdős-Rényi and the SBM

Inference in SBM with variational EM

# Stochastic Block Model: definition

Mixture model point of view: mixture of Erdős-Rényi

## Latent structure

Let  $\mathcal{V} = \{1, \dots, n\}$  be a fixed set of vertices. We give each  $i \in \mathcal{V}$  a **latent label** among a set  $\mathcal{Q} = \{1, \dots, Q\}$  such that

- $\alpha_q = \mathbb{P}(i \in q), \quad \sum_q \alpha_q = 1;$
- $Z_{iq} = \mathbf{1}_{\{i \in q\}}$  are independent hidden variables.

## The conditional distribution of the edges

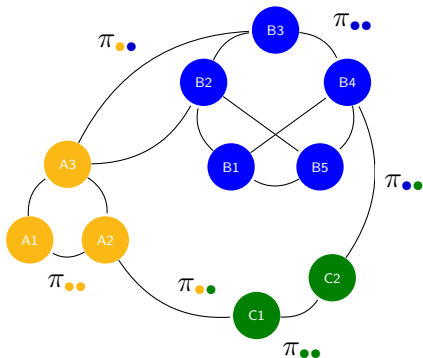
Connexion probabilities depend on the node class belonging:

$$X_{ij} | \{i \in q, j \in \ell\} \sim \mathcal{B}(\pi_{q\ell}) \quad \left( \Leftrightarrow X_{ij} | \{Z_{iq}Z_{j\ell} = 1\} \sim \mathcal{B}(\pi_{q\ell}). \right)$$

The  $Q \times Q$  matrix  $\pi$  gives for all couple of labels

$$\pi_{q\ell} = \mathbb{P}(X_{ij} = 1 | i \in q, j \in \ell).$$

# Stochastic Block Model: the big picture



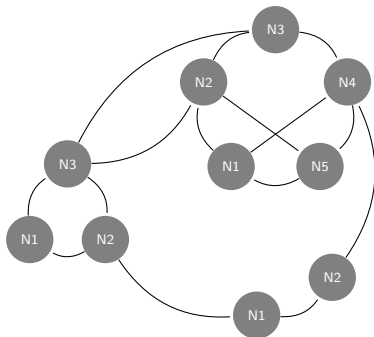
## Stochastic Block Model

Let  $n$  nodes divided into

- $\mathcal{Q} = \{\bullet, \bullet, \bullet\}$  classes
- $\alpha_{\bullet} = \mathbb{P}(i \in \bullet), \bullet \in \mathcal{Q}, i = 1, \dots, n$
- $\pi_{\bullet\bullet} = \mathbb{P}(i \leftrightarrow j | i \in \bullet, j \in \bullet)$

$$Z_i = \mathbf{1}_{\{i \in \bullet\}} \sim^{\text{iid}} \mathcal{M}(1, \alpha), \quad \forall \bullet \in \mathcal{Q},$$
$$X_{ij} | \{i \in \bullet, j \in \bullet\} \sim^{\text{ind}} \mathcal{B}(\pi_{\bullet\bullet})$$

# Stochastic Block Model: unknown parameters



## Stochastic Block Model

Let  $n$  nodes divided into

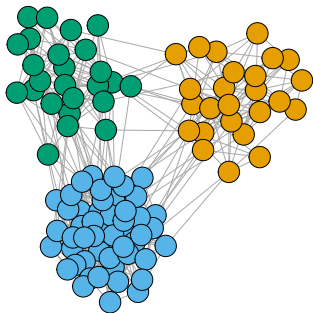
- $\mathcal{Q} = \{\bullet, \bullet, \bullet\}$ ,  $\text{card}(\mathcal{Q})$  known
- $\alpha_{\bullet} = ?$ ,
- $\pi_{\bullet\bullet} = ?$

$$Z_i = \mathbf{1}_{\{i \in \bullet\}} \sim^{\text{iid}} \mathcal{M}(1, \alpha), \quad \forall \bullet \in \mathcal{Q},$$
$$X_{ij} \mid \{i \in \bullet, j \in \bullet\} \sim^{\text{ind}} \mathcal{B}(\pi_{\bullet\bullet})$$

# Stochastic block models – examples of topology

## Community network

```
pi <- matrix(c(0.3,0.02,0.02,0.02,0.3,0.02,0.02,0.02,0.3),3,3)
communities <- igraph::sample_sbm(100, pi, c(25, 50, 25))
plot(communities, vertex.label=NA, vertex.color = rep(1:3,c(25, 50, 25)))
```

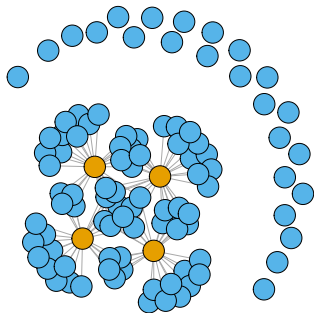




# Stochastic block models – examples of topology

## Star network

```
pi <- matrix(c(0.05,0.3,0.3,0),2,2)
star <- igraph::sample_sbm(100, pi, c(4, 96))
plot(star, vertex.label=NA, vertex.color = rep(1:2,c(4,96)))
```



# Degree distributions

## Conditional degree distribution

The conditional degree distribution of a node  $i \in q$  is

$$D_i | i \in q \sim \text{b}(n-1, \bar{\pi}) \approx \mathcal{P}(\lambda_q), \quad \bar{\pi}_q = \sum_{\ell=1}^Q \alpha_\ell \pi_{q\ell} \quad \lambda_q = (n-1) \bar{\pi}_q$$

## Conditional degree distribution

The degree distribution of a node  $i$  can be approximated by a mixture of Poisson distributions:

$$\mathbb{P}(D_i = k) = \sum_{q=1}^Q \alpha_q \exp\{-\lambda_q\} \frac{\lambda_q^k}{k!}$$

# Likelihoods

## Complete-data loglikelihood

$$\log L(\mathbf{X}, \mathbf{Z}) = \sum_{i,q} Z_{iq} \log \alpha_q + \sum_{i < j, q, \ell} Z_{iq} Z_{j\ell} \log \pi_{q\ell}^{X_{ij}} (1 - \pi_{q\ell})^{1-X_{ij}}.$$

## Conditional expectation of the complete-data loglikelihood

$$\mathbb{E}_{\mathbf{Z}|\mathbf{X}}[\log L(\boldsymbol{\theta}; \mathbf{X}, \mathbf{Z})] = \sum_{i,q} \tau_{iq} \log \alpha_q + \sum_{i < j, q, \ell} \eta_{ijq\ell} \log \pi_{q\ell}^{X_{ij}} (1 - \pi_{q\ell})^{1-X_{ij}}$$

where  $\tau_{iq}, \eta_{ijq\ell}$  are the posterior probabilities:

- $\tau_{iq} = \mathbb{P}(Z_{iq} = 1 | \mathbf{X}) = \mathbb{E}[Z_{iq} | \mathbf{X}]$ .
- $\eta_{ijq\ell} = \mathbb{P}(Z_{iq} Z_{j\ell} = 1 | \mathbf{X}) = \mathbb{E}[Z_{iq} Z_{j\ell} | \mathbf{X}]$ .

# Outline

① Background: mixture models and EM

② The Stochastic Block Model (SBM)

Some Graphs Models and their limitations

Mixture of Erdős-Rényi and the SBM

Inference in SBM with variational EM

# The EM strategy does not apply directly for SBM

Ouch: another intractability problem

- the  $Z_{iq}$  are **not independent** in the SBM framework. . .
- we cannot compute  $\eta_{ijql} = \mathbb{P}(Z_{iq}Z_{jl} = 1|\mathbf{X}) = \mathbb{E}[Z_{iq}Z_{jl}|\mathbf{X}]$ ,
- the conditional expectation  $Q(\boldsymbol{\theta})$ , i.e. the main EM ingredient, is **intractable**.

Solution: mean field approximation

Approximate  $\eta_{ijql}$  by  $\tau_{iq}\tau_{jl}$ , i.e., **assume independence between  $Z_{iq}$**

$\rightsquigarrow$  This can be formalized in the variational framework

# Revisiting the EM algorithm I

## Proposition

*Consider a distribution  $\mathbb{Q}$  for the  $\{Z_{iq}\}$ . We have*

$$\log L(\boldsymbol{\theta}; \mathbf{X}) = \mathbb{E}_{\mathbb{Q}}[\log L(\boldsymbol{\theta}, \mathbf{X}, \mathbf{Z})] + \mathcal{H}(\mathbb{Q}) + \text{KL}(\mathbb{Q} \mid \mathbb{P}(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta})),$$

*where  $\mathcal{H}$  is the entropy and  $\text{KL}(\cdot|\cdot)$  is the Kullback-Leibler divergence:*

$$\mathcal{H}(\mathbb{Q}) = - \sum_z \mathbb{Q}(z) \log \mathbb{Q}(z) = -\mathbb{E}_{\mathbb{Q}}[\log \mathbb{Q}(Z)]$$

$$\mathcal{KL}(\mathbb{Q} \mid \mathbb{P}(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta})) = \sum_z \mathbb{Q}(z) \log \frac{\mathbb{Q}(z)}{\mathbb{P}(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta})} = \mathbb{E}_{\mathbb{Q}} \left[ \log \frac{\mathbb{Q}(z)}{\mathbb{P}(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta})} \right]$$

# Revisiting the EM algorithm II

Let

$$J(\mathbb{Q}, \boldsymbol{\theta}) \triangleq \mathbb{E}_{\mathbb{Q}} (\log L(\boldsymbol{\theta}; \mathbf{X}, \mathbf{Z})) + \mathcal{H}(\mathbb{Q})$$

The steps in the EM algorithm may be viewed as:

Expectation step : choose  $\mathbb{Q}$  to maximize  $J(\mathbb{Q}; \boldsymbol{\theta}^{(t)})$

The solution is  $\mathbb{P}(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}^{(t)})$

Maximization step : choose  $\boldsymbol{\theta}$  to maximize  $J(\mathbb{Q}^{(t)}; \boldsymbol{\theta})$

The solution maximizes  $\mathbb{E}_{\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}^{(t)}} (\log L(\boldsymbol{\theta}; \mathbf{X}, \mathbf{Z}))$

# Variational approximation for SBM

## Problem for SBM

$\mathbb{P}(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}^{(t)})$  cannot be computed thus the E-step cannot be solved.

## Idea

Choose  $\mathbb{Q}$  in a class of function so that the E-step can be solved.

## Family of distribution that factorizes

We chose  $\mathbb{Q}$  so as the  $Z_{iq}$  are marginally independents:

$$\mathbb{Q}(\mathbf{Z}) = \prod_{i=1}^n \mathbb{Q}_i(Z_i) = \prod_{i=1}^n \prod_{q=1}^Q \tau_{iq}^{Z_{iq}},$$

where  $\tau_{iq} = \mathbb{Q}_i(Z_i = q) = \mathbb{E}Q(Z_{iq})$ , with  $\sum_q \tau_{iq} = 1$  for all  $i = 1, \dots, n$ .



# Variational EM for SBM: the criterion

## Lower bound of the loglikelihood

Since  $\mathbb{Q}$  is an approximation of  $\mathbb{P}(\mathbf{Z}|\mathbf{X})$ , the Kullback-Leibler divergence is non-negative and

$$\log L(\boldsymbol{\theta}; \mathbf{X}) \geq \mathbb{E}_{\mathbb{Q}}[\log L(\boldsymbol{\theta}, \mathbf{X}, \mathbf{Z})] + \mathcal{H}(\mathbb{Q}) = J(\mathbb{Q}, \boldsymbol{\theta}).$$

For the SBM,

$$J(\mathbb{Q}, \boldsymbol{\theta}) = \sum_{i,q} \tau_{iq} \log \alpha_q + \sum_{i < j, q, \ell} \tau_{iq} \tau_{j\ell} \log b(X_{ij}; \pi_{q\ell}) - \sum_{i,q} \tau_{iq} \log(\tau_{iq}),$$

$\rightsquigarrow$  we optimize the loglikelihood lower bound  $J(\mathbb{Q}, \boldsymbol{\theta}) = J(\boldsymbol{\tau}, \boldsymbol{\theta})$  in  $(\boldsymbol{\tau}, \boldsymbol{\theta})$ .

# E and M steps for SBM

## Variational E-step

Maximizing  $J(\boldsymbol{\tau})$  for fixed  $\boldsymbol{\theta}$ , we find a fixed-point relationship:

$$\hat{\tau}_{iq} \propto \alpha_q \prod_j \prod_\ell b(X_{ij}, \pi_{q\ell})^{\hat{\tau}_{j\ell}} \quad (1)$$

## M-step

Maximizing  $J(\boldsymbol{\theta})$  for fixed  $\boldsymbol{\tau}$ , we find,

$$\hat{\alpha}_q = \frac{1}{n} \sum_i \hat{\tau}_{iq}, \quad \hat{\pi}_{q\ell} = \frac{\sum_{i \neq j} \hat{\tau}_{iq} \hat{\tau}_{j\ell} X_{ij}}{\sum_{i \neq j} \hat{\tau}_{iq} \hat{\tau}_{j\ell}}. \quad (2)$$

## Model selection

We use our lower bound of the loglikelihood to compute an approximation of the ICL

$$\begin{aligned} \text{vICL}(Q) = \mathbb{E}_{\hat{\mathbb{Q}}}[\log L(\hat{\boldsymbol{\theta}}; \mathbf{X}, \mathbf{Z})] \\ - \frac{1}{2} \left( \frac{Q(Q+1)}{2} \log \frac{n(n-1)}{2} + (Q-1) \log(n) \right), \end{aligned}$$

where

$$\mathbb{E}_{\hat{\mathbb{Q}}}[\log L(\hat{\boldsymbol{\theta}}; \mathbf{X}, \mathbf{Z})] = J(\hat{\boldsymbol{\tau}}, \hat{\boldsymbol{\theta}}) - \mathcal{H}(\hat{\mathbb{Q}}).$$

The variational BIC is just

$$\text{vBIC}(Q) = J(\hat{\boldsymbol{\tau}}, \hat{\boldsymbol{\theta}}) - \frac{1}{2} \left( \frac{Q(Q+1)}{2} \log \frac{n(n-1)}{2} + (Q-1) \log(n) \right).$$

# Example on the French blogosphere I

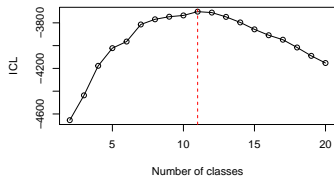
```
library(mixer)
data(blog)
mix.blog <- mixer(x=blog$links,qmin=2,qmax=20)

## Mixer: the adjacency matrix has been transformed in a undirected edge list

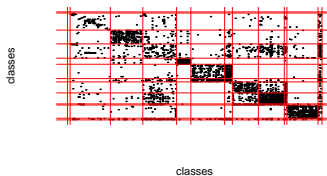
plot(mix.blog)
```

# Example on the French blogosphere II

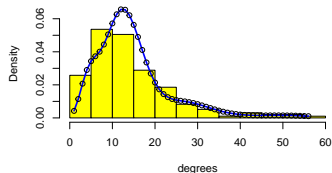
**Integrated Classification Likelihood**



**Reorganized Adjacency matrix**



**Degree distribution**



**Inter/intra class probabilities**

