

Mixture Models

Václav Šmíd

March 20, 2018

Mixture of Gaussians

Probability distribution:

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

where μ_k, Σ_k are mean and covariance matrix of Gaussians weighted by α_k .

Mixture of Gaussians

Probability distribution:

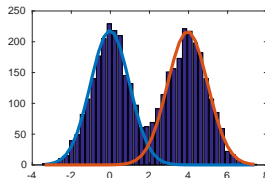
$$p(x) = \sum_{k=1}^K \alpha_k \mathcal{N}(\mu_k, \Sigma_k),$$

where μ_k, Σ_k are mean and covariance matrix of Gaussians weighted by α_k .

- ▶ universal approximation property
- ▶ non-uniqueness
 - ▶ combinatorial,
 - ▶ additive
- ▶ sampling,

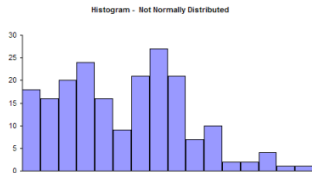
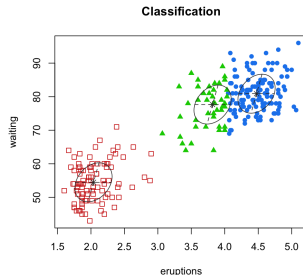
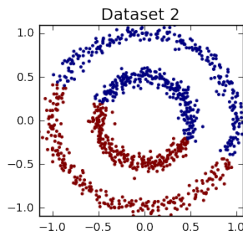
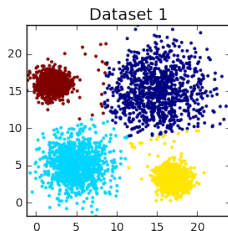
$$p(x) = 0.5\mathcal{N}(0, 1) + 0.5\mathcal{N}(4, 1),$$

- ▶ maximum likelihood,



Uses of mixtures

Clustering (supervised, unsupervised, semi-supervised):



Density representation:

Classification:

Mixture estimation

Consider latent variable $l \in \{\epsilon_1, \dots, \epsilon_K\}$,
 $\epsilon_k = [0, 0, \dots, 1 \dots 0]$. (1-of-n).

$$p(x, l) = p(x|l)p(l),$$

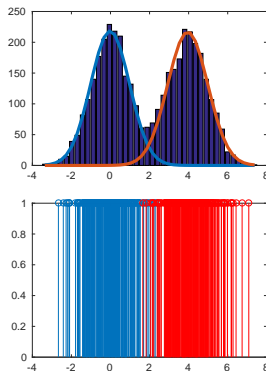
$$p(x|l) = \prod_k \mathcal{N}(\mu_k, \Sigma_k)^{l_k},$$

$$p(l_k = 1) = \alpha_k, \sum_k \alpha_k = 1.$$

$$p(x) = \sum_k p(x|l = \epsilon_k)p(l = \epsilon_k).$$

Multinomial (Bernouli) distribution.

- ▶ Each data point has a label from which component is generated.
- ▶ Estimation of the joint distribution is easier.



Expectation maximization (EM) algorithm

Joint distribution:

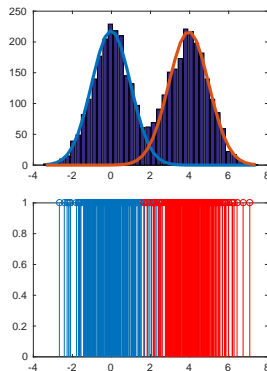
$$p(x, l) = p(x|l)p(l),$$

$$p(x|l) = \prod_k \mathcal{N}(\mu_k, \Sigma_k)^{l_k},$$

$$p(l_k = 1) = \alpha_k, \sum_k \alpha_k = 1.$$

Conditional distribution

$$p(l = \epsilon_k | x) = \frac{p(x, l)}{p(x)} = \frac{\mathcal{N}(\mu_k, \Sigma_k) \alpha_k}{\sum_k \mathcal{N}(\mu_k, \Sigma_k) \alpha_k}$$



Maximum likelihood

Maximizing log-likelihood

$$\log p(x) = \log \left(\sum_{k=1}^K \mathcal{N}(\mu_k, \Sigma_k) \alpha_k \right)$$

$$\begin{aligned} \frac{d}{d\mu_k} \log p(x) &= \frac{1}{\sum_{k=1}^K \mathcal{N}(\mu_k, \Sigma_k) \alpha_k} \alpha_k \mathcal{N}(\mu_k, \Sigma_k) (-\Sigma_k^{-1}(\mu_k - x)), \\ &= p(l = \epsilon_k | x) (-\Sigma_k^{-1}(\mu_k - x)) \end{aligned}$$

After observing N points

$$\begin{aligned} \frac{d}{d\mu_k} \log p(x) &= \sum_{i=1}^n p(l = \epsilon_k | x_i) (-\Sigma_k^{-1}(\mu_k - x_i)) \\ \hat{\mu}_k &= \frac{1}{N_k} \sum_i p(l = \epsilon_k | x_i) x_i \\ N_k &= \sum_i p(l = \epsilon_k | x_i) \end{aligned}$$

$p(l = \epsilon_k | x_i)$ is a function of $\mu_k, \Sigma_k, \forall k$.

Expectation Maximization (EM) algorithm

[Dempster, Laird, Rubin, 1977]

Initialize: choose $\alpha_k^{(0)}, \mu_k^{(0)}, \Sigma_k^{(0)}, \forall k$

Iterate:

1. Compute expected labels:

$$p(I = \epsilon_k | x_i) = \hat{l}_{k,i} = \frac{\mathcal{N}(\mu_k, \Sigma_k) \alpha_k}{\sum_k \mathcal{N}(\mu_k, \Sigma_k) \alpha_k}$$

2. Recompute the component parameters

$$\hat{\mu}_k = \frac{1}{N_k} \sum_i \hat{l}_{k,i} x_i,$$

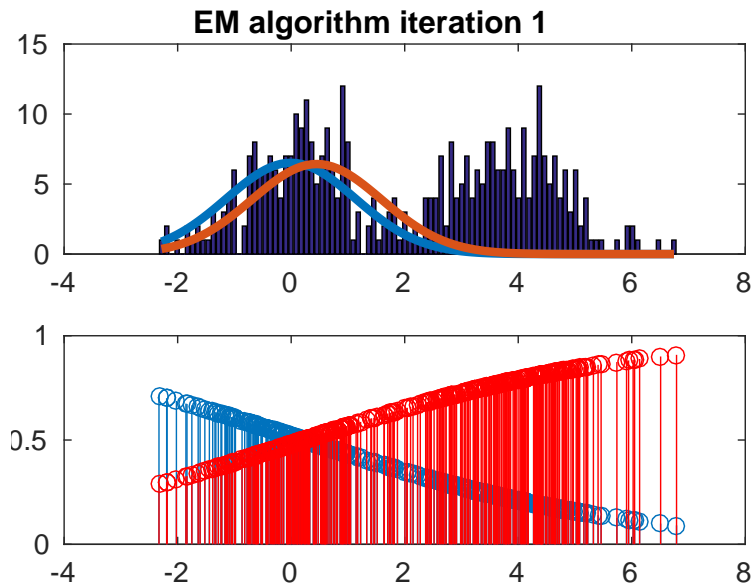
$$\hat{\Sigma}_k = \frac{1}{N_k} \sum_i \hat{l}_{k,i} (x_i - \hat{\mu}_k) (x_i - \hat{\mu}_k)^T,$$

$$\hat{\alpha}_k = \frac{N_k}{N},$$

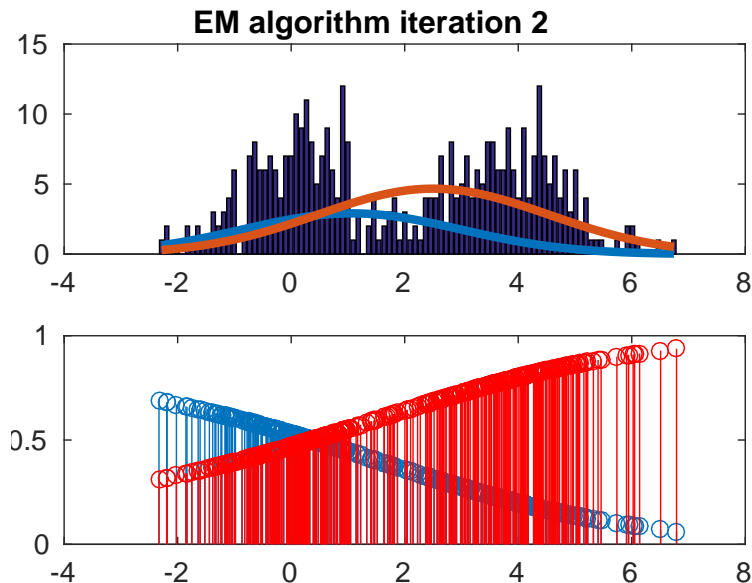
3. (Evaluate log-likelihood)

$$\log p(x) = \sum_i \log \left(\sum_{k=1}^K \mathcal{N}(\hat{\mu}_k, \hat{\Sigma}_k) \hat{\alpha}_k \right)$$

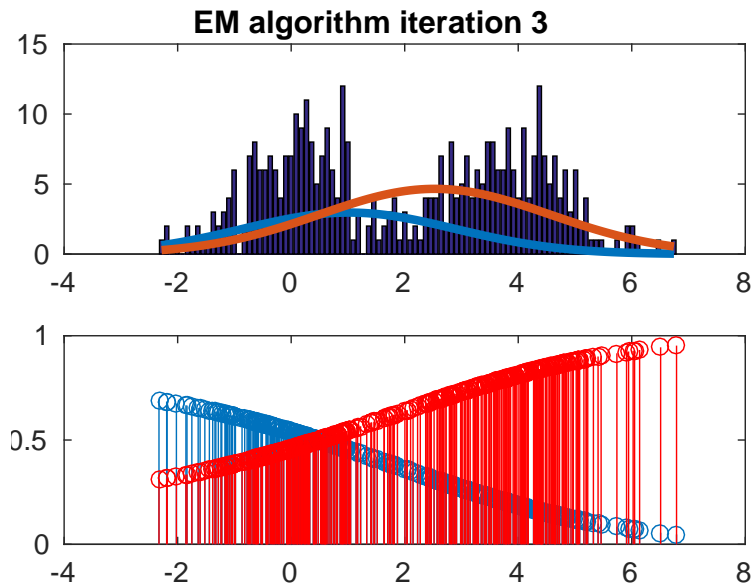
Expectation Maximization (EM) algorithm



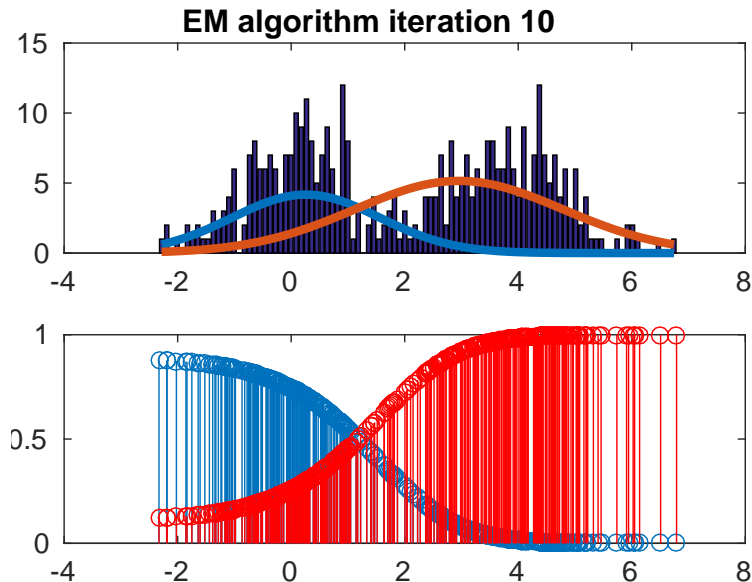
Expectation Maximization (EM) algorithm



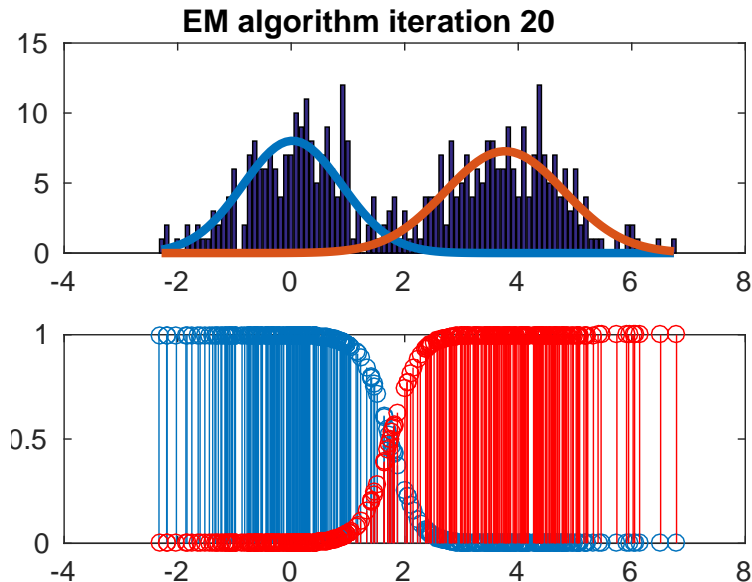
Expectation Maximization (EM) algorithm



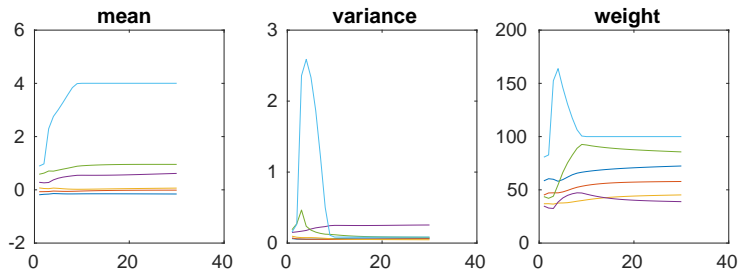
Expectation Maximization (EM) algorithm



Expectation Maximization (EM) algorithm



Expectation Maximization (EM) algorithm



Bayesian treatment

Joint distribution:

$$p(x, l | \alpha) = p(x | l) p(l),$$

$$p(x | l) = \prod_k \mathcal{N}(\mu_k, \omega_k)^{l_k},$$

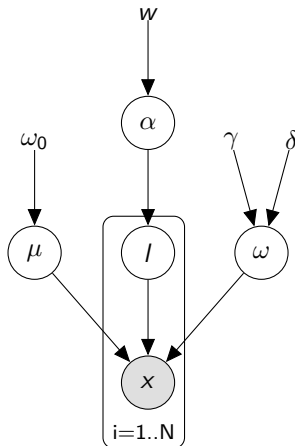
$$p(l_k = 1 | \alpha) = \alpha_k, \sum_k \alpha_k = 1.$$

Priors

$$p(\mu_k | \omega_k) = \mathcal{N}(0, \infty),$$

$$p(\omega_k) = G(0, 0),$$

$$p(\alpha_k) = \text{Di}(w_k) = \frac{\Gamma(\sum_k w_k)}{\prod_k \Gamma(w_k)} \prod_k \alpha_k^{w_{0,k}-1},$$



Variational Bayes for Mixtures

Joint likelihood:

$$\log p(l, \omega, \mu, w) \propto \sum_{i,k} l_{i,k} \left[\frac{1}{2} (\log \omega_k - (x_i - \mu_k) \omega_k (x_i - \mu_k)) \right. \\ \left. + \log \alpha_k \right] - \log \omega_k + (w_{0,k} - 1) \log \alpha_k,$$

Factors, with $n_k = \sum_i \hat{l}_{i,k}$

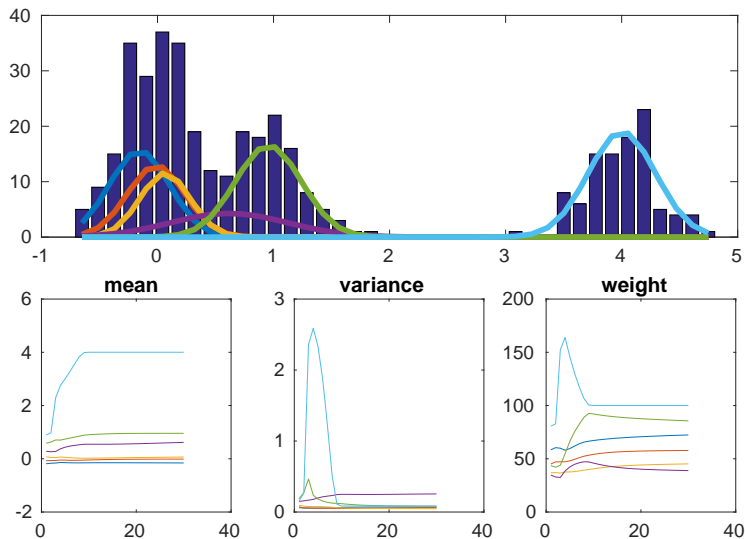
$$q(\omega_k | x) = G(\gamma, \delta), \quad \gamma = n_k \quad \delta = \sum_i \hat{l}_{i,k} (x_i - \hat{\mu}_k)^2 + \textcolor{red}{n_k \hat{\sigma}_k},$$

$$q(\mu_k | x) = \mathcal{N}(\hat{\mu}_k, \hat{\sigma}_k), \quad \hat{\sigma}_k = (n_k \hat{\omega}_k)^{-1} \quad \hat{\mu}_k = \frac{1}{n_k} \sum_i l_{i,k} x_i.$$

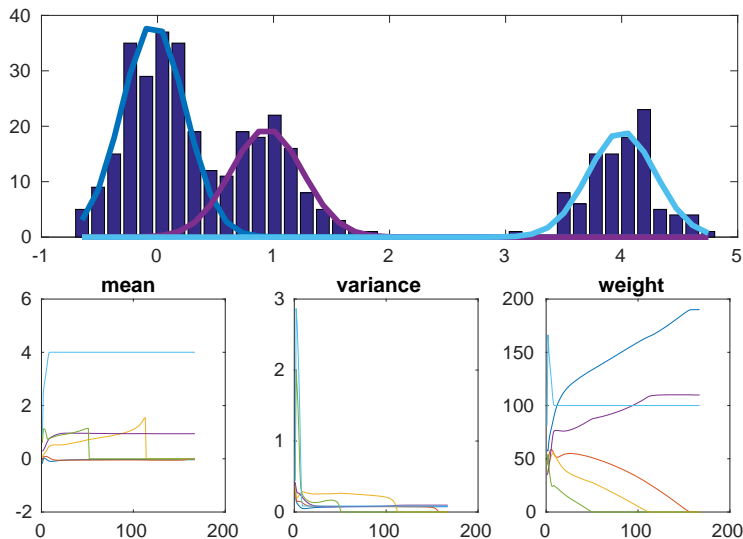
$$q(w | x) = Di(w) \quad w_k = n_k$$

$$q(l_i | x) = Mu(\lambda) \quad \hat{l}_{i,k} = \frac{\lambda_{i,k}}{\sum_i \lambda_{i,k}} \quad \lambda_{i,k} = \exp \frac{1}{2} [-(x_i - \hat{\mu}_k) \hat{\omega}_k (x_i - \hat{\mu}_k) \\ \textcolor{red}{\langle \log \omega_k \rangle - \sigma_\mu \hat{\omega}_k}]$$

EM: $\mu_{true} = \{0, 1, 4\}$, fit $K = 6$ components



VB: $\mu_{true} = \{0, 1, 4\}$ fit $K = 6$ components



Mixture of Gaussians in higher dimensions

Multivariate Gaussians in dimension d :

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Omega}^{-1}),$$

$$\boldsymbol{\mu} \sim \mathcal{N}(\boldsymbol{\mu}_0, (\tau\boldsymbol{\Omega})^{-1})$$

$$\boldsymbol{\Omega} \sim \mathcal{W}(\mathbf{V}, \nu),$$

where \mathcal{W} is the Wishart distribution with ν degrees of freedom.

Covariance matrix:

full covariance: effective number of data $n_k > d$, $O(d^2)$,

scaled identity: homogenous noise σI , (k-means),

diagonal: ignoring rotation of ellipses,

low rank: only selected principal components,

...

Mixture of Gaussians in higher dimensions

Initialization:

random: over what space? cubic...

LHS: latin hypercube sampling

Number of component:

very many: slow convergence

birth and death: random generation

split and merge: evaluate which component to split and/or which two components join into one.

problematic.

Challenge: Patlak Rutland plot

Sequence of scintigraphic images of kidneys.

filename: drsprg_023

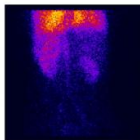
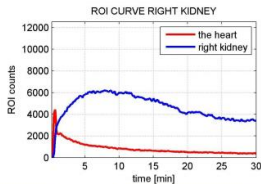
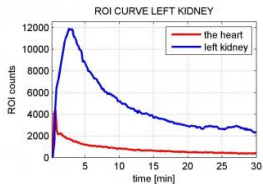
gender = F, age = 31 yrs

CKD stage = 2, LK = 77 %

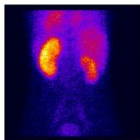
serum Cr - 0, Cr clearance - 0

99mTc-MAG3 - 0, 51Cr-EDTA - 0

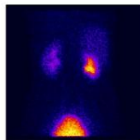
57Co-FLOOD - 1



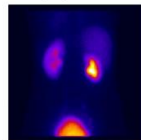
0 - 1 min



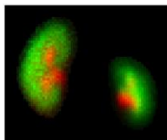
1 - 2 min



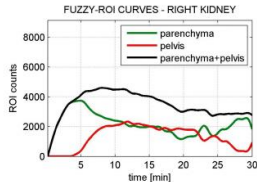
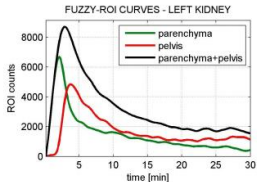
29 - 30 min



MEAN IMAGE



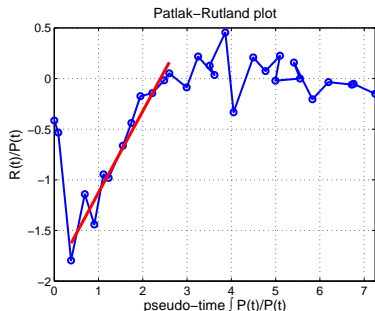
FUZZY ROIS



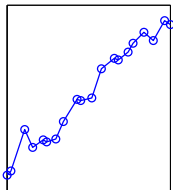
Challenge: Patlak Rutland plot

Patlak Rutland plot is a ratio of parenchyma curve over integral of heart curve.

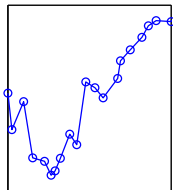
- ▶ typically starts around 1min
- ▶ typically ends around 3min
- ▶ with outliers
- ▶ the slope is a diagnostically important



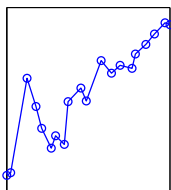
dtpa_1i1c.crv



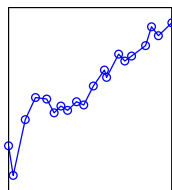
dtpa_2i1c.crv



dtpa_3i1c.crv



dtpa_4i1c.crv



Mixtures

- ▶ Mixture of linear regressions

$$p(y|x) = \sum_k \alpha_k \mathcal{N}(X\theta_k, \sigma_k)$$

- ▶ Mixture of Gamma, Beta distributions – for positive support,
- ▶ Mixture of factor analyzers,
- ▶ Mixture of dynamic models,

Same basic principle:

- ▶ define latent variable with indicator of x being generated from each component.

Assignment

Load data Patlak.mat

35 studies with:

`xpr` x axis

`ypr` y axis

`name` name of the study

`int_start` index where the linear part can start

`int_end` index where the linear part should end

Assignment	points
find slope of linear part for all 35 studies	
a) built-in function	10
b) own code	25