

Paper Reading: Outlier Detection via Parsimonious Mixtures of Contaminated Gaussian Distributions

Nguyen D. Pham

Department of Computer Science
University of Houston

September 15, 2015

Introduction

- » Parametric Density Estimation: Maximum Likelihood
- » Application:
 - Alternative to non-parametric method.
 - Clustering and classification.

Discussion

Introduction

Methodology

- EM Algorithm

- EM Variants

Main Contribution

Experiment Result

Further Reading

The Model

- » Multivariate random variable \mathbf{X} as a mixture model of k components.

$$p(x; \Psi) = \sum_{j=1}^k \pi_j f(x; \vartheta_j)$$

weights $\{\pi_j\}_{j=1}^k$, $\pi_j > 0$, $\sum_{j=1}^k \pi_j = 1$

parameters of the density functions: $\{\vartheta_j\}_{j=1}^k$

parameters set: $\Psi = \{\boldsymbol{\pi}, \boldsymbol{\vartheta}\}$

- » Contaminated Gaussian distribution:

$$f(x; \vartheta_j) = \alpha_j \phi(x; \mu_j, \Sigma_j) + (1 - \alpha_j) \phi(x; \mu_j, \eta_j \Sigma_j)$$

$\alpha_j \in [0, 1]$, $\eta_j > 0$, $\vartheta_j = \{\alpha_j, \mu_j, \Sigma_j, \eta_j\}$

Multivariate Gaussian:

$$\phi(x; \mu, \Sigma) = (2\pi)^{-\frac{p}{2}} |\Sigma|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (x - \mu) \Sigma^{-1} (x - \mu)^T \right\}$$

EM Algorithm

- » Set of observation $\mathbf{X} = \{x_1, x_2, \dots, x_n\}$. The density function defined by set of parameters Θ , $p(x|\Theta)$. The density of the observations:

$$p(\mathbf{X}|\Theta) = \prod_{i=1}^n p(x_i|\Theta) = \mathcal{L}(\Theta|\mathbf{X})$$

Finding Θ that maximizes the likelihood function \mathcal{L}

$$\Theta^* = \arg \max_{\Theta} \mathcal{L}(\Theta|\mathbf{X})$$

- » Assume the observation \mathbf{X} is incomplete, and latent variable \mathbf{Y} . The “true” density function:

$$p(x, y|\Theta) = p(y|x, \Theta)p(x|\Theta)$$

$\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$ is the *complete* data set. Define the *complete-data* likelihood function:

$$\mathcal{L}(\Theta|\mathbf{Z}) = \mathcal{L}(\Theta|\mathbf{X}, \mathbf{Y})$$

- » The likelihood is a function of the random variable \mathbf{Y} , thus we define its expectation over the domain of \mathbf{Y}

$$E[\mathcal{L}(\Theta|\mathbf{Z})] = E[p(x, y|\Theta) | x, \Theta]$$

and take logarithm, the problem becomes:

$$\Theta^* = \arg \max_{\Theta} E[\log p(x, y|\Theta) | x, \Theta]$$

EM Algorithm

» In very simple case, we can solve the problem analytically. In most cases, we need numerical method: EM algorithm, an interactive process where each iteration includes two steps: E-step and M-step.

- E-step: fix the parameters, calculate the expectation.

$$Q\left(\Theta, \Theta^{(i-1)}\right) = E[\log p(x, y | \Theta) | x, \Theta] \quad (1)$$

- M-step: find the parameters to maximize E-step result.

$$\Theta^{(i)} = \arg \max_{\Theta} Q\left(\Theta, \Theta^{(i-1)}\right)$$

» Generalized EM (GEM): relaxing the M-step, only need to find $\Theta^{(i)}$ to increase Q

EM for Mixture Model

- » Assume a mixture of density functions:

$$p(x|\Theta) = \sum_{i=1}^k \omega_i f_i(x|\theta_i)$$

The parameters set $\Theta = (\omega_1, \dots, \omega_k, \theta_1, \dots, \theta_k)$, $\sum_{i=1}^k \omega_i = 1$ and each f_i is a density function.

- » Incomplete-data log-likelihood from n observations \mathbf{X} :

$$\log(\mathcal{L}(\Theta|\mathbf{X})) = \log \prod_{i=1}^n p(x_i, \Theta) = \sum_{i=1}^n \log \left(\sum_{j=1}^k \omega_j p_j(x_i|\theta_j) \right)$$

log of sum is difficult to optimize.

- » Introduce the latent variable $\mathbf{Y} = \{y_i\}_{i=1}^n$ where $y_i \in 1, \dots, k$ indicating the membership of an observation in k components.

EM for Mixture Model

- » If we can observe \mathbf{Y} , the complete-data log-likelihood is:

$$\begin{aligned}\log(\mathcal{L}(\Theta|\mathbf{X}, \mathbf{Y})) &= \log(p(\mathbf{X}, \mathbf{Y}|\Theta)) = \sum_{i=1}^n \log(p(x_i|y_i) p(y_i)) \\ &= \sum_{i=1}^k \log(\omega_{y_i} p_{y_i}(x_i|\theta_{y_i}))\end{aligned}$$

which is easier to process.

- » However \mathbf{Y} is a random variable, we must compute the expectation as shown in equation (1). This is more involved...

EM for Mixture Model

- » Skip to the results: Mixture of Gaussians, the parameters set is $\Theta = \{\omega, \mu, \Sigma\}$ which are the weights, the means, and the covariances.

$$\omega_j^{new} = \frac{1}{n} \sum_{i=1}^n p(j|x_i, \hat{\Theta})$$
$$\mu_j^{new} = \frac{\sum_{i=1}^n x_i p(j|x_i, \hat{\Theta})}{\sum_{i=1}^n p(j|x_i, \hat{\Theta})}$$
$$\Sigma_j^{new} = \frac{\sum_{i=1}^n p(j|x_i, \hat{\Theta}) (x_i - \mu_j^{new}) (x_i - \mu_j^{new})^T}{\sum_{i=1}^n p(j|x_i, \hat{\Theta})}$$

where $j = 1..k$ (number of mixture components).

Model Selection

- » Bayesian information criterion (BIC): Maximum log-likelihood with minimal complexity:

$$BIC = -2 \ln \left(p \left(\mathbf{X} | \hat{\Theta} \right) \right) + \rho \ln(n)$$

ρ is the number of free parameters.

- » Others: ICL, DIC, AIC, all are related to BIC by some approximation.

ECM Algorithm

- » Expectation Conditional Maximization (ECM): A subclass of GEM algorithm. M-step is replaced by multiple CM-steps.
- » Idea: finding multivariate Θ is difficult, it is easier to maximize by one parameter, assuming the others are constants. In general, divide Θ into subsets, find the parameters in each subset while putting constrain on the rest.
- » How many CM-steps? Depend on the analysis of the log-likelihood function.

Parsimonious Variants

- » p -variate domain $\rightarrow p(p+1)/2$ free parameters for each covariance matrix \rightarrow parsimonious models.
- » Eigen decomposition of the covariance matrix:

$$\Sigma_j = \lambda_j \Gamma_j \Delta_j \Gamma_j^T$$

λ_j : Volume.

Γ_j : matrix where columns are normalized eigenvectors: Orientation.

Δ_j : diagonal matrix of eigenvalues in decreasing order: Shape.

- » Constraints on the three components yield fourteen parsimonious models, grouped into three categories: *spherical*, *diagonal*, and *general*.

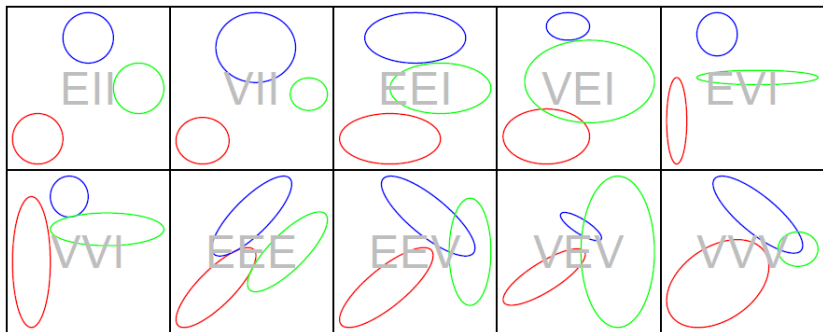
Parsimonious Mixtures of Contaminated Gaussian Distribution models

Table 1: Nomenclature, covariance structure, type of ML solution in the first CM-step of the ECM algorithm (CF=closed form and IP=iterative procedure), and number of free covariance parameters for each member of the PMCGD family.

Family	Model	Volume	Shape	Orientation	Σ_j	ML	Free covariance parameters
Spherical	EII	Equal	Spherical	-	$\lambda \mathbf{I}$	CF	1
	VII	Variable	Spherical	-	$\lambda_j \mathbf{I}$	CF	k
Diagonal	E EI	Equal	Equal	Axis-Aligned	$\lambda \Delta$	CF	p
	VEI	Variable	Equal	Axis-Aligned	$\lambda_j \Delta$	IP	$k + p - 1$
	EVI	Equal	Variable	Axis-Aligned	$\lambda \Delta_j$	CF	$1 + k(p - 1)$
	VVI	Variable	Variable	Axis-Aligned	$\lambda_j \Delta_j$	CF	kp
General	EEE	Equal	Equal	Equal	$\lambda \Delta \Gamma \Delta'$	CF	$p(p + 1) / 2$
	VEE	Variable	Equal	Equal	$\lambda_j \Delta \Gamma \Delta'$	IP	$k + p - 1 + p(p - 1) / 2$
	EVE	Equal	Variable	Equal	$\lambda \Delta_j \Gamma \Delta_j'$	IP	$1 + k(p - 1) + p(p - 1) / 2$
	EEV	Equal	Equal	Variable	$\lambda \Delta \Gamma_j \Delta_j'$	CF	$p + kp(p - 1) / 2$
	VVE	Variable	Variable	Equal	$\lambda_j \Delta_j \Gamma \Delta_j'$	IP	$kp + p(p - 1) / 2$
	VEV	Variable	Equal	Variable	$\lambda_j \Delta_j \Gamma_j \Delta_j'$	IP	$k + p - 1 + kp(p - 1) / 2$
	EVV	Equal	Variable	Variable	$\lambda \Delta_j \Gamma_j \Delta_j'$	CF	$1 + k(p - 1) + kp(p - 1) / 2$
	VVV	Variable	Variable	Variable	$\lambda_j \Delta_j \Gamma_j \Delta_j'$	CF	$kp(p + 1) / 2$

(Punzo and McNicholas)

Bivariate Model Illustration



(Brendan Murphy)

Parsimonious Mixture of Contaminated Gaussian Distributions

- » The paper introduces the model as described in slide 4
- » The parameters set partitions: $\Psi = \{\Psi_1, \Psi_2\}$ where $\Psi_1 = \{\boldsymbol{\pi}, \boldsymbol{\alpha}, \boldsymbol{\mu}, \boldsymbol{\Sigma}\}$ and $\Psi_2 = \{\boldsymbol{\eta}\}$.
- » The authors use **mixture**, a **R** package implementing the parsimonious mixture Gaussians.
- » They implement the 2-step ECM, and give the analytic equations for the VVV model (in section 3.1).
- » They detail the initialization of the contaminated components: $\alpha_j = \eta_j = 1$, meaning: no contamination.
- » The discussion on convergence is somewhat unnecessary, as it is a known result.

Outlier Detection

» Scan through the data set, for each x_i , perform outlier detection in two steps.

- Determine group membership: MAP (maximal a posteriori).

$$j^* = \arg \max_j \pi_j (\alpha_j \phi(x_j; \mu_j, \Sigma_j) + (1 - \alpha_j) \phi(x_j; \mu_j, \eta_j \Sigma_j))$$

- Determine if x_i is an outlier:

$$\max(\alpha_j \phi(x_j; \mu_j, \Sigma_j), (1 - \alpha_j) \phi(x_j; \mu_j, \eta_j \Sigma_j))$$

» The parameter α_j sets the a priori for the good parameters, and is allowed to define with a lower bound or a constant.

Experiment Result

- » Compare BIC and ICL, using synthesized data sets: same performance.
- » On two real data sets: clustering and detect outliers.
- » The results show the strength of the proposed model to detect outliers. But only with 2-d data sets (somewhat quite simple).

Further Reading

- » Dempster A.P., 1976. Maximum Likelihood from Incomplete Data via the EM Algorithm.
- » Donald B.R., 1993. Maximum Likelihood Estimation via the ECM Algorithm: A General Framework.
- » Celeux G., 1993. Gaussian Parsimonious Clustering Models.
- » Steele R.J., 2009. Performance of Bayesian Model Selection Criteria for Gaussian Mixture Models.