

Variational Mixture of Gaussians

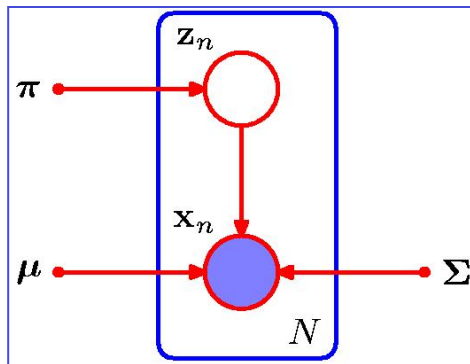
Sargur Srihari
srihari@cedar.buffalo.edu

Objective

- Apply variational inference machinery to Gaussian Mixture Models
- Demonstrates how Bayesian treatment elegantly resolves difficulties with maximum likelihood issues
- Many more complex distributions can be solved using straightforward extensions of this analysis

Graphical Model for GMM

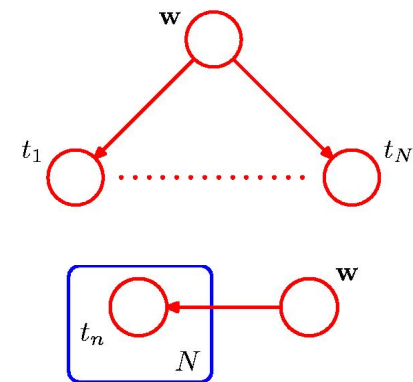
- Graphical model corresponding to likelihood function of standard GMM:



Directed acyclic graph
Representing mixture

- For each observation x_n we have a corresponding latent variable z_n
 - A 1-of-K binary vector with elements z_{nk} for $k=1, \dots, K$
- Denote observed data by $X = \{x_1, \dots, x_N\}$
- Latent variables by $Z = \{z_1, \dots, z_N\}$

Plate Notation:
Equivalent networks



Likelihood Function for GMM

Mixture density function is

x

Since z has values $\{z_k\}$
with probabilities π_k

Therefore Likelihood function is

$$p(X | \pi, \mu, \Sigma) = \prod_{n=1}^N \left\{ \sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \right\}$$

Product is over the N
i.i.d. samples

Therefore log-likelihood function is

$$\ln p(X | \pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \right\}$$

Find parameters π , μ and Σ that maximize log-likelihood

A more difficult problem than for a single Gaussian

GMM m.l.e. expressions

- Obtained using derivatives of log-likelihood

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n$$

Parameters (means)

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (x_n - \mu_k)(x_n - \mu_k)^T$$

Parameters(covariance matrices)

$$\pi_k = \frac{N_k}{N} \quad N_k = \sum_{n=1}^N \gamma(z_{nk})$$

Parameters
(Mixing Coefficients)

All three are in terms of responsibilities

- Not closed form solutions for the parameters
 - Since $\gamma(z_{nk})$ the responsibilities depend on those parameters in a complex way

EM For GMM

- E step

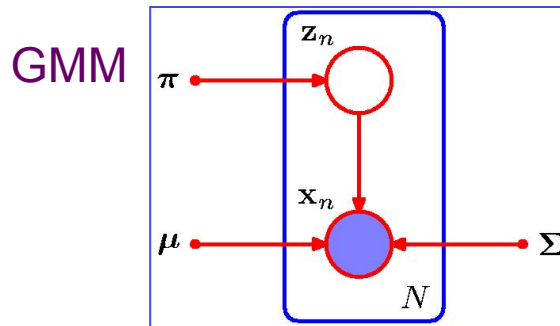
- use current value of parameters

- μ_k, Σ_k, π_k to evaluate posterior probabilities $p(Z/X)$,
i.e., responsibilities $\gamma(z_{nk})$

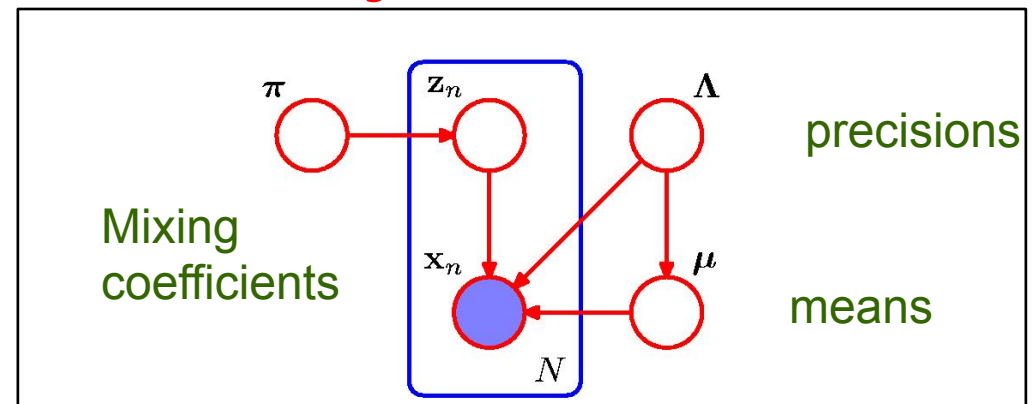
- M step

- use these posterior probabilities to re-estimate $p(X,Z)$: means, covariances and mixing coefficients
wrt $p(Z/X)$

Graphical model for Bayesian GMM



Bayesian
GMM



- To specify model we need these conditional probabilities:
 - $p(Z|\pi)$: conditional distribution of Z given mixing coeffs
 - $p(X|Z, \mu, \Lambda)$:
 - $p(\pi)$: distribution of mixing coefficients
 - $p(\mu, \Lambda)$: prior governing mean and precision of each component

Conditional Distribution Expressions

1. Conditional distribution of $Z=\{z_1,..,z_N\}$ given mix coefficients π

Since components are mutually exclusive $p(z) = \prod_{k=1}^K \pi_k^{z_k}$

$$p(Z | \pi) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}}$$

2. Conditional distribution of observed data $X=\{x_1,..,x_N\}$ given latent variables and component parameters $p(X|Z, \mu, \Lambda)$

– Since components are Gaussian $p(x | z) = \prod_{k=1}^K N(x | \mu_k, \Sigma_k)^{z_k}$

$$p(X | Z, \mu, \Lambda) = \prod \prod N(x_n | \mu_k, \Lambda_k^{-1})^{z_{nk}}$$

– where $\mu = \{\mu_k\}$ and $\Lambda = \{\Lambda_k\}$

- use of precision matrix simplifies further analysis

Parameter Priors: Mixing Coefficients

3. Distribution of mixing coefficients $p(\pi)$

- Conjugate priors simplify analysis
- Dirichlet distribution over π

$$p(\pi) = \text{Dir}(\pi | \alpha_0) = C(\alpha_0) \prod_{k=1}^K \pi_k^{\alpha_0 - 1}$$

- We have chosen the same parameter α_0 for each of the components
- $C(\alpha_0)$ is the normalization constant for the Dirichlet distribution

Parameter Priors: Mean, Precision

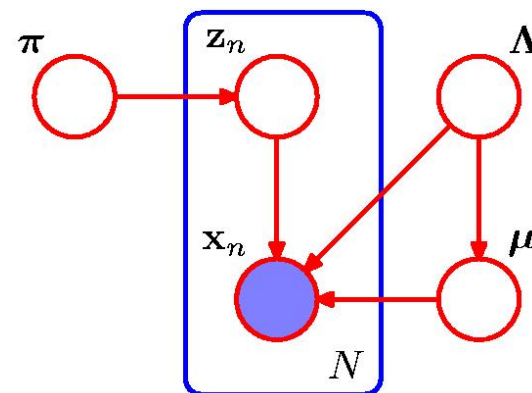
4. Distribution of Mean and Precision of Gaussian components $p(\mu, \Lambda)$

- Gaussian-Wishart prior is

$$\begin{aligned} p(\mu, \Lambda) &= p(\mu | \Lambda) p(\Lambda) \\ &= \prod_{k=1}^K \mathcal{N}(\mu_k | m_0 (\beta_0 \Lambda_k)^{-1}) W(\Lambda_k | W_0, \nu_0) \end{aligned}$$

- Which represents the conjugate prior when both mean and precision are unknown

- Resulting model has:
 - Link between Λ and μ
 - Due to distribution (4) above

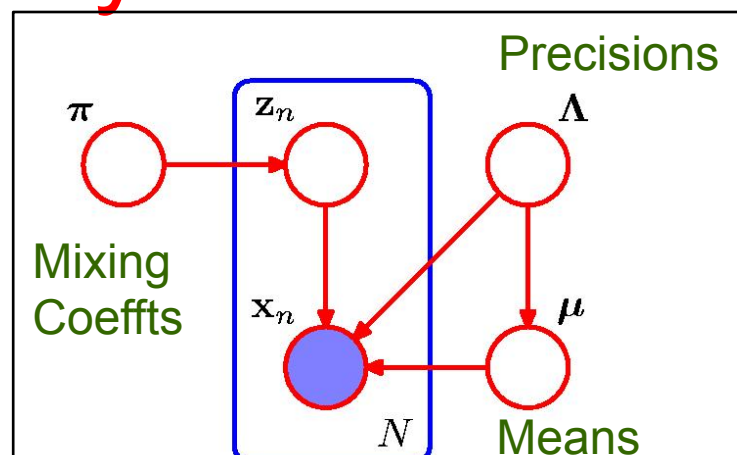


Bayesian Network for Bayesian GMM

- Joint of all random variables:

$$p(X, Z, \pi, \mu, \Lambda) = p(X | Z, \mu, \Lambda) p(Z | \pi) p(\mu | \Lambda) p(\Lambda)$$

- All the factors were given earlier
- Only $X = \{x_1, \dots, x_N\}$ are observed



- This BN provides a nice distinction between latent variables and parameters
 - Variables such as z_n that appear inside the plate are latent variables
 - No of such variables grows with data set
 - Variables outside the plate are parameters
 - Fixed in no. and outside of data set
- From viewpoint of PGMs no fundamental difference¹¹

The variational approach

- Recall GMM $p(\mathbf{x}) = \sum_z p(z)p(\mathbf{x} | z) = \sum_{k=1}^K \pi_k N(\mathbf{x} | \mu_k, \Sigma_k)$
- The EM approach:
 1. Evaluation of posterior distribution $p(Z|X)$
 2. Evaluation of expectation of $p(X,Z)$ wrt to $p(Z|X)$
- Our goal is to specify the variational distribution $q(Z, \pi, \mu, \Lambda)$ which will specify $p(Z, \pi, \mu, \Lambda | X)$

Here $p(z)$ has parameter π
with distribution $p(\pi)$

– Recall $\ln p(X) = L(q) + KL(q \parallel p)$

where

$$L(q) = \int q(Z) \ln \left\{ \frac{p(X, Z)}{q(Z)} \right\} dZ$$

and

$$KL\{q \parallel p\} = - \int q(Z) \ln \left\{ \frac{p(Z | X)}{q(Z)} \right\} dZ$$

Variational Distribution

- In variational inference we can specify q by using a factorized distribution $q(Z) = \prod_{i=1}^M q_i(Z_i)$
 - For Bayesian GMM the latent variables and parameters are Z , π , μ and Λ .

- So we consider the *variational distribution*

$$q(Z, \pi, \mu, \Lambda) = q(Z)q(\pi, \mu, \Lambda)$$

Subscripts for q 's omitted

- Remarkably, this is the only assumption needed for a tractable solution to a Bayesian Mixture Model
- Functional forms of both $q(Z)$ and $q(\pi, \mu, \Lambda)$ are determined automatically by optimizing the variational distribution

Sequential update equations

- Using general result of factorized distributions
 - When $L(q)$ is defined as $L(q) = \int q(Z) \ln \left\{ \frac{p(X, Z)}{q(Z)} \right\} dZ = \int \prod_i q_i \left\{ \ln p(X, Z) - \sum_i \ln q_i \right\} dZ$
 - the q that makes the functional $L(q)$ largest is

$$\ln q_j^*(Z_j) = E_{i \neq j} [\ln p(X, Z)] + \text{const}$$

- For Bayesian GMM log of optimized factor is

$$\ln q^*(Z) = E_{\pi, \mu, \Lambda} [\ln p(X, Z, \pi, \mu, \Lambda)] + \text{const}$$

- Since $p(X, Z, \pi, \mu, \Lambda) = p(X | Z, \mu, \Lambda) p(Z | \pi) p(\mu | \Lambda) p(\Lambda)$ we have

$$\ln q^*(Z) = E_{\pi} [\ln p(Z | \pi)] + E_{\mu, \Lambda} [\ln p(X | Z, \mu, \Lambda)] + \text{const}$$

- Note: Expectations are just weighted sums

Simplification of $q^*(Z)$

- Expression for factor $q^*(Z)$

$$\ln q^*(Z) = E_{\pi} [\ln p(Z | \pi)] + E_{\mu, \Lambda} [\ln p(X | Z, \mu, \Lambda)] + \text{const}$$

- Absorbing terms not depending on Z into constant

$$\ln q^*(Z) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \ln \rho_{nk} + \text{const}$$

$$\text{where } \ln \rho_{nk} = E[\ln \pi_k] + \frac{1}{2} E[\ln |\Lambda_k|] - \frac{D}{2} \ln(2\pi) - \frac{1}{2} E_{\mu_k, \Lambda_k} [(\mathbf{x}_n - \mu_k)^T \Lambda_k (\mathbf{x}_n - \mu_k)]$$

- where D is dimensionality of data variable \mathbf{x}
- Taking exponentials on both sides
- Normalized distribution is

$$q^*(Z) \propto \prod_{n=1}^N \prod_{k=1}^K \rho_{nk}^{z_{nk}}$$

$$q^*(Z) = \prod_{n=1}^N \prod_{k=1}^K r_{nk}^{z_{nk}}$$

where

$$r_{nk} = \frac{\rho_{nk}}{\sum_{j=1}^K \rho_{nj}}$$

r_{nk} are positive since ρ_{nk} are exponentials of real nos. and will sum to one as required

Factor $q^*(Z)$ has same form as prior

- Normalized distribution is
$$q^*(Z) = \prod_{n=1}^N \prod_{k=1}^K r_{nk}^{z_{nk}}$$
- We have found form of q^* to maximize the functional $L(q)$
 - It has same form as prior
$$p(Z | \pi) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}}$$
- Distribution $q^*(Z)$ is discrete and has the standard result $E[z_{nk}] = r_{nk}$,
 - which play the role of responsibilities
- Since equations for $q^*(Z)$ depend on moments of other variables
 - They are coupled and solved iteratively

Variational EM

- Variational E-step: determine responsibilities r_{nk}
- Variational M-step:
 1. determine statistics of data set

$$\begin{aligned} N_k &= \sum_{n=1}^N r_{nk} \\ \bar{\mathbf{x}}_k &= \frac{1}{N_k} \sum_{n=1}^N r_{nk} \mathbf{x}_n \\ S_k &= \frac{1}{N_k} \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \bar{\mathbf{x}}_k)(\mathbf{x}_n - \bar{\mathbf{x}}_k)^T \end{aligned}$$

Responsibility of k^{th} component

Mean of k^{th} component

Covariance matrix of k^{th} component

and

2. find optimal solution for the factor $q(\pi, \mu, \Lambda)$

Factorization of $q(\pi, \mu, \Lambda)$

- Using general result of factorized distributions

$$\ln q_j^*(Z_j) = E_{i \neq j} [\ln p(X, Z)] + \text{const}$$

– We can write

$$\ln q^*(\pi, \mu, \Lambda) = \ln p(\pi) + \sum_{k=1}^K \ln p(\mu_k, \Lambda_k) + E_Z [\ln p(Z | \pi)] + \sum_{n=1}^N \sum_{k=1}^K E[z_{nk}] \ln N(x_k | \mu_k, \Lambda_k^{-1}) + \text{const}$$

- which decomposes into terms involving π and only μ, Λ
- The terms involving μ and Λ comprise sum of terms involving μ_k and Λ_k leading to factorization

$$q(\pi, \mu, \Lambda) = q(\pi) \prod_{k=1}^K q(\mu_k, \Lambda_k)$$

Factor $q(\pi)$ is a Dirichlet

- Given the factorization $q(\pi, \mu, \Lambda) = q(\pi) \prod_{k=1}^K q(\mu_k, \Lambda_k)$
- Consider each factor in turn: $q(\pi)$ and $q(\mu_k, \Lambda_k)$
- (2a) Identifying terms depending on π , $q(\pi)$ has the solution

$$\ln q^*(\pi) = (\alpha_0 - 1) \sum_{k=1}^K \ln z_k + \sum_{k=1}^K \sum_{n=1}^N r_{nk} \ln \pi_k + \text{const}$$

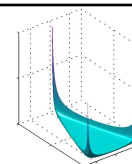
- Taking exponential on both sides we get $q^*(\pi)$ as a Dirichlet

$$q^*(\pi) = \text{Dir}(\pi \mid \alpha)$$

where α has the components $\alpha_k = \alpha_0 + N_k$

Dirichlet:

$$\text{Dir}(\mu \mid \alpha) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_K)} \prod_{k=1}^K \mu_k^{\alpha_k - 1} \text{ where } \alpha_0 = \sum_{k=1}^K \alpha_k$$



$K=3$
 $\alpha_k=0.1$

Factor $q^*(\mu_k, \Lambda_k)$ is a Gaussian-Wishart

(2b) Variational posterior for $q^*(\mu_k, \Lambda_k)$

- Does not further factorize into marginals
- It is a Gaussian-Wishart distribution

$$q^*(\mu_k, \Lambda_k) = \mathcal{N}\left(\mu_k m_k (\beta_k \Lambda_k)^{-1}\right) \mathcal{W}(\Lambda_k | W_0, \nu_0)$$

– \mathcal{W} is the Wishart distribution

- It has the form

$$\mathcal{W}(\Lambda | W, \nu) = B | \Lambda |^{(\nu-D-1)/2} \exp[-1/2 \text{Tr}(W^{-1} \Lambda)]$$

where ν is the no. of degrees of freedom, W is a $D \times D$ scale matrix and Tr is the trace. $B(W, \nu)$ is a normalization constant

- It is the conjugate prior for a Gaussian with known mean and unknown precision matrix Λ

Parameters of $q^*(\mu_k, \Lambda_k)$

- Gaussian-Wishart

$$q^*(\mu_k, \Lambda_k) = \mathcal{N}(\mu_k | m_k, (\beta_k \Lambda_k)^{-1}) \mathcal{W}(\Lambda_k | W_0, \nu_0)$$

– where we have defined

$$\begin{aligned}\beta_k &= \beta_0 + N_k \\ m_k &= \frac{1}{\beta_k} (\beta_0 m_0 + N_k \bar{x}_k) \\ W_k^{-1} &= W_0^{-1} + N_k S_k + \frac{\beta_0 N_k}{\beta_0 + N_k} (\bar{x}_k - m_0)(\bar{x}_k - m_0)^T \\ \nu_k &= \nu_0 + N_k + 1\end{aligned}$$

- These update equations are analogous to M-step of EM for m.l. solution of GMM
 - Involve evaluation of same sums as EM over the data set

Expression for Responsibilities

- For the M step we need expectations $E[z_{nk}] = r_{nk}$
 - Which are obtained by normalizing ρ_{nk}

- Since $r_{nk} = \frac{\rho_{nk}}{\sum_{j=1}^K \rho_{nj}}$ where

$$\ln \rho_{nk} = E[\ln \pi_k] + \frac{1}{2} E[\ln |\lambda_k|] - \frac{D}{2} \ln(2\pi) - \frac{1}{2} E_{\mu_k \Lambda_k} [(\mathbf{x}_n - \mu_k)^T \Lambda_k (\mathbf{x}_n - \mu_k)]$$

- The three expectations wrt variational distribution of parameters are easily evaluated to give

$$\ln \bar{\pi}_k \equiv E[\ln \pi_k] = \psi(\alpha_k) - \nu$$

$$\ln \bar{\Lambda}_k \equiv \frac{1}{2} E[\ln |\lambda_k|] = \sum_{i=1}^D \psi\left(\frac{\nu_k + 1 - i}{2}\right) + D \ln 2 + \ln |W_k|$$

$$E_{\mu_k \Lambda_k} [(\mathbf{x}_n - \mu_k)^T \Lambda_k (\mathbf{x}_n - \mu_k)] = D \beta_k^{-1} + \nu_k (\mathbf{x}_n - \mu_k)^T \Lambda_k (\mathbf{x}_n - \mu_k)$$

ν is the no. of degrees of freedom of Wishart

- ψ is the digamma function with $\hat{\alpha} = \sum_k \alpha_k$

- Digamma $\psi(a) = \frac{d}{da} \ln \Gamma(a)$ appears in the definition of Dirichlet

Evaluation of Responsibilities

- Substituting the three expectations into $\ln \rho_{nk}$

$$r_{nk} \propto \tilde{\pi}_k \tilde{\Lambda}^{1/2} \exp \left\{ -\frac{D}{2\beta_k} - \frac{v_k}{2} (\mathbf{x}_n - \mathbf{m}_k) W_k (\mathbf{x}_n - \mathbf{m}_k) \right\}$$

– This is similar to responsibilities for mle for EM

$$\begin{aligned} \gamma(z_k) \equiv p(z_k = 1 | \mathbf{x}) &= \frac{p(z_k = 1) p(\mathbf{x} | z_k = 1)}{\sum_{j=1}^K p(z_j = 1) p(\mathbf{x} | z_j = 1)} \\ &= \frac{\pi_k N(\mathbf{x} | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(\mathbf{x} | \mu_j, \Sigma_j)} \end{aligned}$$

- which can be written in the form

$$r_{nk} \propto \pi_k |\Lambda_k|^{1/2} \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \mu_k) \Lambda_k (\mathbf{x}_n - \mu_k) \right\}$$

- where we have used precision Λ_k instead of covariance Σ_k to highlight similarity

Summary of Optimization

- Optimization of variational posterior distribution involves cycling between two stages
 - Analogous to E and M steps on m.l. EM
- Variational E-step:
 - Use current distribution over model parameters to evaluate moments and hence evaluate $E[z_{nk}] = r_{nk}$
- Variational M step
 - keep responsibilities fixed; use them to recompute variational distribution over the parameters using

$$\boxed{q^*(\pi) = \text{Dir}(\pi \mid \alpha)}$$
 and $\boxed{q^*(\mu_k, \Lambda_k) = \text{N}(\mu_k m_k (\beta_k \Lambda_k)^{-1}) \mathcal{W}(\Lambda_k \mid W_0, \nu_0)}$

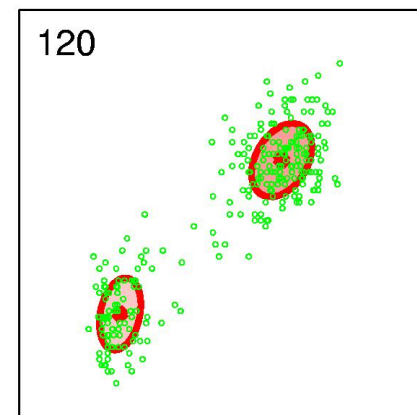
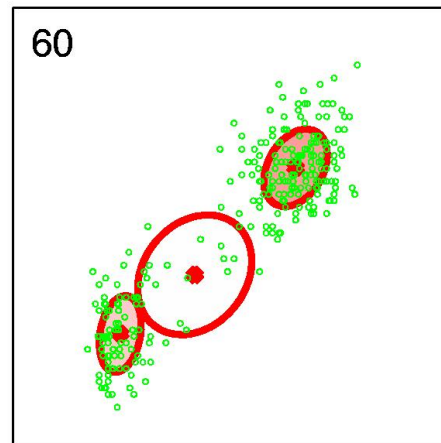
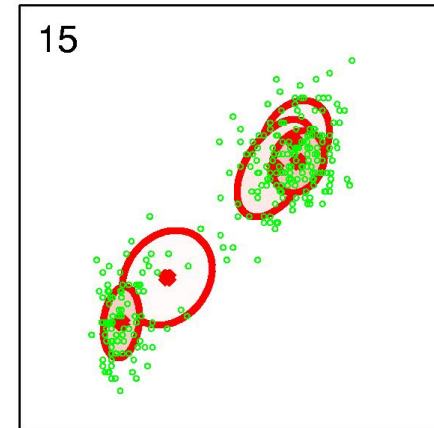
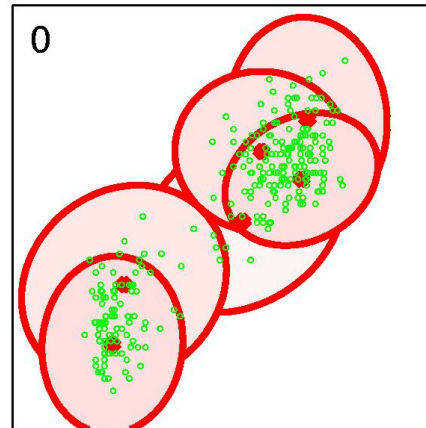
Variational Bayesian GMM

Old Faithful data set

K=6 components

After convergence there are only two components

Density of red ink inside each ellipse shows Mean value of Mixing coefficients



Similarity of Variational Bayes and EM

- Close similarity between variational solution for the Bayesian mixture of Gaussians and the EM algorithm for maximum likelihood
- In the limit as $N \rightarrow \infty$, the Bayesian treatment converges to the maximum likelihood EM
- Variational algorithm is more expensive but problem of singularity is eliminated

Variational Lower Bound

- We can straight-forwardly evaluate the lower bound $L(q)$ for this model

- Recall $\ln p(X) = L(q) + KL(q \parallel p)$

where

$$L(q) = \int q(Z) \ln \left\{ \frac{p(X, Z)}{q(Z)} \right\} dZ$$

and

$$KL\{q \parallel p\} = - \int q(Z) \ln \left\{ \frac{p(Z \mid X)}{q(Z)} \right\} dZ$$

- The lower bound is used to monitor re-estimation to test for convergence

Predictive Density

- In using a Bayesian GMM we will be interested in the predictive density for a new value $\hat{\mathbf{x}}$ of the observed variable
- Assuming corresponding latent variable we can show that

$$p(\hat{\mathbf{x}} | X) = \frac{1}{\hat{\alpha}} \sum_{k=1}^N \alpha_k St(\hat{\mathbf{x}} | \mathbf{m}_k, \mathbf{L}_k, \nu_k + 1 - D)$$

where the k^{th} component has mean \mathbf{m}_k and precision

$$\mathbf{L}_k = \frac{(\nu_k + 1 - D) \beta_k}{(1 + \beta_k)} \mathbf{W}_k$$

- The mixture of Student's T becomes a GMM as $N \rightarrow \infty$

Determining no. of components

- Plot of variational lower bound \mathcal{L} versus no. of components K
- Distinct peak at $K=2$
- For each K model is trained from 100 starts
 - Results shown as +

