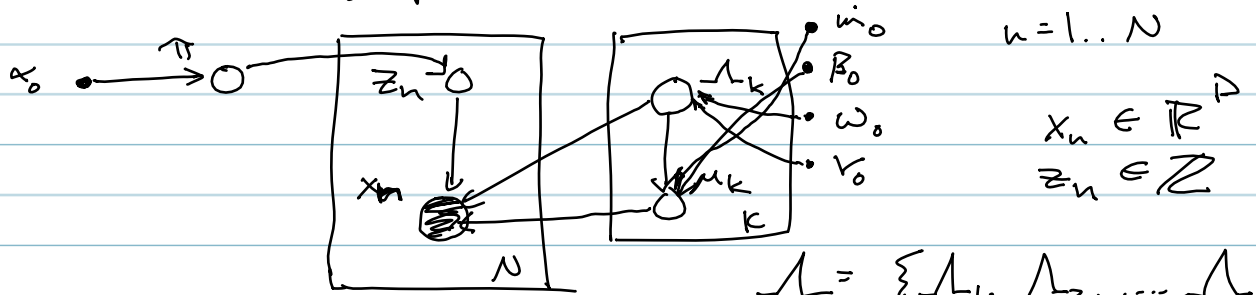


Variational Bayes Updates for a ^{Bayesian} Gaussian Mixture Model

B-GMM graphical model



Latent vars:

$$\pi, z, \lambda, \mu$$

Observed vars:

$$x$$

Parameters:

$$\alpha_0, m_0, B_0, W_0, v_0$$

hyperprior for class labels, hyperprior for class means & covariances

$$\lambda = \{\lambda_1, \lambda_2, \dots, \lambda_K\}$$

K is # of classes

$$\lambda_k \in \mathbb{R}^{D \times D}, \lambda_k \text{ pos-semidef.}$$

$$\mu = \{\mu_1, \mu_2, \dots, \mu_K\}$$

$$\mu_k \in \mathbb{R}^D$$

$$\pi = [\pi_1, \dots, \pi_K]$$

class prior probabilities

Likelihood of $z | \pi$, discrete

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

z_{nk} is vector with one non-zero element equal to 1

Conditional likelihood of $X | z, \mu, \lambda$

$$p(X | z, \mu, \lambda) = \prod_{n=1}^N \prod_{k=1}^K \mathcal{N}(x_n | \mu_k, \lambda_k^{-1})^{z_{nk}}$$

Priors of π, μ, λ

$$p(\vec{\pi}) = \text{Dir}(\vec{\pi} | \vec{\alpha}_0) = C(\vec{\alpha}_0) \cdot \prod_{k=1}^K \pi_k^{\alpha_0 - 1}$$

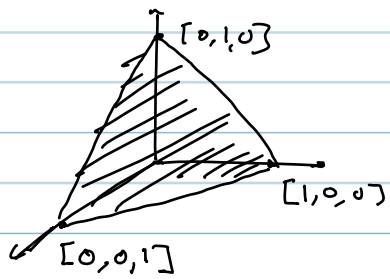
$$\vec{\alpha}_0 = [\alpha_0, \alpha_0, \dots, \alpha_0]$$

$$C(\vec{\alpha}_0) = \frac{\Gamma(K\alpha_0)}{\Gamma(\alpha_0)\Gamma(\alpha_0)\dots\Gamma(\alpha_0)}_K$$

pg 687
(B.23)

Dirichlet Dist.

Dirichlet distribution is a distribution over distributions



distribution is a vector of numbers btw 0 and 1 whose sum = 1

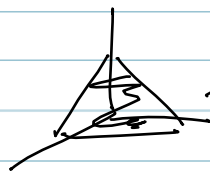
the set of all such vectors is called the simplex

3-dim prob. vectors, corners of simplex are

$$[1, 0, 0], [0, 1, 0], [0, 0, 1]$$

Simplex is a 2^d plane in 3^d

Dir. dist.



Example dist's on the simplex



$$P([1, 0, 0]) > P([0.33, 0.33, 0.33])$$

this dist sparsity encourages

Prior on class means & covariances

$$P(\mu, \Lambda) = P(\mu | \Lambda) P(\Lambda)$$

choose a Gaussian-Wishart prior governing the mean & precision matrix for each class

$$= \prod_{k=1}^K N(\mu_k | \mu_0, (\beta_0 \Lambda_k)^{-1}) W(\Lambda_k | W_0, \nu_0)$$

* Note - $\nu_0 > \text{~~3~~ } D-1$

pg 693 (bottom)

B GMM Joint dist'n

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

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

Note: only X is observed. Everything else is latent

Interpretation of Model: Unsupervised clustering / learning, discovering latent "structure" in the data (clusters), multi-modal density estimation where the latent density is unknown and complicated.

Variational Approximation ~~to GMM Posterior~~

Choose a "magic" factorization (this choice is not always obvious a priori).

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

Only necessary assumption to achieve practical results.

The functional forms of $q(Z)$ and $q(\pi, \mu, \Lambda)$ will be determined automatically.

~~Remember~~ coupled "sequential/iterative" updates can be derived by applying update eqn. for each factor

Remember 10.9

$$\ln q_j^*(z_j) = \mathbb{E}_{z_{i \neq j}} [\ln p(X, Z)] + \text{const}$$

Recipe

Apply this rule to BGMM with given factors and do lots of algebra.

Start with update rule for factor $q(z)$.

- Want to discover $q^*(z)$'s functional / distributional form.
- Want to be able to compute the parameters of this distribution by using expectations computed from other factors

$$\ln q^*(z) = \mathbb{E}_{\pi, \mu, \Lambda} [\ln p(X, z, \mu, \Lambda)] + \text{const}$$

- note as we go along, all factors that are not a function of z will be absorbed into the const.

- plug in joint dist'n

$$\begin{aligned} \ln q^*(z) &= \mathbb{E}_{\mu, \pi, \Lambda} [\ln p(z|\pi) + \ln p(X|z, \mu, \Lambda)] + \text{const} \\ &= \mathbb{E}_{\pi} [\ln p(z|\pi)] + \mathbb{E}_{\mu, \Lambda} [\ln p(X|z, \mu, \Lambda)] + \text{const} \end{aligned}$$

↑
other terms
in the joint

plug in def's of each

$$= \mathbb{E}_{\pi} \left[\sum_{n=1}^N \sum_{k=1}^K z_{nk} \ln \pi_{1k} \right] + \mathbb{E}_{\mu, \Lambda} \left[\sum_{n=1}^N \sum_{k=1}^K z_{nk} \cdot \ln \mathcal{N}(x_n | \mu_k, \Lambda_k^{-1}) \right]$$

evaluate the log under of Normal distribution

Normal dist.

$$\mathcal{N}(x | \mu, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \cdot \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

a-d (from C.13) $|A^{-1}| = \frac{1}{|A|}$

$$\ln q^*(z)$$

copy from previous page

$$= \sum_{n=1}^N \sum_{k=1}^K z_{nk} \mathbb{E}_{\pi} [\ln \pi_k] + \sum_{n=1}^N \sum_{k=1}^K z_{nk} \cdot \ln \mathcal{N}(x_n | \mu_k, \Sigma_k^{-1})$$

$$= \sum_{n=1}^N \sum_{k=1}^K z_{nk} \mathbb{E}_{\pi} [\ln \pi_k] + \text{const}$$

$$+ \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left(-\frac{D}{2} \ln(z\pi) - \frac{1}{2} \mathbb{E}_{\Sigma} [\ln |\Sigma_k^{-1}|] - \frac{1}{2} \mathbb{E}_{\mu, \Sigma} \left[(x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) \right] \right)$$

let's define

$$\ln p_{nk} = \mathbb{E}_{\pi} [\ln \pi_k] + \frac{1}{2} \mathbb{E} [\ln |\Sigma_k^{-1}|] - \frac{D}{2} \ln(z\pi) - \frac{1}{2} \mathbb{E}_{\mu, \Sigma} \left[(x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) \right]$$

← note sign change

then we can write

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

and by exponentiating

$$q^*(z) = C \cdot \prod_{n=1}^N \prod_{k=1}^K p_{nk}^{z_{nk}}$$

- How do we normalize this distribution?
- How do we compute the terms in (expectations?)

1) Normalization: $q^*(z)$ (unnormalized version) is N iid R.V.'s that are binary indicators over K states

- this means that for each n $q^*(z_n)$ can be normalized by summing over all K states, i.e.

$$\Gamma_{nk} = \sum_j p_{nj}$$

Then $q^*(z) = \prod_n \prod_k \Gamma_{nk}^{z_{nk}}$

← "distributional form"
2) computation?

Note: p_{nk} is a w exponent of a real value, so always non-negative

corollary - r_{nk} also non-negative and will sum to 1

Note: $E[z_{nk}] = 0 \cdot (1 - r_{nk}) + 1 \cdot (r_{nk}) = r_{nk}$

r_{nk} can be interpreted as a "responsibility"

- We have $q^*(z)$ which is a function of z but requires expectations of various quantities taken w.r.t. the other factors.

Other factors

Because we already know the answer, let's write down the following statistics of the observed data

$$N_k = \sum_{n=1}^N r_{nk}$$
$$\bar{x}_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n$$
$$S_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \bar{x}_k)(x_n - \bar{x}_k)^T$$

- Note these look like the per-class updates for the EM GMM algorithm.

Next $q(\pi, \mu, \Sigma)$ and remember - use general result

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

- + 1 induced factorization - will turn out to factor nicely
- 2 choice of conjugate priors will help tremendously

$$\ln q_j^*(z_j) = \mathbb{E}_{i \neq j} [\ln p(x, z)] + \text{const.}$$

remember

$$p(x, z, \pi, \mu, \Lambda) = p(x|z, \mu, \Lambda) p(z|\pi) p(\pi) p(\mu, \Lambda) p(\Lambda)$$

Applying this for $q^*(\pi, \mu, \Lambda)$

$$\begin{aligned} \ln q^*(\pi, \mu, \Lambda) &= \mathbb{E}_{q(z)} [\ln (p(x|z, \mu, \Lambda) p(z|\pi) p(\pi) p(\mu, \Lambda))] + \text{const.} \\ &= \underbrace{\mathbb{E}_z [\ln (p(z|x, \mu, \Lambda))]}_{\text{involve only } z, \mu, \Lambda} + \underbrace{\mathbb{E}_z [\ln p(z|\pi)] + \mathbb{E}_z [\ln p(\pi)] + \mathbb{E}_z [\ln p(\mu, \Lambda)]}_{\text{involve only } \pi \neq z} + \text{const.} \end{aligned}$$

$$\Rightarrow q(\pi, \mu, \Lambda) = q(\pi) q(\mu, \Lambda) \quad (\text{further factorization})$$

and noting that $p(\mu, \Lambda) = \prod_k p(\mu_k, \Lambda_k)$

$$= p(z|x, \mu, \Lambda) = \prod_k \prod_{\lambda} \pi_{\lambda k} z_{\lambda k}$$

$$\begin{aligned} \Rightarrow q(\pi, \mu, \Lambda) &= q(\pi) q(\mu, \Lambda) \\ &= q(\pi) \prod_{k=1}^K q(\mu_k, \Lambda_k) \end{aligned}$$

- net effect: can optimize each of these factors independently

To start:

$$q^*(\pi) = \mathbb{E}_z [\ln p(z|\pi)] + \ln p(\pi) + \text{const}$$

Remember Dir distribution

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

and

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

Plug these into $\ln q^*(\pi) = \mathbb{E}_z [\ln p(z | \pi)] + \ln p(\pi) + \text{const}$

$$\ln q^*(\pi) = (\alpha_0 - 1) \cdot \sum_{k=1}^K \ln \pi_k + \sum_{n=1}^N \sum_{k=1}^K \mathbb{E}[z_{nk}] \cdot \ln \pi_k + \text{const}$$

rearrange

$$\ln q^*(\pi) = \sum_{k=1}^K \left(\alpha_0 - 1 + \sum_{n=1}^N \mathbb{E}[z_{nk}] \right) \ln \pi_k + \text{const}$$

now remember that $\mathbb{E}[z_{nk}] = r_{nk}$

and that $N_k = \sum_{n=1}^N r_{nk}$

so

$$\ln q^*(\pi) = \sum_{k=1}^K (\alpha_0 - 1 + N_k) \ln \pi_k + \text{const}$$

$$\Rightarrow q^*(\pi) = C \cdot \prod_{k=1}^K \pi_k^{(\alpha_0 - 1 + N_k)}$$

$$\Rightarrow q^*(\pi) = \text{Dir}(\pi | \vec{\alpha})$$

$$\vec{\alpha} = [\alpha_0 + N_1, \alpha_0 + N_2, \dots, \alpha_0 + N_K]$$

examine carefully,
this is a Dir dist
- hint: this is true
because of co-j. choice

To compute the responsibilities we need to compute

$$\mathbb{E}_{q^*(\pi)} [\ln \pi_k] = ?$$

This is given in PRML in B page 687

$$\mathbb{E}[\ln \pi_k] = \Psi(\alpha_0 + N_k) - \Psi\left(\alpha_0 + \sum_{k=1}^K N_k\right)$$

Remember $\sum_{k=1}^K N_k = N$

aside
 $\psi(a) = \frac{d}{da} \Gamma(a)$ is called the "digamma" function (known as psi in matlab)
 a lot like $\log(x)$

In order to complete the VB approximate inference algorithm for GMMs we need the distributional form of $q^*(\mu, \Sigma) = \prod_k q^*(\mu_k, \Sigma_k)$.

This involves a large amount of algebra, but is similar in spirit to the last two examples.

To make further factorization obvious

$$\ln(q^*(\mu, \Sigma)) = \sum_k \mathbb{E}_z \left[\mathbb{E}_x \left[\ln p(\mu_k, \Sigma_k) \right] \right] + \sum_k \sum_n \mathbb{E}_z \left[\mathbb{E}_x \left[\ln N(x_n | \mu_k, \Sigma_k^{-1})^{z_{nk}} \right] \right] + \text{const}$$

$$= \underbrace{\sum_k \ln p(\mu_k, \Sigma_k) + \sum_k \sum_n \mathbb{E}_z \left[z_{nk} \right] \ln N(x_n | \mu_k, \Sigma_k^{-1})}_{\text{sums over } k \text{ i.i.d. products and } k \text{ independent factors } q(\mu_k, \Sigma_k)}$$

$$\ln(q^*(\mu_k, \Sigma_k)) = \ln P(\mu_k, \Sigma_k) + \sum_n \mathbb{E}_{q^*(z)} [z_{nk}] \ln N(x_n | \mu_k, \Sigma_k^{-1})$$

is nasty algebraically, but with experience in conjugate families, one can almost guess the answer

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

where

$$\beta_k = \beta_0 + N_k$$

$$m_k = \frac{1}{\beta_k} (\beta_0 m_0 + N_k \bar{x}_k)$$

$$W_k^{-1} = W_0^{-1} + N_k S_k + \frac{\beta_0 N_k}{\beta_0 + N_k} (\bar{x}_k - m_0)(\bar{x}_k - m_0)^T$$

$$\nu_k = \nu_0 + N_k + 1$$

Now we have (i-boxes) all distributional forms for VB updates factors.

VB updates to factor parameters require deriving or computing

$$\mathbb{E}_{\Lambda_k} [\ln | \Lambda_k |] \neq \mathbb{E}_{\mu_k, \Lambda_k} [(x_n - \mu_k)^T \Lambda_k (x_n - \mu_k)]$$

$$\downarrow 10.65 = \sum_{i=1}^D 4 \left(\frac{\nu_k + 1 - i}{2} \right) + D \cdot \ln 2 + \ln |W_k|$$

$$\downarrow 10.64 = D \beta_k^{-1} + \nu_k (x_n - m_k)^T W_k (x_n - m_k)$$

Now we have everything to compute the $q^*(z)$, and given $q^*(z)$ we can update the parameters for $q^*(\pi) = q^*(\mu_k, \Lambda_k) \forall k$.

Running this until convergence will produce an approximate posterior

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

with parameters given by formulae in boxes.

VB for GMM's

- Note: - Coupled update equations roughly correspond to the E-M steps of EM for GMM's
- Functional form of variational factors is a consequence of the choice of conjugate priors

Interesting aside: Model selection: α_0 controls "sparsity" of resulting model, $\alpha_0 < 1$ prefers "sparse" (low component count) models (i.e. posterior distribution over models places high score on models with few components)

Under the VB GMM the expected value of the class probabilities is

$$\mathbb{E}[\pi_k] = \frac{\alpha_k + N_k}{K\alpha_0 + N}$$

if $N_k \approx 0$ then as $N \rightarrow \infty$ $\mathbb{E}[\pi_k] \rightarrow 0$ if α_0 is small. If $\alpha_0 \rightarrow \infty$ then $\alpha_k \approx \alpha_0$ and as N gets big $\mathbb{E}[\pi_k] \rightarrow \frac{1}{K}$

Bayesian treatment	ML, EM
No singularities	Singularities
A <u>little</u> computational overhead	Computationally minimal

VB
model selection w/o cross validation