

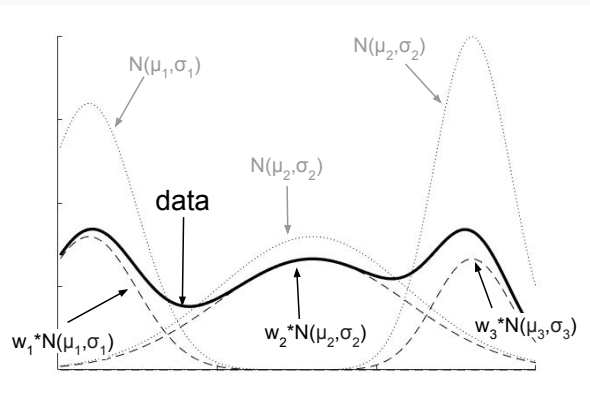
K-means, Gaussian Mixture Models, UBM-GMM

Pattern Recognition, Jul-Nov 2019

Indian Institute of Technology Madras

September 6, 2019

Gaussian mixture models (GMM)



- ▶ A GMM is the weighted sum of individual Gaussian distributions

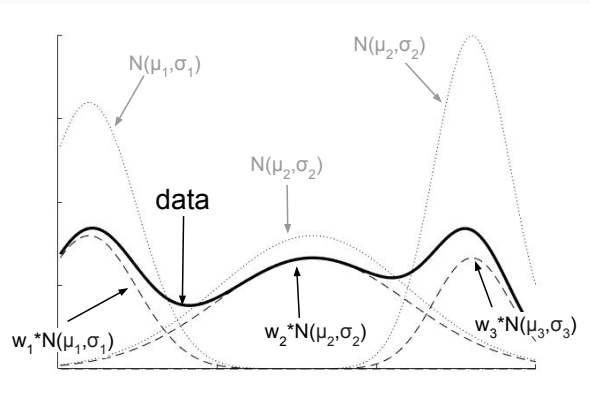
$$P(\vec{x}|\lambda_s) = \sum_{k=1}^M w_k \mathcal{N}(\vec{x}|\vec{\mu}_k, \Sigma_k)$$

$$\lambda_s = \{w_k, \vec{\mu}_k, \Sigma_k\}_{k=1}^M$$

$$\sum_{k=1}^K w_k = 1$$

$$\mathcal{N}(\vec{x}|\vec{\mu}_k, \Sigma_k) = (2\pi)^{-d/2} |\Sigma_k|^{-1/2} \exp \left\{ -\frac{1}{2} (\vec{x} - \vec{\mu}_k)^T \Sigma_k^{-1} (\vec{x} - \vec{\mu}_k) \right\}$$

Gaussian mixture models (GMM)



- ▶ The problem of fitting a GMM is an incomplete data problem. Hence, the mixture needs to be estimated iteratively using Expectation Maximization (E-M) algorithm.

Parameter estimation of GMM using E-M algorithm

- ▶ Estimate parameters $(\bar{\mu}_k, \Sigma_k)$.
- ▶ Define one more quantity 'responsibility'.

$$P(z_k = 1|x) = \frac{P(z_k = 1)P(\bar{x}|z_k = 1)}{\sum_{j=1}^K P(\bar{x}|z_j)P(z_j = 1)} \quad (1)$$

$P(z_k = 1|x) = \gamma_k$, which is the *responsibility* of k_{th} mixture in describing a point x .

$$\begin{aligned} P(z_k = 1|x) &= \frac{P(z_k = 1)P(\bar{x}|z_k = 1)}{\sum_{j=1}^K P(\bar{x}|z_j)P(z_j = 1)} \\ &= \frac{\pi_k N(\bar{x}; \bar{\mu}_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(\bar{x}; \bar{\mu}_j, \Sigma_j)} \end{aligned}$$

$$N(\bar{x}, \bar{\mu}_j, \Sigma_j) = \frac{-1}{(\sqrt[2]{2\pi})^d |\Sigma|} e^{(\frac{-1}{2}(\bar{x} - \bar{\mu}_j)^t \Sigma^{-1}(\bar{x} - \bar{\mu}_j))}$$

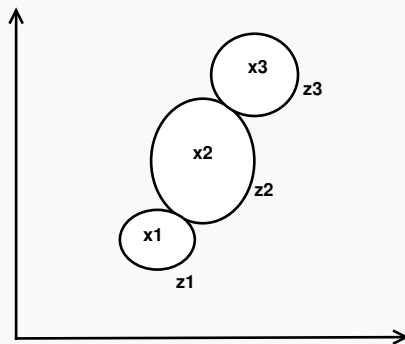
A single point will be described completely by all the mixtures together.

$$\sum_{k=1}^K \gamma_k = 1$$

$$\gamma_{nk} = P(z_k = 1 | \bar{x}_n)$$

$$\sum_{n=1}^N \gamma_{nk} = N_k$$

N_k is the effective number of points that belongs to k^{th} cluster



$$\gamma(z_1 = 1 | \bar{x}_2) = .125, \quad \gamma(z_1 = 1 | \bar{x}_1) = .75, \quad \gamma(z_1 = 1 | \bar{x}_3) = .05$$

$$\gamma(z_2 = 1 | \bar{x}_2) = .75, \quad \gamma(z_1 = 1 | \bar{x}_1) = .20, \quad \gamma(z_2 = 1 | \bar{x}_3) = .20$$

$$\gamma(z_3 = 1 | \bar{x}_2) = .125, \quad \gamma(z_1 = 1 | \bar{x}_1) = .05, \quad \gamma(z_3 = 1 | \bar{x}_3) = .75$$

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^{N_k} \gamma_{nk} \bar{x}_n$$

Estimation of parameters

$$\theta_{ML} = \arg \max_{\theta} \log p(D|\theta)$$

$$\bar{\theta} = [\bar{\mu}_1, \bar{\mu}_2, \dots, \bar{\mu}_k, \Sigma_1, \Sigma_2, \dots, \Sigma_k, \pi_1, \pi_2, \dots, \pi_k]$$

$$\ln P(D|\theta) = \ln \prod_{n=1}^N P(\bar{x}_n|\theta) = l(\theta)$$

$$\ln P(D|\theta) = \sum_{n=1}^N \ln \sum_{j=1}^K \pi_j N(\bar{x}_n; \bar{\mu}_j, \Sigma_j)$$

$$\frac{\partial l(\bar{\theta})}{\partial \bar{\theta}} = 0$$

$$\frac{\partial l(\bar{\theta})}{\partial \bar{\mu}_k} = 0; \frac{\partial l(\bar{\theta})}{\partial \bar{\Sigma}_k} = 0$$

π_k requires a constraint $\sum_{k=1}^K \pi_k = 1$

$$\begin{aligned}\frac{\partial l(\bar{\theta})}{\partial \bar{\mu}_k} &= \sum_{n=1}^N \frac{\partial}{\partial \bar{\mu}_k} \ln \sum_{j=1}^K \pi_j N(\bar{x}_n; \bar{\mu}_j, \Sigma_j) \\ &= \sum_{n=1}^N \frac{\pi_k}{\sum_{j=1}^K \pi_j N(\bar{x}_n; \bar{\mu}_j, \Sigma_j)} \frac{\partial}{\partial \bar{\mu}_k} N(\bar{x}_n; \bar{\mu}_k, \Sigma_k) \\ \frac{\partial}{\partial \bar{\mu}_k} N(\bar{x}_n; \bar{\mu}_k, \Sigma_k) &= \frac{-1}{(\sqrt[2]{2\pi})^d |\Sigma|^{1/2}} e^{(\frac{-1}{2}(\bar{x}_n - \bar{\mu}_k)^t \Sigma_k^{-1} (\bar{x}_n - \bar{\mu}_k))}\end{aligned}$$

$$\frac{\partial}{\partial \bar{x}} \bar{x}^t M \bar{x} = 2M \bar{x}$$

$$\begin{aligned}\frac{\partial l(\bar{\theta})}{\partial \bar{\mu}_k} &= -N(\bar{x}_n; \bar{\mu}_k, \Sigma_k) \Sigma_k^{-1} (\bar{x}_n - \bar{\mu}_k) \\ &= \sum_{n=1}^N \frac{\pi_k N(\bar{x}; \bar{\mu}_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(\bar{x}; \bar{\mu}_j, \Sigma_j)} \\ \frac{\partial}{\partial \bar{\mu}_k} = 0 &\implies \sum_{n=1}^N \gamma_{nk} \Sigma_k^{-1} (\bar{x}_n - \bar{\mu}_k) = 0 \\ \sum_{n=1}^N \gamma_{nk} \bar{x}_n &= \sum_{n=1}^N \gamma_{nk} \bar{\mu}_k\end{aligned}$$

$$\bar{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} \bar{x}_n$$

To find Σ_k , take derivative with respect to Σ_k and equate it to 0.

$$\bar{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} (\bar{x}_n - \hat{\mu}_k)(\bar{x}_n - \hat{\mu}_k)^t$$

Identities

- ▶ $\frac{\partial |A|}{\partial A} = |A| A^{-1}$
- ▶ $\frac{\partial \bar{\mu}^t A^{-1} \bar{\mu}}{\partial A} = A^{-1} \bar{\mu} \bar{\mu}^t A^{-1}$
- ▶ $\frac{\partial (\bar{u} \cdot \bar{v})}{\partial \bar{x}} = \bar{v}^t \frac{d \bar{u}}{d \bar{x}} + \bar{u}^t \frac{d \bar{v}}{d \bar{x}}$

Estimation of π_k such that $\sum_{j=1}^K \pi_j = 1$

This is a constraint optimization problem.

$$L(\pi_k, \lambda) = \sum_{n=1}^K \ln \sum_{j=1}^K \pi_j N(\bar{x}_n; \bar{\mu}_j, \Sigma_j, \pi_j) - \lambda \left(\sum_{j=1}^K \pi_j - 1 \right)$$

$$\frac{\partial L(\pi_k, \lambda)}{\partial \pi_k} = \sum_{n=1}^N \frac{\pi_k N(\bar{x}; \bar{\mu}_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(\bar{x}; \bar{\mu}_j, \Sigma_j)}$$

$$\lambda \pi_k = \sum_{n=1}^N \gamma_{nk}$$

$$\lambda \sum_{k=1}^K \pi_k = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk}$$

$$\pi_k = \frac{N_k}{N}$$

GMM : EM algorithm (Recap)

1) Random initialization of $\vec{\mu}_k, \Sigma_k$ and w_k

2) Expectation-Step

Align vectors to model

$$\gamma_{nk} = \frac{w_k \mathcal{N}(\vec{x}_n | \vec{\mu}_k, \Sigma_k)}{\sum_{j=1}^M w_j \mathcal{N}(\vec{x}_n | \vec{\mu}_j, \Sigma_j)}$$

$$N_k = \sum_{n=1}^N \gamma_{nk}$$

→ γ_{nk} is the responsibility of k -th component towards n -th feature vector

3) Maximization-Step

Update model parameters by maximum likelihood estimation (MLE)

$$\hat{\vec{\mu}}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} \vec{x}_n$$

$$\hat{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} (\vec{x}_n - \hat{\vec{\mu}}_k)(\vec{x}_n - \hat{\vec{\mu}}_k)^T$$

$$\hat{w}_k = \frac{N_k}{N}$$

4) Repeat step 2 and 3 until convergence.

GMM-proof of convergence

- ▶ Let $\bar{\Theta}^{\text{old}}$ be parameters used at the start of any EM-iteration and $\bar{\Theta}$ be the updated parameters.
- ▶ Let \mathcal{D} be the given set of data points and \bar{z} be the latent variable defined by $\bar{\Theta}^{\text{old}}$.
- ▶ We need to prove that $P(\mathcal{D} \mid \bar{\Theta}) \geq P(\mathcal{D} \mid \bar{\Theta}^{\text{old}})$ for every iteration.
- ▶ In every iteration we need to maximize $E[\log P(\mathcal{D}, \bar{z} \mid \bar{\Theta}) \mid \mathcal{D}, \bar{\Theta}^{\text{old}}]$
- ▶ Lets define the auxiliary function as follows:

$$\begin{aligned} A(\bar{\Theta}, \bar{\Theta}^{\text{old}}) &= \sum_z P(\bar{z} \mid \mathcal{D}, \bar{\Theta}^{\text{old}}) \log P(\mathcal{D}, \bar{z} \mid \bar{\Theta}) \\ &= \sum_z P(\bar{z} \mid \mathcal{D}, \bar{\Theta}^{\text{old}}) \log (P(\bar{z} \mid \mathcal{D}, \bar{\Theta}) P(\mathcal{D} \mid \bar{\Theta})) \\ &= \sum_z P(\bar{z} \mid \mathcal{D}, \bar{\Theta}^{\text{old}}) \log (P(\bar{z} \mid \mathcal{D}, \bar{\Theta})) \\ &\quad + \sum_z P(\bar{z} \mid \mathcal{D}, \bar{\Theta}^{\text{old}}) \log (P(\mathcal{D} \mid \bar{\Theta})) \end{aligned}$$

GMM-proof of convergence (Contd..)

- ▶ Substituting $\bar{\Theta}^{\text{old}}$ for $\bar{\Theta}$ in the previous equation we get

$$\begin{aligned} A(\bar{\Theta}^{\text{old}}, \bar{\Theta}^{\text{old}}) &= \sum_z P(\bar{z} | \mathcal{D}, \bar{\Theta}^{\text{old}}) \log(P(\bar{z} | \mathcal{D}, \bar{\Theta}^{\text{old}})) \\ &\quad + \sum_z P(\bar{z} | \mathcal{D}, \bar{\Theta}^{\text{old}}) \log(P(\mathcal{D} | \bar{\Theta}^{\text{old}})) \end{aligned}$$

- ▶ Now $\log P(\mathcal{D} | \bar{\Theta}) - \log P(\mathcal{D} | \bar{\Theta}^{\text{old}})$ can be written as:

$$\begin{aligned} \log P(\mathcal{D} | \bar{\Theta}) - \log P(\mathcal{D} | \bar{\Theta}^{\text{old}}) &= A(\bar{\Theta}, \bar{\Theta}^{\text{old}}) - A(\bar{\Theta}^{\text{old}}, \bar{\Theta}^{\text{old}}) \\ &\quad + \sum_z P(\bar{z} | \mathcal{D}, \bar{\Theta}^{\text{old}}) \log \left(\frac{P(\bar{z} | \mathcal{D}, \bar{\Theta}^{\text{old}})}{P(\bar{z} | \mathcal{D}, \bar{\Theta})} \right) \end{aligned} \quad (2)$$

GMM-proof of convergence (Contd..)

In Equation 2

- ▶ The last part of the equation is *Kullback-Leibler* divergence which is always positive or null.
- ▶ The update equation are derived such that the A increases (See <https://courses.iitm.ac.in/mod/resource/view.php?id=789> for derivation based on auxiliary function)

Therefore, the likelihood always increases from iteration to iteration.

UBM-GMM

Universal background model (UBM)

- ▶ UBM is a GMM trained on huge data pooled together from all the available classes.
- ▶ To overcome the huge data requirement for training GMM for individual classes.
- ▶ UBM is a GMM built with all class data using MLE algorithm.
- ▶ Maximum A-Posteriori (MAP) adaptation is used to train target models.
- ▶ UBM does not discriminate different speaker but acts as a reference for all class models.

UBM-GMM : MAP adaptation (contd...)

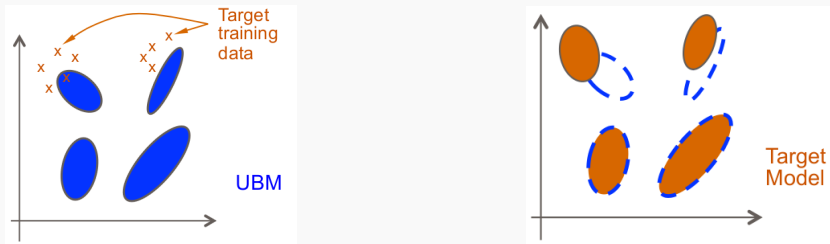


Figure 1: Adaptation using MAP

- Target model mean is updated using sufficient statistics, mixing co-efficient (α_k) and relevance factor r as

$$\alpha_k = \frac{N_k}{N_k + r}$$

$$\hat{\vec{\mu}}_k^{new} = \alpha_k \hat{\vec{\mu}}_k + (1 - \alpha_k) \vec{\mu}_k^{ubm}$$

UBM-GMM : MAP adaptation (contd...)

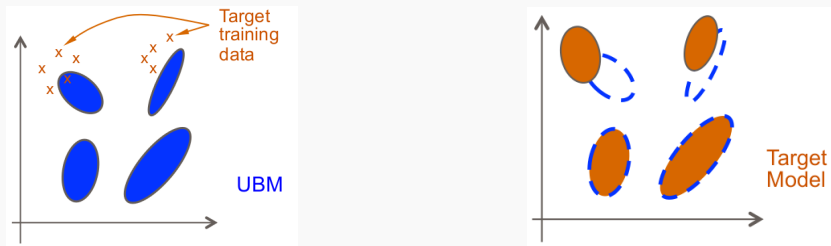


Figure 2: Adaptation using MAP

- The covariance can also be adapted using the same formule. But for most of the application we limit ourselves to adapting the mean.

Scoring in UBM-GMM

UBM acts as a reference for all class models. Hence the scores are calculated with reference to the UBM.

Likelihood scoring (in the context of speech)

The average log-likelihood ratio score for a test utterance $\mathcal{X} = \{\bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2, \dots, \bar{\mathbf{x}}_T\}$ with claim is calculated as

$$\text{LogLikelihood}_{\text{avg}}(\mathcal{X}, \lambda_{\text{claim}}, \lambda_{\text{UBM}}) = \frac{1}{T} \sum_{t=1}^T \{ \log P(\bar{\mathbf{x}}_t | \lambda_{\text{claim}}) - \log P(\bar{\mathbf{x}}_t | \lambda_{\text{UBM}}) \}$$

Advantages of MAP in UBM-GMM

- ▶ Mixtures of UBM and target GMM have one-one correspondence
- ▶ Only few mixture components contribute to the a particular class's feature vectors
- ▶ $\text{LogLikelihood}_{avg}$ can be computed with maximum contributing C mixture components of the UBM.