

A tutorial of Hidden Markov Model (HMM)

Nov. 2015

Qiuqiang Kong

q.kong@qmul.ac.uk

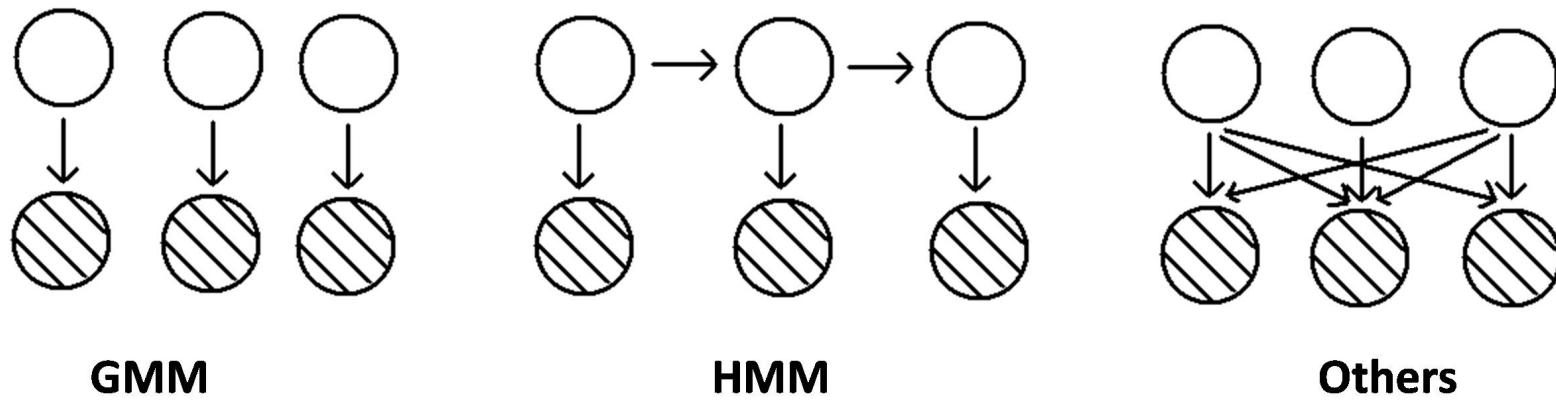
Application of HMM

- Stock price analysis
- Auto speech recognition
- Character recognition
- Sequence alignment
- etc.

Probabilistic Graphical Model (PGM)

PGM is a probabilistic model. The graph expresses the conditional dependence structure between random variables.

Example:



White circle represents latent variable.
Solid circle represents observed variable.
All of GMM, HMM, Others belongs to PGM.

EM algorithm for PGM

EM algorithm is used to estimate parameters in probabilistic graphic model (PGM) with latent variables*.

EM algorithm for PGM with latent variables

1. Init parameters
2. E step: $q(\mathbf{Z}) = p(\mathbf{Z} | \mathbf{X}\theta^{\text{old}})$
3. M step: $\hat{\theta} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}})$
where $Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}\theta^{\text{old}}) \ln p(\mathbf{XZ} | \theta)$
4. If converge then stop, otherwise goto 2

* More details can be seen in GMM Tutorial.

Generally, estimation of $p(\mathbf{Z}|\mathbf{X}\theta^{\text{old}})$ is difficult.

However, the independence of PGM will simplify the model.

D-separation property^[1]

Let A, B, C be set of nodes. We can check whether A and B are conditional independent on C in belowing way.

Check all the paths from node in A to node C. The path is said to be blocked if

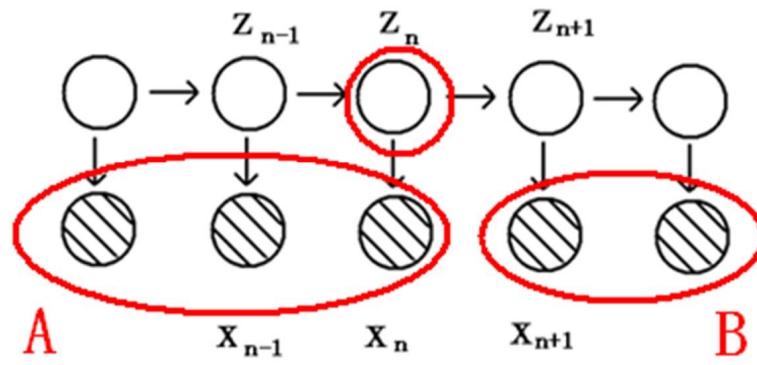
- (a) \exists head-to-tail or tail-to-tail nodes on the path, and the node is in C.
- (b) \forall head-to-head nodes, neither the node, nor any of its descendants is in C

If all paths are blocked, then A is said to be d-separated from B by C.

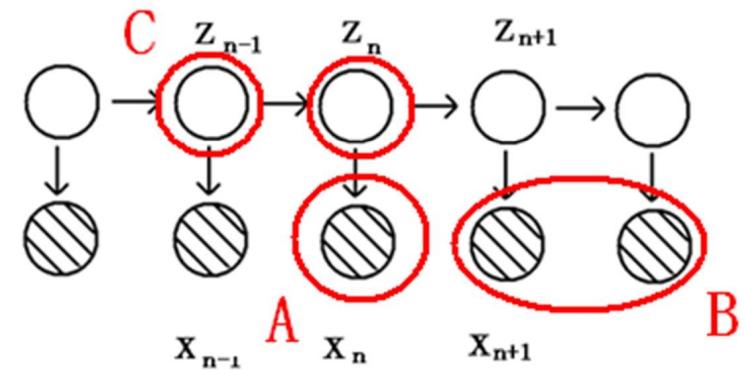
That is $P(AB|C) = P(A|C)P(B|C)$.

[1] Bishop, Christopher M. *Pattern recognition and machine learning*. Springer, 2006. Chap. 8.

Exercises



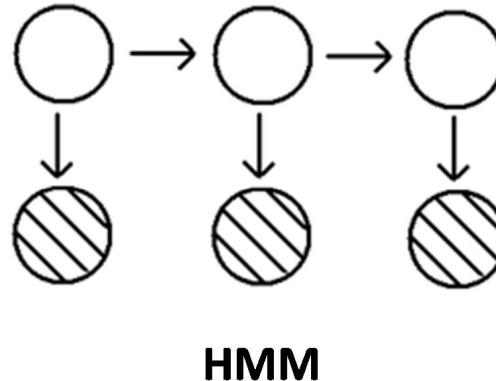
$$p(AB|\mathbf{z}_n) = p(A|\mathbf{z}_n)p(B|\mathbf{z}_n)$$



$$p(ABC|\mathbf{z}_n) = p(A|\mathbf{z}_n)p(B|\mathbf{z}_n)p(C|\mathbf{z}_n)$$

HMM model

Hidden Markov Model (HMM) is a kind of PGM with latent variables.



HMM's **joint probability distribution** over both latent and observed variables is

$$p(\mathbf{XZ}|\theta) = p(\mathbf{z}_1|\boldsymbol{\pi}) \prod_{n=2}^N p(\mathbf{z}_n | \mathbf{z}_{n-1}, \mathbf{A}) \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n, \phi) \quad (1)$$

where the parameters are $\theta = \{\boldsymbol{\pi}, \mathbf{A}, \phi\}$

$\boldsymbol{\pi}$ is the start probability of each state.

\mathbf{A} is the transition matrix.

ϕ is the parameter associated with emission distribution (can be multinomial, Gauss, GMM, etc.)

Three basic problem of HMM

- **How to train a HMM model?**
- **Given a sequence X, how to get the likelihood $p(X)$ from HMM model?**
- **How to find the best decoding path of HMM model?**

How to train HMM model?

HMM is a kind of probabilistic graphic model (PGM) with latent variables.

Apply EM algorithm to HMM, $Q(\theta, \theta^{\text{old}})$ is

$$\begin{aligned}
 Q(\theta, \theta^{\text{old}}) &= \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}\theta^{\text{old}}) \ln p(\mathbf{XZ} | \theta) \\
 &= \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}\theta^{\text{old}}) \left[\ln p(\mathbf{z}_1) + \sum_{n=2}^N \ln p(\mathbf{z}_n | \mathbf{z}_{n-1}) + \sum_{n=1}^N \ln p(\mathbf{x}_n | \mathbf{z}_n) \right] \\
 &= \sum_{\mathbf{z}_1} p(\mathbf{z}_1 | \mathbf{X}\theta^{\text{old}}) \ln p(\mathbf{z}_1) + \sum_{n=2}^N \sum_{\mathbf{z}_{n-1}, \mathbf{z}_n} p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}\theta^{\text{old}}) \ln p(\mathbf{z}_n | \mathbf{z}_{n-1}) + \sum_{n=1}^N \sum_{\mathbf{z}_n} p(\mathbf{z}_n | \mathbf{X}\theta^{\text{old}}) \ln p(\mathbf{x}_n | \mathbf{z}_n) \\
 &= \underbrace{\sum_{k=1}^K \gamma(z_{nk}) \ln \pi_k}_{\textcircled{1}} + \underbrace{\sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \xi(z_{n-1,j} z_{n,k}) \ln A_{jk}}_{\textcircled{2}} + \underbrace{\sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \ln p(\mathbf{x}_n | z_{nk} \phi)}_{\textcircled{3}} \quad (2)
 \end{aligned}$$

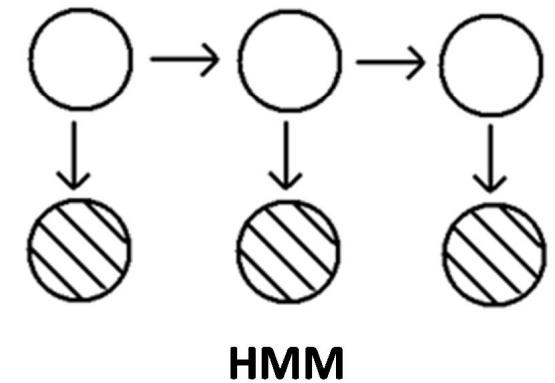
$$\text{Where } \gamma(z_{nk}) = p(z_{nk} | \mathbf{X}\theta^{\text{old}}) \quad (3)$$

$$\xi(z_{n-1,j} z_{n,k}) = p(z_{n-1,j} z_{n,k} | \mathbf{X}\theta^{\text{old}}) \quad (4)$$

π is only associated with $\textcircled{1}$

A is only associated with $\textcircled{2}$

ϕ is only associated with $\textcircled{3}$



E step

Evaluate $p(\mathbf{Z} | \mathbf{X}\theta^{\text{old}})$.

From the form of $Q(\theta, \theta^{\text{old}})$ of HMM, there is no need to calculate $p(\mathbf{Z} | \mathbf{X}\theta^{\text{old}})$ for all \mathbf{Z} . Many terms vanish. So just need to calculate

$$\gamma(z_{nk}) = p(z_{nk} | \mathbf{X}\theta^{\text{old}})$$

$$\xi(z_{n-1,j} z_{n,k}) = p(z_{n-1,j} z_{n,k} | \mathbf{X}\theta^{\text{old}})$$

Generally, $p(\mathbf{Z} | \mathbf{X}\theta^{\text{old}})$ is difficult to estimate in PGM. While using dependency of HMM, we can simplify the computation.

We will show $\gamma(z_{nk})$ can be decomposed to forward and backward term.

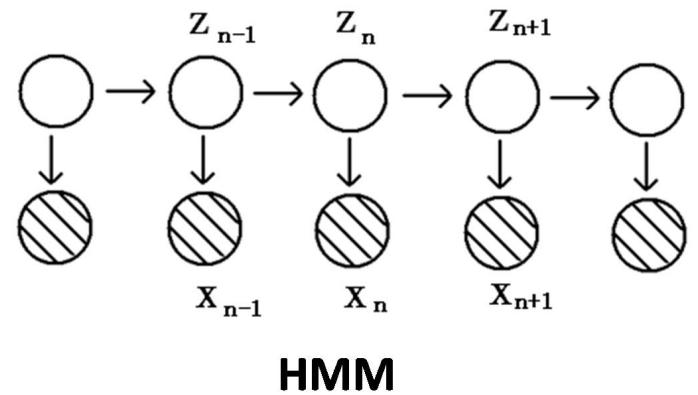
$$\gamma(\mathbf{z}_n) = \frac{\hat{\alpha}(\mathbf{z}_n)}{\text{forward}} \frac{\hat{\beta}(\mathbf{z}_n)}{\text{backward}} \quad (5)$$

Both of forward and backward term can be computed efficiently.

Forward Backward algorithm

Using independence of PGM, we can decompose $\gamma(\mathbf{z}_n | \mathbf{X})$ into forward and backward term.

$$\begin{aligned}
 \gamma(\mathbf{z}_n) &= p(\mathbf{z}_n | \mathbf{X}) \\
 &= \frac{p(\mathbf{X} | \mathbf{z}_n) p(\mathbf{z}_n)}{p(\mathbf{X})} \\
 &= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n | \mathbf{z}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n) p(\mathbf{z}_n)}{p(\mathbf{x}_1, \dots, \mathbf{x}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{x}_1, \dots, \mathbf{x}_n)} \\
 &= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n)}{p(\mathbf{x}_1, \dots, \mathbf{x}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{x}_1, \dots, \mathbf{x}_n)} \\
 &= \hat{\alpha}(\mathbf{z}_n) \hat{\beta}(\mathbf{z}_n)
 \end{aligned} \tag{6}$$



where $\hat{\alpha}(\mathbf{z}_n) = \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n)}{p(\mathbf{x}_1, \dots, \mathbf{x}_n)}$ (Forward) (7)

$$\hat{\beta}(\mathbf{z}_n) = \frac{p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n)}{p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{x}_1, \dots, \mathbf{x}_n)} \quad (\text{Backward}) \quad (8)$$

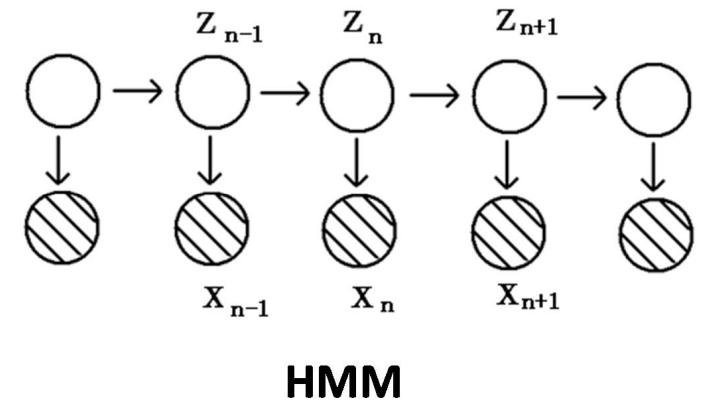
Compute $\hat{\alpha}(\mathbf{z}_n)$

$\hat{\alpha}(\mathbf{z}_n)$ Can be computed efficiently using dependency of HMM

$$\text{denote } c_n = p(\mathbf{x}_n | \mathbf{x}_1, \dots, \mathbf{x}_{n-1}) \quad (9)$$

$$\hat{\alpha}(\mathbf{z}_1) = p(\mathbf{z}_1 | \mathbf{x}_1) = \frac{p(\mathbf{x}_1 | \mathbf{z}_1)p(\mathbf{z}_1)}{c_1} \quad (10)$$

$$\begin{aligned} \hat{\alpha}(\mathbf{z}_n) &= p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n) / p(\mathbf{x}_1, \dots, \mathbf{x}_n) \\ &= p(\mathbf{x}_1, \dots, \mathbf{x}_n | \mathbf{z}_n)p(\mathbf{z}_n) / p(\mathbf{x}_1, \dots, \mathbf{x}_n) \\ &= p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_n)p(\mathbf{x}_n | \mathbf{z}_n) / p(\mathbf{x}_1, \dots, \mathbf{x}_n) \\ &= \sum_{\mathbf{z}_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_{n-1}, \mathbf{z}_n)p(\mathbf{x}_n | \mathbf{z}_n) / p(\mathbf{x}_1, \dots, \mathbf{x}_n) \\ &= \sum_{\mathbf{z}_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_n | \mathbf{z}_{n-1})p(\mathbf{z}_{n-1})p(\mathbf{x}_n | \mathbf{z}_n) / p(\mathbf{x}_1, \dots, \mathbf{x}_n) \\ &= \sum_{\mathbf{z}_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{z}_{n-1})p(\mathbf{z}_n | \mathbf{z}_{n-1})p(\mathbf{z}_{n-1})p(\mathbf{x}_n | \mathbf{z}_n) / (p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1})p(\mathbf{x}_n | \mathbf{x}_1, \dots, \mathbf{x}_{n-1})) \\ &= p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{z}_{n-1})p(\mathbf{z}_n | \mathbf{z}_{n-1}) / (p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1})c_n) \\ &= \frac{1}{c_n} p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} \hat{\alpha}(\mathbf{z}_{n-1})p(\mathbf{z}_n | \mathbf{z}_{n-1}) \end{aligned} \quad (11)$$



where $c_1 = \sum_{\mathbf{z}_1} p(\mathbf{x}_1 | \mathbf{z}_1) p(\mathbf{z}_1)$ (12)

$$c_n = \sum_{\mathbf{z}_n} \left[p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} \hat{\alpha}(\mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1}) \right] \quad (\text{integrate both side of (11)}) \quad (13)$$

*Byproduct: the likelihood of a sequence \mathbf{X} : $p(\mathbf{X}) = \prod_{n=1}^N c_n$ (14)

Compute $\hat{\beta}(\mathbf{z}_n)$

$$\begin{aligned} \hat{\beta}(\mathbf{z}_n) &= p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n) / p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{x}_1, \dots, \mathbf{x}_n) \\ &= \sum_{\mathbf{z}_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_{n+1}) / \left[p(\mathbf{x}_{n+2}, \dots, \mathbf{x}_N | \mathbf{x}_1, \dots, \mathbf{x}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{x}_1, \dots, \mathbf{x}_n) \right] \\ &= \sum_{\mathbf{z}_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1}) / \left[p(\mathbf{z}_n) / \left[p(\mathbf{x}_{n+2}, \dots, \mathbf{x}_N | \mathbf{x}_1, \dots, \mathbf{x}_{n+1}) c_{n+1} \right] \right] \\ &= \sum_{\mathbf{z}_{n+1}} p(\mathbf{x}_{n+2}, \dots, \mathbf{x}_N | \mathbf{z}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}) p(\mathbf{z}_n | \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1}) / \left[p(\mathbf{z}_n) / \left[p(\mathbf{x}_{n+2}, \dots, \mathbf{x}_N | \mathbf{x}_1, \dots, \mathbf{x}_{n+1}) c_{n+1} \right] \right] \\ &= \frac{1}{c_{n+1}} \sum_{\mathbf{z}_{n+1}} \hat{\beta}(\mathbf{z}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n) \end{aligned} \quad (15)$$

where $\hat{\beta}(\mathbf{z}_n) = \frac{\gamma(\mathbf{z}_n)}{\hat{\alpha}(\mathbf{z}_n)} = \frac{p(\mathbf{z}_n | \mathbf{X})}{p(\mathbf{z}_n | \mathbf{X})} = 1$ (16)

Compute $\xi(\mathbf{z}_{n-1}\mathbf{z}_n)$

Using conditional dependency of HMM

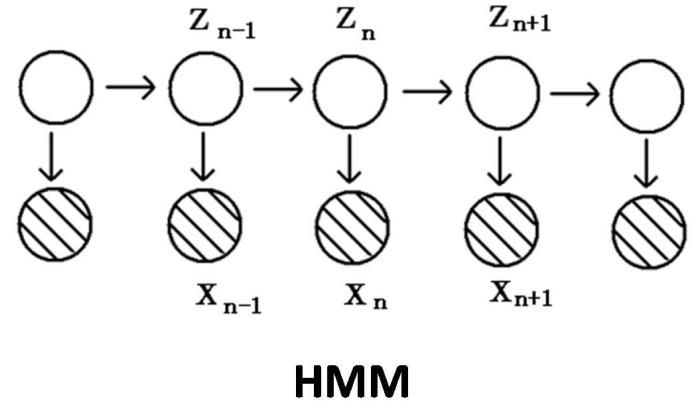
$$\xi(\mathbf{z}_{n-1}\mathbf{z}_n) = p(\mathbf{z}_{n-1}\mathbf{z}_n | \mathbf{X})$$

$$= \frac{p(\mathbf{X} | \mathbf{z}_{n-1}\mathbf{z}_n) p(\mathbf{z}_{n-1}\mathbf{z}_n)}{p(\mathbf{X})}$$

$$= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{z}_n) p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n) p(\mathbf{z}_{n-1}\mathbf{z}_n)}{p(\mathbf{X})}$$

$$= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_n) / p(\mathbf{z}_n) p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n) p(\mathbf{z}_{n-1}\mathbf{z}_n)}{p(\mathbf{X})}$$

$$= \hat{\alpha}(\mathbf{z}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n | \mathbf{z}_{n-1}) \hat{\beta}(\mathbf{z}_n) \quad (17)$$



M step

In M step, need to optimize $Q(\theta, \theta^{\text{old}})$ with respect to $\theta = \{\pi, \mathbf{A}, \phi\}$

$$Q(\theta, \theta^{\text{old}}) = \sum_{k=1}^K \gamma(z_{nk}) \ln \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \xi(z_{n-1,j} z_{n,k}) \ln A_{jk} + \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \ln p(\mathbf{x}_n | z_{nk} \phi) \quad (18)$$



π is only associated with ①

\mathbf{A} is only associated with ②

ϕ is only associated with ③

π, \mathbf{A}, ϕ are independent, so they can be optimized respectively.

① maximize π

Take constraint $\sum_{j=1}^K \pi_j = 1$ into consideration, introduce Lagrange function

$$R_1(\boldsymbol{\pi}, \lambda) = \sum_{k=1}^K \gamma(z_{nk}) \ln \pi_k + \lambda \left(\sum_{j=1}^K \pi_j - 1 \right) \quad (19)$$

Let $\frac{\partial R_1(\boldsymbol{\pi}, \lambda)}{\partial \pi_k} = 0 \quad k = 1, \dots, K$ $\rightarrow \pi_k = \frac{\gamma_{1k}}{\sum_{j=1}^K \gamma_{1j}}, \quad k = 1, \dots, K$ (20)

② maximize A

Take constraint $\sum_{l=1}^K A_{jl} = 1 \quad l = 1, \dots, K$ into consideration, introduce Lagrange function

$$R_2(\mathbf{A}, \lambda_1, \dots, \lambda_K) = \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \xi(z_{n-1,j} z_{n,k}) \ln A_{jk} + \lambda_1 \left(\sum_{l=1}^K A_{1l} - 1 \right) + \dots + \lambda_K \left(\sum_{l=1}^K A_{Kl} - 1 \right) \quad (21)$$

Let $\frac{\partial R_2(\mathbf{A}, \lambda_1, \dots, \lambda_K)}{\partial A_{jk}} = 0 \quad k = 1, \dots, K$ $\rightarrow A_{jk} = \frac{\sum_{n=2}^N \xi(z_{n-1,j} z_{n,k})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j} z_{n,l})} \quad j, k = 1, \dots, K$ (22)

This shows the elements which are zero in A_{jk} will keep zero all the time.
So if you want to get left-right HMM just initialize A as upper diagonal matrix.

③ Maximize ϕ

the emission distribution can be multinomial, Gaussian, etc. We will discuss separately.

I . Emission distribution is discrete multinomial distribution

probability density function

$$p(\mathbf{x}_n | z_{nk} \phi) = \prod_{m=1}^M B_{km}^{x_{nm}} \quad s.t. \sum_{l=1}^M B_{kl} = 1, \quad k = 1, \dots, K \quad (23)$$

where $\phi = \{B_{km}\}$ $k = 1, \dots, K, m = 1, \dots, M$

B_{km} represents the probability of m-th event at k-th state.

To estimate B_{km} , introduce Lagrange function

$$R_3(\phi, \phi^{\text{old}}, \lambda_1, \dots, \lambda_K) = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \ln p(\mathbf{x}_n | z_{nk} \phi) + \lambda_1 \left(\sum_{l=1}^M B_{1l} - 1 \right) + \lambda_K \left(\sum_{l=1}^M B_{Kl} - 1 \right) \quad (24)$$

$$\frac{\partial R_3(\phi, \phi^{\text{old}}, \lambda_1, \dots, \lambda_K)}{\partial B_{km}} = 0 \quad \rightarrow \quad B_{km} = \frac{\sum_{n=1}^N \gamma_{nk} x_{nm}}{\sum_{n=1}^N \gamma_{nk} \sum_{m=1}^M x_{nm}} \quad k = 1, \dots, K; m = 1, \dots, M \quad (25)$$

EM algorithm for multinomial-HMM

1. Init parameters $\theta = \{\pi, \mathbf{A}, \phi\}$

where $\phi = \{B_{km}\} \quad k = 1, \dots, K; m = 1, \dots, M$

2. E step

Calculate $\gamma(\mathbf{z}_n)$ using (6)

Calculate $\xi(\mathbf{z}_{n-1} \mathbf{z}_n)$ using (17)

3. M step

$$\pi_k = \frac{\gamma_{1k}}{\sum_{j=1}^K \gamma_{1j}}, \quad k = 1, \dots, K \quad A_{jk} = \frac{\sum_{n=2}^N \xi(z_{n-1,j} z_{n,k})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j} z_{n,l})} \quad j, k = 1, \dots, K$$

$$B_{km} = \frac{\sum_{n=1}^N \gamma_{nk} x_{nm}}{\sum_{n=1}^N \gamma_{nk} \sum_{m=1}^M x_{nm}} \quad k = 1, \dots, K; m = 1, \dots, M$$

4. If converge then stop, otherwise goto 2.

Example

Generate Multinomial-HMM for belowing sequences

Data{1} = [1 1 1 4 1 1 1 2 2 2 2 2 1 2 2 2 2 3 3 3 3 3 1 3 3 3]

Data{2} = [1 1 2 1 1 1 1 2 2 2 3 2 2 2 2 2 3 3 3 3 3 4 3 3 3]

State num: 3, multinomial num: 4

Output

$$\pi = [1, 0, 0]$$

$$\mathbf{A} = \begin{bmatrix} 0.87 & 0.13 & 0.00 \\ 0.00 & 0.90 & 0.10 \\ 0.00 & 0.00 & 0.10 \end{bmatrix}$$

$$\mathbf{B} = \begin{bmatrix} 0.85 & 0.06 & 0.05 \\ 0.08 & 0.89 & 0.00 \\ 0.00 & 0.05 & 0.89 \\ 0.07 & 0.00 & 0.06 \end{bmatrix}$$

II . Emission distribution is Gaussian distribution

$$\text{probability density function } p(\mathbf{x}_n | z_{nk}, \phi) = N(\mathbf{x}_n | \boldsymbol{\mu}_k \Sigma_k) \quad (26)$$

where $\phi = \{\boldsymbol{\mu}_k, \Sigma_k\} \quad k = 1, \dots, K$

$$\begin{aligned} \text{Define } R_3(\phi, \phi^{old}) &= \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \ln p(\mathbf{x}_n | z_{nk}, \phi) \\ &= \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \left[-\frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\Sigma_k| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right] \end{aligned} \quad (27)$$

To estimate parameters, let derivative of (27) with respect to

$\phi = \{\boldsymbol{\mu}_k, \Sigma_k\} \quad k = 1, \dots, K$ be zero

$$\begin{aligned} \frac{\partial R_3(\phi, \phi^{old})}{\partial \boldsymbol{\mu}_k} = 0 \quad k = 1, \dots, K &\quad \rightarrow \quad \boldsymbol{\mu}_k = \frac{\sum_{n=1}^N \gamma_{nk} \mathbf{x}_n}{\sum_{n=1}^N \gamma_{nk}} \quad k = 1, \dots, K \\ \frac{\partial R_3(\phi, \phi^{old})}{\partial \Sigma_k} = 0 \quad k = 1, \dots, K &\quad \rightarrow \quad \Sigma_k = \frac{\sum_{n=1}^N \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{\sum_{n=1}^N \gamma_{nk}} \quad k = 1, \dots, K \end{aligned} \quad (28)$$

EM algorithm for Gaussian-HMM

1. Init parameters $\theta = \{\pi, \mathbf{A}, \phi\}$

where $\phi = \{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K\} \quad k = 1, \dots, K$

2. E step

Calculate $\gamma(\mathbf{z}_n)$ using (6)

Calculate $\xi(\mathbf{z}_{n-1}\mathbf{z}_n)$ using (17)

3. M step

$$\pi_k = \frac{\gamma_{1k}}{\sum_{j=1}^K \gamma_{1j}}, \quad k = 1, \dots, K$$

$$\boldsymbol{\mu}_k = \frac{\sum_{n=1}^N \gamma_{nk} \mathbf{x}_n}{\sum_{n=1}^N \gamma_{nk}} \quad k = 1, \dots, K$$

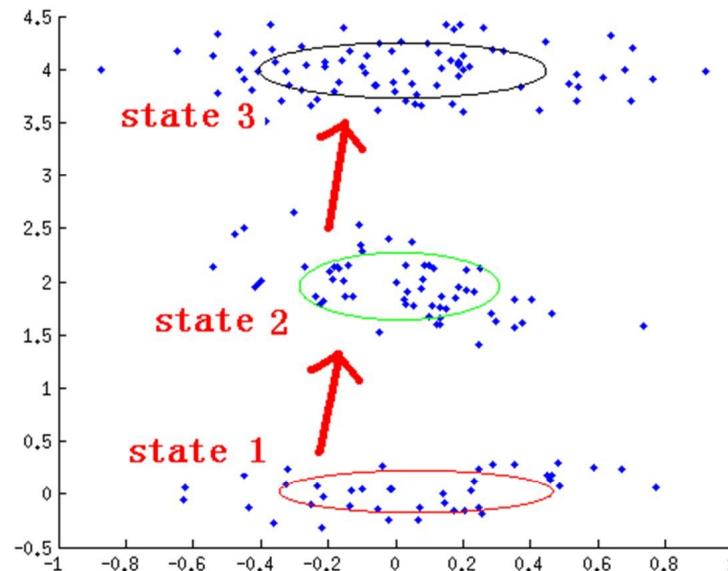
$$A_{jk} = \frac{\sum_{n=2}^N \xi(z_{n-1,j} z_{n,k})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j} z_{n,l})} \quad j, k = 1, \dots, K$$

$$\boldsymbol{\Sigma}_k = \frac{\sum_{n=1}^N \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{\sum_{n=1}^N \gamma_{nk}} \quad k = 1, \dots, K$$

4. If converge then stop, otherwise goto 2.

Example

The graph below shows the trained Gaussian-HMM model using created data.



State num: 3

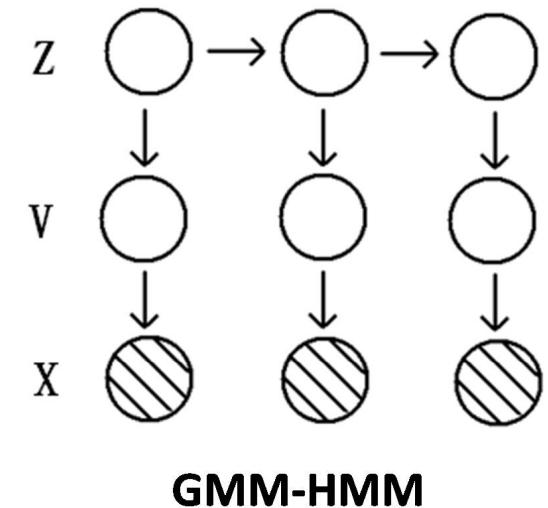
$$\boldsymbol{\pi} = [1, 0, 0] \quad \mathbf{A} = [0.95, 0.05, 0.00 \\ 0.00, 0.97, 0.03 \\ 0.00, 0.00, 1.00]$$

$\phi = \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\} \quad k=1, \dots, K$ is shown in graph

GMM-HMM

Gaussian distribution is not able to capture distributions with many centers. Especially in ASR, where the timbre differentiates from person to person.

In GMM-HMM, every single state is a GMM model. Define v as 1-of- K^*M random variable, where K is the number of state. M is the mixture number of GMM. $v_{km}=1$ represents the k -th state, m -th mixture occurs.



Parameters are $\theta = \{\pi, \mathbf{A}, \phi\}$

where $\phi = \{B_{km}, \boldsymbol{\mu}_{km}, \boldsymbol{\Sigma}_{km}\} \quad k = 1, \dots, K; m = 1, \dots, M$

$B_{km} = p(v_{km} | z_k)$ is the probability of m -th mixture under k -th state

$\boldsymbol{\mu}_{km}, \boldsymbol{\Sigma}_{km}$ is mean and covariance of the k -th state, m -mixture, respectively.

Probability density function

$$p(\mathbf{x} | z_k \phi) = \sum_{m=1}^M p(v_{km} | z_k \phi) p(\mathbf{x} | v_{km} \phi) = \sum_{m=1}^M B_{km} N(\mathbf{x} | \boldsymbol{\mu}_{km}, \boldsymbol{\Sigma}_{km}) \quad (29)$$

Apply EM algorithm to GMM-HMM

Compared with GMM-HMM, there are other latent variables \mathbf{v}_n , which dominates the emission GMM

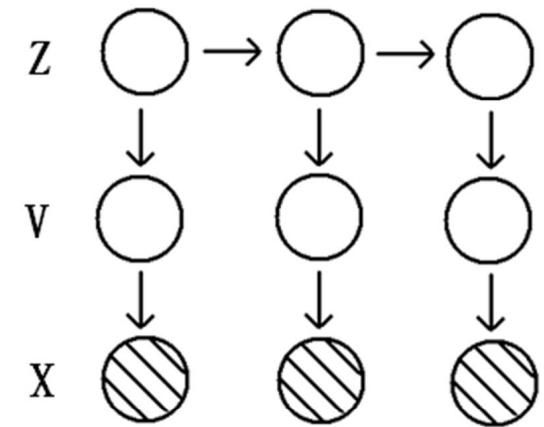
denote $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
 $\mathbf{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_N\}$
 $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_N\}$

(30)

Using independence of GMM-HMM, calculate $Q(\theta, \theta^{\text{old}})$

GMM-HMM

$$\begin{aligned}
 Q(\theta, \theta^{\text{old}}) &= \sum_{\mathbf{VZ}} p(\mathbf{VZ} | \mathbf{X}\theta^{\text{old}}) \ln p(\mathbf{VXZ} | \theta) \\
 &= \sum_{\mathbf{VZ}} p(\mathbf{VZ} | \mathbf{X}\theta^{\text{old}}) \left[\ln p(\mathbf{z}_1) + \sum_{n=2}^N \ln p(\mathbf{z}_n | \mathbf{z}_{n-1}) + \sum_{n=1}^N \ln p(\mathbf{v}_n | \mathbf{z}_n) + \sum_{n=1}^N \ln p(\mathbf{x}_n | \mathbf{v}_n) \right] \\
 &= \sum_{\mathbf{z}_1} p(\mathbf{z}_1 | \mathbf{X}\theta^{\text{old}}) \ln p(\mathbf{z}_1) + \sum_{n=2}^N \sum_{\mathbf{z}_{n-1}\mathbf{z}_n} p(\mathbf{z}_{n-1}\mathbf{z}_n | \mathbf{X}\theta^{\text{old}}) \ln p(\mathbf{z}_n | \mathbf{z}_{n-1}) + \sum_{n=1}^N \sum_{\mathbf{v}_n\mathbf{z}_n} p(\mathbf{v}_n\mathbf{z}_n | \mathbf{X}\theta^{\text{old}}) [\ln p(\mathbf{v}_n | \mathbf{z}_n) + \ln p(\mathbf{x}_n | \mathbf{v}_n)] \\
 &= \sum_{k=1}^K \gamma(z_{1k}) \ln \pi_k + \underbrace{\sum_{n=1}^N \sum_{j=1}^K \sum_{k=1}^K \xi(z_{n-1,j} z_{n,k}) \ln A_{jk}}_{\textcircled{1}} + \underbrace{\sum_{n=1}^N \sum_{k=1}^K \sum_{m=1}^M \eta_{nkm} \left[\ln B_{km} - \frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\Sigma_{km}| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \Sigma_{km}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right]}_{\textcircled{2}} \quad (31)
 \end{aligned}$$



E step

Calculate $p(\mathbf{VZ} | \mathbf{X}\theta^{\text{old}})$

From (31), no need to calculate all $p(\mathbf{VZ} | \mathbf{X}\theta^{\text{old}})$, only the terms below is needed

$$\gamma(z_{nk}) = p(z_{nk} | \mathbf{X}\theta^{\text{old}}) \quad (32)$$

$$\xi(z_{n-1,j} z_{n,k}) = p(z_{n-1,j} z_{n,k} | \mathbf{X}\theta^{\text{old}}) \quad (33)$$

$$\eta_{nkm} = p(v_{nkm} z_{nk} | \mathbf{X}\theta^{\text{old}}) \quad (34)$$

(32), (33) can be calculated in same way as (6), (17)

To calculate (34), just use forward-backward algorithm similar to (6)

$$\eta(\mathbf{v}_n \mathbf{z}_n) = p(\mathbf{v}_n \mathbf{z}_n | \mathbf{X}) = \underbrace{\hat{\alpha}(\mathbf{v}_n \mathbf{z}_n)}_{\text{forward}} \underbrace{\hat{\beta}(\mathbf{v}_n \mathbf{z}_n)}_{\text{backward}} \quad (35)$$

$$\text{where } \hat{\alpha}(\mathbf{v}_n \mathbf{z}_n) = \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{v}_n \mathbf{z}_n)}{p(\mathbf{x}_1, \dots, \mathbf{x}_1)} \quad (36)$$

$$\hat{\beta}(\mathbf{v}_n \mathbf{z}_n) = \frac{p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{v}_n \mathbf{z}_n)}{p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{x}_1, \dots, \mathbf{x}_n)} \quad (37)$$

denote $c_n = p(\mathbf{x}_n | \mathbf{x}_1, \dots, \mathbf{x}_{n-1})$

$\hat{\alpha}(\mathbf{v}_n \mathbf{z}_n)$, $\hat{\beta}(\mathbf{v}_n \mathbf{z}_n)$, c_n can be computed similar to (10) – (16)

Results are given without detailed deduction

$$\hat{\alpha}(\mathbf{v}_1 \mathbf{z}_1) = p(\mathbf{v}_1 \mathbf{z}_1 | \mathbf{x}_1) = \frac{p(\mathbf{x}_1 | \mathbf{v}_1 \mathbf{z}_1) p(\mathbf{v}_1 \mathbf{z}_1)}{p(\mathbf{x}_1)} = \frac{p(\mathbf{x}_1 | \mathbf{v}_1) p(\mathbf{v}_1 | \mathbf{z}_1) p(\mathbf{z}_1)}{c_1} \quad (38)$$

$$\hat{\alpha}(\mathbf{v}_n \mathbf{z}_n) = \frac{1}{c_n} p(\mathbf{x}_n | \mathbf{v}_n) p(\mathbf{v}_n | \mathbf{z}_n) \sum_{\mathbf{v}_{n-1} \mathbf{z}_{n-1}} \hat{\alpha}(\mathbf{v}_{n-1} \mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1}) \quad (39)$$

$$c_1 = p(\mathbf{x}_1) = \sum_{\mathbf{v}_1 \mathbf{z}_1} p(\mathbf{v}_1 \mathbf{z}_1 | \mathbf{x}_1) = \sum_{\mathbf{v}_1 \mathbf{z}_1} p(\mathbf{z}_1) p(\mathbf{v}_1 | \mathbf{z}_1) p(\mathbf{x}_1 | \mathbf{v}_1) \quad (40)$$

$$c_n = \sum_{\mathbf{v}_n \mathbf{z}_n} \left[p(\mathbf{x}_n | \mathbf{v}_n) p(\mathbf{v}_n | \mathbf{z}_n) \sum_{\mathbf{v}_{n-1} \mathbf{z}_{n-1}} \hat{\alpha}(\mathbf{v}_{n-1} \mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1}) \right] \quad (41)$$

$$\hat{\beta}(\mathbf{v}_N \mathbf{z}_N) = 1 \quad (42)$$

$$\hat{\beta}(\mathbf{v}_n \mathbf{z}_n) = \frac{1}{c_{n+1}} \sum_{\mathbf{v}_{n+1} \mathbf{z}_{n+1}} \hat{\beta}(\mathbf{v}_{n+1} \mathbf{z}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{v}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n) p(\mathbf{v}_{n+1} | \mathbf{z}_{n+1}) \quad (43)$$

M step

Maximize π, A

The maximization of π, A is same as (20), (22).

$$\pi_k = \frac{\gamma_{1k}}{\sum_{j=1}^K \gamma_{1j}}, \quad k = 1, \dots, K \quad (44)$$

$$A_{jk} = \frac{\sum_{n=2}^N \xi(z_{n-1,j} z_{n,k})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j} z_{n,l})} \quad j, k = 1, \dots, K \quad (45)$$

Maximize ϕ

To maximize $\phi = \{B_{km}, \mu_{km}, \Sigma_{km}\}$ $k = 1, \dots, K; m = 1, \dots, M$

Taking constraint $\sum_{l=1}^M B_{kl} = 1, \quad k = 1, \dots, K$ into account

Introduce Lagrange function

$$R_3(\phi, \phi^{old}, \lambda_1, \dots, \lambda_K) = \sum_{n=1}^N \sum_{k=1}^K \sum_{m=1}^M \eta_{nkm} \left[\ln B_{km} - \frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\Sigma_{km}| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_{km}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right] + \lambda_1 \left(\sum_{l=1}^M B_{1l} - 1 \right) + \lambda_K \left(\sum_{l=1}^M B_{Kl} - 1 \right) \quad (46)$$

Let the derivative of $R_3(\phi, \phi^{\text{old}}, \lambda_1, \dots, \lambda_K)$ With respect to
 $\phi = \{B_{km}, \mu_{km}, \Sigma_{km}\} \quad k = 1, \dots, K; m = 1, \dots, M$ be zero

$$\frac{\partial R_3(\phi, \phi^{\text{old}}, \lambda_1, \dots, \lambda_K)}{\partial B_{km}} = 0$$

$$B_{km} = \frac{\sum_{n=1}^N \eta_{nkm}}{\sum_{n=1}^N \sum_{j=1}^M \eta_{nj}}$$

$$\frac{\partial R_3(\phi, \phi^{\text{old}}, \lambda_1, \dots, \lambda_K)}{\partial \mu_{km}} = 0$$



$$\mu_{km} = \frac{\sum_{n=1}^N \eta_{nkm} \mathbf{x}_n}{\sum_{n=1}^N \eta_{nkm}} \quad k = 1, \dots, K; m = 1, \dots, M \quad (47)$$

$$\frac{\partial R_3(\phi, \phi^{\text{old}}, \lambda_1, \dots, \lambda_K)}{\partial \Sigma_{km}} = 0$$

$$\Sigma_{km} = \frac{\sum_{n=1}^N \eta_{nkm} (\mathbf{x}_n - \mu_{km})(\mathbf{x}_n - \mu_{km})^T}{\sum_{n=1}^N \eta_{nkm}}$$

EM algorithm for GMM-HMM

1. Init parameters $\theta = \{\pi, \mathbf{A}, \phi\}$

where $\phi = \{B_{km}, \boldsymbol{\mu}_{km}, \boldsymbol{\Sigma}_{km}\} \quad k = 1, \dots, K; m = 1, \dots, M$

2. E step

Calculate $\gamma(\mathbf{z}_n)$, $\xi(\mathbf{z}_{n-1}\mathbf{z}_n)$, $\eta(\mathbf{v}_n \mathbf{z}_n)$ from (32), (33), (34)

3. M step

$$\pi_k = \frac{\gamma_{1k}}{\sum_{j=1}^K \gamma_{1j}}, \quad k = 1, \dots, K$$

$$A_{jk} = \frac{\sum_{n=2}^N \xi(z_{n-1,j} z_{n,k})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j} z_{n,l})} \quad j, k = 1, \dots, K$$

$$B_{km} = \frac{\sum_{n=1}^N \eta_{nkm}}{\sum_{n=1}^N \sum_{j=1}^M \eta_{nj}}$$

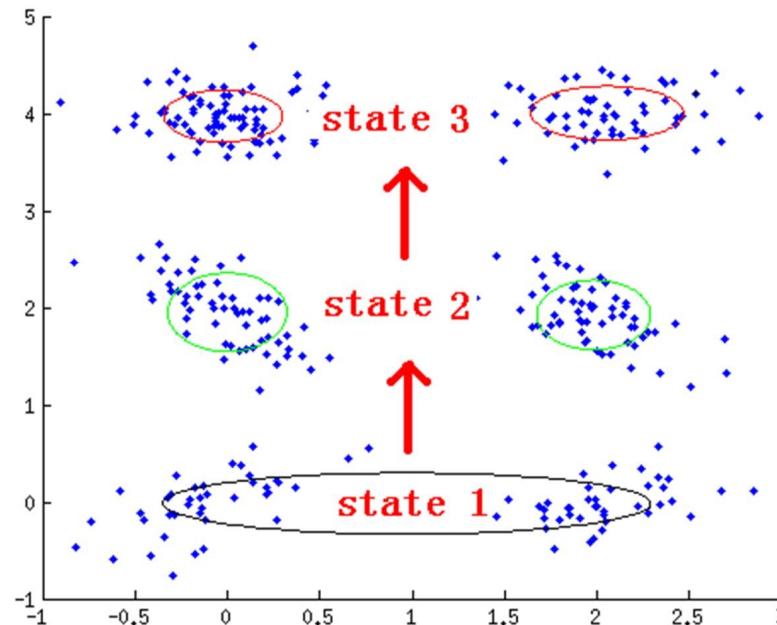
$$\boldsymbol{\mu}_{km} = \frac{\sum_{n=1}^N \eta_{nkm} \mathbf{x}_n}{\sum_{n=1}^N \eta_{nkm}}$$

$$\boldsymbol{\Sigma}_{km} = \frac{\sum_{n=1}^N \eta_{nkm} (\mathbf{x}_n - \boldsymbol{\mu}_{km})(\mathbf{x}_n - \boldsymbol{\mu}_{km})^T}{\sum_{n=1}^N \eta_{nkm}}$$

4. If converge then stop, otherwise goto 2.

Example

A GMM-HMM with state num:3, mix num:2



$$\pi: [0, 0, 1]$$

$$A: [1.00, 0.00, 0.00, \\ 0.03, 0.97, 0.00, \\ 0.00, 0.05, 0.95]$$

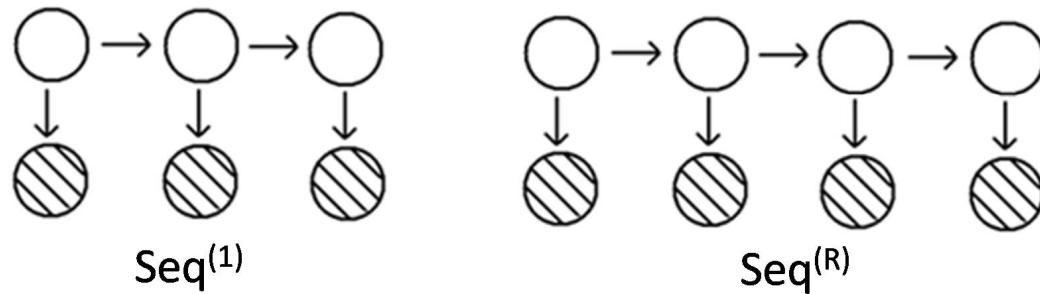
$$\phi = \{B_{km}, \mu_{km}, \Sigma_{km}\} \quad k=1, \dots, K; m=1, \dots, M$$

is shown on the graph

Multi Sequence Training

Till now, we train HMM only use one sequence.

In ASR, we have many utterance to train one HMM model for each phoneme / word.
This part will show how to train HMM using multi sequences.



Using independency of PGM, $Q(\theta, \theta^{\text{old}})$

$$\begin{aligned} Q(\theta, \theta^{\text{old}}) &= \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}\theta^{\text{old}}) \ln p(\mathbf{XZ} | \theta) \\ &= \sum_{r=1}^R \left(\sum_{\mathbf{Z}^{(r)}} p(\mathbf{Z}^{(r)} | \mathbf{X}^{(r)}\theta^{\text{old}}) \ln p(\mathbf{X}^{(r)}\mathbf{Z}^{(r)} | \theta) \right) \\ &= \sum_{r=1}^R \sum_{k=1}^K \gamma_{1k}^{(r)} \ln \pi_k + \sum_{r=1}^R \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \xi^{(r)}(z_{n-1,j} z_{n,k}) \ln A_{jk} + \sum_{l=r}^R \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk}^{(r)} \ln p(\mathbf{x}_n | z_{nk}) \quad (48) \end{aligned}$$

Where $\mathbf{X}^{(r)}$ is the observation variables of r-th seq.
 $\mathbf{Z}^{(r)}$ is the latent variables of r-th seq.

E STEP

Estimate γ , ξ , η for each sequences separately.

$$\gamma_{nk}^{(r)} = p(z_{nk}^{(r)} \mid \mathbf{X}^{(r)} \theta^{old}) \quad (49)$$

$$\xi^{(r)}(z_{n-1,j} z_{n,k}) = p(z_{n-1,j}^{(r)} z_{n,k}^{(r)} \mid \mathbf{X}^{(r)} \theta^{old}) \quad (50)$$

$$\eta_{nkm}^{(r)} = p(v_{nkm}^{(r)} z_{nk}^{(r)} \mid \mathbf{X} \theta^{old}) \quad (\text{for GMM only}) \quad (51)$$

M STEP

Optimize (48) with respect to $\theta = \{\pi, \mathbf{A}, \phi\}$

$$\pi_k = \frac{\sum_{r=1}^R \gamma_{1k}^{(r)}}{\sum_{r=1}^R \sum_{j=1}^K \gamma_{1j}^{(r)}} \quad (52)$$

$$A_{jk} = \frac{\sum_{r=1}^R \sum_{n=2}^N \xi^{(r)}(z_{n-1,j} z_{n,k})}{\sum_{r=1}^R \sum_{l=1}^K \sum_{n=2}^N \xi^{(r)}(z_{n-1,j} z_{n,l})} \quad (53)$$

I . Multinomial Distribution

$$B_{km} = \frac{\sum_{r=1}^R \sum_{n=1}^{N^{(r)}} \gamma_{nk}^{(r)} x_{nm}^{(r)}}{\sum_{r=1}^R \sum_{n=1}^{N^{(r)}} \gamma_{nk}^{(r)} \sum_{m=1}^M x_{nm}^{(r)}} \quad (54)$$

II . Gauss Distribution

$$\begin{aligned} \boldsymbol{\mu}_k &= \frac{\sum_{r=1}^R \sum_{n=1}^{N^{(r)}} \gamma_{nk}^{(r)} \mathbf{x}_n^{(r)}}{\sum_{r=1}^R \sum_{n=1}^{N^{(r)}} \gamma_{nk}^{(r)}} \\ \boldsymbol{\Sigma}_k &= \frac{\sum_{r=1}^R \sum_{n=1}^{N^{(r)}} \gamma_{nk}^{(r)} (\mathbf{x}_n^{(r)} - \boldsymbol{\mu}_k)(\mathbf{x}_n^{(r)} - \boldsymbol{\mu}_k)^T}{\sum_{r=1}^R \sum_{n=1}^{N^{(r)}} \gamma_{nk}^{(r)}} \end{aligned} \quad (55)$$

III. GMM

$$\begin{aligned} B_{km} &= \frac{\sum_{r=1}^R \sum_{n=1}^{N^{(r)}} \eta_{nkm}^{(r)}}{\sum_{r=1}^R \sum_{n=1}^{N^{(r)}} \sum_{j=1}^M \eta_{nj}^{(r)}} \\ \boldsymbol{\mu}_{km} &= \frac{\sum_{r=1}^R \sum_{n=1}^{N^{(r)}} \eta_{nkm}^{(r)} \mathbf{x}_n}{\sum_{r=1}^R \sum_{n=1}^{N^{(r)}} \eta_{nkm}^{(r)}} \\ \boldsymbol{\Sigma}_{km} &= \frac{\sum_{r=1}^R \sum_{n=1}^{N^{(r)}} \eta_{nkm}^{(r)} (\mathbf{x}_n - \boldsymbol{\mu}_{km})(\mathbf{x}_n - \boldsymbol{\mu}_{km})^T}{\sum_{r=1}^R \sum_{n=1}^{N^{(r)}} \eta_{nkm}^{(r)}} \end{aligned} \quad (56)$$

Decoding of HMM

Q: How to find the best decoding path?

$$\begin{aligned}\mathbf{A:} \quad \mathbf{Z}_{opt} &= \max_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}) = \max_{\mathbf{Z}} p(\mathbf{ZX}) \\ &= \max_{\mathbf{Z}} p(\mathbf{z}_1) \prod_{n=2}^N p(\mathbf{z}_n | \mathbf{z}_{n-1}) \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n)\end{aligned}\tag{57}$$

However, we need to evaluate all possible \mathbf{Z} which is K^N times to get accurate solution. This infeasible.

Viterbi Algorithm

Use greedy algorithm to estimate optimized path step by step.

By discarding paths with low probability and storing previous step,
the computation complexity decreased to K*N

Viterbi Algorithm

$$V_{1k} = p(\mathbf{x}_1 | z_{1k}) \times p(z_{1k})$$

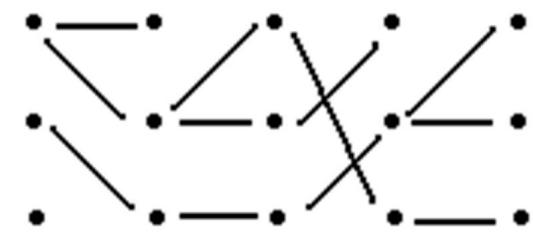
$$V_{1k} = p(\mathbf{x}_1 | z_{1k}) \times p(z_{1k})$$

for $n = 2 : N$

$$V_{nk} = \max_j (V_{n-1,j} \times p(z_{nk} | z_{n-1,j}) \times p(\mathbf{x}_n | z_{nj}))$$

$$\text{path}(n-1) = j$$

$$\text{path}(n) = \arg \max_j V_{nj}$$

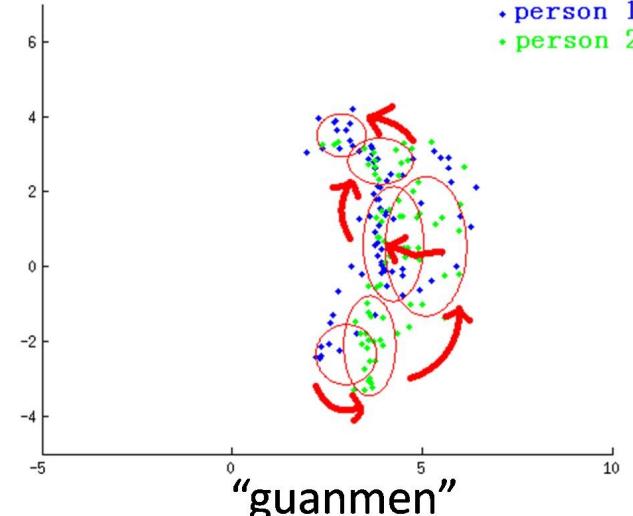
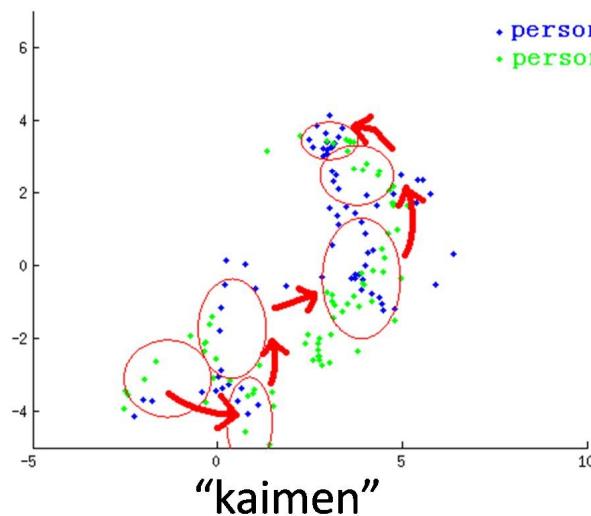


Path stored

Experiments in small vocabulary ASR

12 Mfcc feature (only choose the 1st and 2nd dimension to plot)

1. Gaussian-HMM, state num: 6



$$\pi = [1, 0, 0, 0, 0, 0]$$

$$A = \begin{matrix} 0.86 & 0.14 & 0 & 0 & 0 & 0 \\ 0 & 0.90 & 0.10 & 0 & 0 & 0 \\ 0 & 0 & 0.87 & 0.13 & 0 & 0 \\ 0 & 0 & 0 & 0.97 & 0.03 & 0 \\ 0 & 0 & 0 & 0 & 0.91 & 0.09 \\ 0 & 0 & 0 & 0 & 0 & 1.00 \end{matrix}$$

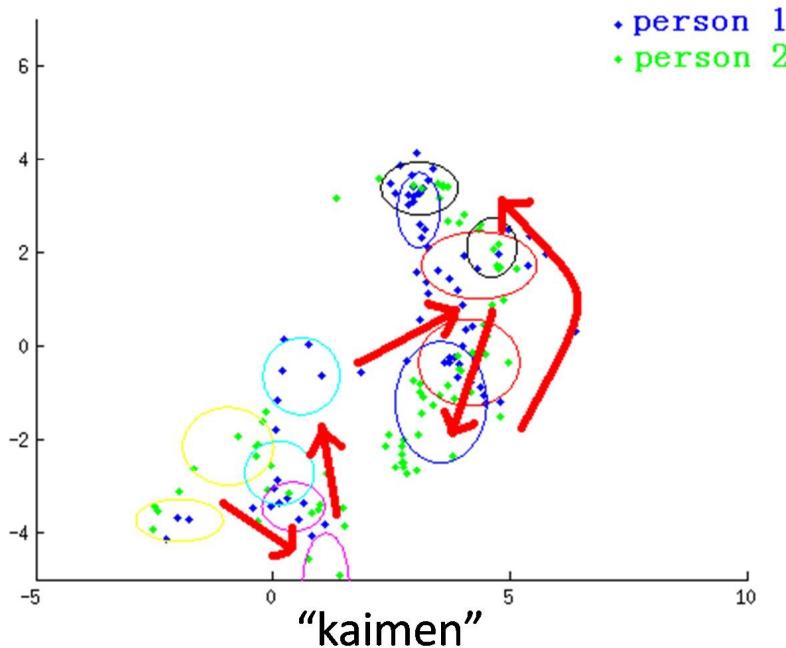
$$\pi = [1, 0, 0, 0, 0, 0]$$

$$A = \begin{matrix} 0.89 & 0.11 & 0 & 0 & 0 & 0 \\ 0 & 0.87 & 0.13 & 0 & 0 & 0 \\ 0 & 0 & 0.77 & 0.23 & 0 & 0 \\ 0 & 0 & 0 & 0.97 & 0.03 & 0 \\ 0 & 0 & 0 & 0 & 0.95 & 0.05 \\ 0 & 0 & 0 & 0 & 0 & 1.00 \end{matrix}$$

2. GMM-HMM (6 states, 2 mixutre)

With the increase of dataset, Gaussian-HMM is not able to capture the timbre of different people, gender, age.

We can use GMM-HMM model instead.



However, GMM-HMM is sensitive to the initial parameters. The tricks we are using in this example is:

1. Use Gaussian-HMM to train different people separately.
2. Combine the data point which are in the same state. Use GMM to initialize the parameters.
3. Run GMM-HMM to fine-tune the model.

Results

Dataset: 20 Isolated Chinese words. 11 male + 9 female. Altogether 800 pronunciations.

10 male and 9 female for training. 10 male and 10 female for testing.

Feature: 12 dimension MFCC

Model: Gaussian-HMM, GMM-HMM

Model	Accuracy
Gaussian-HMM	82.25%
GMM-HMM	84.00%

Weakness of HMM

1. Markov assumption

The next state is only dependent upon the current state. So is poor at capturing long-range correlations between the observed variables.

$$p(\mathbf{z}_{n+1} | \mathbf{z}_1, \dots, \mathbf{z}_n) = p(\mathbf{z}_{n+1} | \mathbf{z}_n) \quad (58)$$

2. Stationary assumption

$$p(\mathbf{z}_{n+1} | \mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{z}_{n-1}) \quad (59)$$

3. Output independence assumption

The current output is conditionally independent of the previous output.

$$p(\mathbf{X} | \mathbf{Z}) = \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{Z}) \quad (60)$$

Tricks of HMM

1. HMM is more sensitive to the initial parameters than GMM. So it is easy to get into local minimum.

Solve: Use GMM or other methods to initialize parameters.

Initialize parameters randomly and run HMM separately for several times.

2. For Gaussian-HMM & GMM-HMM, if $\text{eig}(\Sigma)$ is too small. Then Σ^{-1} will be unstable.

Solve: if $\text{eig}(\Sigma) < \epsilon$ then $\Sigma = \Sigma + \sigma I$

3. If $p(\mathbf{x}|\mathbf{z})$ is Gaussian or GMM pdf, underflow may occurs.

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

if this term is too small,
after exp it will underflow

Solve: use $\ln p(\mathbf{x}|\mathbf{z})$ to instead $p(\mathbf{x}|\mathbf{z})$ in code implementation.

For Gaussian

$$\ln p(\mathbf{x}|\mathbf{z}) = \ln \pi - \frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

For GMM

$$\begin{aligned} \ln p(\mathbf{x}|\mathbf{z}) &= \ln \sum_{m=1}^m \pi_m N(\mathbf{x}|\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m) \\ &= \ln \sum_{m=1}^m \exp\left(\ln \pi_m - \frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\boldsymbol{\Sigma}_m| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_m)^T \boldsymbol{\Sigma}_m^{-1} (\mathbf{x} - \boldsymbol{\mu}_m)\right) \\ &= \left[\ln \sum_{m=1}^m \exp\left(\ln \pi_m - \frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\boldsymbol{\Sigma}_m| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_m)^T \boldsymbol{\Sigma}_m^{-1} (\mathbf{x} - \boldsymbol{\mu}_m) - U\right) \right] + U \end{aligned}$$

if this term is too small,
after exp it will underflow

Normalization factor,
To avoid underflow

4. According to (11), (15), $\hat{\alpha}(\mathbf{z}_n), \hat{\beta}(\mathbf{z}_n)$ may underflow, and $\gamma(\mathbf{z}_n)$ will be unstable

Solve: Use $\ln \gamma(\mathbf{z}_n), \ln c_n, \ln \hat{\alpha}(\mathbf{z}_n), \ln \hat{\beta}(\mathbf{z}_n)$ replace $\gamma(\mathbf{z}_n), c_n, \hat{\alpha}(\mathbf{z}_n), \hat{\beta}(\mathbf{z}_n)$

$$\begin{aligned}\ln \hat{\alpha}(\mathbf{z}_n) &= -\ln c_n + \ln p(\mathbf{x}_n | \mathbf{z}_n) + \ln \sum_{\mathbf{z}_{n-1}} \exp \left(\ln \hat{\alpha}(\mathbf{z}_{n-1}) + \ln p(\mathbf{z}_n | \mathbf{z}_{n-1}) \right) \\ &= -\ln c_n + \ln p(\mathbf{x}_n | \mathbf{z}_n) + \left[\ln \sum_{\mathbf{z}_{n-1}} \exp \left(\ln \hat{\alpha}(\mathbf{z}_{n-1}) + \ln p(\mathbf{z}_n | \mathbf{z}_{n-1}) \right) - U \right] + U\end{aligned}$$

if this term is too small,
after exp it will underflow

Normalization factor,
To avoid underflow

The same strategy can be applied to $\ln c_n, \ln \hat{\beta}(\mathbf{z}_n)$

$\gamma(\mathbf{z}_n) = \hat{\alpha}(\mathbf{z}_n) \hat{\beta}(\mathbf{z}_n)$ will turn to

$$\ln \gamma(\mathbf{z}_n) = \ln \hat{\alpha}(\mathbf{z}_n) + \ln \hat{\beta}(\mathbf{z}_n)$$

Furthermore, for parameter estimation, such as (28) will turn to

$$\boldsymbol{\mu}_k = \frac{\sum_{n=1}^N \gamma_{nk} \mathbf{x}_n}{\sum_{n=1}^N \gamma_{nk}} = \frac{\sum_{n=1}^N \exp[\ln \gamma_{nk}] \mathbf{x}_n}{\sum_{n=1}^N \exp[\ln \gamma_{nk}]} = \frac{\sum_{n=1}^N \exp[(\ln \gamma_{nk}) - U] \mathbf{x}_n}{\sum_{n=1}^N \exp[(\ln \gamma_{nk}) - U]}$$

**if this term is too small,
after exp it will underflow**

**Normalization factor,
To avoid underflow**

Matlab Code

<https://github.com/qiuqiangkong/matlab-hmm>

THANK YOU!