

Hidden Markov Models

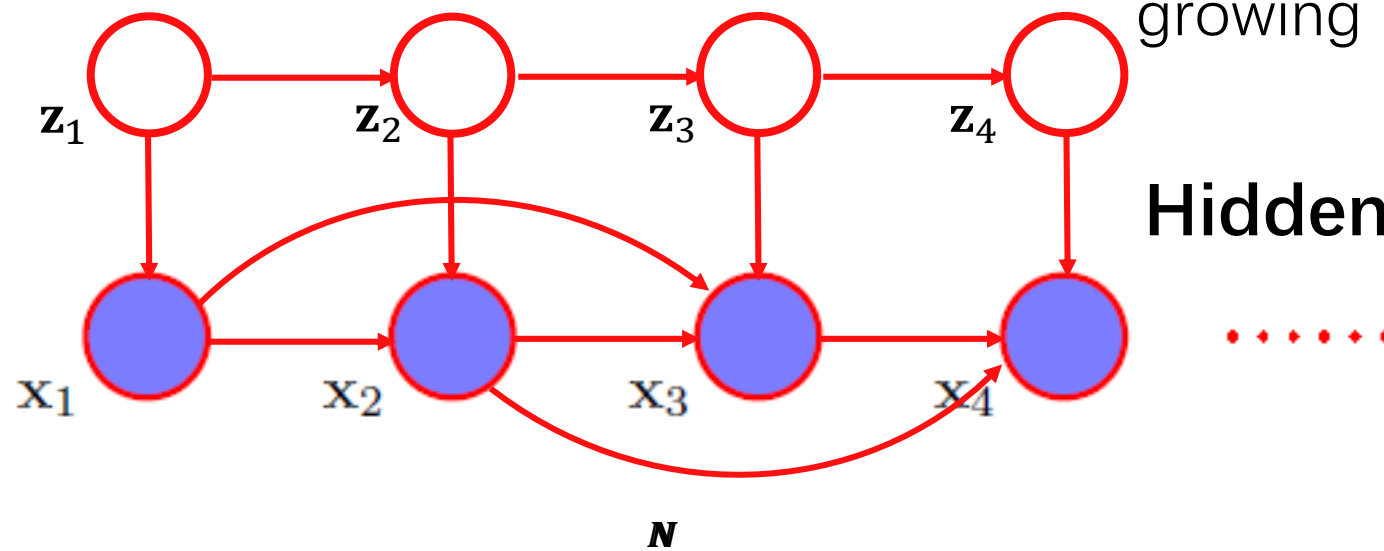
2019.01.25

Contents

- Introduction
- Algorithms
 - Sum-Product
 - Forward-Backward
 - Viterbi
- Application in KG
- Conclusion

Introduction

For large order Markov Chain, it's hard to evaluate the exponentially growing parameters.

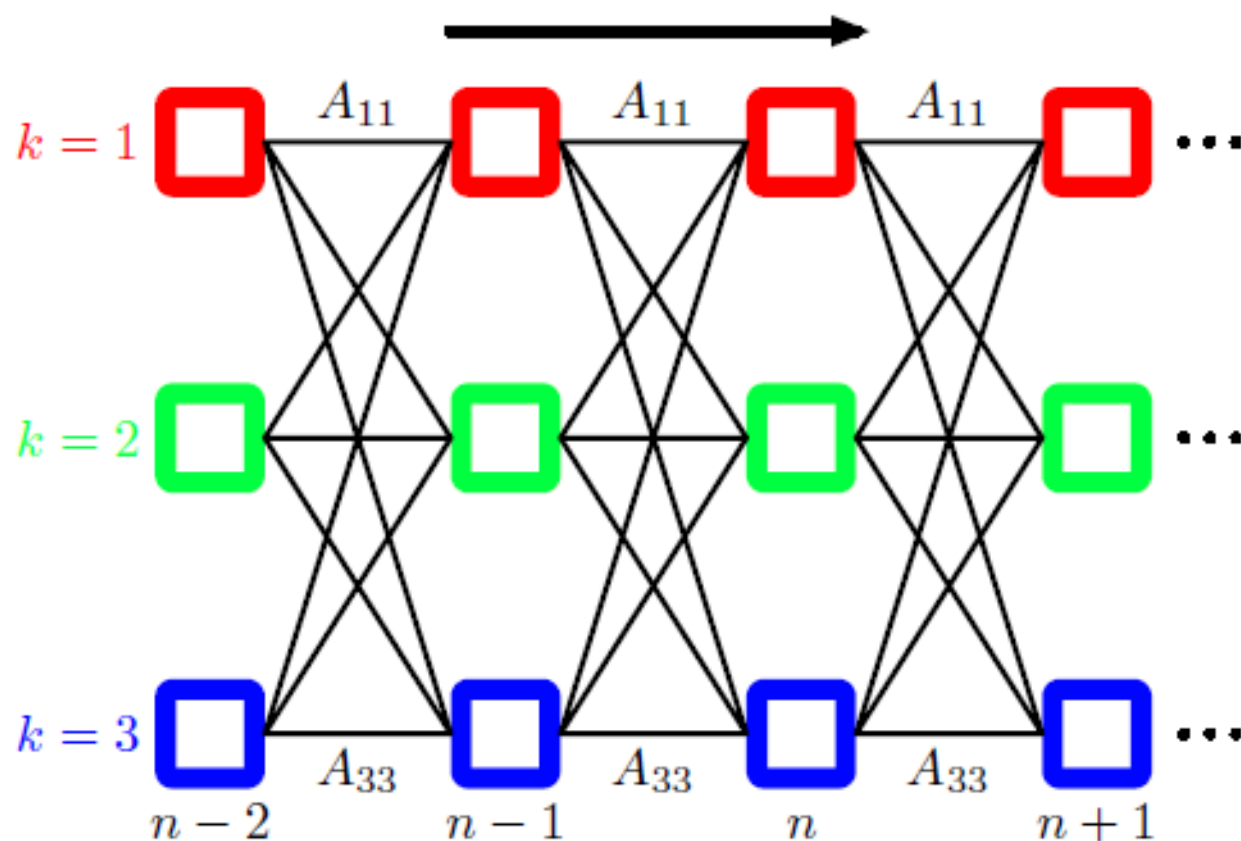
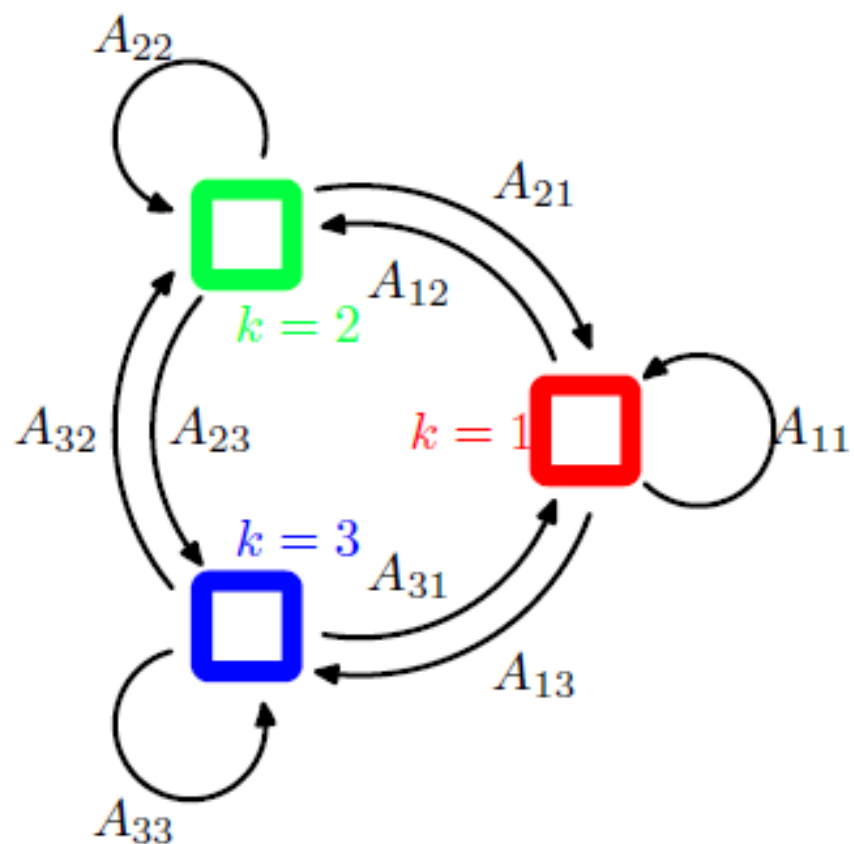


Hidden Markov Model(HMM)

$$p(\mathbf{x}) = p(x_1) \prod_{n=2}^N p(x_n | z_{n-1}, x_{n-2}) \quad \text{First order Markov Chain}$$

$$p(\mathbf{x}, \mathbf{z}) = p(z_1) \prod_{n=2}^N \underbrace{p(z_n | z_{n-1})}_{\text{transition}} \prod_{n=2}^N \underbrace{p(x_n | z_n)}_{\text{emission}}$$

Introduction – Transition Probability



Contents

- Introduction
- Algorithms
 - Sum-Product
 - Forward-Backward
 - Viterbi
- Application in KG
- Conclusion

Algorithms

1

Backward-Forward (Sum-product)

2

Maximum Likelihood ---- EM

3

Viterbi Algorithm

Algorithms – EM framework

Likelihood Function:

$$p(\mathbf{X}|\boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}).$$

In E step, we estimate $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old})$. Considering the relationship between \mathbf{z}_n and \mathbf{z}_{n-1} , we define:

$$\gamma(z_{nk}) = \mathbb{E}[z_{nk}] = \sum_{\mathbf{Z}} \gamma(\mathbf{Z}) z_{nk}$$

$$\xi(z_{n-1,j}, z_{nk}) = \mathbb{E}[z_{n-1,j} z_{nk}] = \sum_{\mathbf{Z}} \gamma(\mathbf{Z}) z_{n-1,j} z_{nk}$$

In M step, we maximize:

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{old}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$$

Here $\boldsymbol{\theta}$ represents parameters.

Algorithms – Forward-Backward

To seek an efficient procedure for evaluating the quantities $\gamma(z_{nk})$ and $\xi(z_{n-1,j}, z_{nk})$

$$\begin{aligned}p(\mathbf{X}|\mathbf{z}_n) &= p(\mathbf{x}_1, \dots, \mathbf{x}_n | \mathbf{z}_n) \\ &\quad p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n) \\ p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{x}_n, \mathbf{z}_n) &= p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{z}_n) \\ p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{z}_{n-1}, \mathbf{z}_n) &= p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{z}_{n-1}) \\ p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n, \mathbf{z}_{n+1}) &= p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_{n+1}) \\ p(\mathbf{x}_{n+2}, \dots, \mathbf{x}_N | \mathbf{z}_{n+1}, \mathbf{x}_{n+1}) &= p(\mathbf{x}_{n+2}, \dots, \mathbf{x}_N | \mathbf{z}_{n+1}) \\ p(\mathbf{X} | \mathbf{z}_{n-1}, \mathbf{z}_n) &= p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{z}_{n-1}) \\ &\quad p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n) \\ p(\mathbf{x}_{N+1} | \mathbf{X}, \mathbf{z}_{N+1}) &= p(\mathbf{x}_{N+1} | \mathbf{z}_{N+1}) \\ p(\mathbf{z}_{N+1} | \mathbf{z}_N, \mathbf{X}) &= p(\mathbf{z}_{N+1} | \mathbf{z}_N)\end{aligned}$$

Algorithms – Forward-Backward

Using Bayes' theorem, and we have:

$$\gamma(\mathbf{z}_n) = p(\mathbf{z}_n|\mathbf{X}) = \frac{p(\mathbf{X}|\mathbf{z}_n)p(\mathbf{z}_n)}{p(\mathbf{X})}$$

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

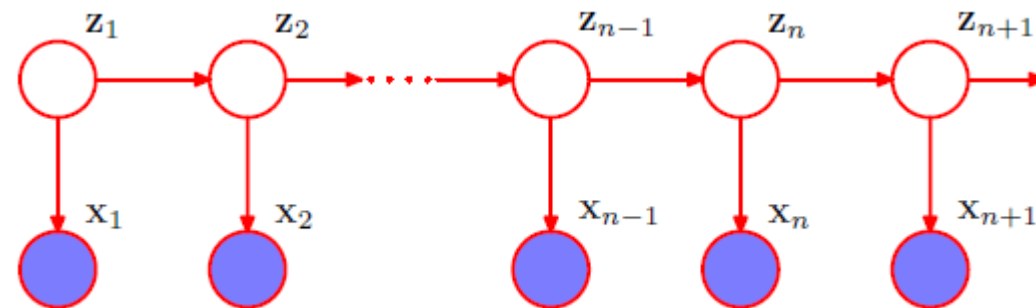
Where we defined:

$$\alpha(\mathbf{z}_n) \equiv p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n)$$

$$\beta(\mathbf{z}_n) \equiv p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N|\mathbf{z}_n)$$

Algorithms – Forward-Backward

$$\begin{aligned}
 \alpha(z_n) &= p(x_1, \dots, x_n, z_n) \\
 &= p(x_1, \dots, x_n | z_n) p(z_n) \\
 &= p(x_n | z_n) p(x_1, \dots, x_{n-1} | z_n) p(z_n) \\
 &= p(x_n | z_n) p(x_1, \dots, x_{n-1}, z_n) \\
 &= p(x_n | z_n) \sum_{z_{n-1}} p(x_1, \dots, x_{n-1}, z_{n-1}, z_n) \\
 &= p(x_n | z_n) \sum_{z_{n-1}} p(x_1, \dots, x_{n-1}, z_n | z_{n-1}) p(z_{n-1}) \\
 &= p(x_n | z_n) \sum_{z_{n-1}} p(x_1, \dots, x_{n-1} | z_{n-1}) p(z_n | z_{n-1}) p(z_{n-1}) \\
 &= p(x_n | z_n) \sum_{z_{n-1}} p(x_1, \dots, x_{n-1}, z_{n-1}) p(z_n | z_{n-1})
 \end{aligned}$$

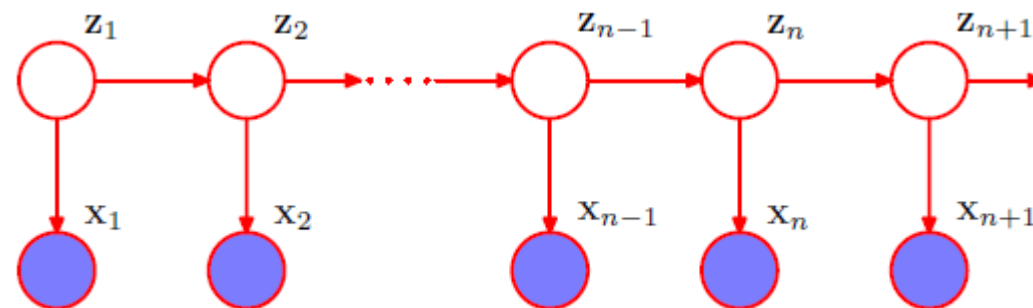


Then we have:

$$\alpha(z_n) = p(x_n | z_n) \sum_{z_{n-1}} \alpha(z_{n-1}) p(z_n | z_{n-1})$$

Algorithms – Forward-Backward

$$\begin{aligned}\beta(\mathbf{z}_n) &= p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n) \\ &= \sum_{\mathbf{z}_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N, \mathbf{z}_{n+1} | \mathbf{z}_n) \\ &= \sum_{\mathbf{z}_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n, \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n) \\ &= \sum_{\mathbf{z}_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n) \\ &= \sum_{\mathbf{z}_{n+1}} p(\mathbf{x}_{n+2}, \dots, \mathbf{x}_N | \mathbf{z}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n)\end{aligned}$$



Then we have:

$$\beta(\mathbf{z}_n) = \sum_{\mathbf{z}_{n+1}} \beta(\mathbf{z}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n)$$



Algorithms – Forward-Backward

Then we estimate $\xi(z_{n-1,j}, z_{nk})$

$$\begin{aligned}\xi(z_{n-1}, z_n) &= p(z_{n-1}, z_n | \mathbf{X}) \\ &= \frac{p(\mathbf{X} | z_{n-1}, z_n) p(z_{n-1}, z_n)}{p(\mathbf{X})} \\ &= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | z_{n-1}) p(\mathbf{x}_n | z_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | z_n) p(z_n | z_{n-1}) p(z_{n-1})}{p(\mathbf{X})} \\ &= \frac{\alpha(z_{n-1}) p(\mathbf{x}_n | z_n) p(z_n | z_{n-1}) \beta(z_n)}{p(\mathbf{X})}\end{aligned}\tag{13.43}$$

Then $\xi(z_{n-1,j}, z_{nk})$ can be estimated by α and β recursion.



Algorithms – Forward-Backward

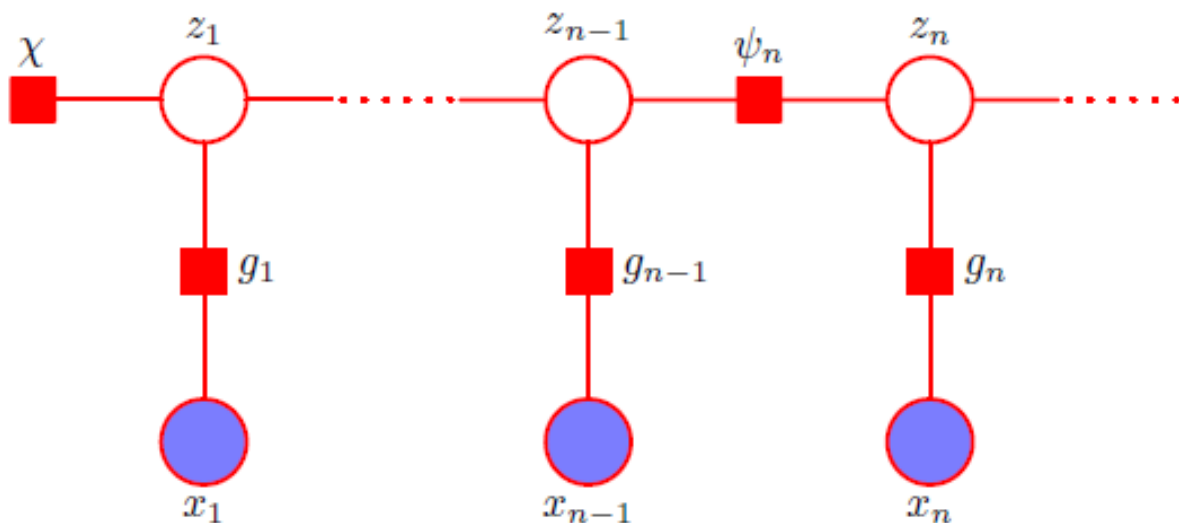
For prediction, we get:

$$\begin{aligned} p(\mathbf{x}_{N+1}|\mathbf{X}) &= \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1}, \mathbf{z}_{N+1}|\mathbf{X}) \\ &= \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1}|\mathbf{z}_{N+1})p(\mathbf{z}_{N+1}|\mathbf{X}) \\ &= \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1}|\mathbf{z}_{N+1}) \sum_{\mathbf{z}_N} p(\mathbf{z}_{N+1}, \mathbf{z}_N|\mathbf{X}) \\ &= \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1}|\mathbf{z}_{N+1}) \sum_{\mathbf{z}_N} p(\mathbf{z}_{N+1}|\mathbf{z}_N)p(\mathbf{z}_N|\mathbf{X}) \\ &= \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1}|\mathbf{z}_{N+1}) \sum_{\mathbf{z}_N} p(\mathbf{z}_{N+1}|\mathbf{z}_N) \frac{p(\mathbf{z}_N, \mathbf{X})}{p(\mathbf{X})} \\ &= \frac{1}{p(\mathbf{X})} \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1}|\mathbf{z}_{N+1}) \sum_{\mathbf{z}_N} p(\mathbf{z}_{N+1}|\mathbf{z}_N)\alpha(\mathbf{z}_N) \end{aligned}$$



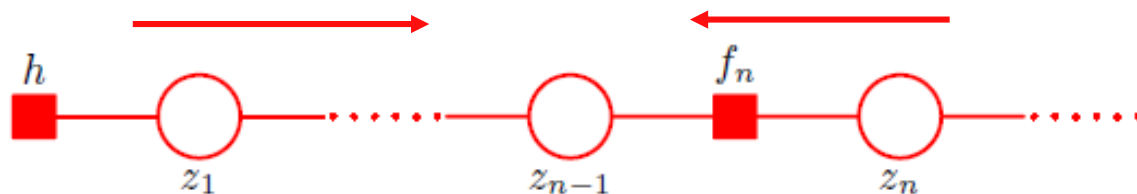
Algorithms – Sum-Product

In the perspective of belief propagation:



$$h(z_1) = p(z_1)p(x_1|z_1)$$

$$f_n(z_{n-1}, z_n) = p(z_n|z_{n-1})p(x_n|z_n)$$



$$\mu_{f_{n+1} \rightarrow f_n}(z_n) = \sum_{z_{n+1}} f_{n+1}(z_n, z_{n+1}) \mu_{f_{n+2} \rightarrow f_{n+1}}(z_{n+1})$$

$\alpha(z_n)$



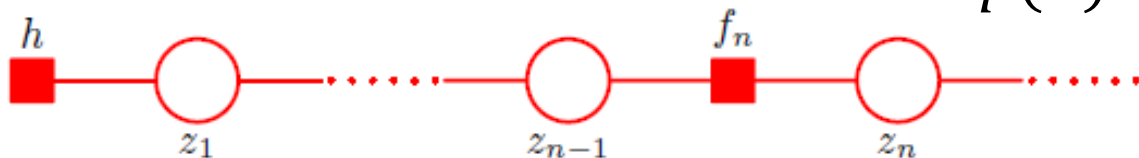
Algorithms – Sum-Product

Estimation of $\gamma(\mathbf{z}_n)$ and $\xi(\mathbf{z}_{n-1}, \mathbf{z}_n)$:

$$p(\mathbf{z}_n, \mathbf{X}) = \mu_{f_n \rightarrow \mathbf{z}_n}(\mathbf{z}_n) \mu_{f_{n+1} \rightarrow \mathbf{z}_n}(\mathbf{z}_n) = \alpha(\mathbf{z}_n) \beta(\mathbf{z}_n)$$

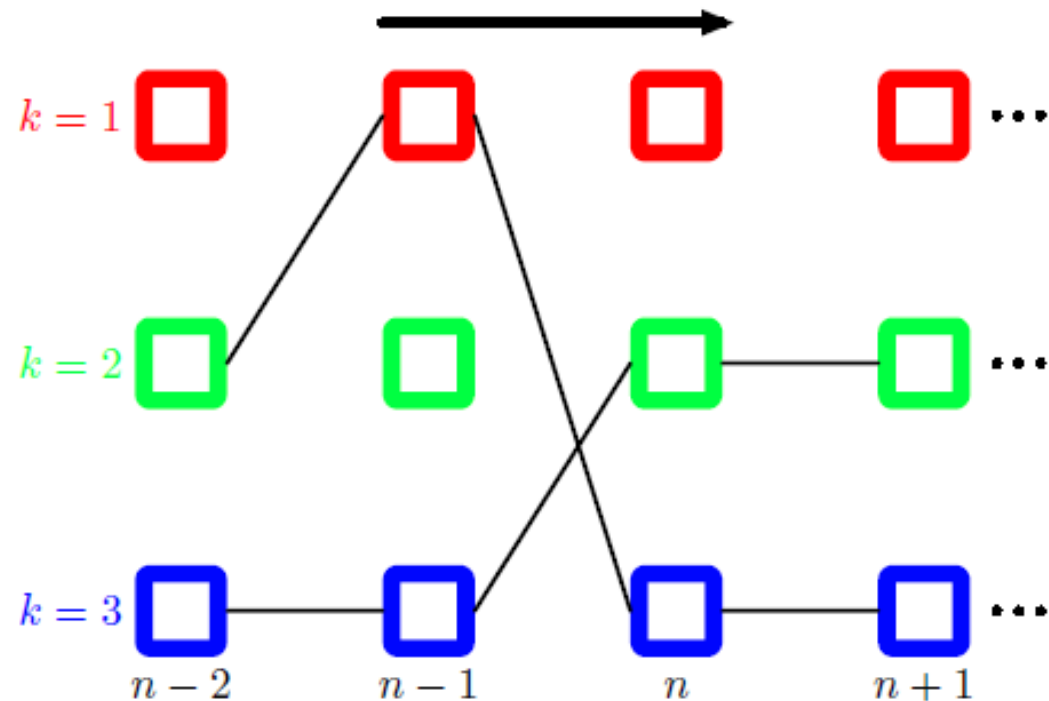
$$\gamma(\mathbf{z}_n) = \frac{p(\mathbf{z}_n, \mathbf{X})}{p(\mathbf{X})} = \frac{\alpha(\mathbf{z}_n) \beta(\mathbf{z}_n)}{p(\mathbf{X})}$$

$$\begin{aligned} \xi(\mathbf{z}_{n-1}, \mathbf{z}_n) &= p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}) \\ &= f(\mathbf{z}_{n-1}, \mathbf{z}_n) \mu_{\mathbf{z}_n \rightarrow f_n}(\mathbf{z}_n) \mu_{\mathbf{z}_{n-1} \rightarrow f_n}(\mathbf{z}_{n-1}) \\ &= f(\mathbf{z}_{n-1}, \mathbf{z}_n) \mu_{f_n \rightarrow \mathbf{z}_n}(\mathbf{z}_n) \mu_{f_{n-1} \rightarrow \mathbf{z}_{n-1}}(\mathbf{z}_{n-1}) \\ &= f(\mathbf{z}_{n-1}, \mathbf{z}_n) \alpha(\mathbf{z}_{n-1}) \beta(\mathbf{z}_n) \\ &= \frac{\alpha(\mathbf{z}_{n-1}) \beta(\mathbf{z}_n) p(\mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n)}{p(\mathbf{X})} \end{aligned}$$



Algorithms – Viterbi

In many applications of hidden Markov models, the latent variables have some meaningful interpretation, and so it is often of interest to find the most probable sequence of hidden states for a given observation sequence.



Algorithms – Viterbi

$$\begin{aligned}\mu_{\mathbf{z}_n \rightarrow f_{n+1}}(\mathbf{z}_n) &= \mu_{f_n \rightarrow \mathbf{z}_n}(\mathbf{z}_n) \\ \mu_{f_{n+1} \rightarrow \mathbf{z}_{n+1}}(\mathbf{z}_{n+1}) &= \max_{\mathbf{z}_n} \left\{ \ln f_{n+1}(\mathbf{z}_n, \mathbf{z}_{n+1}) + \mu_{\mathbf{z}_n \rightarrow f_{n+1}}(\mathbf{z}_n) \right\} \\ \omega(\mathbf{z}_{n+1}) &= \ln p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}) + \max_{\mathbf{z}_n} \{ \ln p(\mathbf{z}_{n+1} | \mathbf{z}_n) + \omega(\mathbf{z}_n) \} \\ \omega(\mathbf{z}_n) &\equiv \mu_{f_n \rightarrow \mathbf{z}_n}(\mathbf{z}_n)\end{aligned}$$

Where

$$\omega(\mathbf{z}_1) = \ln p(\mathbf{z}_1) + \ln p(\mathbf{x}_1 | \mathbf{z}_1)$$

By taking the logarithm and then exchanging maximizations and summations,

$$\omega(\mathbf{z}_n) = \max_{\mathbf{z}_1, \dots, \mathbf{z}_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_1, \dots, \mathbf{z}_n)$$

Contents

- Introduction
- Algorithms
 - Sum-Product
 - Forward-Backward
 - Veterbi
- Application: HMM in KG
- Conclusion

Application – HMM in KG

问题：给定训练好的模型，给定一句话，预测每个词对应的实体标签

输入：模型 $\lambda=(A,B,\Pi)$ ，观测序列 $O=(\text{浙}, \text{江}, \text{大}, \text{学}, \text{位}, \text{于}, \text{杭}, \text{州})$

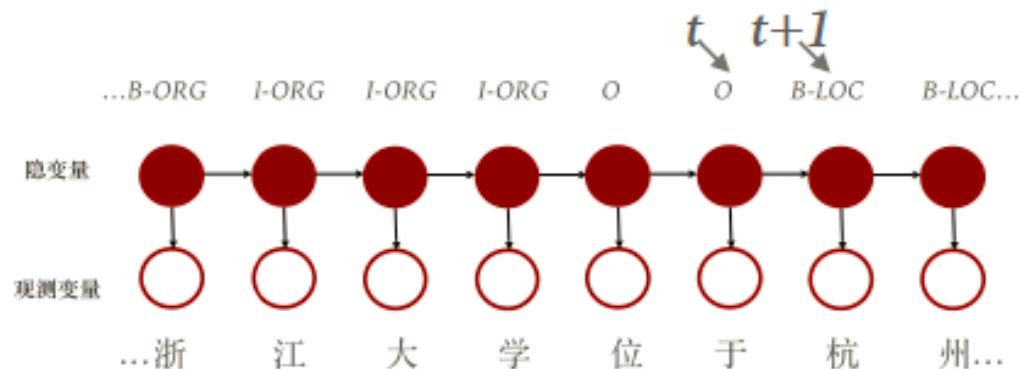
输出：最有可能的隐藏状态序列 $I=\{i_1, i_2, \dots, i_T\}$ ，即实体标签序列

动态规划算法的局部状态：在时刻 t 隐藏状态为 i 所有可能的状态转移路径 i_1, i_2, \dots, i_t 中的概率最大值

$$\delta_t(i) = \max_{i_1, i_2, \dots, i_{t-1}} P(i_t = i, i_1, i_2, \dots, i_{t-1}, o_t, o_{t-1}, \dots, o_1 | \lambda), \quad i = 1, 2, \dots, N$$

核心递推式：

$$\begin{aligned} \delta_{t+1}(i) &= \max_{i_1, i_2, \dots, i_t} P(i_{t+1} = i, i_1, i_2, \dots, i_t, o_{t+1}, o_t, \dots, o_1 | \lambda) \\ &= \max_{1 \leq j \leq N} [\delta_t(j) a_{ji}] b_i(o_{t+1}) \end{aligned}$$



Application – HMM in KG



问题：给定训练好的模型，给定一句话，预测每个词对应的实体标签

输入：模型 $\lambda=(A,B,\Pi)$ ，观测序列 $O=(\text{浙}, \text{江}, \text{大}, \text{学}, \text{位}, \text{于}, \text{杭}, \text{州})$

输出：最有可能的隐藏状态序列 $I=\{i_1, i_2, \dots, i_T\}$ ，即实体标签序列

1. 初始化局部状态

$$\delta_1(i) = \pi_i b_i(o_1), i = 1, 2 \dots N$$

$$\Psi_1(i) = 0, i = 1, 2 \dots N$$

时刻1，输出为 o_1 时，各个隐藏状态的可能性。

2. 进行动态规划递推时刻 $t=2, 3, \dots, T$ 时刻的局部状态

$$\delta_t(i) = \max_{1 \leq j \leq N} [\delta_{t-1}(j) a_{ji}] b_i(o_t), i = 1, 2 \dots N$$

$$\Psi_t(i) = \arg \max_{1 \leq j \leq N} [\delta_{t-1}(j) a_{ji}], i = 1, 2 \dots N$$

在 t 时刻，所有从 $t-1$ 时刻的状态 j 中，取最大概率。

从 $t-1$ 时刻的状态中，选择使 t 时刻概率最大的那个隐藏状态的编号

3. 如此递推，可计算最后时刻 T 最大的 $\delta_T(i)$ ，即为最可能隐藏状态序列出现的概率

$$P^* = \max_{1 \leq j \leq N} \delta_T(i)$$

4. 计算时刻 T 最大的 $\Psi_t(i)$ ，即为时刻 T 最可能的隐藏状态。

$$i_T^* = \arg \max_{1 \leq j \leq N} [\delta_T(i)]$$

5. 利用局部状态 $\Psi(i)$ 开始回溯，最终得到解码的序列，如：“...B-ORG, I-ORG, I-ORG, I-ORG, O, O, B-LOC, B-LOC...”。

Application

- Text summarization[1]
- Named Entity Recognition[2]
- Spectral Algorithm[3]
- ...

[1]Conroy J M, O'leary D P. Text summarization via hidden markov models[C]//Proceedings of the 24th annual international ACM SIGIR conference on Research and development in information retrieval. ACM, 2001: 406-407.

[2] Morwal S, Jahan N, Chopra D. Named entity recognition using hidden Markov model (HMM)[J]. International Journal on Natural Language Computing (IJNLC), 2012, 1(4): 15-23.

[3] Hsu D, Kakade S M, Zhang T. A spectral algorithm for learning hidden Markov models[J]. Journal of Computer and System Sciences, 2012, 78(5): 1460-1480.

Contents

- Introduction
- Algorithms
 - Sum-Product
 - Forward-Backward
 - Veterbi
- Application: HMM in KG
- Conclusion

Conclusion

- A method for learning sequential observations
- An Extraordinary way to extract latent features
- A fascinating example for belief propagation