

Machine Learning for Graphs and Sequential Data Exercise Sheet 06

Autoregressive Models, Markov Chains, Hidden Markov Models

Exercises marked with a (*) will be discussed in the in-person exercise session.

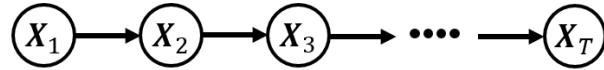
Problem 1: Consider the stationary AR(p) process $X_t = c + \sum_{i=1}^p \phi_i X_{t-i} + \epsilon$ with $\epsilon \sim \mathcal{N}(0, \sigma^2)$. We denote by μ the mean $E[X_t]$ and by γ_i the autocovariance $Cov(X_t, X_{t-i})$. Show:

1. $\mu = \frac{c}{1 - \sum_{i=1}^p \phi_i}$, for all t
2. $\gamma_0 = \sum_{j=1}^p \phi_j \gamma_{-j} + \sigma^2$
3. $\gamma_i = \sum_{j=1}^p \phi_j \gamma_{i-j}$, for all $t, i \in [1, p]$

Problem 2: (*) Let \mathbf{X}_t be a 2-D random vector:

$$\mathbf{X}_t = \begin{bmatrix} u_t \\ v_t \end{bmatrix} \quad \text{where } u_t, v_t \in \{1, 2, \dots, K\}. \quad (1)$$

Consider the following Markov chain.



Model parameters are as follows:

- initial distribution $\boldsymbol{\pi}_x \in \mathbb{R}^{K \times K}$ that parametrizes $\Pr(\mathbf{X}_1)$:

$$\Pr \left(\mathbf{X}_1 = \begin{bmatrix} i \\ j \end{bmatrix} \right) = \boldsymbol{\pi}_x(i, j). \quad (2)$$

- transition probability matrix $\mathbf{A}_x \in \mathbb{R}^{K \times K \times K \times K}$ that parametrizes $\Pr(\mathbf{X}_{t+1} | \mathbf{X}_t)$:

$$\Pr \left(\mathbf{X}_{t+1} = \begin{bmatrix} i_{t+1} \\ j_{t+1} \end{bmatrix} \mid \mathbf{X}_t = \begin{bmatrix} i_t \\ j_t \end{bmatrix} \right) = \mathbf{A}_x(i_t, j_t, i_{t+1}, j_{t+1}). \quad (3)$$

Because of the Markov property of \mathbf{X}_t , the joint probability can be factorized as

$$\Pr(\mathbf{X}_1, \dots, \mathbf{X}_T) = \Pr(\mathbf{X}_1) \prod_{t=1}^{T-1} \Pr(\mathbf{X}_{t+1} | \mathbf{X}_t).$$

In this task, we refer to this model as “2-D first-order Markov chain”.

- a) Does the sequence $[u_1, \dots, u_T]$ (where $u_t \in \{1, 2, \dots, K\}$) defined in Eq. (1) have the first-order Markov property? Why or why not?
 - b) Let $[Y_1, \dots, Y_T] \in \{1, 2\}^T$ be a first-order Markov chain with initial probability distribution $\boldsymbol{\pi}_y \in \mathbb{R}^2$ and transition probabilities $\mathbf{A}_y \in \mathbb{R}^{2 \times 2}$.
-

Recap

AR: Learn using Pseudo - In // Yule-walke

ML: - Markov property $p(z_t | z_{1:t-1}) = p(z_t | z_{t-1})$

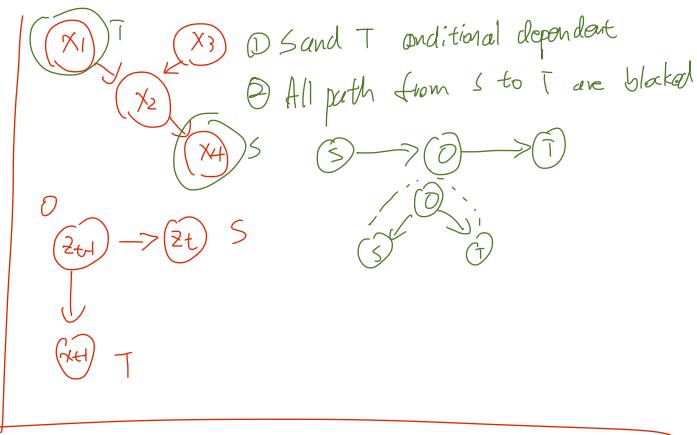
- MLE counting Transition frequentia!

HMM: - Filtering $p(z_t | x_{1:t})$ = forward

- Smoothing $p(z_t | x_{1:T})$ = forward - backward

- MAP $\underset{z_{1:T}}{\operatorname{argmax}} p(z_{1:T} | x_{1:T})$ vertebral

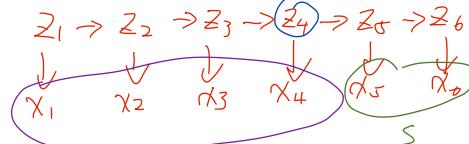
- Learning \Rightarrow EM



$$\text{Why } p(x_{t+1:T} | z_t = k, x_{1:t})$$

$$= p(x_{t+1:T} | z_t = k)$$

$$\neq p(x_{t+1:T} | x_{1:t})$$



Problem 1: Consider the stationary AR(p) process $X_t = c + \sum_{i=1}^p \phi_i X_{t-i} + \epsilon$ with $\epsilon \sim \mathcal{N}(0, \sigma^2)$. We denote by μ the mean $E[X_t]$ and by γ_i the autocovariance $Cov(X_t, X_{t-i})$. Show:

$$1. \mu = \frac{c}{1 - \sum_{i=1}^p \phi_i}, \text{ for all } t$$

$$2. \gamma_0 = \sum_{j=1}^p \phi_j \gamma_{-j} + \sigma^2$$

$$3. \gamma_i = \sum_{j=1}^p \phi_j \gamma_{i-j}, \text{ for all } t, i \in [1, p]$$

$$4. \text{For all } t, E(X_t) = E(X_{t-i})$$

$$E(X_t) = E(c) + \sum_{i=1}^p \phi_i E(X_{t-i}) + \cancel{E(\epsilon)} \\ = E(X_t)$$

$$E(X_t) \left[1 - \sum_{i=1}^p \phi_i \right] = c$$

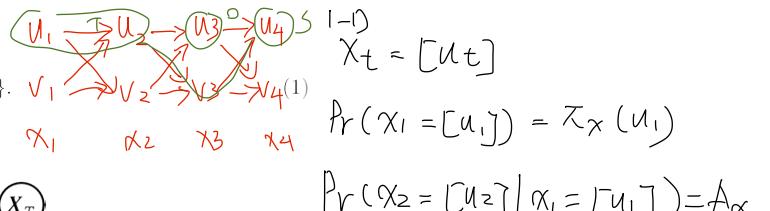
$$E(X_t) = \mu = \frac{c}{1 - \sum_{i=1}^p \phi_i}$$

$$2. \gamma_0 = Cov(X_t, X_t) = Var(X_t) = \underbrace{Var(c)}_{\sigma^2} + \underbrace{Var(\epsilon)}_{\sigma^2} + \underbrace{Var\left(\sum_{i=1}^p \phi_i X_{t-i}\right)}_{\sum_{i=1}^p \phi_i Cov(X_i, X_{t-i})} \\ = \sigma^2 + \sum_{j=1}^p \phi_j \cancel{\phi_j} \gamma_j = \gamma_j$$

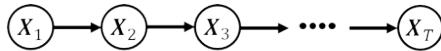
$$3. \gamma_i = Cov(X_t, X_{t-i}) = \underbrace{Cov(c, X_{t-i})}_{0} + \underbrace{Cov(\epsilon, X_{t-i})}_{0} + \sum_{j=1}^p \phi_j \underbrace{Cov(X_{t-j}, X_{t-i})}_{\gamma_{i-j}} \\ = \sum_{j=1}^p \phi_j \gamma_{i-j}$$

Problem 2: (*) Let \mathbf{X}_t be a 2-D random vector:

$$\mathbf{X}_t = \begin{bmatrix} u_t \\ v_t \end{bmatrix} \quad \text{where } u_t, v_t \in \{1, 2, \dots, K\}.$$



Consider the following Markov chain.



Model parameters are as follows:

- initial distribution $\pi_x \in \mathbb{R}^{K \times K}$ that parametrizes $\Pr(\mathbf{X}_1)$:

$$\Pr\left(\mathbf{X}_1 = \begin{bmatrix} i \\ j \end{bmatrix}\right) = \pi_x(i, j).$$

- transition probability matrix $\mathbf{A}_x \in \mathbb{R}^{K \times K \times K \times K}$ that parametrizes $\Pr(\mathbf{X}_{t+1} | \mathbf{X}_t)$:

$$\Pr\left(\mathbf{X}_{t+1} = \begin{bmatrix} i_{t+1} \\ j_{t+1} \end{bmatrix} \mid \mathbf{X}_t = \begin{bmatrix} i_t \\ j_t \end{bmatrix}\right) = \mathbf{A}_x(i_t, j_t, i_{t+1}, j_{t+1}). \quad (3)$$

Because of the Markov property of \mathbf{X}_t , the joint probability can be factorized as

$$\Pr(\mathbf{X}_1, \dots, \mathbf{X}_T) = \Pr(\mathbf{X}_1) \prod_{t=1}^{T-1} \Pr(\mathbf{X}_{t+1} | \mathbf{X}_t).$$

In this task, we refer to this model as “2-D first-order Markov chain”.

- a) Does the sequence $[u_1, \dots, u_T]$ (where $u_t \in \{1, 2, \dots, K\}$) is defined in Eq. (1) have the first-order Markov property? Why or why not?

- b) Let $[Y_1, \dots, Y_T] \in \{1, 2\}^T$ be a first-order Markov chain with initial probability distribution $\pi_y \in \mathbb{R}^2$ and transition probabilities $\mathbf{A}_y \in \mathbb{R}^{2 \times 2}$.

- Briefly explain why the sequence $\begin{bmatrix} Y_2 \\ Y_1 \end{bmatrix}, \begin{bmatrix} Y_3 \\ Y_2 \end{bmatrix}, \dots, \begin{bmatrix} Y_T \\ Y_{T-1} \end{bmatrix}$ is a 2-D first-order Markov chain.

- Compute initial and transition probabilities, π_y and \mathbf{A}_y (defined in Eqs. (2) and (3)) for the sequence $\begin{bmatrix} Y_2 \\ Y_1 \end{bmatrix}, \begin{bmatrix} Y_3 \\ Y_2 \end{bmatrix}, \dots, \begin{bmatrix} Y_T \\ Y_{T-1} \end{bmatrix}$. $Y_1 \rightarrow Y_2 \rightarrow Y_3 \rightarrow Y_4 \rightarrow Y_5$

$Y_{i=1}^T$ is first-order Markov chain

$$p(Y_T | Y_{T-1}, \dots, Y_1) = p(Y_T | Y_{T-1})$$

$$p\left(\begin{bmatrix} Y_1 \\ Y_{T-1} \end{bmatrix} \mid \begin{bmatrix} Y_{T-2} \\ Y_{T-1} \end{bmatrix}, \dots, \begin{bmatrix} Y_2 \\ Y_1 \end{bmatrix}\right)$$

because Y_T is only dependent on Y_{T-1}

$$= p\left(\begin{bmatrix} Y_1 \\ Y_{T-1} \end{bmatrix} \mid \begin{bmatrix} Y_{T-1} \\ Y_{T-2} \end{bmatrix}\right) \Leftrightarrow \text{2D-Markov chain}$$

$$\Pr(\mathbf{X}_1 = \begin{bmatrix} Y_2 \\ Y_1 \end{bmatrix} = \begin{bmatrix} i \\ j \end{bmatrix}) = \pi_x(Y_2, Y_1) = \pi_x(i, j)$$

$$? = \Pr(Y_2 = i | Y_1 = j) \cdot p(Y_1 = j)$$

$$= A_y(i, j) \cdot \pi_y(j)$$

$$\Pr(\mathbf{X}_{t+1} = \begin{bmatrix} Y_1 \\ Y_{T-1} \end{bmatrix} = \begin{bmatrix} i \\ j \end{bmatrix} \mid \mathbf{X}_t = \begin{bmatrix} Y_{T-1} \\ Y_{T-2} \end{bmatrix} = \begin{bmatrix} j \\ k \end{bmatrix}) = \mathbf{A}_y$$

$$= \Pr(Y_T = i \mid Y_{T-1} = j, Y_{T-2} = k) \Pr(Y_{T-1} = j \mid Y_{T-1} = j, Y_{T-2} = k)$$

$$= A_y(i, j) \cdot \cancel{A_y(j, k)} = 1$$

- Briefly explain why the sequence $\begin{bmatrix} Y_2 \\ Y_1 \end{bmatrix}, \begin{bmatrix} Y_3 \\ Y_2 \end{bmatrix}, \dots, \begin{bmatrix} Y_T \\ Y_{T-1} \end{bmatrix}$ is a 2-D first-order Markov chain.
- Compute initial and transition probabilities, π_x and \mathbf{A}_x (defined in Eqs. (2) and (3)) for the sequence $\begin{bmatrix} Y_2 \\ Y_1 \end{bmatrix}, \begin{bmatrix} Y_3 \\ Y_2 \end{bmatrix}, \dots, \begin{bmatrix} Y_T \\ Y_{T-1} \end{bmatrix}$.

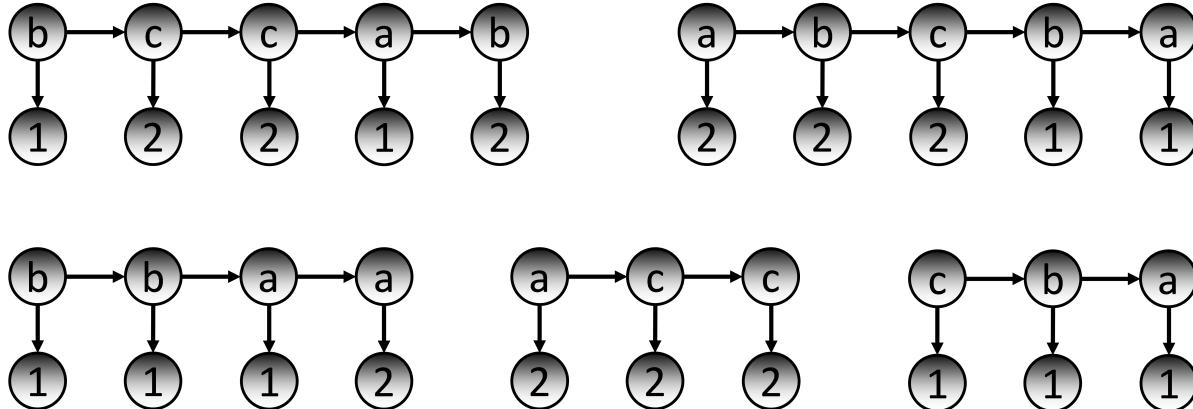
Problem 3: (*) Consider an HMM where hidden variables are in $\{1, 2\}$ and observed variables are in $\{a, b, c\}$. Let the model parameters be as follows:

$$A = \begin{matrix} & \begin{matrix} 1 & 2 \end{matrix} \\ \begin{matrix} 1 \\ 2 \end{matrix} & \begin{bmatrix} 0.2 & 0.8 \\ 0.5 & 0.5 \end{bmatrix} \end{matrix} \quad B = \begin{matrix} & \begin{matrix} a & b & c \end{matrix} \\ \begin{matrix} 1 \\ 2 \end{matrix} & \begin{bmatrix} 0.2 & 0 & 0.8 \\ 0.4 & 0.6 & 0 \end{bmatrix} \end{matrix} \quad \pi = \begin{matrix} & \begin{matrix} 1 \\ 2 \end{matrix} \\ & \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \end{matrix}$$

Assume that the sequence $X_{1:5} = [cabac]$ is observed.

1. Filtering: find the distribution $P(Z_3|X_{1:3})$.
2. Smoothing: find the distribution $P(Z_3|X_{1:5})$.
3. Viterbi algorithm: find the most probable sequence $[Z_1, \dots, Z_5]$.

Problem 4: Consider an HMM where states Z_t are in $\{a, b, c\}$ and emissions X_t are in $\{1, 2\}$. Given is the following set of fully-observed instances (two sequences of length 5, one sequence of length 4, and two sequences of length 3):



Learn the parameters of the HMM (i.e. $\pi \in \mathbb{R}^3$, $\mathbf{A} \in \mathbb{R}^{3 \times 3}$, and $\mathbf{B} \in \mathbb{R}^{3 \times 2}$) using maximum-likelihood estimation.

Problem 3: (*) Consider an HMM where hidden variables are in $\{1, 2\}$ and observed variables are in $\{a, b, c\}$. Let the model parameters be as follows:

$$A = \begin{bmatrix} 1 & 2 \\ 1 & 0.2 & 0.8 \\ 2 & 0.5 & 0.5 \end{bmatrix} \quad B = \begin{bmatrix} a & b & c \\ 1 & 0.2 & 0 & 0.8 \\ 2 & 0.4 & 0.6 & 0 \end{bmatrix} \quad \pi = \begin{bmatrix} 1 & 0.5 \\ 2 & 0.5 \end{bmatrix}$$

$Z_i \in \{1, 2\} \Leftrightarrow$ Markov property

$X_i \in \{a, b, c\} \Leftrightarrow$ time only depend

Assume that the sequence $X_{1:5} = [cabac]$ is observed.

1. Filtering: find the distribution $P(Z_3|X_{1:3})$.
2. Smoothing: find the distribution $P(Z_3|X_{1:5})$.
3. Viterbi algorithm: find the most probable sequence $[Z_1, \dots, Z_5]$.

$$1. P(Z_3|X_{1:3}) = \frac{P(Z_3, X_{1:3})}{\sum_{j=1}^k P(X_{1:3}, Z_3=j)}$$

$$\alpha_1(1) = P(Z_1=1, X_1) = P(X_1|Z_1=1) \cdot P(Z_1=1) = 0.8 \cdot 0.5 = 0.4$$

$$\begin{aligned} \alpha_2(1) &= P(Z_2=1, X_{1:2}) \\ &= P(X_2|Z_2=1, X_1) \cdot P(Z_2=1, X_1) \\ &= P(X_2=a|Z_2=1) \cdot \sum_{j=1}^k P(Z_2=1, Z_1=j, X_1) \\ &= 0.2 \cdot \sum_{j=1}^k P(Z_2=1|Z_1=j, X_1) \cdot P(Z_1=j, X_1) \\ &= 0.2 \cdot [P(Z_2=1|Z_1=1) \cdot P(Z_1=1, X_1=c) + P(Z_2=1|Z_1=2) \cdot P(Z_1=2, X_1=c)] \\ &\approx 0.2 \cdot [0.2 \cdot 0.4 + 0] \\ &= 0.016 \end{aligned}$$

$$\begin{aligned} \alpha_2(2) &= P(X_2=a|Z_2=2) \sum_j P(Z_2=2|Z_1=j) P(Z_1=j, X_1) \\ &= 0.4 [0.8 \cdot 0.4 + 0] \\ &\approx 0.128 \end{aligned}$$

$$\begin{aligned} \alpha_3 &= \beta_3(X_3) \odot A' \alpha_2 \\ &= \begin{bmatrix} 0 \\ 0.6 \end{bmatrix} \odot \left(\begin{bmatrix} 0.2 & 0.8 \\ 0.4 & 0.5 \end{bmatrix} \begin{bmatrix} 0.016 \\ 0.128 \end{bmatrix} \right) = \begin{bmatrix} 0 \\ 0.6 \end{bmatrix} \odot \begin{bmatrix} 0.672 \\ 0.768 \end{bmatrix} = \begin{bmatrix} 0 \\ 0.4608 \end{bmatrix} \end{aligned}$$

$$P(Z_3|X_{1:3}) = \frac{\alpha_3(1)}{\alpha_3(1) + \alpha_3(2)} = 0 // \frac{\alpha_3(2)}{\alpha_3(1) + \alpha_3(2)} = 1$$

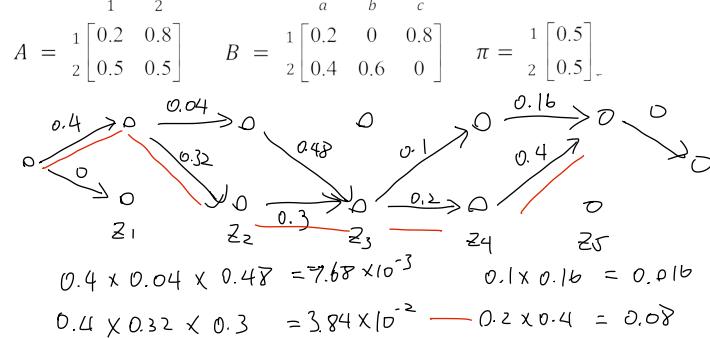
$$\begin{aligned} 2. P(Z_3|X_{1:5}) &= \beta_5(X_{1:5}) \\ \beta_5 &= 1 \\ \beta_4 &= \begin{bmatrix} 0.2 & 0.8 \\ 0.5 & 0.5 \end{bmatrix} \left(\begin{bmatrix} 0.8 \\ 0 \end{bmatrix} \odot \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right) \\ &= \begin{bmatrix} -0.16 \\ 0.4 \end{bmatrix} \\ \beta_3 &= \begin{bmatrix} 0.2 & 0.8 \\ 0.5 & 0.5 \end{bmatrix} \left(\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} \odot \begin{bmatrix} 0.16 \\ 0.4 \end{bmatrix} \right) \\ &= \begin{bmatrix} 0.2 & 0.8 \\ 0.5 & 0.5 \end{bmatrix} \begin{bmatrix} 0.032 \\ 0.016 \end{bmatrix} \\ &= \begin{bmatrix} 0.012 \\ 0.024 \end{bmatrix} \end{aligned}$$

cabac

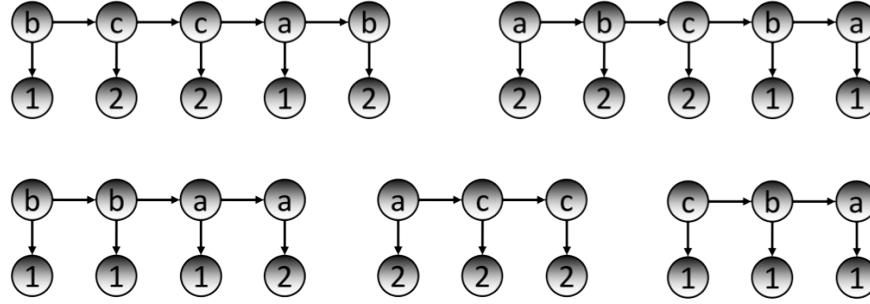
$$A = \begin{bmatrix} 1 & 2 \\ 1 & 0.2 & 0.8 \\ 2 & 0.5 & 0.5 \end{bmatrix} \quad B = \begin{bmatrix} a & b & c \\ 1 & 0.2 & 0 & 0.8 \\ 2 & 0.4 & 0.6 & 0 \end{bmatrix} \quad \pi = \begin{bmatrix} 1 & 0.5 \\ 2 & 0.5 \end{bmatrix}$$

$$P(Z_3|X_{1:5}) = \frac{\alpha_3(k)\beta_3(k)}{\sum \alpha_t(k)\beta_t(k)} = \frac{0}{0} + \frac{0.04608 \cdot 0.024}{0 + \dots} = 1$$

$$\begin{aligned}
3. \underset{z}{\operatorname{argmax}} p(z_{1:T} | x_{1:T}) &= \underset{z}{\operatorname{argmax}} \log p(z_{1:T}, x_{1:T}) \sim \underset{z}{\operatorname{argmax}} (\log p(x_1|z_1) p(z_1) + \sum \log p(x_t|z_t) p(z_t|z_{t-1})) \\
&\sim -\log p(z_1=j) p(x_1|z_1=j) \\
&= -\log \left[\frac{1}{2} \cdot 0.8 \right] \\
&\sim -\log p(z_t=j | z_{t-1}=i) p(x_t|z_t=j) \\
&\quad \vdots
\end{aligned}$$



Problem 4: Consider an HMM where states Z_t are in $\{a, b, c\}$ and emissions X_t are in $\{1, 2\}$. Given is the following set of fully-observed instances (two sequences of length 5, one sequence of length 4, and two sequences of length 3):



Learn the parameters of the HMM (i.e. $\pi \in \mathbb{R}^3$, $A \in \mathbb{R}^{3 \times 3}$, and $B \in \mathbb{R}^{3 \times 2}$) using maximum-likelihood estimation.

$$\begin{aligned}
&\max p(z_{1:T}, x_{1:T}) \\
&\max \log p(z_{1:T}, x_{1:T}) \\
&= \log p(z_1) + \sum \log p(z_t | z_{t-1}) + \sum \log p(x_t | z_t) \\
&= \sum_{k=1}^5 \sum_k \mathbb{I}(z_1=k) \log \pi_k + \sum_{t=2}^T \sum_{i,j} \mathbb{I}(z_t=j, z_{t-1}=i) \cdot \log A_{ij} + \sum_{t=1}^T \sum_{i,j} \mathbb{I}(x_t=j, z_t=i) \cdot \log B_{ij}
\end{aligned}$$

$$\sum_{K=1}^5 \sum_{k=1}^5 \mathbb{I}(z_1=k) \cdot \log \pi_k + \dots$$

$$\textcircled{1} \quad \text{s.t. } \pi_k \geq 0, \quad \sum_{k=1}^5 \pi_k = 1$$

$$\max_{\pi} \sum_{k=1}^5 \mathbb{I}(z_1=k) \log \pi_k + \lambda (\sum_{k=1}^5 \pi_k - 1)$$

$$\begin{aligned}
\frac{\partial}{\partial \pi_k} &= \sum_{n=1}^5 \mathbb{I}(z_1=k) \cdot \frac{1}{\pi_k} + \lambda = 0 \Rightarrow \pi_k = -\frac{\frac{5}{\lambda} \mathbb{I}(z_1=k)}{\lambda} \propto \begin{pmatrix} \sum \mathbb{I}(z_1=1)/N \\ \sum \mathbb{I}(z_1=2)/N \\ \sum \mathbb{I}(z_1=3)/N \end{pmatrix} \\
\textcircled{2} \quad \sum_j A(i,j) &= 1
\end{aligned}$$

$$\textcircled{3} \quad \sum_{j'} B(i', j') = 1$$