

Exercises on Hidden Markov Model

immediate

April 9, 2018

1 Demonstration of KL Divergence

- Convex function

$$f(\lambda a + (1 - \lambda)b) \leq \lambda f(a) + (1 - \lambda)f(b) \quad (\lambda \geq 0)$$

- By induction, we obtain the [Jensen's inequality](#)

$$f\left(\sum_{i=1}^M \lambda_i x_i\right) \leq \sum_{i=1}^M \lambda_i f(x_i)$$

where $\lambda_i \geq 0$ and $\sum_i \lambda_i = 1$ for any set of points $\{x_i\}$.

- Interpret λ_i as the probability distribution over a discrete variable x taking the values $\{x_i\}$

$$f(\mathbb{E}[x]) \leq \mathbb{E}[f(x)]$$

- For continuous variables, Jensen's inequality

$$f\left(\int \mathbf{x} p(\mathbf{x}) d\mathbf{x}\right) \leq \int f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

- $-\ln x$ is a convex function
- Apply Jensen's inequality to the KL divergence

$$KL(p||q) = - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x} \geq - \ln \int q(\mathbf{x}) d\mathbf{x} = 0$$

2 Sample

Q: Verify the conditional distribution for \mathbf{x}_n given all of the observations up to time n is given by

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

A: We first of all find the joint distribution $p(\mathbf{x}_1, \dots, \mathbf{x}_n)$ by marginalizing over the variables $\mathbf{x}_{n+1}, \dots, \mathbf{x}_N$, to give

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_n) &= \sum_{\mathbf{x}_{n+1}} \cdots \sum_{\mathbf{x}_N} p(\mathbf{x}_1, \dots, \mathbf{x}_N) \\ &= \sum_{\mathbf{x}_{n+1}} \cdots \sum_{\mathbf{x}_N} p(\mathbf{x}_1) \prod_{m=2}^N p(\mathbf{x}_m | \mathbf{x}_{m-1}) \\ &= p(\mathbf{x}_1) \prod_{m=2}^n p(\mathbf{x}_m | \mathbf{x}_{m-1}) \end{aligned}$$

Now we evaluate the required conditional distribution

$$p(\mathbf{x}_n | \mathbf{x}_1, \dots, \mathbf{x}_{n-1}) = \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n)}{\sum_{\mathbf{x}_n} p(\mathbf{x}_1, \dots, \mathbf{x}_n)} = \frac{p(\mathbf{x}_1) \prod_{m=2}^n p(\mathbf{x}_m | \mathbf{x}_{m-1})}{\sum_{\mathbf{x}_n} p(\mathbf{x}_1) \prod_{m=2}^n p(\mathbf{x}_m | \mathbf{x}_{m-1})}$$

Note that any factors which do not depend on \mathbf{x}_n will cancel between numerator and denominator, giving

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

3 Exercise 1

Q: Show that second-order Markov chain described by the joint distribution

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N) = p(\mathbf{x}_1) p(\mathbf{x}_2 | \mathbf{x}_1) \prod_{n=3}^N p(\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{x}_{n-2})$$

satisfies the conditional independent property

$$p(\mathbf{x}_n | \mathbf{x}_1, \dots, \mathbf{x}_{n-1}) = p(\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{x}_{n-2})$$

A: The marginal distribution over the variables $p(\mathbf{x}_1, \dots, \mathbf{x}_n)$ is given by

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_n) &= \sum_{\mathbf{x}_{n+1}} \cdots \sum_{\mathbf{x}_N} p(\mathbf{x}_1, \dots, \mathbf{x}_N) \\ &= \sum_{\mathbf{x}_{n+1}} \cdots \sum_{\mathbf{x}_N} p(\mathbf{x}_1) p(\mathbf{x}_2 | \mathbf{x}_1) \prod_{m=3}^N p(\mathbf{x}_m | \mathbf{x}_{m-1}, \mathbf{x}_{m-2}) \\ &= p(\mathbf{x}_1) p(\mathbf{x}_2 | \mathbf{x}_1) \prod_{m=3}^n p(\mathbf{x}_m | \mathbf{x}_{m-1}, \mathbf{x}_{m-2}) \end{aligned}$$

The required conditional distribution is then given by

$$p(\mathbf{x}_n | \mathbf{x}_1, \dots, \mathbf{x}_{n-1}) = \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n)}{\sum_{\mathbf{x}_n} p(\mathbf{x}_1, \dots, \mathbf{x}_n)} = \frac{p(\mathbf{x}_1) p(\mathbf{x}_2 | \mathbf{x}_1) \prod_{m=3}^n p(\mathbf{x}_m | \mathbf{x}_{m-1}, \mathbf{x}_{m-2})}{\sum_{\mathbf{x}_n} p(\mathbf{x}_1) p(\mathbf{x}_2 | \mathbf{x}_1) \prod_{m=3}^n p(\mathbf{x}_m | \mathbf{x}_{m-1}, \mathbf{x}_{m-2})}$$

Again, cancelling factors independent of \mathbf{x}_n between numerator and denominator we obtain

$$p(\mathbf{x}_n | \mathbf{x}_1, \dots, \mathbf{x}_{n-1}) = \frac{p(\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{x}_{n-2})}{\sum_{\mathbf{x}_n} p(\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{x}_{n-2})} = p(\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{x}_{n-2})$$

4 Exercise 2

Q: The joint probability distribution over both latent and observed variables of HMM is

$$p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) = p(\mathbf{z}_1 | \boldsymbol{\pi}) \left[\prod_{t=2}^T p(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{A}) \right] \prod_{t=1}^T p(\mathbf{x}_t | \mathbf{z}_t, \boldsymbol{\phi})$$

where $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$, $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_T\}$, and $\boldsymbol{\theta} = \{\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\phi}\}$ denotes the set of parameters governing the model. Given the marginal distribution

$$\begin{aligned} \gamma(z_{tk}) &= \mathbb{E}[z_{tk}] = \sum_{\mathbf{z}} \gamma(\mathbf{z}) z_{tk} \\ \xi(z_{t-1,i}, z_{tj}) &= \mathbb{E}[z_{t-1,i}, z_{tj}] = \sum_{\mathbf{z}} \gamma(\mathbf{z}) z_{t-1,i} z_{tj} \end{aligned}$$

The expectation of the logarithm of the complete-data likelihood function defined by

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{old}) &= \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}^{old}) \ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) \\ &= \sum_{k=1}^K \gamma(z_{1k}) \ln \pi_k + \sum_{t=2}^T \sum_{i=1}^K \sum_{j=1}^K \xi(z_{t-1,i}, z_{tj}) \ln A_{ij} + \sum_{t=1}^T \sum_{k=1}^K \gamma(z_{tk}) \ln p(\mathbf{x}_t | \boldsymbol{\phi}_k) \end{aligned}$$

A: In the M step, we maximize $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{old})$ with respect to the parameters $\boldsymbol{\theta} = \{\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\phi}\}$ in which we treat $\gamma(\mathbf{z}_t)$ and $\xi(\mathbf{z}_{t-1}, \mathbf{z}_t)$ as constant. Maximization with respect to $\boldsymbol{\phi}$ and \mathbf{A} is easily achieved using appropriate Lagrange multipliers.

Maximization with respects to $\boldsymbol{\pi}$: Take account of the summation constraint

$$\sum_{k=1}^K \pi_k = 1 \quad (1)$$

We there first omit terms from $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{old})$ which are independent of $\boldsymbol{\pi}$, and then add a Lagrange multiplier term to enforce the constraint, giving the following function to be maximized

$$\tilde{Q} = \sum_{k=1}^K \gamma(z_{1k}) \ln \pi_k + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) \quad (2)$$

Setting the derivative with respect to π_k equal to zero we obtain

$$0 = \gamma(z_{1k}) \frac{1}{\pi_k} + \lambda \quad (3)$$

We now multiply through by π_k and then sum over k and make use of the summation constraint to give

$$\lambda = - \sum_{k=1}^K \gamma(z_{1k}) \quad (4)$$

Substituting back into (2)

$$\tilde{Q} = \sum_{k=1}^K \gamma(z_{1k}) \ln \pi_k - \sum_{k=1}^K \gamma(z_{1k}) \left(\sum_{k=1}^K \pi_k - 1 \right) \quad (5)$$

and solving for π_k we obtain

$$\pi_k = \frac{\gamma(z_{1k})}{\sum_{i=1}^K \gamma(z_{1i})} \quad (6)$$

Maximization with respect to \mathbf{A} : For the maximization with respect to \mathbf{A} we follow the same steps and first omit terms from $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{old})$ which are independent of \mathbf{A} , and then add appropriate Lagrange multiplier terms to enforce the summation constraints. In this case there are K constraints to be satisfied since we must have

$$\sum_{j=1}^K A_{ij} = 1 \quad (7)$$

for $i = 1, \dots, K$. We introduce K Lagrange multipliers λ_i for $i = 1, \dots, K$, and maximize the following function

$$\hat{Q} = \sum_{t=2}^T \sum_{i=1}^K \sum_{j=1}^K \xi(z_{t-1,i}, z_{tj}) \ln A_{ij} + \sum_{i=1}^K \lambda_i \left(\sum_{j=1}^K A_{ij} - 1 \right) \quad (8)$$

Setting the derivative of \hat{Q} with respect to A_{ij} to zero we obtain

$$0 = \sum_{t=2}^T \xi(z_{t-1,i}, z_{tj}) \frac{1}{A_{ij}} + \lambda_i \quad (9)$$

Again we multiply through by A_{ij} and then sum over j and make use of the summation constraint to give

$$\lambda_i = - \sum_{t=2}^T \sum_{j=1}^K \xi(z_{t-1,i}, z_{tj}) \quad (10)$$

Substituting for λ_i in (8) and solving for A_{ij} we obtain

$$A_{ij} = \frac{\sum_{t=2}^T \xi(z_{t-1,i}, z_{tj})}{\sum_{k=1}^K \sum_{t=2}^T \xi(z_{t-1,i}, z_{tk})} \quad (11)$$

Maximization with respect to $\boldsymbol{\phi}$: To maximize $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{old})$ with respect to $\boldsymbol{\phi}_k$, we notice that only the final term depends on $\boldsymbol{\phi}_k$.

In the case of Gaussian emission densities we have $p(\mathbf{x}|\boldsymbol{\phi}_k) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$, and maximize the following function

$$\bar{Q} = \sum_{t=1}^T \sum_{k=1}^K \gamma(z_{tk}) \ln \mathcal{N}(\mathbf{x}_t | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (12)$$

$$= \sum_{t=1}^T \sum_{k=1}^K \gamma(z_{tk}) \left(\frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x}_t - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_t - \boldsymbol{\mu}_k) \right) \quad (13)$$

Setting the derivative of \bar{Q} with respect to $\boldsymbol{\mu}_k$ to zero we obtain

$$\sum_{t=1}^T \gamma(z_{tk}) \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_t - \boldsymbol{\mu}_k) = 0 \quad (14)$$

$$\Rightarrow \sum_{t=1}^T \gamma(z_{tk}) \mathbf{x}_t = \boldsymbol{\mu}_k \sum_{t=1}^T \gamma(z_{tk}) \quad (15)$$

$$\Rightarrow \boldsymbol{\mu}_k = \frac{\sum_{t=1}^T \gamma(z_{tk}) \mathbf{x}_t}{\sum_{t=1}^T \gamma(z_{tk})} \quad (16)$$

Next if we define

$$N_k = \sum_{t=1}^T \gamma(z_{tk}) \quad (17)$$

$$\hat{\mathbf{S}}_k = \sum_{t=1}^T \gamma(z_{tk}) (\mathbf{x}_t - \boldsymbol{\mu}_k) (\mathbf{x}_t - \boldsymbol{\mu}_k)^T \quad (18)$$

then we can rewrite \bar{Q} in the form

$$\bar{Q} = -\frac{N_k D}{2} \ln(2\pi) - \frac{N_k}{2} \ln |\boldsymbol{\Sigma}_k| - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_k^{-1} \hat{\mathbf{S}}_k) \quad (19)$$

Setting the derivative of \bar{Q} with respect to $\boldsymbol{\Sigma}_k^{-1}$ to zero we obtain

$$\frac{N_k}{2} \boldsymbol{\Sigma}_k^T - \frac{1}{2} \hat{\mathbf{S}}_k^T = 0 \quad (20)$$

$$\Rightarrow \boldsymbol{\Sigma}_k = \frac{\hat{\mathbf{S}}_k}{N_k} = \frac{\sum_{t=1}^T \gamma(z_{tk}) (\mathbf{x}_t - \boldsymbol{\mu}_k) (\mathbf{x}_t - \boldsymbol{\mu}_k)^T}{\sum_{t=1}^T \gamma(z_{tk})} \quad (21)$$