

# 隐马尔可夫模型

## 1 隐马尔可夫模型

### 1.1 模型定义

**定义：** 隐马尔可夫模型是关于时序的概率模型，描述一个由隐藏的马尔可夫链生成不可观测的状态随机序列，再有各个状态生成一个观测随机序列的过程。

其中，隐藏的马尔可夫链随机生成状态序列；每个状态生成一个观测，由此产生的观测随机序列，称为观测序列。

隐马尔可夫模型可由初始概率分布、状态转移概率分布和观测概率分布确定。可按如下数学形式进行严格定义：

设  $Q = \{q_1, q_2, \dots, q_N\}$  是所有状态的集合， $V = \{v_1, v_2, \dots, v_M\}$  是所有观测的集合。 $I = (i_1, i_2, \dots, i_T)$  是转态序列， $O = (o_1, o_2, \dots, o_T)$  是观测序列。

状态的转移矩阵为  $A = [a_{ij}]_{N \times N}$ ，其中  $a_{ij} = P(i_{t+1} = q_j | i_t = q_i), i = 1, 2, \dots, N, j = 1, 2, \dots, N$  表示在时刻  $t$  下不可观测的状态序列从  $q_i$  转移到  $q_j$  的概率。

状态观测转移的观测概率矩阵为  $B = [b_j(k)]_{N \times M}, j = 1, 2, \dots, N, k = 1, 2, \dots, M$ ，其中  $b_j(k) = P(o_t = v_k | i_t = q_j)$  表示在时刻  $t$  下不可观测的状态  $q_j$  转化为可观测的观测  $v_k$  的概率。

状态的初始概率向量为  $\pi = (\pi_i)$ ，其中  $\pi_i = P(i_1 = q_i), i = 1, 2, \dots, N$  表示  $t_1$  时刻下处于状态  $q_i$  的概率。

隐马尔可夫模型由初始状态概率向量  $\pi$ ，状态转移概率矩阵  $A$  和观测概率矩阵  $B$  决定。其中状态序列由  $\pi, A$  决定，观测序列由  $B$  决定。因此，隐马尔可夫模型可表示为三元组： $\lambda = (A, B, \pi)$ 。

隐马尔可夫模型具有以下两个基本假设：

1. 马尔可夫假设： $P(i_t | i_{t-1}, o_{t-1}, \dots, i_1, o_1) = P(i_t | i_{t-1})$
2. 观测独立性假设： $P(o_t | o_T, i_T, o_{T-1}, i_{T-1}, \dots, o_1, i_1) = P(o_t | i_t)$

隐马尔可夫模型具有以下三个基本问题：

1. **概率计算问题：** 给定模型  $\lambda = (A, B, \pi)$  和观测序列  $O = (o_1, o_2, \dots, o_T)$ ，计算在  $\lambda$  下观测序列  $O$  出现的概率  $P(O | \lambda)$

2. **学习问题**：已知观测序列  $O = (o_1, o_2, \dots, o_T)$ ，估计模型  $\lambda = (A, B, \pi)$  的参数
3. **预测问题**：也称为解码问题。已知模型  $\lambda = (A, B, \pi)$  和观测序列  $O = (o_1, o_2, \dots, o_T)$ ，求最有可能的状态序列  $I = (i_1, i_2, \dots, i_T)$

## 1.2 概率计算算法

### 1.2.1 直接计算法

对于给定的状态序列  $I = (i_1, i_2, \dots, i_T)$ ，其出现概率为

$$P(I|\lambda) = \pi_{i_1} a_{i_1 i_2} a_{i_2 i_3} \dots a_{i_{T-1} i_T}$$

对于固定的状态序列  $I = (i_1, i_2, \dots, i_T)$  和观测序列  $O = (o_1, o_2, \dots, o_T)$  出现的概率为

$$P(O|I) = b_{i_1}(o_1) b_{i_2}(o_2) \dots b_{i_T}(o_T)$$

因此， $O, I$  的联合概率为

$$P(O, I|\lambda) = P(O|I, \lambda) P(I|\lambda) = \pi_{i_1} b_{i_1}(o_1) a_{i_1 i_2} b_{i_2}(o_2) \dots a_{i_{T-1} i_T} b_{i_T}(o_T)$$

对所有  $I$  求和即可得到观测序列  $O$  出现的概率：

$$\begin{aligned} P(O|\lambda) &= \sum_I P(O, I|\lambda) \\ &= \sum_{i_1, i_2, \dots, i_T} \pi_{i_1} b_{i_1}(o_1) a_{i_1 i_2} b_{i_2}(o_2) \dots a_{i_{T-1} i_T} b_{i_T}(o_T) \end{aligned}$$

此计算方法的计算复杂度为  $O(TN^T)$ 。因此在实际操作中计算量过大，并不可行，需要使用前向-后向算法减少计算量。

### 1.2.2 前向算法

**定义**：给定隐马尔可夫模型  $\lambda$ ，定义到时刻  $t$  部分观测序列为  $o_1, o_2, \dots, o_t$  且状态为  $q_i$  的概率为**前向概率**，记为：

$$\alpha_t(i) = P(o_1, o_2, \dots, o_t, i_t = q_i | \lambda)$$

因此，可以获得以下递推式：

$$\begin{aligned}
\alpha_{t+1}(i) &= P(o_1, o_2, \dots, o_{t+1}, i_{t+1} = q_i | \lambda) \\
&= \sum_{j=1}^N P(o_1, o_2, \dots, o_{t+1}, i_{t+1} = q_i, i_t = q_j | \lambda) \\
&= \sum_{j=1}^N P(o_1, o_2, \dots, o_t, i_t = q_j | \lambda) P(i_{t+1} = q_i | i_t = q_j, \lambda) P(o_{t+1} | i_{t+1} = q_i, \lambda) \\
&= [\sum_{j=1}^N \alpha_t(j) a_{ji}] b_i(o_{t+1})
\end{aligned}$$

因此，可按如下步骤计算前向概率  $\alpha_t(i)$ ：

1. 初始化  $\alpha_1(i) = \pi_i b_i(o_1), i = 1, 2, \dots, N$
2. 对  $t = 1, 2, \dots, T-1$  递推：

$$\alpha_{t+1}(i) = [\sum_{j=1}^N \alpha_t(j) a_{ji}] b_i(o_{t+1}), i = 1, 2, \dots, N$$

3. 计算  $P(O|\lambda) = \sum_{i=1}^N \alpha_T(i)$

前向算法的计算复杂度是  $O(TN^2)$ ，其复杂度大幅度减少的原因是每一次计算直接引用了前一次计算的结果。

### 1.2.3 后向算法

**定义：**给定隐马尔可夫模型  $\lambda$ ，定义在时刻  $t$  状态为  $q_i$  的条件下，从  $t+1$  到  $T$  的部分观测序列  $o_{t+1}, o_{t+2}, \dots, o_T$  的概率为**后向概率**，记为：

$$\beta_t(i) = P(o_{t+1}, o_{t+2}, \dots, o_T | i_t = q_i, \lambda)$$

同样，对  $\beta_t(i)$  也可获得如下递推式：

$$\begin{aligned}
\beta_t(i) &= P(o_{t+1}, o_{t+2}, \dots, o_T | i_t = q_i, \lambda) \\
&= \sum_{j=1}^N P(o_{t+1}, o_{t+2}, \dots, o_T, i_{t+1} = q_j | i_t = q_i, \lambda) \\
&= \sum_{j=1}^N P(o_{t+2}, \dots, o_T, | i_{t+1} = q_j, \lambda) P(i_{t+1} = q_j | i_t = q_i, \lambda) P(o_{t+1} | i_{t+1} = q_j, \lambda) \\
&= \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)
\end{aligned}$$

因此，可按如下步骤计算后向概率  $\beta_t(i)$ ：

1. 初始化  $\beta_T(i) = 1, i = 1, 2, \dots, N$
2. 对  $t = T-1, T-2, \dots, 1$  递推:

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j), i = 1, 2, \dots, N$$

3. 计算  $P(O|\lambda) = \sum_{i=1}^N \pi_i b_i(o_1) \beta_1(i)$

后向算法的时间复杂度同样为  $O(TN^2)$ 。

利用前向算法和后向算法，还可以得到关于单个状态和两个状态的概率计算公式

1. 给定模型  $\lambda$  和观测  $O$ ，模型在时刻  $t$  处于状态  $q_i$  的概率定义为  $\gamma_t(i) = P(i_t = q_i | O, \lambda)$ 。可以得到:

$$\gamma_t(i) = P(i_t = q_i | O, \lambda) = \frac{P(i_t = q_i, O | \lambda)}{P(O | \lambda)} = \frac{P(i_t = q_i, O | \lambda)}{\sum_{j=1}^N P(i_t = q_j, O | \lambda)} = \frac{\alpha_t(i) \beta_t(i)}{\sum_{j=1}^N \alpha_t(j) \beta_t(j)}$$

2. 给定模型  $\lambda$  和观测  $O$ ，模型在时刻  $t$  处于状态  $q_i$  且在时刻  $t+1$  处于状态  $q_j$  的概率定义为  $\xi_t(i, j) = P(i_t = q_i, i_{t+1} = q_j | O, \lambda)$ 。可以得到:

$$\begin{aligned} \xi_t(i, j) &= P(i_t = q_i, i_{t+1} = q_j | O, \lambda) = \frac{P(i_t = q_i, i_{t+1} = q_j, O | \lambda)}{P(O | \lambda)} \\ &= \frac{P(i_t = q_i, i_{t+1} = q_j, O | \lambda)}{\sum_{i=1}^N \sum_{j=1}^N P(i_t = q_i, i_{t+1} = q_j, O | \lambda)} \\ &= \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)} \end{aligned}$$

### 1.3 学习算法

#### 1.3.1 监督学习算法

当状态序列和观测序列均是已知时，通过序列  $\{(O_1, I_1), (O_2, I_2), \dots, (O_S, I_S)\}$  使用极大似然估计可以简单得到  $\pi, A, B$  中的参数。

$$\hat{a}_{ij} = \frac{A_{ij}}{\sum_{i=1}^N \sum_{j=1}^N A_{ij}}, \hat{b}_j(k) = \frac{B_{jk}}{\sum_{k=1}^M B_{jk}}, \hat{\pi}_i = \frac{\Pi_i}{\sum_{j=1}^N \Pi_j}$$

其中， $A_{ij}$  表示从状态  $q_i$  转移到状态  $q_j$  的频数， $B_{jk}$  表示从状态  $b_j$  变化为观测  $o_k$  的频数， $\Pi_i$  表示首个状态为  $q_i$  的频数。

#### 1.3.2 Baum-Welch 算法

当只知道观测序列  $\{O_1, O_2, \dots, O_T\}$ ，而不知道状态序列  $\{I_1, I_2, \dots, I_T\}$  时，需要使用 EM 算法对参数进行估计。此时， $\{O_1, O_2, \dots, O_T\}$  是不完全数据，而  $\{(O_1, I_1), (O_2, I_2), \dots, (O_S, I_S)\}$  是完全数据。

算法的 E 步是确定 Q 函数  $Q(\lambda, \bar{\lambda})$ :

$$\begin{aligned} Q(\lambda, \bar{\lambda}) &= E_I[\log P(O, I|\lambda)|O, \bar{\lambda}] \\ &= \sum_I \log P(O, I|\lambda) P(I|O, \bar{\lambda}) \\ &= \sum_I \log P(O, I|\lambda) \frac{P(O, I|\bar{\lambda})}{P(O|\bar{\lambda})} \\ &\propto \sum_I \log P(O, I|\lambda) P(O, I|\bar{\lambda}) \end{aligned}$$

而  $P(O, I|\lambda) = \pi_{i_1} b_{i_1}(o_1) a_{i_1 i_2} b_{i_2}(o_2) \dots a_{i_{T-1} i_T} b_{i_T}(o_T)$ , 因此:

$$Q(\lambda, \bar{\lambda}) = \sum_I \log \pi_{i_1} P(O, I|\bar{\lambda}) + \sum_I \left( \sum_{t=1}^{T-1} \log a_{i_t i_{t+1}} \right) P(O, I|\bar{\lambda}) + \sum_I \left( \sum_{t=1}^T \log b_{i_t}(o_t) \right) P(O, I|\bar{\lambda})$$

算法的 M 步需要极大化上述 Q 函数  $Q(\lambda, \bar{\lambda})$ 。注意到  $\lambda$  的参数  $A, B, \pi$  分别位于 Q 函数的三项之中, 相互不影响, 因此只要分别极大化每一项即可。注意参数  $A, B, \pi$  满足一定的约束条件:  $\sum_{i=1}^N \pi_i = 1, \sum_{j=1}^N a_{ij} = 1, \sum_{k=1}^M b_j(k) = 1$ , 最大化时需要使用拉格朗日乘数法。得到的结果为:

$$\begin{aligned} \hat{\pi}_i &= \frac{P(O, i_1 = i|\bar{\lambda})}{P(O|\bar{\lambda})} = \gamma_1(i) \\ \hat{a}_{ij} &= \frac{\sum_{t=1}^{T-1} P(O, i_t = i, i_{t+1} = j|\bar{\lambda})}{\sum_{t=1}^{T-1} P(O, i_t = i|\bar{\lambda})} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)} \\ \hat{b}_j(k) &= \frac{\sum_{t=1}^T P(O, i_t = j|\bar{\lambda}) I(o_t = v_k)}{\sum_{t=1}^T P(O, i_t = j|\bar{\lambda})} = \frac{\sum_{t=1}^T \gamma_t(j) I(o_t = v_k)}{\sum_{t=1}^T \gamma_t(j)} \end{aligned}$$

EM 算法作用于隐马尔可夫模型的计算如上所示。其被称为 Baum-Welch 算法, 具体步骤如下:

**输入:** 观测数据  $O = (o_1, o_2, \dots, o_T)$ 。

1. 初始化, 选取  $a_{ij}^{(0)}, b_j(k)^{(0)}, \pi_i^{(0)}$ , 得到模型  $\lambda^{(0)} = (A^{(0)}, B^{(0)}, \pi^{(0)})$
2. 使用前向-后向算法计算  $\alpha_t(i), \beta_t(i)$
3. 计算  $\gamma_t(j), \xi_t(i, j)$ , 并更新参数得到  $\lambda^{(n)} = (A^{(n)}, B^{(n)}, \pi^{(n)})$ :

$$\begin{aligned} \hat{\pi}_i &= \frac{P(O, i_1 = i|\bar{\lambda})}{P(O|\bar{\lambda})} = \gamma_1(i) \\ \hat{a}_{ij} &= \frac{\sum_{t=1}^{T-1} P(O, i_t = i, i_{t+1} = j|\bar{\lambda})}{\sum_{t=1}^{T-1} P(O, i_t = i|\bar{\lambda})} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)} \\ \hat{b}_j(k) &= \frac{\sum_{t=1}^T P(O, i_t = j|\bar{\lambda}) I(o_t = v_k)}{\sum_{t=1}^T P(O, i_t = j|\bar{\lambda})} = \frac{\sum_{t=1}^T \gamma_t(j) I(o_t = v_k)}{\sum_{t=1}^T \gamma_t(j)} \end{aligned}$$

## 4. 重复 2 和 3 直至收敛

输出：隐马尔可夫模型参数  $\hat{A}, \hat{B}, \hat{\pi}$ 。

## 1.4 预测算法

## 1.4.1 近似算法

近似算法在每个时刻  $t$  分别选择该时刻最有可能出现的状态  $i_t^*$ ，从而得到一个状态序列  $I^* = (i_1^*, i_2^*, \dots, i_T^*)$ 。

由于  $\gamma_t(i) = \frac{P(i_t=q_i, O|\lambda)}{P(O|\lambda)}$ ，因此易于发现：

$$i_t^* = \arg \max_{1 \leq i \leq N} \gamma_t(i)$$

近似算法比较简单，但没有从整体上考虑状态序列。

## 1.4.2 维特比算法

维特比算法使用动态规划求概率最大的路径。其思想是，要使从  $i_1^*$  到  $i_T^*$  的路径最优，那么对于任意的  $2 \leq t \leq T-1$ ，从  $i_t^*$  到  $i_T^*$  的路径都是最优的。因此维特比算法从  $t = 1$  开始，递推地向后计算各条路径的最大概率，直到  $t = T$ 。在计算完最大概率后，再从  $i_T^*$  开始，由后向前求解节点  $i_{T-1}^*, \dots, i_1^*$ ，从而得出最优路径  $I^* = (i_1^*, i_2^*, \dots, i_T^*)$ 。

在向前计算最大概率时，需要递推进行计算以提高效率。定义时刻  $t$  时状态为  $i$  的所有单个路径  $(i_1, i_2, \dots, i_t)$  中的概率最大值为：

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

因此有

$$\begin{aligned} \delta_{t+1}(i) &= \max_{i_1, i_2, \dots, i_t} P(i_{t+1} = i, i_t, \dots, i_1, o_{t+1}, o_t, \dots, o_1 | \lambda) \\ &= \max_{1 \leq j \leq N} [\delta_t(j) a_{ji}] b_i(o_{t+1}), \quad i = 1, 2, \dots, N; t = 1, 2, \dots, T-1 \end{aligned}$$

在反向寻找最优节点时，也可定义在时刻  $t$  时状态为  $i$  的所有单个路径  $(i_1, i_2, \dots, i_{t-1}, i)$  中概率最大的路径的第  $t-1$  个节点为：

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

使用函数  $\delta$  和  $\Psi$ ，可以得到维特比算法，其步骤如下：

输入：模型  $\lambda = (A, B, \pi)$  和观测  $O = \{o_1, o_2, \dots, o_T\}$

1. 初始化

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

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

2. 向前递推，对  $t = 1, 2, \dots, T-1$ ，计算

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

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

3. 找出最大概率  $P^* = \max_{1 \leq i \leq N} \delta_T(i)$  并确定最终节点  $i_T^* = \arg \max_{1 \leq i \leq N} \delta_T(i)$

4. 反向找出最优节点，对  $t = T-1, T-2, \dots, 1$ ，计算

$$i_t^* = \Psi_{t+1}(i_{t+1}^*)$$

输出：最优路径  $I^* = (i_1^*, i_2^*, \dots, i_T^*)$ 。

## 2 代码实现

本次考虑的隐马尔可夫模型参数如下：

$$\pi = (0.2, 0.4, 0.4)^T, A = \begin{bmatrix} 0.5 & 0.2 & 0.3 \\ 0.3 & 0.5 & 0.2 \\ 0.2 & 0.2 & 0.6 \end{bmatrix}, B = \begin{bmatrix} 0.5 & 0.5 \\ 0.4 & 0.6 \\ 0.7 & 0.3 \end{bmatrix}$$

从该隐马尔可夫模型中随机生成 50 个长度为 50 的观测序列作为此次的数据。

### 2.1 sklearn 实现

sklearn 中的隐马尔可夫模型的模块已被独立到 hmmlearn 包中。使用 MultinomialHMM 函数可以模拟或者拟合离散的隐马尔可夫模型，该函数的主要属性有：

- startprob\_：初始状态概率向量
- transmat\_：状态转移矩阵
- emissionprob\_：观测概率矩阵

首先，根据隐马尔可夫模型的参数模拟出 50 个长度为 50 的观测序列：

```
import numpy as np
from hmmlearn.hmm import MultinomialHMM

np.random.seed(12345)
model = MultinomialHMM(n_components=3)
model.startprob_ = np.array([0.2, 0.4, 0.4])
model.transmat_ = np.array([[0.5, 0.2, 0.3],
                             [0.3, 0.5, 0.2],
                             [0.2, 0.3, 0.5]])
model.emissionprob_ = np.array([[0.5, 0.5],
                                 [0.4, 0.6],
                                 [0.7, 0.3]])
O = np.concatenate([model.sample(50)[0] for _ in range(50)],axis=0)
```

再创建一个新的类，对其进行拟合并输出结果：

```
remodel = MultinomialHMM(n_components=3)
remodel.fit(O,lengths=[50 for _ in range(50)])

print('初始状态概率向量:\n',remodel.startprob_,
      '\n状态转移矩阵:\n',remodel.transmat_,
      '\n观测概率矩阵:\n',remodel.emissionprob_)

## 初始状态概率向量:
## [0.29689949 0.43891086 0.26418965]
## 状态转移矩阵:
## [[0.34383868 0.29584149 0.36031983]
## [0.34050106 0.3076312 0.35186774]
## [0.34480482 0.29234502 0.36285016]]
## 观测概率矩阵:
## [[0.60855692 0.39144308]
## [0.22427814 0.77572186]
## [0.71926203 0.28073797]]
```