# 隐马尔可夫模型（Hidden Markov Model, HMM）.

SYSU-SPF 机器学习石开修班 讲义.　　　　付星宇. 中山大学数学学院 88 逸仙班

本讲义分为以下几个部分：

① — 例子引入.

② — HMM定义与经典问题.

③ — Forward 算法： $P(O|\lambda)$.

④ — Viterbi 算法： $I^* = \underset{I}{\arg\max}\, P(I|O,\lambda)$

⑤ — Baum-welch 算法： $\lambda^* = \underset{\lambda}{\arg\max}\, P(O|\lambda)$


## ① — 例子引入.

考虑一个人，他有两种状态(State)，状态集合 $I = \{$开心，不开心$\}$.
在每一个状态下，他会做出一些可被人们察觉的行为(Observation)，行为
集 $O = \{$睡觉，玩，工作$\}$. 现在开始，我们每天都观测这个人的行为,
假设观测了 100天，那么我们会得到这个人 100天的行为序列，但
是这100天此人每天的状态是不可直接测得的.
HMM 便是对这个例子的数学建模，我们希望通过对行为序列
$(O_1, O_2 \cdots O_{100})$ 的研究，<u>找出可能的隐藏状态序列.</u>
　　　　　　　　　　↳注：状态序列不可测，故 <u>Hidden</u>
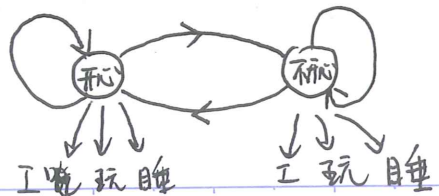
进一步假设：

— <u>齐次马氏性假设</u>. 即 $\forall t \in N^+$, $t \geq 2$，第t天的状态只与第 (t-1)天状态
　有关，i.e.

$$P(i_t = i \mid i_{t-1} = j, i_{t-2} = \cdots, \cdots, i_1 = \cdots)$$
$$= P(i_t = i \mid i_{t-1} = j) = a_{ji} \quad (\text{一个不依赖t的常数}).$$

— <u>发射(emission)无关假设</u>. 即 $\forall t \in N^+$，第t天的行为只与第t天的状态
　有关，i.e.

$$P(O_t = o \mid i_1 O_1\, i_2 O_2 \cdots i_t\, i_{t+1} O_{t+1} \cdots i_T O_T)$$
$$= P(O_t = o \mid i_t = j) = b_{jo} \quad (\text{一个不依赖t的常数}).$$

可以用一张图来描述以上的模型:

HMM 广泛运用在 交易择时 , 语音识别 , 基因测段等方向 . 据说 HMM
是美国著名对冲基金 文也复兴科技公司 的 大奖章基金背后的模型 .

② 一 HMM定义 与经典问题 .

设 $\{I_t \mid t \in \mathbb{N}^+\}$ 是一离散时间的马尔可夫过程 , $I_t \in I$ (状态空间) .
且 $|I| = n < +\infty$ . 该马尔可夫过程的转移矩阵为 $A$ , i.e.

$$A_{ij} = P(I_t = j \mid I_{t-1} = i) , \quad \forall i, j \in I .$$

且初始分布为 $\pi$ . i.e.

$$\pi_i = P(I_1 = i) \qquad \forall i \in I .$$

由随机过程的知识知 $(A, \pi)$ 完全决定了 $\{I_t \mid t \in \mathbb{N}^+\}$ .
此外在每一个时刻 $t$ 下 , 状态 $I_t$ 会发射 (emit) 一个观测 (observation) $O_t$ .
设 $O_t \in O$ ( 观测空间 ) , 且有发射矩阵 $B$ , i.e.

$$|O| = m .$$

$$B_{io} = P(O_t = o \mid I_t = i)$$

其中 $\{O_t \mid t \in \mathbb{N}^+\}$ 观测序列可见 , 而 $\{I_t \mid t \in \mathbb{N}^+\}$ 是 Hidden 的 .
令 $\lambda = (A, B, \pi)$ , 则 $\lambda$ 完全决定一个 HMM .

HMM 有以下三个重要问题:

— 给定 $\lambda = (A, B, \pi)$ , 给定观测 $O_T = (O_1, O_2 \cdots O_T)$ , 计算
这一观测的可能性 . i.e. $P(O_T \mid \lambda)$ . (用 Forward 算法求解) .

— 给定 $\lambda = (A, B, \pi)$ , 给定 观测序列 $O_T = (O_1, \cdots O_T)$ 计算
最可能导致这种观测的 状态链 $I_T^*$ . i.e. $I_T^* = \underset{I}{\arg\max} P(I \mid \lambda, O_T)$ .
(用 Viterbi 算法) .

— 给定 $O_T = (O_1, O_2, \cdots, O_T)$ , 给定 $|I| = n$ , $|O| = m$ , 找最可能的
系统参数 $\lambda = (A, B, \pi)$ . i.e. $\lambda^* = \underset{\lambda}{\arg\max} P(O \mid \lambda)$ .
(用 Baum-welch 算法) .

2.

③- Forward 算法.

给定 $\lambda = (A, B, \pi)$    给定 $O_T = (o_1, \cdots, O_T)$.

令 $\alpha_t(i) = P((o_1 o_2 \cdots o_t), i_t = q_i \mid \lambda)$    (称为 前向概率).

由全概率公式 $P(O_T \mid \lambda) = \sum_{i_T \in I} P(O_T, i_T = i \mid \lambda)$

$$= \sum_{i_T \in I} \alpha_T(i)$$

且注意到 $\alpha_{t+1}(i) = P((o_1, o_2 \cdots o_{t+1}), i_{t+1} = i \mid \lambda)$

$$= P(O_{t+1} \mid (o_1, o_2 \cdots o_t), i_{t+1} = i, \lambda) \times P((o_1, \cdots, o_t), i_{t+1} = i \mid \lambda)$$

(发射独立性) $= b_i O_{t+1} \times P((o_1 \cdots o_t) i_{t+1} = i \mid \lambda)$.

(全概) $= b_i O_{t+1} \times \sum_{j=1}^{N} P((o_1 \cdots o_t), i_t = j, i_{t+1} = i \mid \lambda)$

$$= b_i O_{t+1} \sum_{j=1}^{N} P(i_{t+1} = i \mid i_t = j, (o_1 \cdots o_t), \lambda) \times P(i_t = j, (o_1 \cdots o_t) \mid \lambda)$$

(马氏性) $= b_i O_{t+1} \sum_{j=1}^{N} a_{ji} \alpha_t(j)$.

故 我们可以一层一层, 从下往上, 递推计算 $\alpha$ 矩阵 $= \left[ \alpha_t(i) \right]_{T \times N}$.

其中第一行 $\alpha_1(i) = P(O_1, i_1 = i \mid \lambda)$

$$= P(O_1 \mid i_1 = i, \lambda) \times P(i_1 = i \mid \lambda)$$

$$= b_i O_1 \times \pi_i.$$

这便是 Forward 算法. 比传统的计算方法快很多.

3.

④- Viterbi 算法

给定 $\lambda = (A, B, \pi)$    给定 $O_T = (o_1 \cdots O_T)$.

定义 $\delta_t(i) = \max\limits_{(i_1 \cdots i_{t-1})} P\left((o_1 \cdots o_t), i_t = i, (i_1 \cdots i_{t-1}) \mid \lambda\right)$.

显然有 $\max\limits_{I} P(I \mid O_T, \lambda) = \max\limits_{I} \dfrac{P(I, O_T \mid \lambda)}{P(O_T \mid \lambda)}$

$= \dfrac{\max\limits_{I} P(I, O_T \mid \lambda)}{P(O_T \mid \lambda)} = \dfrac{\max\limits_{j \in I} \delta_T(j)}{P(O_T \mid \lambda)}$    (动态规划)

而 $\delta_{t+1}(i) = \max\limits_{(i_1 \cdots i_t)} P\left((o_1 \cdots o_{t+1}) \text{ 且 } i_{t+1} = i \text{ 且 } (i_1 \cdots i_t) \mid \lambda\right)$

$= \max\limits_{(i_1 \cdots i_t)} b_{i o_{t+1}} \times P\left((o_1 \cdots o_t)\ i_{t+1} = i\ (i_1 \cdots i_t) \mid \lambda\right)$

$= \max\limits_{(i_1 \cdots i_t)} b_{i o_{t+1}} \times a_{i_t i} P\left((o_1 \cdots o_t)\ (i_1 \cdots i_t) \mid \lambda\right)$

(动态规划) $= \max\limits_{j}\left(\max\limits_{(i_1 \cdots i_{t-1})} a_{ji} \times b_{i o_{t+1}} P\left((o_1 \cdots o_t)\ (i_1 \cdots i_{t-1} j) \mid \lambda\right)\right)$

$= \max\limits_{j} b_{i o_{t+1}} a_{ji} \delta_t(j)$

故我们可以 (递推) 地 计算 最大似然 状态列.

4.

① — Baum-welch算法.

We want to maximize the log-likelihood function

$l(\lambda) = \log p(O|\lambda)$. By EM Algorithm, we can

maximize over $Q(\lambda, \bar{\lambda}) = \sum_I p(O, I|\bar{\lambda}) \times \log p(O, I|\lambda)$

while $p(O, I|\lambda) = \pi_{i_1} b_{i_1}(O_1) a_{i_1 i_2} b_{i_2}(O_2) \cdots a_{i_{T-1} i_T} b_{i_T}(O_T)$.

E Step:

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

$$+ \sum_I \left( \sum_{t=1}^{T} \log b_{i_t}(O_t) \right) p(O, \lambda, |\bar{\lambda})$$

M step:

Use Lagrange multiplier and set partial derivative

to zero.

Details in P182 of '统计学习方法'.