

隐马尔可夫模型

机器学习研究室

计算机科学与技术学院
吉林大学

目 录

- HMM的由来
- 马尔可夫性和马尔可夫链
- HMM实例
- HMM的三个基本算法
- 主要参考文献

HMM的由来

- 俄国数学家A.A.马尔可夫于1907年第一次提出马尔科夫模型
 - 马尔可夫模型
 - 马尔可夫链
 - 隐马尔可夫模型

Andrei A Markov





- **Gifted Russian mathematician**
- **June 14, 1856 - July 20, 1922**
- **St Petersburg University**
- **Doctoral Advisor: Pafnuty Chebyshev**

http://en.wikipedia.org/wiki/Andrey_Markov

Andrei A Markov

- **Research field:**
 - number theory
 - continuous fraction theory
 - differential equations
 - probability theory
 - statistics
- **Markov is best known for his work in probability and for stochastic processes especially **Markov chains**.**

- 
- 
- At the age of 30 Markov became a professor at St. Petersburg University and a member of St. Petersburg Academy of Sciences
 - More than 120 scientific papers
 - Classical textbook “**Calculus of Probabilities**”

Pafnuty Chebyshev

- One of the founding fathers of Russian mathematics
- His famous students :
 - Dmitry Grave (格雷夫)
 - Aleksandr Korkin,
 - Aleksandr Lyapunov
 - Andrei Markov
 -
- According to the Mathematics Genealogy Project, Chebyshev has about 4 000 mathematical descendants.
- His famous work in the field of probability, statistics and number theory



独立事件概率

- 设想我们再作一连串的实验，而每次实验所可能发生的结果定为 $E_1, E_2, \dots, E_n, \dots$ 。（可能是有限，也可能是无限）。每一个结果 E_k ，我们若能给定一个出现的可能性 p_k （即概率），则对某一特定样本之序列 $E_{j1} E_{j2} \dots E_{jn}$ ，我们知道它出现的概率是 $p(E_{j1} E_{j2} \dots E_{jn}) = p_{j1} \dots p_{jn}$ 。

非独立事件概率

- 常用的统计中，彼此相互「独立」是最有用的一个概念。简单直白地说，互相「独立」就是彼此毫不相干，一点牵涉都没有。
- 但是实际生活中很多事件是相互关联的。
- [非独立]也就是说互相关联的意思，但是要怎样相关呢？如何在相关中作一些简单的分类呢？
- 马可夫连就是要描述在「相关」这个概念中最简单的一种。但即使如此，有关马可夫链的理论已经相当丰富了。在概率理论中，它占了很大一部分内容。

非独立事件概率

- 在马可夫链中我们考虑最简单的「相关」性。在在这种情况下，我们不能给任一个事件 E_j 一个概率 p_j 。但我们给一对事件 (E_j, E_k) 一个概率 p_{jk} ，这个时候 p_{jk} 的解释是一种条件概率，就是假设在某次实验中 E_j 已经出现，而在下一次实验中 E_k 出现的概率。除了 p_{jk} 之外，我们还需要知道第一次实验中 E_j 出现的机率 a_j 。有了这些资料后，一个样本序列 $E_{j0} E_{j1} \dots E_{jn}$ （也就是说第零次实验结果是 E_{j0} ，第一次是 E_{j1} ……第 n 次实验是 E_{jn} ）的概率就很清楚的是 $P(E_{j0}, E_{j1}, E_{jn})$
$$= a_j p_{j0j1} p_{j1j2} \dots p_{jn-1jn}.$$

马尔可夫性

- 如果一个过程的“将来”仅依赖“现在”而不依赖“过去”，则此过程具有**马尔可夫性**，或称此过程为**马尔可夫过程**
- $X(t+1) = f(X(t))$

马尔科夫链

- 时间和状态都离散的马尔科夫过程称为马尔科夫链
- 记作 $\{X_n = X(n), n = 0, 1, 2, \dots\}$
 - 在时间集 $T = \{0, 1, 2, \dots\}$ 上对离散状态的过程相继观察的结果
 - 链的状态空间记做 $I = \{a_1, a_2, \dots\}, a_i \in R.$
- 条件概率 $P_{ij}(m, m+n) = P\{X_{m+n} = a_j | X_m = a_i\}$ 称为马氏链在时刻 m 处于状态 a_i 条件下, 而在时刻 $m+n$ 转移到状态 a_j 的转移概率。

Example of Markov Chain

Weather: A Markov Model



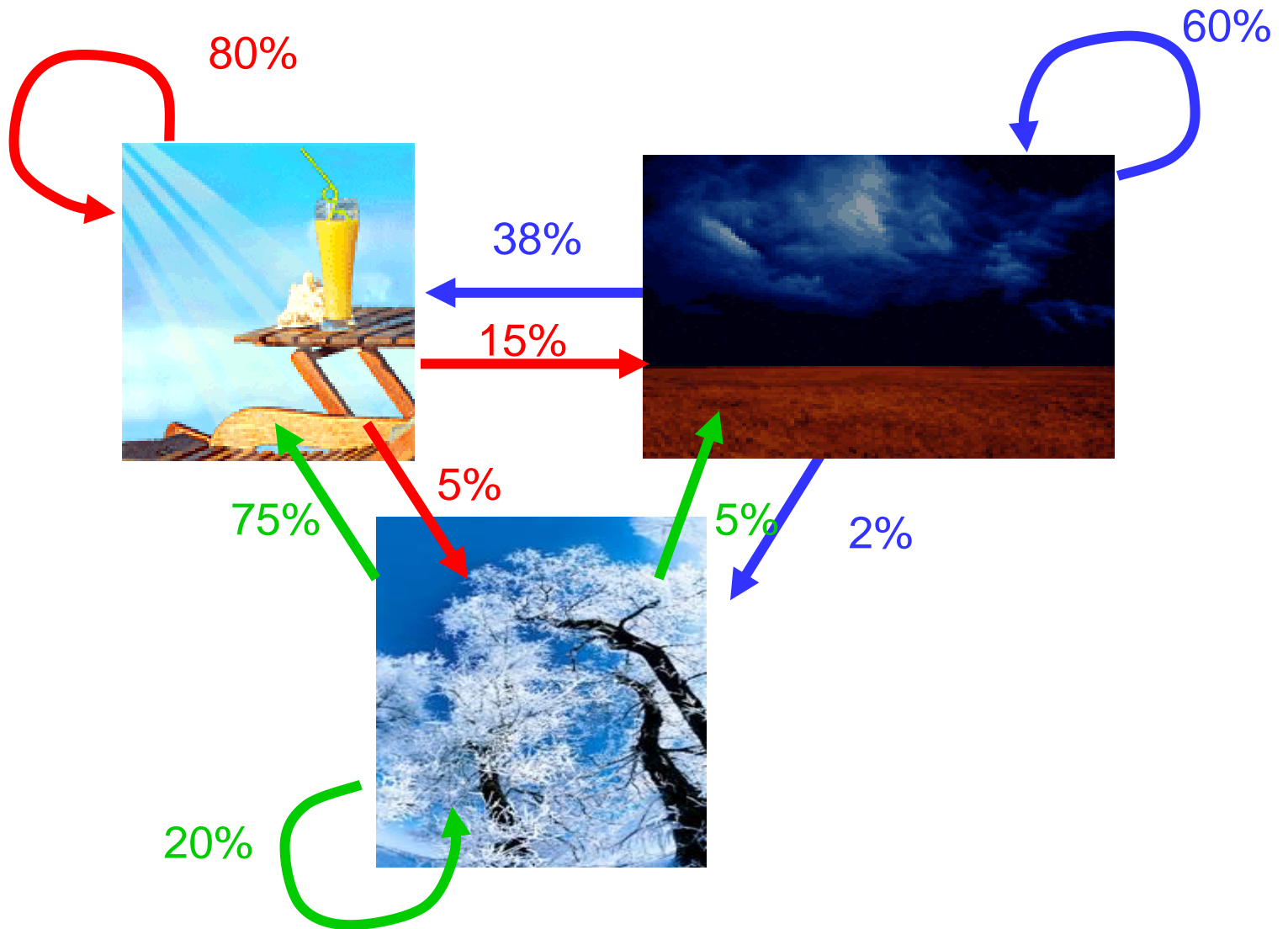
Sunny



Rainy



Snowy



Markov Model

- **States:**

$$\{S_1, S_2, \dots, S_N\}$$

- **State transition probabilities:**

$$a_{ij} = P(q_{t+1} = S_j \mid q_t = S_i)$$

represents the probability of moving from state i to state j

$A = \{a_{ij}\}$: transition probability matrix

$$0 \leq a_{ij} \leq 1; \sum_{j=1}^N a_{ij} = 1, 1 \leq i \leq N$$

Markov Model

- Initial state distribution:

$$\pi_i = P(q_1 = S_i)$$

$$0 \leq \pi_i \leq 1, 1 \leq i \leq N; \sum_{j=1}^N \pi_i = 1$$

Markov Model

- Markov Model= (π, A)
 - Begins ($t=1$) in some initial state (s).
 - At each time step ($t=1, 2, \dots$) the system moves from **current** to **next** state (possibly the same as the current state) according to **transition probabilities** associated with **current** state.

HMM的应用领域

- 语音识别
- 机器视觉
 - 人脸检测
 - 机器人足球
- 图像处理
 - 图像去噪
 - 图像识别
- 生物医学分析
 - DNA/蛋白质序列分析

Weather Markov Model

- **States:**

$$\{S_{\text{sunny}}, S_{\text{rainy}}, S_{\text{snowy}}\}$$


- **State transition probabilities:**

$$A = \begin{pmatrix} .8 & .15 & .05 \\ .38 & .6 & .02 \\ .75 & .05 & .2 \end{pmatrix}$$

- **Initial state distribution:**

$$\pi = (.7 \quad .25 \quad .05)$$

Basic Calculations-1

- Given the weather on the first day  , what is the probability that the weather for consecutive six days



t=1

Sunny



t=2

Rainy



t=3

Rainy



t=4

Rainy



t=5

Snowy



t=6

Snowy

Basic Calculations-1

$$\begin{aligned} &P(S_{\text{sunny}}, S_{\text{rainy}}, S_{\text{rainy}}, S_{\text{rainy}}, S_{\text{snowy}}, S_{\text{snowy}} \mid \text{Model}) \\ &= P(S_{\text{sunny}}) \cdot P(S_{\text{rainy}} \mid S_{\text{sunny}}) \cdot P(S_{\text{rainy}} \mid S_{\text{rainy}}) \cdot P(S_{\text{rainy}} \mid S_{\text{rainy}}) \\ &\quad \cdot P(S_{\text{snowy}} \mid S_{\text{rainy}}) \cdot P(S_{\text{snowy}} \mid S_{\text{snowy}}) \\ &= 0.7 \cdot 0.15 \cdot 0.6 \cdot 0.6 \cdot 0.02 \cdot 0.2 = 0.0001512 \end{aligned}$$

$$\pi = (.7 \quad .25 \quad .05)$$

$$A = \begin{pmatrix} .8 & .15 & .05 \\ .38 & .6 & .02 \\ .75 & .05 & .2 \end{pmatrix}$$

Basic Calculations-2



- Given that the system is in a known weather S_i , what is the probability that it stays in the same weather for consecutive d days: e.g.

$$Q = \{ \underbrace{s_i, s_i, s_i, \dots, s_i}_d, s_j, i \neq j \}$$

$$p(Q | Model, q_1 = s_i) = p(q_1 = s_i, Q | Model) / p(q_1 = s_i)$$

$$= \sum_{j=1, j \neq i}^N p(q_1 = s_i, \{ \underbrace{s_i, s_i, s_i, \dots, s_i}_d, s_j \} | Model) / p(q_1 = s_i)$$

$$= p(q_1 = s_i) (p(s_i | s_i))^{d-1} \sum_{j=1, j \neq i}^N p(s_j | s_i) / p(q_1 = s_i)$$




$$= a_{ii}^{d-1} (1 - a_{ii}) = p_i(d)$$

Basic Calculations-3

- Conditioned on starting the weather , compute the expected number of duration in weather

无穷级数的求解

$$\overline{d}_i = \sum_{d=1}^{\infty} dp_i(d) = \sum_{d=1}^{\infty} da_{ii}^{d-1} (1 - a_{ii}) = \frac{1}{1 - a_{ii}}$$

- Expected number of consecutive  days is 5.
- Expected number of consecutive  days is 2.5.
- Expected number of consecutive  days is 1.25.

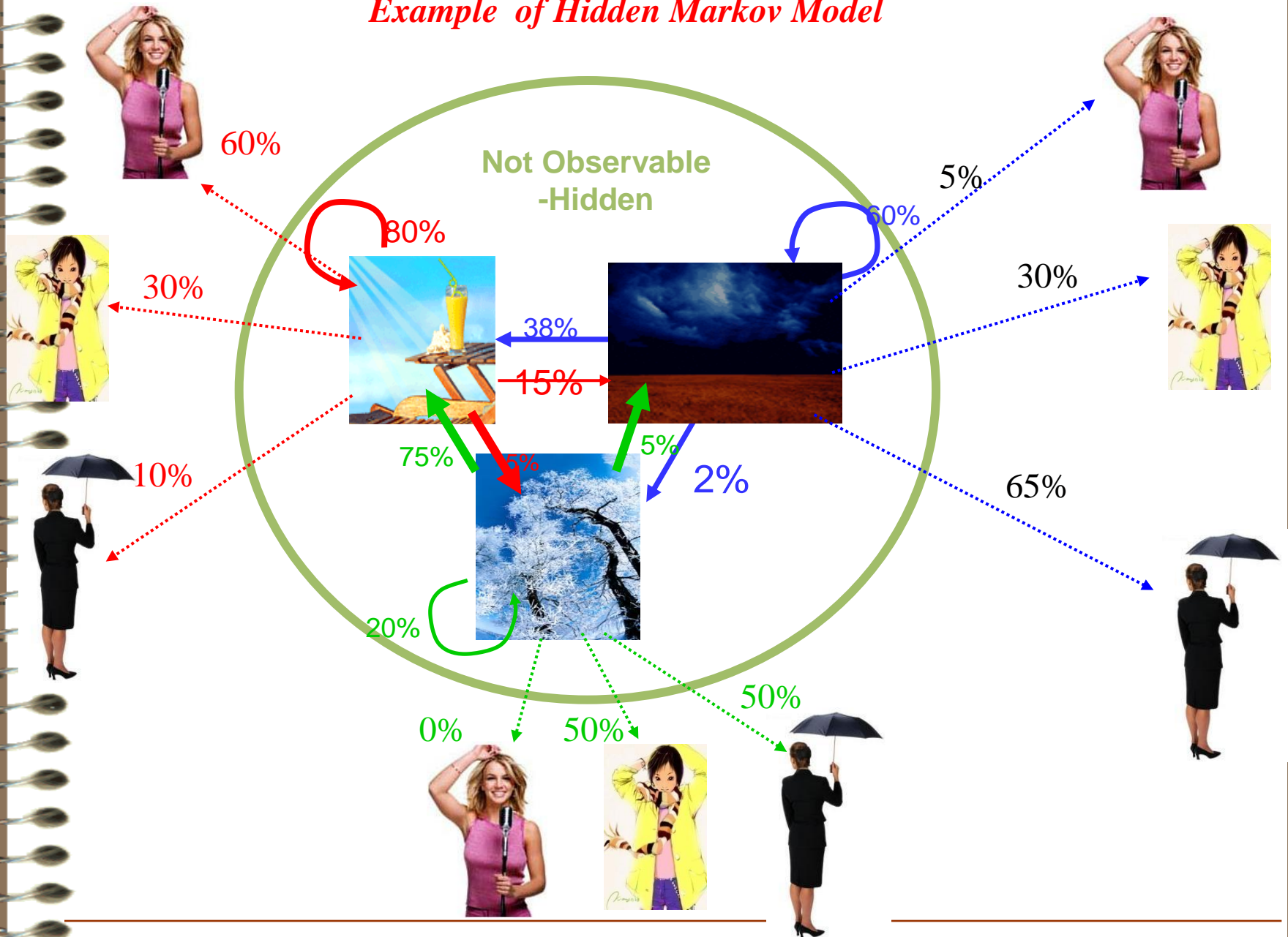
注：等比数列通项 $A_n = A_1 \cdot q^{(n-1)}$ ，前n项和 $S_n = [A_1(1 - q^n)] / (1 - q)$

A spiral-bound notebook with a brown cover is shown from a top-down perspective. A bright yellow sticky note is attached to the right side of the notebook. The sticky note has rounded corners and a small tab on the right edge. On the left side of the sticky note, there are two small, light green circular marks, one near the top-left corner and one near the bottom-left corner, resembling punch holes or decorative elements. The main text on the sticky note is in a bold, red, sans-serif font, and the subtitle is in a blue, sans-serif font.

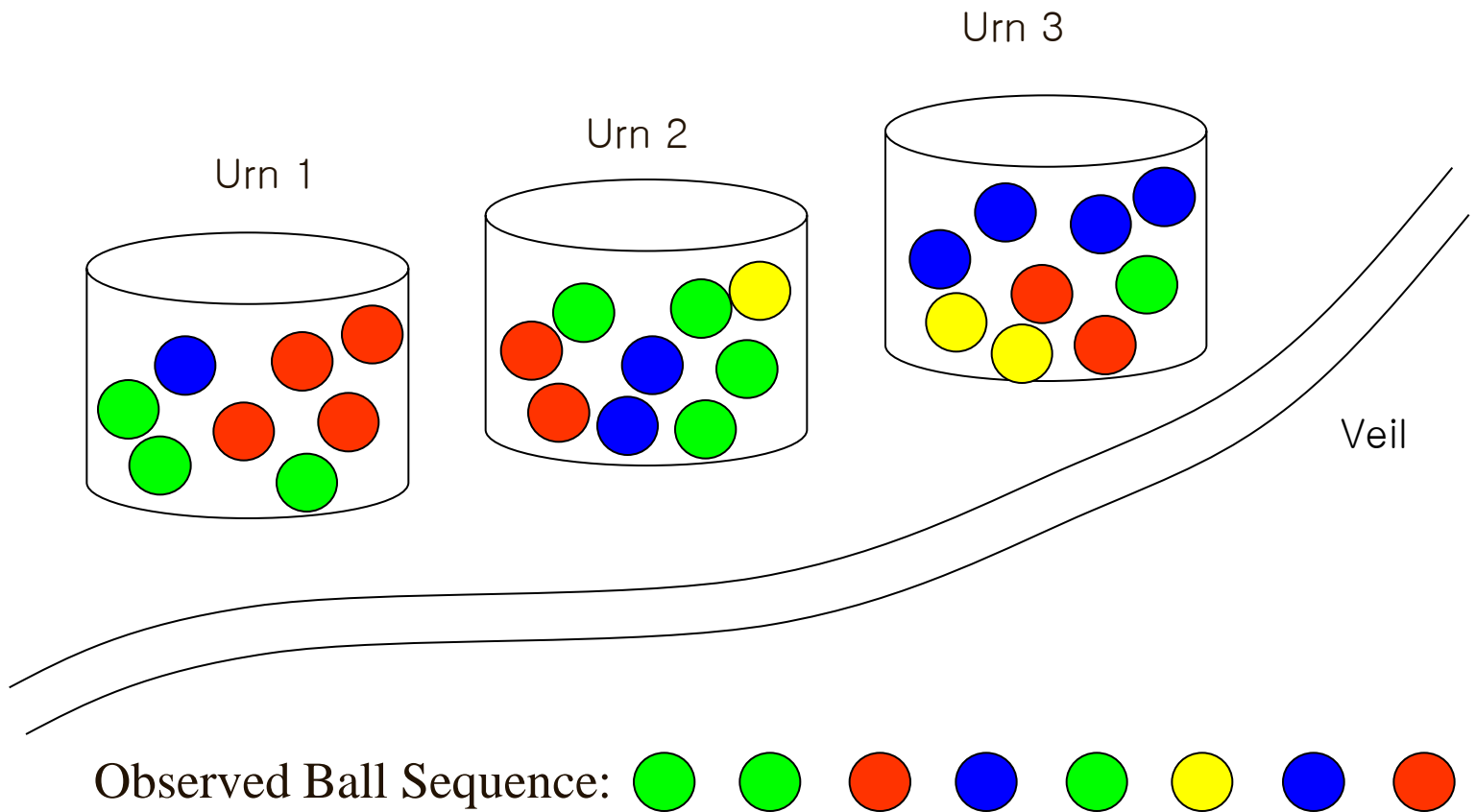
**The World is not
fully observable!!!**

Hidden Markov Model

Example of Hidden Markov Model

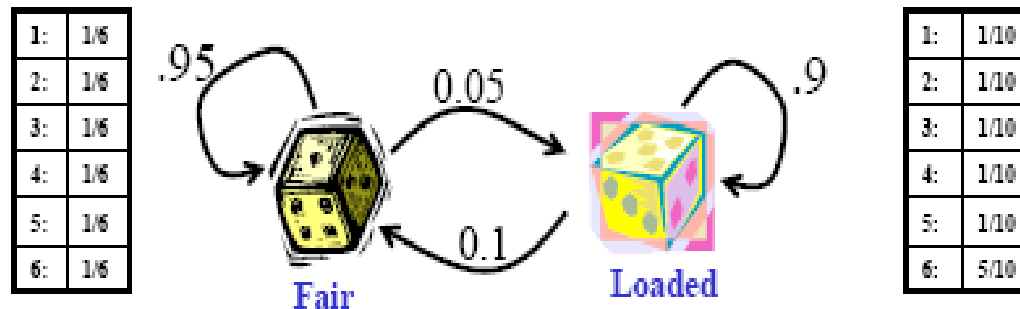


HMM实例



又一个HMM实例

- Occasionally dishonest casino



Simulted a sequence of 100 “die usages”

[illegible]

Question: what's the long-run fraction of each type?

However, we only observe the result of the rolls of a die: can we tell?

HMM实例—描述

- 设有 N 个缸，每个缸中装有很多彩球，球的颜色由一组概率分布描述。实验进行方式如下
 - 根据初始概率分布，随机选择 N 个缸中的一个开始实验
 - 根据缸中球颜色的概率分布，随机选择一个球，记球的颜色为 O_1 ，并把球放回缸中
 - 根据描述缸的转移的概率分布，随机选择下一口缸，重复以上步骤。
- 最后得到一个描述球的颜色序列 O_1, O_2, \dots ，称为观察值序列 O 。

HMM实例——约束

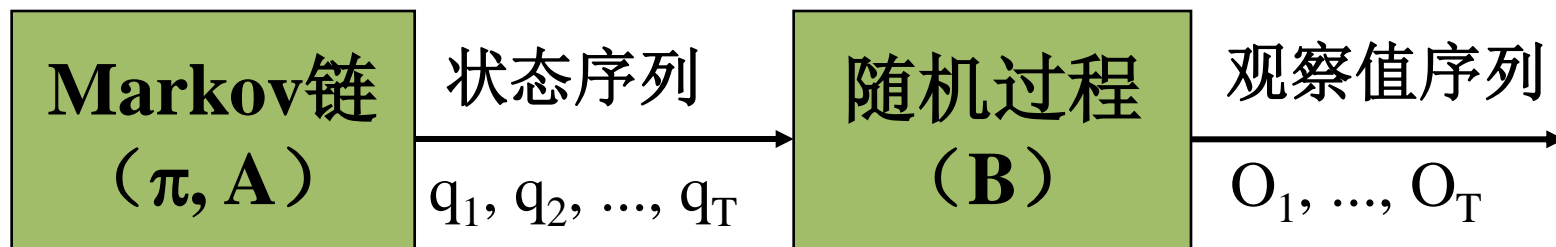
在上述实验中，有几个要点需要注意：

- 不能被直接观察缸间的转移
- 从缸中所选取的球的颜色和缸并不是一一对应的
- 每次选取哪个缸由一组转移概率决定

HMM概念

- HMM的状态是不确定或不可见的，只有通过观测序列的随机过程才能表现出来
- 观察到的事件与状态并不是一一对应，而是通过一组概率分布相联系
- HMM是一个双重随机过程，两个组成部分：
 - **马尔可夫链**：描述状态的转移，用**转移概率**描述。
 - **一般随机过程**：描述状态与观察序列间的关系，用**观察值概率**描述。

HMM组成



HMM的组成示意图

HMM的基本要素

- 用五元组 $\lambda = (N, M, \pi, A, B)$ 来描述HMM，或简写为 $\lambda = (\pi, A, B)$

| 参数 | 含义 | 实例 |
|-------|----------------|----------------------|
| N | 状态数目 | 缸的数目 |
| M | 每个状态可能的观察值数目 | 彩球颜色数目 |
| A | 与时间无关的状态转移概率矩阵 | 在选定某个缸的情况下，选择另一个缸的概率 |
| B | 给定状态下，观察值概率分布 | 每个缸中的颜色分布 |
| π | 初始状态空间的概率分布 | 初始时选择某口缸的概率 |

HMM可解决的问题

- 问题1：给定观察序列 $O=O_1, O_2, \dots, O_T$, 以及模型 $\lambda=(\pi, A, B)$, 如何计算 $P(O|\lambda)$? 计算问题
- 问题2：给定观察序列 $O=O_1, O_2, \dots, O_T$ 以及模型 λ , 如何选择一个对应的状态序列 $S = q_1, q_2, \dots, q_T$, 使得 S 能够最为合理的解释观察序列 O ? 解码问题
- 问题3：如何调整模型参数 $\lambda=(\pi, A, B)$, 使得 $P(O|\lambda)$ 最大? 学习问题

解决问题1 基础方法

- 给定一个固定的状态序列 $S=(q_1, q_2, q_3 \dots)$

$$P(O | S, \lambda) = \prod_{t=1}^T P(O_t | q_t, \lambda) = b_{q_1}(O_1) b_{q_2}(O_2) \cdots b_{q_T}(O_T)$$

其中, $b_{q_t}(O_t)$ 表示在 q_t 状态下观测到 O_t 的概率

$$P(O | \lambda) = \sum_{\text{所有 } S} P(O | S, \lambda) P(S | \lambda)$$

- $N=5, M=100, \Rightarrow$ 计算量 10^{72}

计算量过大

解决问题1 前向法

- 动态规划
- 定义前向变量

第t个状态为i

$$\alpha_t(i) = P(O_1, O_2, \dots, O_t, q_t = \theta_i \mid \lambda) \quad 1 \leq t \leq T$$

– 初始化:

$$\alpha_1(i) = \pi_i b_i(O_1) \quad 1 \leq i \leq N$$

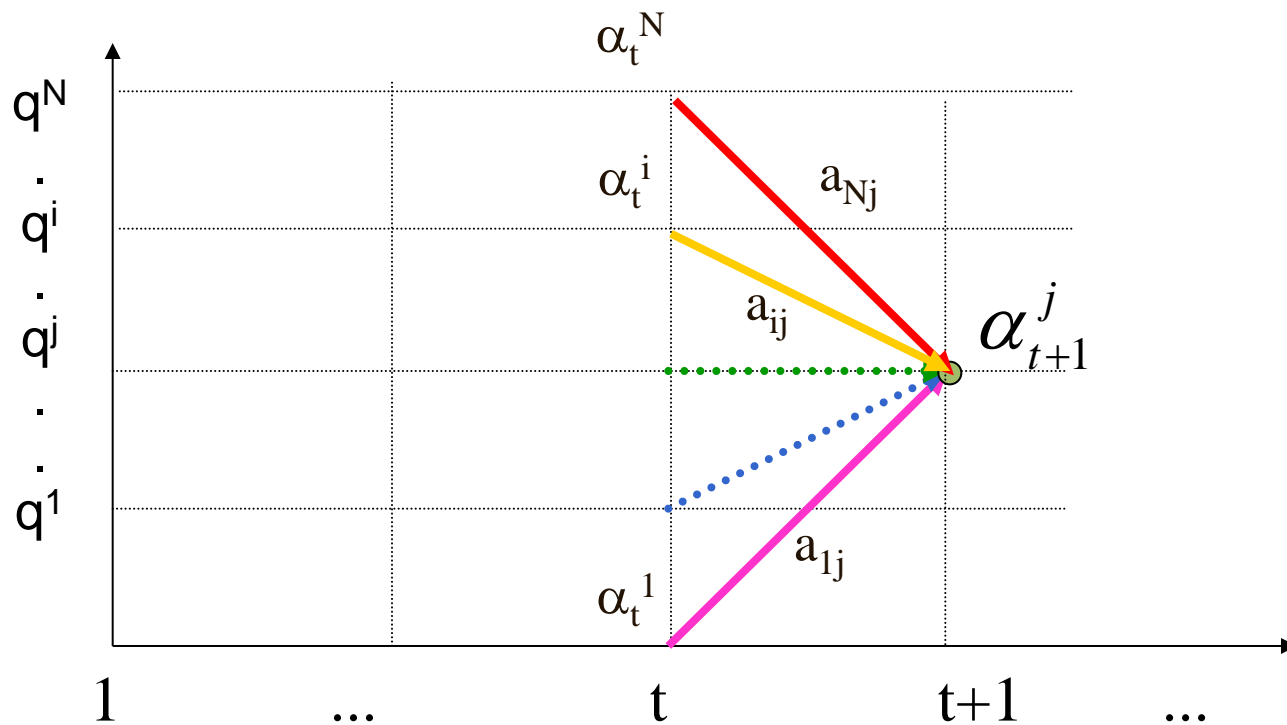
– 递归:

$$\alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(O_{t+1}) \quad 1 \leq t \leq T-1, 1 \leq j \leq N$$

– 终结:

$$P(O \mid \lambda) = \sum_{i=1}^N \alpha_T(i)$$

前向法示意图



$N=5, M=100, \Rightarrow$ 计算量3000

解决问题1 后向法

观测规律:

$$\begin{aligned}
 p(o_1 \dots o_T | \lambda) &= \sum_{i=1}^N \sum_{q_2 \dots q_T} p(o_1 \dots o_T, q_1 = s_i, q_2 \dots q_T | \lambda) \\
 &= \sum_{i=1}^N \pi_i b_i(o_1) \sum_{q_2 \dots q_T} p(o_2 \dots o_T, q_2 \dots q_T | q_1 = s_i, \lambda)
 \end{aligned}$$

(o1...ot already generated)

Starting from at time t state s_i
Generating $o_{t+1} \dots o_T$

算法设计:

$$\beta_t(i) = \sum_{q_{t+1} \dots q_T} p(o_{t+1} \dots o_T, q_{t+1} \dots q_T | q_t = s_i, \lambda)$$

$$\begin{aligned}
 &= \sum_{q_{t+1} \dots q_T} p(o_{t+2} \dots o_T, q_{t+2} \dots q_T | q_{t+1}, \lambda) p(q_{t+1} | q_t = s_i, \lambda) p(o_{t+1} | q_{t+1}, \lambda) \\
 &= \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \sum_{q_{t+2} \dots q_T} p(o_{t+2} \dots o_T, q_{t+2} \dots q_T | q_{t+1} = s_j, \lambda) \\
 &= \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)
 \end{aligned}$$

解决问题1 后向法

- The data likelihood is

$$\begin{aligned} p(o_1 \dots o_T \mid \lambda) &= \sum_{i=1}^N \pi_i b_i(o_1) \beta_1(i) \\ &= \sum_{i=1}^N \alpha_1(i) \beta_1(i) = \sum_{i=1}^N \alpha_t(i) \beta_t(i) \quad \text{for any } t \end{aligned}$$

Complexity:

$$O(TN^2)$$

解决问题1 后向法

- 与前向法类似
- 定义后向变量

$$\beta_t(i) = P(O_{t+1}, O_{t+2}, \dots, O_T, q_t = \theta_i | \lambda) \quad 1 \leq t \leq T-1$$

– 初始化:

$$\beta_T(i) = 1 \quad 1 \leq i \leq N$$

– 递归:

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(O_{t+1}) \beta_{t+1}(j) \quad t = T-1, T-2, \dots, 1, 1 \leq i \leq N$$

– 终结:

$$P(O | \lambda) = \sum_{i=1}^N \pi_i b_i(O_1) \beta_1(i)$$

解决问题2 Viterbi算法

- 目的：给定观察序列O以及模型 λ , 如何选择一个对应的状态序列S, 使得S能够最为合理的解释观察序列O?
 - N和T分别为状态个数和序列长度
- 定义：

$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} P(q_1 q_2 \cdots q_{t-1}, q_t = i, O_1, O_2, \dots, O_t / \lambda)$$

我们所要找的，就是T时刻最大的 $\delta_T(i)$ 所代表的那个状态序列

Viterbi算法(续)

- 初始化:

$$\delta_1(i) = \pi_i b_i(O_1), \quad 1 \leq i \leq N$$

$$\varphi_1(i) = 0, \quad 1 \leq i \leq N$$

- 递归:

$$\delta_t(j) = \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] b_j(O_t), \quad 2 \leq t \leq T, 1 \leq j \leq N$$

$$\varphi_t(j) = \arg \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}], \quad 2 \leq t \leq T, 1 \leq j \leq N$$

- 终结:

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

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

- 求S序列:

$$q_t^* = \varphi_{t+1}(q_{t+1}^*), \quad t = T-1, T-2, \dots, 1$$

Baum-Welch算法(模型训练算法)

- 目的：给定观察值序列 O ，通过计算确定一个模型 λ ，使得 $P(O|\lambda)$ 最大。
- 算法步骤：
 1. 初始模型（待训练模型） λ_0 ，
 2. 基于 λ_0 以及观察值序列 O ，训练新模型 λ ；
 3. 如果 $\log P(X|\lambda) - \log(P(X|\lambda_0)) < \text{Delta}$ ，说明训练已经达到预期效果， 算法结束。
 4. 否则，令 $\lambda_0 = \lambda$ ，继续第2步工作

Baum-Welch算法(续)

- 定义:

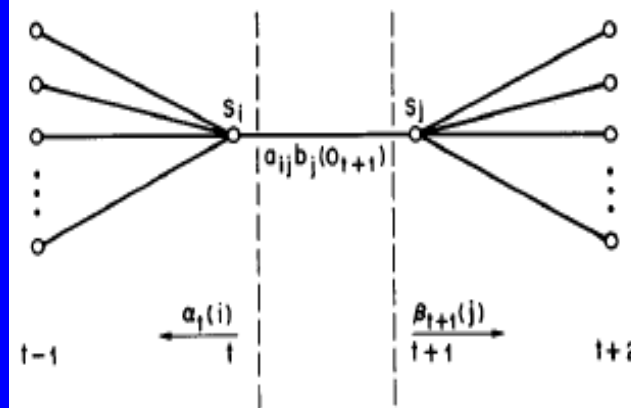
给定模型 λ 和观察序列条件下, 从 i 到 j 的转移概率定义为 $\xi_t(i, j)$

$$\begin{aligned}\xi_t(i, j) &= P(s_t = i, s_{t+1} = j | X, \lambda) \\ &= \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(x_{t+1}) \beta_{t+1}(j)}\end{aligned}$$

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j) \quad t\text{时刻处于状态} S_i \text{的概率}$$

$$\sum_{t=1}^{T-1} \gamma_t(i) = \text{整个过程中从状态} S_i \text{转出的次数 (number of time) 的预期}$$

$$\sum_{t=1}^{T-1} \xi_t(i, j) = \text{从} S_i \text{跳转到} S_j \text{次数的预期}$$



Baum-Welch算法(续2)

- 参数估计:

Reestimate :

$$\begin{aligned}\hat{a}_{ij} &= \frac{\text{expected count of transitions from } i \text{ to } j}{\text{expected count of stays at } i} \\ &= \frac{\sum_t \xi_t(i, j)}{\sum_t \sum_j \xi_t(i, j)}\end{aligned}$$

$$\begin{aligned}\hat{b}_j(k) &= \frac{\text{expected number of times in state } j \text{ and observing symbol } k}{\text{expected number of times in state } j} \\ &= \frac{\sum_{t, O_t=k} \gamma_t(j)}{\sum_t \gamma_t(j)}\end{aligned}$$

$$\pi_i = \text{当 } t=1 \text{ 时处于 } S_i \text{ 的概率} = \gamma_1(i)$$

基于HMM的CpG岛识别

- 指DNA上一个区域，此区域含有大量相联的胞嘧啶（C）、鸟嘌呤（G），以及使两者相连的磷酸酯键（p）。哺乳类基因中的启动子上，含有约40%的CpG岛（人类约70%）。一般CpG岛的长度约300到3000个碱基对（bp）。
- 在许多基因的启动子（promotor）或“起始”区域周围，甲基化经常被抑制。这些区域包含浓度相对较高的CpG对，与此段区域对应的染色体区段一起被称作CpG岛
- CpG岛常位于管家基因和其他在细胞中被频繁表达基因的启动子区域

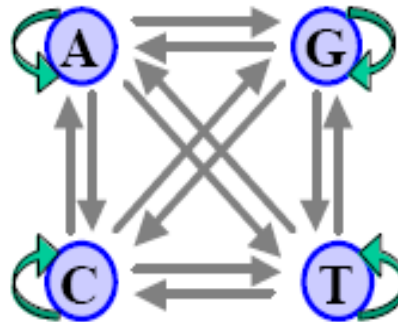
基于HMM的CpG岛识别

- Problem:

Given a short stretch of genomic data, does it come from a CpG island or not?

- States: {A,T,G,C}

- Transition probabilities matrix P



基于HMM的CpG岛识别

● Training the Markov Models

- Transition probabilities matrix P^- based on known non-CpG Island sequences
- Transition probabilities matrix P^+ based on known CpG Island sequences

● “+”model:

- Use transition matrix P^+

● “-”model:

- Use transition matrix P^-

P^+

| + | A | C | G | T |
|---|-------|-------|-------|-------|
| A | 0.180 | 0.274 | 0.426 | 0.120 |
| C | 0.171 | 0.368 | 0.274 | 0.188 |
| G | 0.161 | 0.339 | 0.375 | 0.125 |
| T | 0.079 | 0.355 | 0.384 | 0.182 |

P^-

| - | A | C | G | T |
|---|-------|-------|-------|-------|
| A | 0.300 | 0.205 | 0.285 | 0.210 |
| C | 0.322 | 0.298 | 0.078 | 0.302 |
| G | 0.248 | 0.246 | 0.298 | 0.208 |
| T | 0.177 | 0.239 | 0.292 | 0.292 |

基于HMM的CpG岛识别

□ For a new sequence \mathbf{x} :

$$\mathbf{x} = \mathbf{x}_1 \dots \mathbf{x}_L$$

□ Probability of a sequence

$$P(x) = P(x_L | x_{L-1}) \dots P(x_2 | x_1) P(x_1)$$

$$P(x_1) = P(x_1 | x_0)$$

$$P(x) = \prod_{i=0}^{L-1} P(x_{i+1} | x_i) = \prod_{i=0}^{L-1} a_{x_i, x_{i+1}}$$

基于HMM的CpG岛识别

- New sequence: $x = \text{TGCAGCG}$, x from CpG island regions?

$$\begin{aligned} P(x | +\text{model}) &= P_+(G|C) * P_+(C|G) * P_+(G|A) * P_+(A|C) * P_+(C|G) * P_+(G|T) \\ &= 0.274 * 0.339 * 0.426 * 0.171 * 0.339 * 0.384 \\ &= 0.000880819444 \\ P(x | -\text{model}) &= P_-(G|C) * P_-(C|G) * P_-(G|A) * P_-(A|C) * P_-(C|G) * P_-(G|T) \\ &= 0.078 * 0.246 * 0.285 * 0.322 * 0.246 * 0.292 \\ &= 0.00012648773 \end{aligned}$$

- $\text{Ratio} = P(x | +\text{model}) / P(x | -\text{model}) = 6.96 > 1$
- x more likely comes from CpG island regions
- In fact, take 'log': * to +

P+

| + | A | C | G | T |
|---|-------|-------|-------|-------|
| A | 0.180 | 0.274 | 0.426 | 0.120 |
| C | 0.171 | 0.368 | 0.274 | 0.188 |
| G | 0.161 | 0.339 | 0.375 | 0.125 |
| T | 0.079 | 0.355 | 0.384 | 0.182 |

P-

| - | A | C | G | T |
|---|-------|-------|-------|-------|
| A | 0.300 | 0.205 | 0.285 | 0.210 |
| C | 0.322 | 0.298 | 0.078 | 0.302 |
| G | 0.248 | 0.246 | 0.298 | 0.208 |
| T | 0.177 | 0.239 | 0.292 | 0.292 |

基于HMM的CpG岛识别

- Further Problem:

For a (stretch of a) genomic sequence, where are the CpG islands?

TAAAAAATAAATATGTTTAATTTGTGAACTGATTACCATCAGAAT

- States: {N, C}
- Observations: {A,T,G,C}
- Transition probabilities matrix P: 2×2



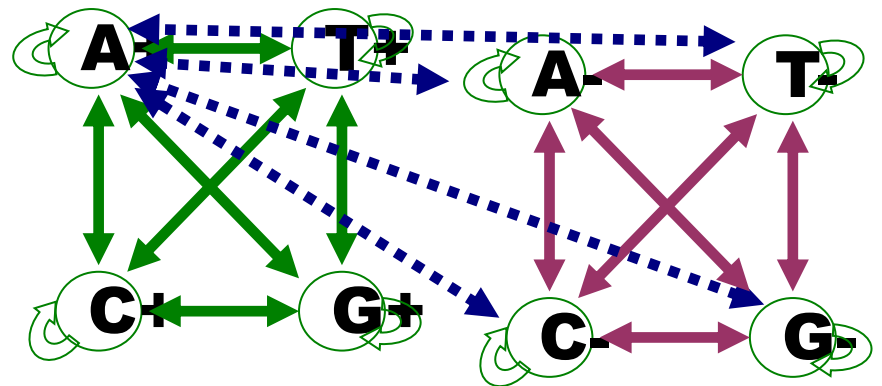
基于HMM的CpG岛识别

- Further Problem:

For a (stretch of a) genomic sequence, where are the CpG islands?

TAAAAAATAAATATGTTTAATTTGTGAAGTACTGATTACCATCAGAAT

- States: 8 States: {A+, C+, G+, T+, A-, C-, G-, T- }
- Observations: {A,T,G,C}
- Transition probabilities matrix P: 8*8



基于HMM的中文分词

- 小明硕士毕业于吉林大学计算机学院
- 小明/硕士/毕业于/吉林大学/计算机/学院
- 假设我们有一个标注好的数据集 $S=\{s_1, s_2, \dots, s_n\}$ ，格式如上。目标是训练一个模型，对于任意一句话，能够给出分词。

基于HMM的中文分词

- 小明硕士毕业于吉林大学计算机学院
- 小明/硕士/毕业/于/吉林大学/计算机/学院
- B E B E B E S B M M E B M E B E
- B表示分词后词语的首字，M表示分词后词语的非首非尾自（如果词语字数 ≥ 3 ），E表示词语的尾字，S表示单字成词
- $\lambda = (N, M, \pi, A, B)$ $N=4, M=|D|$ （汉字总数）

基于HMM的中文分词

- 小明硕士毕业于吉林大学计算机学院
- 小明/硕士/毕业/于/吉林大学/计算机/学院
- B E B E B E S B M M E B M E B E
- B表示分词后词语的首字，M表示分词后词语的非首非尾自（如果词语字数 ≥ 3 ），E表示词语的尾字，S表示单字成词
- $\lambda = (N, M, \pi, A, B)$ $N=4, M=|D|$ （汉字总数）
- 统计初始概率矩阵 π
- 就是统计初始状态B、M、E、S的出现频率

基于HMM的中文分词

- 小明硕士毕业于吉林大学计算机学院
- 小明/硕士/毕业/于/吉林大学/计算机/学院
- B E B E B E S BM ME BM E B E
- B表示分词后词语的首字，M表示分词后词语的非首非尾自（如果词语字数 ≥ 3 ），E表示词语的尾字，S表示单字成词
- $\lambda = (N, M, \pi, A, B)$ $N=4, M=|D|$ （汉字总数）
- 统计状态转移概率矩阵A

| | B | M | E | S |
|---|-----|-----|-----|-----|
| B | b_b | b_m | b_e | b_s |
| M | m_b | m_m | m_e | m_s |
| E | e_b | e_m | e_e | e_s |
| S | s_b | s_m | s_e | s_s |

知乎 @baiziyou

基于HMM的中文分词

- 小明硕士毕业于吉林大学计算机学院
- 小明/硕士/毕业/于/吉林大学/计算机/学院
- B E B E B E S B M M E B M E B E
- B表示分词后词语的首字，M表示分词后词语的非首非尾自（如果词语字数 ≥ 3 ），E表示词语的尾字，S表示单字成词
- $\lambda = (N, M, \pi, A, B)$ $N=4, M=|D|$ （汉字总数）
- 统计发射概率矩阵B

| | B | M | E | S |
|-----|------|------|------|------|
| 字1 | b_c1 | m_c1 | e_c1 | s_c1 |
| 字2 | b_c2 | m_c2 | e_c2 | s_c2 |
| ... | | | | |
| 字m | b_cm | m_cm | e_cm | s_cm |

基于HMM的中文词性标注

- 小明硕士毕业于吉林大学计算机学院
- 小明/硕士/毕业/于/吉林大学/计算机/学院
- 名词 名词 动词介词 名词 名词 名词
- $\lambda = (N, M, \pi, A, B)$ N =词性类别数, $M=|C|$ (词的总数)

作业

- **States:** $\{S_{sunny}, S_{rainy}, S_{snowy}\}$

- **Observations:** $\{O_{skirt}, O_{coat}, O_{umbrella}\}$

- **State transition probabilities:** $A = \begin{pmatrix} .8 & .15 & .05 \\ .38 & .6 & .02 \\ .75 & .05 & .2 \end{pmatrix}$

- **emission probability :**

- **Initial state distribution:**

$$\pi = (.7 \quad .25 \quad .05)$$

$$B = \begin{pmatrix} .6 & .3 & .1 \\ .05 & .3 & .65 \\ 0 & .5 & .5 \end{pmatrix}$$

作业

- 给定下列观察序列O:



- 1. 出现此序列的概率为多少(前向算法)?
- 2. 此序列对应的天气状况序列是什么?

Pafnuty Chebyshev (契比雪夫1821-1894)

