

基于隐式马尔可夫链的基因发现模型和算法

章 铭, 陆菊康

(上海大学计算机学院, 上海 200072)

摘 要: 探讨了隐式马尔可夫链在基因发现中的应用。提出了一个基于GHMM (泛化的隐式马尔可夫链) 的基因发现系统的简化的模型, 论述了用该模型和扩展的Viterbi算法发现基因的方法, 介绍了用于描述编码区和非编码区及信号的模型和实现。

关键词: 基因发现; HMM; GHMM

HMM-based Gene-finding Model and Algorithm

ZHANG Ming, LU Jukang

(College of Computer, Shanghai University, Shanghai 200072)

【Abstract】 This article describes the application of Hidden Markov Chain in the gene-finding system. A simplified model of gene-finding system, which is based on GHMM, is presented. Then it describes the method of finding gene with this model and modified Viterbi algorithm. Also it introduces the models for coding and non-coding regions and signals.

【Key words】 Gene finding; HMM; GHMM

随着人类基因组工程的完成以及越来越多的生物体的DNA被测序出来, 从大量的序列中获取有价值的信息变得日益重要。而其中最具有重要意义的是识别出蛋白质的编码区域, 并由此推测出完整的基因结构和相应蛋白质的结构。在过去十几年的时间里, 研究人员已开发出了许多基因发现(识别)系统, 其中的一些得到了广泛的应用, 如GeneMark^[1]、GenScan^[2]、Genie^[3]等。这些系统采用了各种各样的技术, 包括隐式马尔可夫模型、神经网络、动态编程和决策树等。其中基于隐式马尔可夫模型(HMM)的算法得到了广泛的研究和应用。本文论述基于HMM的扩展形式、GHMM的基因发现系统的框架和实现中的若干问题。

1 模型

1.1 HMM和GHMM基础

HMM是马尔可夫链的推广, 和马尔可夫链一样可看作是一状态机, 由一组状态及相关的转移组成。其与普通的马尔可夫链的不同之处在于, 它的每一个状态也是一个随机过程, 用以描述状态和输出的观察值之间的统计对应关系。因此, 马尔可夫链输出的是一状态序列, 而HMM输出的则是一观察值序列, 产生这一观察值序列的路径(状态序列)是不可见的(隐含的), 因此称其为“隐式的”。

HMM可表示为 $\lambda = \{A, B, \pi\}$ 。其中 π 表示初始状态概率分布矢量。 A 代表状态转移概率分布矩阵 $\{a_{ij}\}$ ($1 \leq i, j \leq N, N$ 是模型中状态的数目), $a_{ij} = P[q_{t+1}=S_j | q_t=S_i]$ 表示从状态 S_i 转移到 S_j 概率。 B 代表观察值概率矩阵 $\{b_j(k)\}$ ($1 \leq j \leq N, 1 \leq k \leq M, M$ 是观察值的数目, 即离散字母表的大小), $b_j(k) = P[v_k \text{ at } t | q_t=S_j]$ 表示在时刻 t , 状态 S_j 生成观察值 v_k 的概率。

对于一观察值序列 $O = o_1 o_2 \dots o_T$, 假设有状态序列 $Q = q_1 q_2 \dots q_T$, 则在模型 λ 中, O 由 Q 生成的概率是

$$P(O|Q, \lambda) = \prod_{t=1}^T P(o_t | q_t, \lambda) = b_{q_1}(o_1) * b_{q_2}(o_2) * \dots * b_{q_T}(o_T)$$

而状态序列 Q 的概率是

$$P(Q|\lambda) = \pi_{q_1} * a_{1q_2} * a_{2q_3} * \dots * a_{(T-1)q_T}$$

由于 $P(O|\lambda) = P(O|Q, \lambda)P(Q|\lambda)$, 因此在给定模型 λ 的

条件下, 生成观察值序列 O 的概率就是

$$P(O|\lambda) = \sum_{\text{all } Q} P(O|Q, \lambda)P(Q|\lambda) = \sum_{q_1, q_2, \dots, q_T} \pi_{q_1} * b_{q_1}(o_1) * a_{1q_2} * b_{q_2}(o_2) * \dots * a_{(T-1)q_T} * b_{q_T}(o_T)$$

通常对于某一观察值序列 O , 有多条路径可生成它。而生成 O 的最优的(也就是能最好地“解释” O)的状态序列 Q 往往揭示了模型隐含的部分, 解释了被模型所描述的事物的结构和特性。一种广泛使用的确定最优路径的办法是选择单个最佳状态序列(single best state sequence), 也就是使 $P(Q|O, \lambda)$ 达到最大值的状态序列。这可由Viterbi算法计算得到。在基因发现领域, 观察值为4种脱氧核苷酸(A, T, C, G), 而HMM模型的状态则对应于基因的各个功能序列段。而基因发现的方法就是对于一条DNA序列, 通过Viterbi算法计算出其最优的状态序列(路径)。最终得到的路径就描述了该DNA序列中的预测的基因结构。关于HMM更详细的论述见参考文献[1]。

GHMM (Generalized HMM) 是HMM的扩展形式。标准的HMM的每个状态产生一个观察值, 而GHMM的每个状态则可产生0或多个观察值。这使得GHMM能更直观地描述基因的结构, 并使框架有更好的可扩展性。

1.2 用于基因发现的GHMM框架

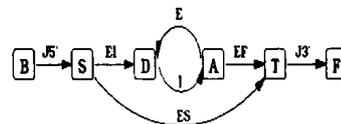


图1 多外显子基因的GHMM框架

图1是用于多外显子基因的GHMM框架。与一般的状态机不同, 这里代表了状态机中的状态, 节点代表了状态之间的转移, 那么在一个分析中, 将各边(状态)连接起来就可得到一条完整的基因序列。在GHMM中, 状态和转移都可具体地对应到基因的各个组成部分。其状态的具体含义是:

作者简介: 章 铭 (1977—), 男, 硕士生, 研究方向: 生物信息学, 数据仓库; 陆菊康, 副教授

收稿日期: 2002-09-23

J5'-5'端不翻译区, J3'-3'端不翻译区, E1-初始外显子, E-内部的外显子, I-内含子, EF-末尾的外显子, ES-单外显子。节点对应于基因中的信号: D-5'拼接点, A-3'拼接点, S-翻译开始区, T-翻译结束区。而B、F则是整个状态机的起始和终止位点。

当给定了一条候选的DNA序列X, 找到一路径(称为分析, parse) Φ , 使该分析 Φ 相对于其他分析来说, 生成X的概率最大。该 Φ 中的各个状态即对应于基因的各个功能子序列, 而 Φ 就描述了预测的基因的结构。可将这些概念形式化如下:

设M为一GHMM模型。由于GHMM的每一状态可生成多个符号(观察值), 因此, 定义X为一有序的子序列集合 $\{x_1, x_2, \dots, x_k\}$, 使得X是各子序列的连接:

$$X = x_1 x_2 \dots x_k$$

并定义分析为一有序的状态/序列对的集合:

$$\Phi = \{ \{q_1, x_1\}, \{q_2, x_2\}, \dots, \{q_k, x_k\} \}$$

其中 q_i ($1 \leq i \leq k$) 是分析 Φ 的第i个状态(即模型中的边), 对应于基因的功能单元。而 x_i 就是X中的第i个子序列, 是与该功能单元 q_i 相对应的子序列。这里规定 q_1 是离开开始节点B的边, q_k 是指向结束节点F的边。可计算 $P[X, \Phi]$ 如下(隐含的以M为条件):

$$P[X, \Phi] = P(q_1|B) \prod_{i=1}^{k-1} P(x_i|q_i) \prod_{i=1}^{k-1} P(q_{i+1}|\text{node}(q_i))$$

$\text{node}(q_i)$ 是边 q_i 指向的节点。每个 $P(x_i|q_i)$ 项都称为“内容检测器”(content sensor), 是按 q_i 的模型生成 x_i 的概率。每个状态自身可以是任意的模型, 如HMM、GHMM或神经网络。GHMM将系统层的模型与每个状态的模型(“内容探测器”)的实现细节相分离, 使系统模块化, 提高了整个系统的可扩展性, 也使得能够用任意合适的模型来描述基因各个功能单元。分析序列X, 就是要找到一个 Φ , 使得 $P[X, \Phi|M]$ 最大。这可用Viterbi算法来实现。

2 系统设计和实现

2.1 状态的模型

在GHMM中, 状态对应于基因的功能子序列, 如内含子, 外显子和5'、3'端的不翻译区。在基因发现领域, 用于计算这些序列的模型称为“内容检测器”。这些模型本身也可以是HMM、GHMM。这里以编码区(包括E1、EF、E、ES)和非编码区(包括I、5'3'UTR)为例, 介绍相应状态所采用的模型。

编码区就是基因中编码蛋白质的部分, 每3个核苷酸(密码子、codon)编码一个氨基酸。而非编码区则是基因中不编码蛋白质的部分。在进化选择的作用下, 各个子序列区域就有了不同的核苷酸使用特性。就可以利用不同类型的马尔可夫随机过程(或其他的模型, 如神经网络)来描述这些子序列。

非编码区域一般使用同构马尔可夫模型(homogeneous Markov model)。一阶的同构马尔可夫模型包括:(1)初始概率向量 PN_0 ;(2)转移矩阵PN, 这些参数可由训练集合中非编码区的单核苷酸和二核苷酸的数量来计算。

由于在编码区中每3个核苷酸(密码子)编码一个氨基酸, 随机过程就需要捕捉密码字每个位置的信息, 因此编码区使用异构马尔可夫模型(inhomogeneous Markov model)。一阶的异构马尔可夫模型包括:(1)3个初始状态概率向量: $P1_0, P2_0, P3_0$, 各自对应密码子中的一个位置。每个向量包括元素 $P1_{0i}, P2_{0i}, P3_{0i}, i=1, 2, 3, 4$, 分别代表了4种核

酸;(2)3个转移矩阵 $P1, P2, P3$, 包括元素 $P1_{ij}, P2_{ij}, P3_{ij}, i, j=1, 2, 3, 4$ 。

当给定一段DNA序列 $S = s_1 s_2 \dots s_n$, s_i 是核苷酸, 并设n是3的倍数。则S在非编码区的概率是

$$P[S|NON] = PN_0(s_1) * PN(s_2|s_1) * \dots * PN(s_n|s_{n-1})$$

S在编码区的概率根据S的第一个核苷酸在密码子中的位置(1, 2, 3), 可分成3种互斥的结果。

$$P[S|COD1] = P1_0(s_1) * P1(s_2|s_1) * P2(s_3|s_2) * P3(s_4|s_3) * \dots * P2(s_n|s_{n-1})$$

$$P[S|COD2] = P2_0(s_1) * P2(s_2|s_1) * P3(s_3|s_2) * P1(s_4|s_3) * \dots * P3(s_n|s_{n-1})$$

$$P[S|COD3] = P3_0(s_1) * P3(s_2|s_1) * P1(s_3|s_2) * P2(s_4|s_3) * \dots * P1(s_n|s_{n-1})$$

由于密码子是由3个核苷酸组成的, 因此二阶的马尔可夫模型是描述编码区的模型的最小阶数。但更高阶的模型需要更多的训练数据。我们采用了五阶的马尔可夫模型, 它可以捕捉到两个相邻密码字的序列信息。通过对上述一阶的异构马尔可夫模型进行扩充, 即可获得五阶马尔可夫模型。

2.2 信号的检测

信号是基因中的功能位点, 如剪切位点、启动子。它通常是由几个和几十个核苷酸组成的序列片断(或由多个这样的片断组成), 是基因不同功能单元序列之间的分界, 标志着模型中状态之间的转移的发生, 如剪切位点是内含子和外显子的分界, 而启动子则标志了整个基因的起始。因此, 在基因发现的过程中, 正确的检测信号是非常重要的。

和检测功能子序列的“内容检测器”相对应, 用以检测信号的模型称为“信号检测器”(signal sensor)。它既可以作为“内容检测器”的一部分, 也可以通过修改公式, 将其作为独立的项并入公式中, 如Genie。

用于剪切位点的“信号检测器”通常采用权值矩阵方法(WMM, Weight Matrix Method)。在这种方法中, 长度为 λ 的信号中的核苷酸被认为是由与位置相关的概率分布独立生成的。当给定了一个信号模型“+”和一段序列 $X = x_1 \dots x_\lambda$, 则这段序列的概率就是

$$P[X|+] = P'_{WMM}(X) = \prod_{i=1}^{\lambda} P^{(i)}_{i, x_i}$$

其中 $P^{(i)}$ 是在信号位置i上生成核苷酸j的概率。我们采用了对WMM的扩展, 权值阵列模型(WAM, Weight Array Model)。它实际上是一个一阶异构马尔可夫链, 其中位置i的概率分布依赖于位置i-1上的核苷酸种类。

2.3 用Viterbi算法寻找基因

在得到各个子模块之后, 就可利用Viterbi算法在给定的DNA序列中发现基因了。Viterbi算法^[1]是HMM中用以计算“最佳”路径的动态编程算法(Dynamic Programming)。由于这里用到的GHMM模型是对标准HMM模型扩展, 因此对算法的实现过程做了相应的修改。具体的算法如下:

首先给所有的状态(J5', E1, E, I, EF, J3')编号为 $q_1, q_2, q_3, q_4, q_5, q_6$ 。子序列 $s_{i,j}$ 在状态 q_k 下的概率是 $P(s_i, j|q_k)$ 。并设 $\gamma_i(j)$ 为在位置j结束于状态 Q 的子序列 $S_{i,j}$ 的最佳分析的联合概率。并设给定序列的长度为L。

(1)初始化: 由于只寻找完整的基因,即以状态J5'(q₁)开头的分析。因此有

$$\gamma_i(1) = P(q_1|B)P(s_{i,1}|q_1), \gamma_i(l) = 0 \quad (2 \leq i \leq 6)$$

(2)递归:

$$\gamma_i(j) = \max_{1 \leq k \leq 6} [\gamma_i(k) P(q_i|\text{node}(q_k)) P(s_{i,j}|q_i)]$$

$$(1 \leq i \leq 6, 2 \leq j \leq N, 2 \leq s \leq 6)$$

(下转第145页)

