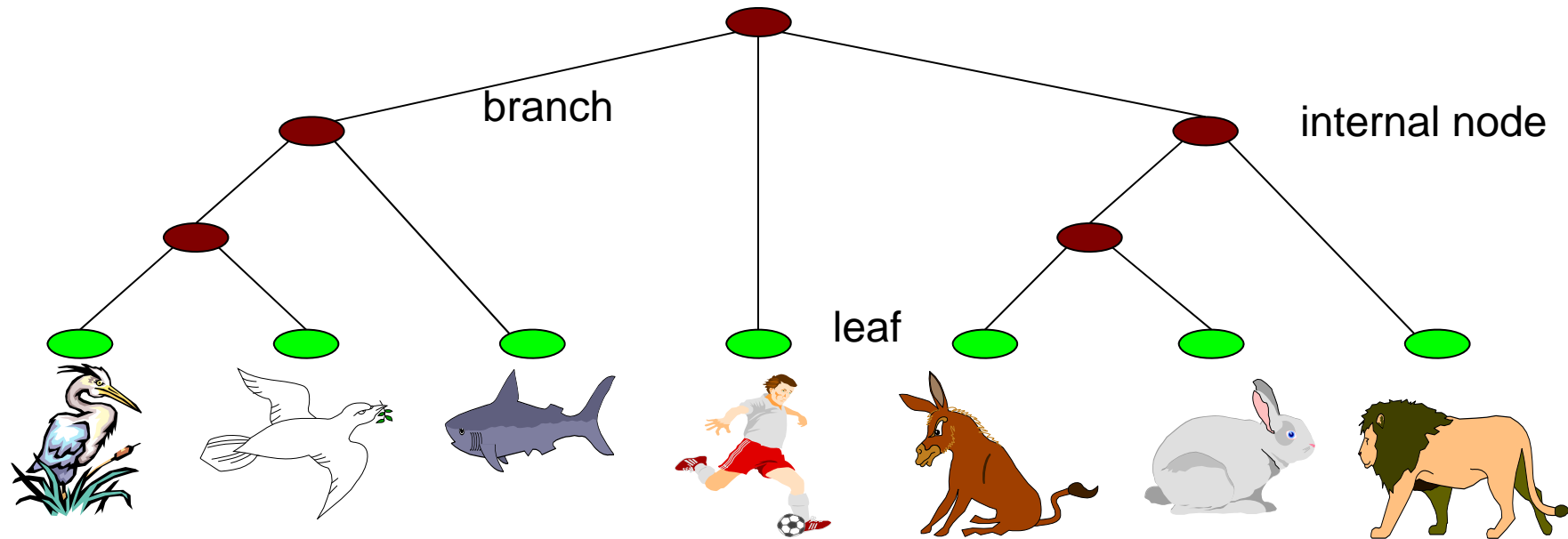


# 第4章：进化树构建的概率方法

- 问题介绍
- 进化树构建方法的概率方法

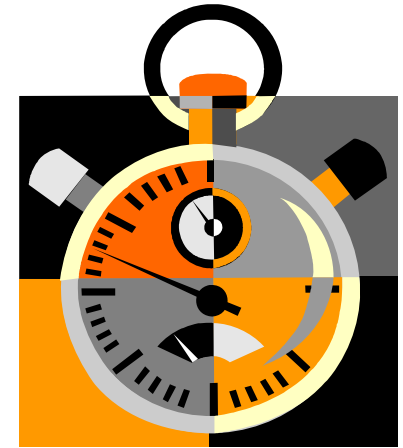
部分Slides修改自University of Basel的Michael Springmann  
课程“CS302 Seminar Life Science Informatics”的讲义

# Phylogenetic Tree



- Topology: bifurcating
  - Leaves -  $1...N$
  - Internal nodes  $N+1...2N-2$
- Branch length

# Molecular Clock Hypothesis



- Amount of genetic difference between sequences is a function of time since separation.
- Rate of molecular change is constant (enough) to predict times of divergence

# Likelihood of a Tree

- Given:
  - $n$  aligned sequences  $M = X_1, \dots, X_n$
  - A tree  $T$ , leaves labeled with  $X_1, \dots, X_n$
- Reconstruction  $t^*$ :
  - Labeling of internal nodes
  - Branch lengths

Goal: Find optimal reconstruction  $t^*$  : One maximizing the likelihood  $P(M|T, t^*)$

# Probabilistic Methods

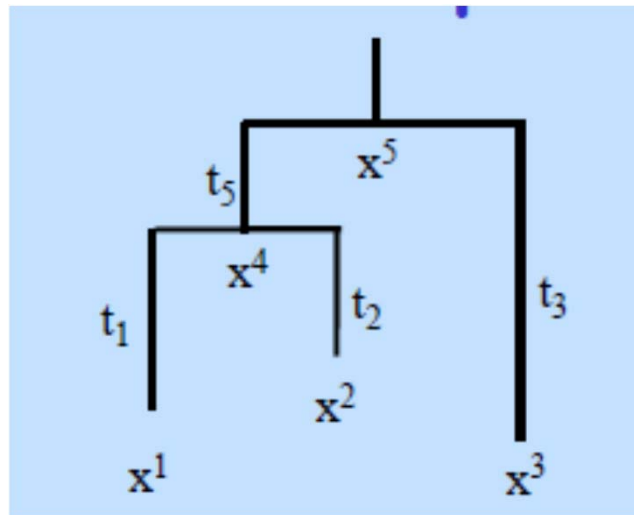
- The phylogenetic tree represents a generative probabilistic model (like HMMs) for the observed sequences.
- Background probabilities:  $q(a)$
- Mutation probabilities:  $P(a|b, t)$
- Models for evolutionary mutations
  - Jukes Cantor
  - Kimura 2-parameter model
- Such models are used to derive the probabilities

# Probabilistic Model

- Assumptions:
  - Each character is independent
  - The branching is a Markov process: The probability that a node  $x$  has a specific label is only a function of the parent node  $y$  and the branch length  $t$  between them
  - The probabilities  $P(x|y,t)$  are known

# Example

- Given then tree



$$\begin{aligned} &P(x_1, x_2, x_3, x_4, x_5 | T, t^*) \\ &= P(x_1 | x_4, t_1) P(x_2 | x_4, t_2) P(x_3 | x_5, t_3) P(x_4 | x_5, t_5) \end{aligned}$$

# 连续时间马氏链

- 随机变量族  $\{\xi_t(\omega), t \geq 0\}$  称为连续时间马氏链, 若  $\xi_n$  是离散的(其状态空间至多是可数集, 即有限或者与自然数一一对应), 而且对于  $\forall m \geq 0, \forall s > s_1 \geq \cdots \geq s_m \geq 0$  及任意状态序列  $i, j, i_1, \cdots, i_m$ , 都有

$$\begin{aligned} & P(\xi_{t+s} = j \mid \xi_s = i, \cdots, \xi_{s_1} = i_1, \cdots, \xi_{s_m} = i_m) \\ &= P(\xi_{t+s} = j \mid \xi_s = i) \end{aligned}$$



# 转移概率矩阵

- 定义：对连续时间马氏链  $\{\xi_t(\omega), t \geq 0\}$ ，记

$$p_{ij}(s, t) = P(\xi_t = j | \xi_s = i), \forall t \geq s$$
$$P(s, t) = (p_{ij}(s, t))$$

称无穷矩阵  $P(s, t)$  为转移概率矩阵

- 易证转移概率满足Chapman-Kolmogorov方程

$$p_{ij}(s, u) = \sum_k p_{ik}(s, t) p_{kj}(t, u)$$

# 时齐的连续时间马氏链

- 定义：如果连续时间马氏链的转移概率矩阵 $P(s, s+t)$ 与起始时间 $s$ 无关，称此时的马氏链为时齐的。那么此时的转移概率矩阵就可用一个时间参数来标度

$$P(t - s) = P(s, t)$$

- 相应地，chapman-kolmogorov方程简化为

$$P(t + s) = P(s)P(t)$$

# 连续时间马氏链的有限维分布

- 类似于离散时间马氏链，连续时间马氏链的有限维分布有初始分布和转移概率唯一确定。设  $\mu_i(0) = P(\xi_0 = i)$ , 转移概率矩阵  $\mathbf{P}(\mathbf{t})$ , 那么对  $\forall 0 < t_1 < t_2 < \cdots < t_n$ ,

$$\begin{aligned} &P(\xi_0 = i_0, \xi_{t_1} = i_1, \cdots, \xi_{t_n} = i_n) \\ &= \mu_{i_0}(0) p_{i_0 i_1}(t_1) p_{i_1 i_2}(t_2 - t_1) \cdots p_{i_{n-1} i_n}(t_n - t_{n-1}) \end{aligned}$$

# 绝对概率

- 令  $\mu_i(t) = P(\xi_t = i)$ , 称向量

$$\mu(t) = \begin{pmatrix} \mu_1(t) \\ \vdots \\ \mu_i(t) \\ \vdots \end{pmatrix}$$

为t时刻的绝对概率

# 转移速率矩阵(Q矩阵)

- 马氏链中时齐马氏链的多步转移矩阵是最小的转移矩阵 $P$ 的次方。
- 连续情况下对应于无穷小生成元，但此时需要加条件
- 定义：设  $P(t) \rightarrow I, t \rightarrow 0$ , 若转移概率矩阵在0点的导数存在，记为  $P'(0) = Q$ , 称之为转移概率矩阵的转移速率矩阵(简称为Q矩阵)。

# 有限状态情形

- 当状态数有限时，如 $\mathbf{Q}$ 矩阵满足  $\sum_j q_{ij} = 0$  (即所谓保守)，此时转移概率矩阵由 $\mathbf{Q}$ 矩阵唯一确定

$$P(t) = e^{\mathbf{Q}t} = I + \frac{\mathbf{Q}t}{1!} + \frac{\mathbf{Q}^2 t^2}{2!} + \cdots + \frac{\mathbf{Q}^k t^k}{k!} + \cdots$$

- 通常转移速率矩阵 $\mathbf{Q}$ 已知，我们可以通过指数矩阵给出转移概率矩阵 $\mathbf{P}(t)$ .

# Kolmogorov方程和Master方程

- Kolmogorov向前方程(Fokker-Planck方程):

$$P'(t) = P(t)Q, P(0) = I$$

- Kolmogorov向后方程:

$$P'(t) = QP(t), P(0) = I$$

- Master方程

$$(\mu'_1(t), \mu'_2(t), \dots, ) = (\mu_1(t), \mu_2(t), \dots, )Q$$

其中  $\mu_k(t) = P(\xi_t = k)$  为绝对概率

# Molecular Evolution

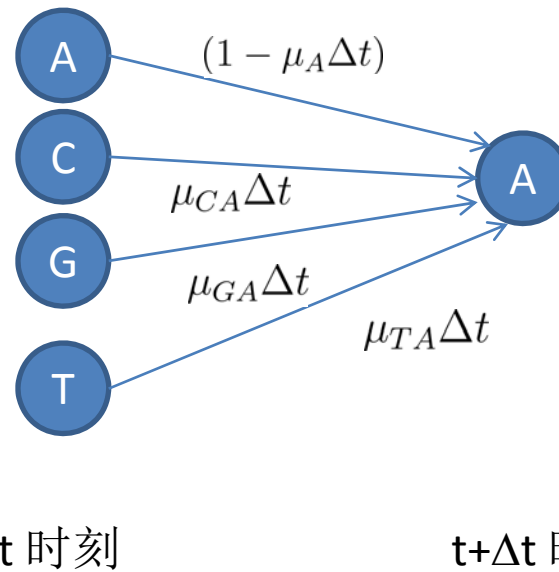
Q: How can we model evolution on nucleotide level? (ignore gaps, focus on substitutions)

A: Consider what happens at a specific position for small time interval  $\Delta t$

- $p(t)$  = vector of probabilities of {A,C,G,T} at time  $t$
- $\mu_{AC}$  = rate of transition from A to C per unit time
- $\mu_A = \mu_{AC} + \mu_{AG} + \mu_{AT}$  rate of transition out of A



# Molecular Evolution



$$p_A(t + \Delta t) = p_A(t)(1 - \mu_A \Delta t) + p_C(t)\mu_{CA} \Delta t + p_G(t)\mu_{GA} \Delta t + p_T(t)\mu_{TA} \Delta t$$

# Molecular Evolution

- 同理可得

$$\begin{aligned} p_C(t + \Delta t) &= p_C(t)(1 - \mu_C \Delta t) \\ &\quad + p_A(t)\mu_{AC}\Delta t + p_G(t)\mu_{GC}\Delta t + p_T(t)\mu_{TC}\Delta t \end{aligned}$$

$$\begin{aligned} p_G(t + \Delta t) &= p_G(t)(1 - \mu_G \Delta t) \\ &\quad + p_A(t)\mu_{AG}\Delta t + p_C(t)\mu_{CG}\Delta t + p_T(t)\mu_{TG}\Delta t \end{aligned}$$

$$\begin{aligned} p_T(t + \Delta t) &= p_T(t)(1 - \mu_T \Delta t) \\ &\quad + p_A(t)\mu_{AT}\Delta t + p_C(t)\mu_{CT}\Delta t + p_G(t)\mu_{GT}\Delta t \end{aligned}$$

# Molecular Evolution

In vector notation, we get

$$p(t + \Delta t) = p(t) + P(t)Q\Delta t$$

where  $Q$  is the substitution rate matrix

$$Q = \begin{pmatrix} -\mu_A & \mu_{AG} & \mu_{AC} & \mu_{AT} \\ \mu_{GA} & -\mu_G & \mu_{GC} & \mu_{GT} \\ \mu_{CA} & \mu_{CG} & -\mu_C & \mu_{CT} \\ \mu_{TA} & \mu_{TG} & \mu_{TC} & -\mu_T \end{pmatrix}$$

# Molecular Evolution

- 写成列向量形式

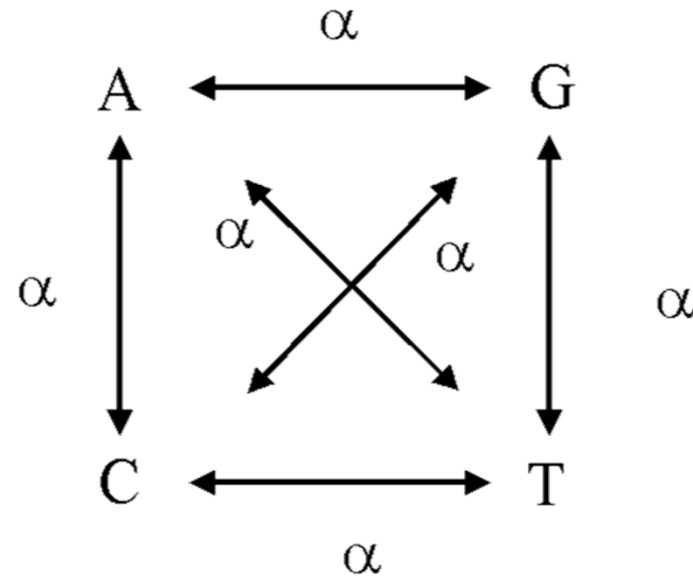
$$\vec{p}(t + \Delta t) = \vec{p}(t) + Q^T \vec{p}(t)$$

- 取转置后即为Master方程
- 其中Q为转移速率矩阵

$$Q = \begin{pmatrix} -\mu_A & \mu_{AG} & \mu_{AC} & \mu_{AT} \\ \mu_{GA} & -\mu_G & \mu_{GC} & \mu_{GT} \\ \mu_{CA} & \mu_{CG} & -\mu_C & \mu_{CT} \\ \mu_{TA} & \mu_{TG} & \mu_{TC} & -\mu_T \end{pmatrix}$$

# Jukes Cantor model

- Mutation occurs at a constant rate
- Each nucleotide is equally likely to mutate into any other nucleotide with rate  $\alpha$ .



$$Q = \begin{pmatrix} -3\alpha & \alpha & \alpha & \alpha \\ \alpha & -3\alpha & \alpha & \alpha \\ \alpha & \alpha & -3\alpha & \alpha \\ \alpha & \alpha & \alpha & -3\alpha \end{pmatrix}$$

# Substitution Matrix

- 由对称性，可设

$$P(t) = \begin{pmatrix} \gamma(t) & s(t) & s(t) & s(t) \\ s(t) & \gamma(t) & s(t) & s(t) \\ s(t) & s(t) & \gamma(t) & s(t) \\ s(t) & s(t) & s(t) & \gamma(t) \end{pmatrix}$$

- 由Kolmogorov后向方程

$$\frac{dP(t)}{dt} = QP(t)$$

# Substitution Matrix

- 可得方程

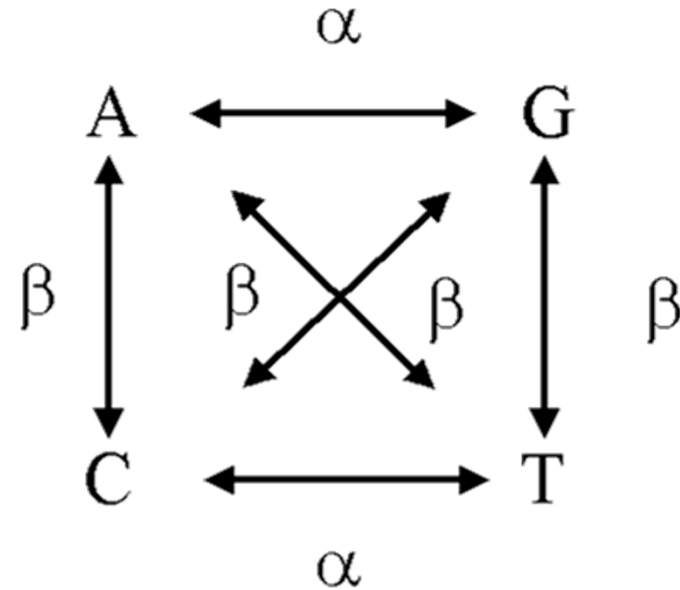
$$\begin{cases} \frac{d\gamma(t)}{dt} = -3\alpha\gamma(t) + 3\alpha s(t) \\ \frac{ds(t)}{dt} = -\alpha s(t) + \alpha\gamma(t) \end{cases}$$

- 容易求得

$$\begin{aligned} \gamma(t) &= \frac{1}{4}(1 + 3e^{-4\alpha t}) \\ s(t) &= \frac{1}{4}(1 - e^{-4\alpha t}) \end{aligned}$$

# Kimura 2-parameter Model

- Allows a different rate for transitions and transversions.



$$Q = \begin{pmatrix} -\alpha - 2\beta & \beta & \alpha & \beta \\ \beta & -\alpha - 2\beta & \beta & \alpha \\ \alpha & \beta & -\alpha - 2\beta & \beta \\ \beta & \alpha & \beta & -\alpha - 2\beta \end{pmatrix}$$



# Substitution Matrix

- 由对称性，可设转移概率矩阵

$$P(t) = \begin{pmatrix} \gamma(t) & s(t) & u(t) & s(t) \\ s(t) & \gamma(t) & s(t) & u(t) \\ u(t) & s(t) & \gamma(t) & s(t) \\ s(t) & u(t) & s(t) & \gamma(t) \end{pmatrix}$$

- 由Kolmogorov后向方程

$$\frac{dP(t)}{dt} = QP(t)$$

# Substitution Matrix

- 可得方程

$$\begin{cases} \frac{d\gamma(t)}{dt} = -(2\beta + \alpha)\gamma(t) + 2\beta s(t) + \alpha u(t) \\ \frac{ds(t)}{dt} = -2\beta s(t) + \beta\gamma(t) + \beta u(t) \\ \frac{du(t)}{dt} = -(2\beta + \alpha)u(t) + 2\beta s(t) + \alpha\gamma(t) \end{cases}$$

- 容易求得

$$\begin{cases} s(t) = \frac{1}{4}(1 - e^{-4\beta t}) \\ s(t) = \frac{1}{4}(1 + e^{-4\beta t} - 2e^{-2(\alpha+\beta)t}) \\ \gamma(t) = 1 - 2s(t) - u(t) \end{cases}$$

# Substitution Matrix: General Case

- 对于对称矩阵 $Q$ 可以对角化, 即存在正交矩阵 $U$ , 和特征值  $\lambda_1 \geq \cdots \geq \lambda_n$  , 使得

$$Q = U^T \text{diag}\{\lambda_1, \cdots, \lambda_n\} U$$

- 于是

$$P(t) = U^T \text{diag}\{e^{\lambda_1 t}, \cdots, e^{\lambda_n t}\} U$$

# PAM矩阵

- Point accepted mutation (Dayhoff et al 1978)
- Given an tree of protein family, the frequency matrix  $A_{ab}$  counting the occurrence of an “a” in the ancestral sequence was replaced by a “b” in the descendant.
- Estimate the conditional probability  $p(b|a)$

$$P(b|a) = B_{a,b} = \frac{A_{ab}}{\sum_c A_{ac}}$$

# PAM矩阵

- Scaling B

$$C_{ab} = \sigma B_{ab}, C_{aa} = \sigma B_{aa} + (1 - \sigma)$$

- Such that the expected number of substitution is 1%, i.e.

$$\sum_{ab} q_a q_b C_{ab} = 0.01$$

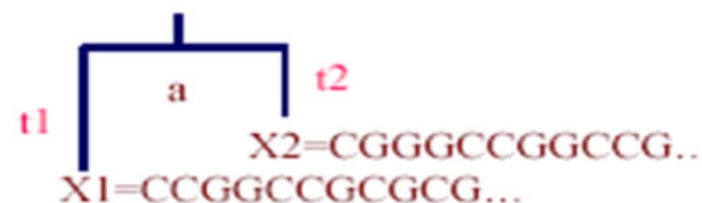
- Then the PAM(1) matrix is given by

$$S(1) = (C_{ab})$$

# Calculating the Likelihood for Ungapped Alignments

$$P(X^1, X^2|T, t_1, t_2) = \prod_{u=1}^N P(X_u^1, X_u^2|T, t_1, t_2)$$

$$P(X_u^1, X_u^2|T, t_1, t_2) = \sum_a q_a P(X_u^1|a, t_1) P(X_u^2|a, t_2)$$



- 假设突变符合JC model, 等初始概率  $q_A = q_C = q_G = q_T = \frac{1}{4}$

$$\begin{aligned} P(C, C|T, t_1, t_2) &= q_C \gamma(t_1) \gamma(t_2) + q_G s(t_1) s(t_2) + q_A s(t_1) s(t_2) + q_T s(t_1) s(t_2) \\ &= \frac{1}{3} (r(t_1) r(t_1) + 3S(t_1) S(t_2)) \end{aligned}$$

$$P(C, G|T, t_1, t_2) = P(G, C|T, t) = \frac{1}{4} (\gamma(t_1) s(t_1) + s(t_1) \gamma(t_2) + 2s(t_1) s(t_2))$$

$$P(X^1, X^2|T, t_1, t_2) = 16^{-(n_1+n_2)} (1 + 3e^{-4\alpha(t_1+t_2)})^{n_1} (1 - e^{-4\alpha(t_1+t_2)})^{n_2}$$

其中n1是匹配数， n2是不匹配数目.

# Calculating the Likelihood for Ungapped Alignments

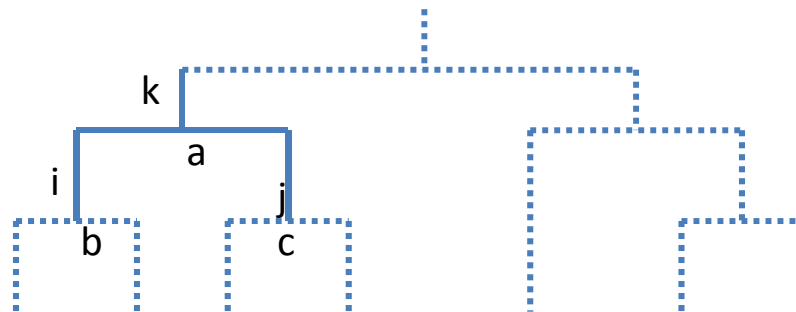
- $n$  sequences of length  $N$ , site  $u=1\dots N$
- Given a rooted tree contains  $2n - 1$  nodes,  $1\dots n$  being the leaf nodes,  $n+1 \dots 2n-1$  non-leaf, tree lengths  $t_1, \dots, t_{2n-1}$ .
- Let  $a(i)$  denote the ancestor of node  $a^i$

$$P(x^1, \dots, x^n | T, t) = \prod_{u=1}^N P(x_u^1, \dots, x_u^n | T, t)$$

$$P(x_u^1, \dots, x_u^n | T, t) = \sum_{a^{n+1}, \dots, a^{2n-1}} q_{a^{2n-1}} \prod_{i=n+1}^{2n-2} P(a^i | a^{\alpha(i)}, t_i) \\ \times \prod_{i=1}^n P(x_u^i | a^{\alpha(i)}, t_i)$$

# Felsenstein's Recursive Algorithm

- Let  $P(L_k | a)$  denote the probability of all the leafs below node  $k$  given that the residue at  $k$  is  $a$ .
- Then we compute  $P(L_k | a)$  from the probabilities  $P(L_i | b)$  and  $P(L_j | c)$  for all  $b$  and  $c$ , where  $i$  and  $j$  are the daughter nodes of  $k$ .





# Felsenstein's Recursive Algorithm

- Initialization: set  $k=2n-1$
- Recursion: Compute  $P(L_k | a)$  for all  $a$  as follows:
  - If  $k$  is leaf node:  $P(L_k | a)=1$  only if  $a = x_u^k$ .
  - If  $k$  is not a leaf node:
    - Compute  $P(L_i | a)$ ,  $P(L_j | a)$  for all  $a$  at the daughter nodes  $i, j$ , and set  $P(L_k | a) = \sum_{bc} P(b|a, t_i) P(L_i | b) P(c|a, t_j) P(L_j | c)$
- Termination: Likelihood at site  $u$ ,

$$P(x_u | T, t) = \sum_a P(L_{2n-1} | a) q_a$$

# Reversibility & Independence of Root Position

- The score of the optimal tree is independent of the root position if and only if:
  - the substitution matrix is **multiplicative**
  - the substitution matrix is **reversible**
- A substitution matrix is reversible if for all a,b and t:

$$P(b|a, t)q_a = P(a|b, t)q_b$$

# Maximum Likelihood (ML)

- Score each tree by
  - Assumption of independent positions “m”
- Branch lengths  $t$  can be optimized
  - Gradient Ascent
  - EM
- We look for the highest scoring tree
  - Exhaustive
  - Sampling methods (Metropolis)

# Computational Problem

- Such procedures are computationally expensive!
- Computation of optimal parameters, per candidate, requires non-trivial optimization step.
- Spend non-negligible computation on a candidate, even if it is a low scoring one.
- In practice, such learning procedures can only consider small sets of candidate structures

# 参考文献

- S. Durbin, S. Eddy, A. Krogh and G. Mitchison. Biological Sequence Analysis—Probabilistic Models of Proteins and Nucleic Acids. 1998, Cambridge University Press.