

Substitution Models

Tomáš Vinař

November 4, 2021

Substitution models, notation

$P(b|a, t)$: probability that if we start with symbol a , after time t we will see symbol b

Transition probability matrix:

$$S(t) = \begin{pmatrix} P(A|A, t) & P(C|A, t) & P(G|A, t) & P(T|A, t) \\ P(A|C, t) & P(C|C, t) & P(G|C, t) & P(T|C, t) \\ P(A|G, t) & P(C|G, t) & P(G|G, t) & P(T|G, t) \\ P(A|T, t) & P(C|T, t) & P(G|T, t) & P(T|T, t) \end{pmatrix}$$

Substitution models, basic properties

- $S(0) = I$

- $\lim_{t \rightarrow \infty} S(t) = \begin{pmatrix} \pi_A & \pi_C & \pi_G & \pi_T \\ \pi_A & \pi_C & \pi_G & \pi_T \\ \pi_A & \pi_C & \pi_G & \pi_T \\ \pi_A & \pi_C & \pi_G & \pi_T \end{pmatrix}$

Distribution π is called stationary (equilibrium)

- $S(t_1 + t_2) = S(t_1)S(t_2)$ (multiplicativity)

- Jukes-Cantor model should also satisfy

$$S(t) = \begin{pmatrix} 1 - 3s(t) & s(t) & s(t) & s(t) \\ s(t) & 1 - 3s(t) & s(t) & s(t) \\ s(t) & s(t) & 1 - 3s(t) & s(t) \\ s(t) & s(t) & s(t) & 1 - 3s(t) \end{pmatrix}$$

$$S(t) = \begin{pmatrix} 1 - 3s(t) & s(t) & s(t) & s(t) \\ s(t) & 1 - 3s(t) & s(t) & s(t) \\ s(t) & s(t) & 1 - 3s(t) & s(t) \\ s(t) & s(t) & s(t) & 1 - 3s(t) \end{pmatrix}$$

$$S(2t) = S(t)^2 =$$

$$= \begin{pmatrix} 1 - 6s(t) + 12s(t)^2 & 2s(t) - 4s(t)^2 & 2s(t) - 4s(t)^2 & 2s(t) - 4s(t)^2 \\ 2s(t) - 4s(t)^2 & 1 - 6s(t) + 12s(t)^2 & 2s(t) - 4s(t)^2 & 2s(t) - 4s(t)^2 \\ 2s(t) - 4s(t)^2 & 2s(t) - 4s(t)^2 & 1 - 6s(t) + 12s(t)^2 & 2s(t) - 4s(t)^2 \\ 2s(t) - 4s(t)^2 & 2s(t) - 4s(t)^2 & 2s(t) - 4s(t)^2 & 1 - 6s(t) + 12s(t)^2 \end{pmatrix}$$

$$\approx \begin{pmatrix} 1 - 6s(t) & 2s(t) & 2s(t) & 2s(t) \\ 2s(t) & 1 - 6s(t) & 2s(t) & 2s(t) \\ 2s(t) & 2s(t) & 1 - 6s(t) & 2s(t) \\ 2s(t) & 2s(t) & 2s(t) & 1 - 6s(t) \end{pmatrix}$$

for $t \rightarrow 0$

Substitution rate matrix (matica rýchlostí, matica intenzít)

- Substitution rate matrix for Jukes-Cantor model:

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

- For very small t we have $S(t) \approx I + Rt$
- Rate α is the probability of a change per unit of time for very small t , or derivative of $s(t)$ with respect to t at $t = 0$
- Solving the differential equation for the Jukes-Cantor model we get $s(t) = (1 - e^{-4\alpha t})/4$

Jukes-Cantor model

$$S(t) = \begin{pmatrix} (1 + 3e^{-4\alpha t})/4 & (1 - e^{-4\alpha t})/4 & (1 - e^{-4\alpha t})/4 & (1 - e^{-4\alpha t})/4 \\ (1 - e^{-4\alpha t})/4 & (1 + 3e^{-4\alpha t})/4 & (1 - e^{-4\alpha t})/4 & (1 - e^{-4\alpha t})/4 \\ (1 - e^{-4\alpha t})/4 & (1 - e^{-4\alpha t})/4 & (1 + 3e^{-4\alpha t})/4 & (1 - e^{-4\alpha t})/4 \\ (1 - e^{-4\alpha t})/4 & (1 - e^{-4\alpha t})/4 & (1 - e^{-4\alpha t})/4 & (1 + 3e^{-4\alpha t})/4 \end{pmatrix}$$

The rate matrix is typically normalized so that there is on average one substitution per unit of time, here $\alpha = 1/3$

Jukes-Cantor model, summary

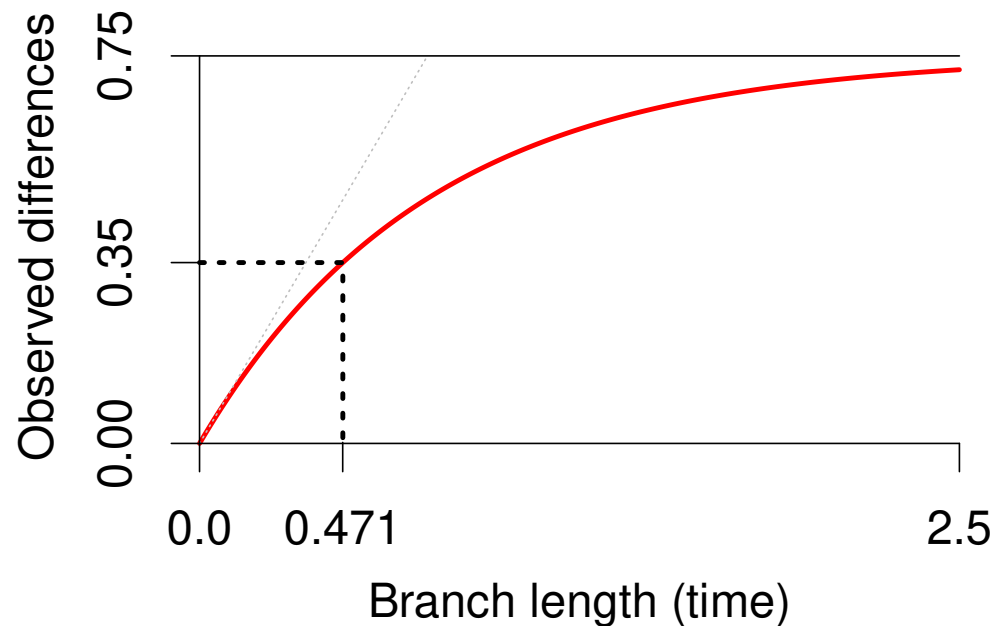
- $S(t)$: matrix 4×4 , where $S(t)_{a,b} = P(b|a, t)$ is the probability that if we start with base a , after time t we have base b .
- Jukes-Cantor model assumes that $P(b|a, t)$ is the same for all $a \neq b$
- For a given time t , off-diagonal elements are $s(t)$, diagonal $1 - 3s(t)$
- Rate matrix R : for J-C off-diagonal α , diagonal -3α
- For very small t we have $S(t) \approx I - Rt$
- Rate α is the probability of a change per unit of time for very small t , or derivative of $s(t)$ with respect to t for $t = 0$
- Solving the differential equation for the Jukes-Cantor model, we get $s(t) = (1 - e^{-4\alpha t})/4$
- The rate matrix is typically normalized so that there is on average one substitution per unit of time, that is, $\alpha = 1/3$

Correction of evolutionary distances

$$\Pr(X_{t_0+t} = C \mid X_{t_0} = A) = \frac{1}{4}(1 - e^{-\frac{4}{3}t})$$

The expected number of observed changes per base in time t :

$$D(t) = \Pr(X_{t_0+t} \neq X_{t_0}) = \frac{3}{4}(1 - e^{-\frac{4}{3}t})$$



Correction of observed distances

$$D = \frac{3}{4} \left(1 - e^{-\frac{4}{3}t} \right) \quad \Rightarrow \quad t = -\frac{3}{4} \ln \left(1 - \frac{4}{3}D \right)$$

More complex models

- General rate matrix R

$$R = \begin{pmatrix} \cdot & \mu_{AC} & \mu_{AG} & \mu_{AT} \\ \mu_{CA} & \cdot & \mu_{CG} & \mu_{CT} \\ \mu_{GA} & \mu_{GC} & \cdot & \mu_{GT} \\ \mu_{TA} & \mu_{TC} & \mu_{TG} & \cdot \end{pmatrix}$$

- μ_{xy} is the rate at which base x changes to a different base y
- Namely, $\mu_{xy} = \lim_{t \rightarrow 0} \frac{\Pr(y | x, t)}{t}$
- The diagonal is added so that the sum of each row is 0
- There are models with a smaller number of parameters (compromise between J-C and an arbitrary matrix)

Kimura model

- A and G are purines, C and T pyrimidines
- Purines more often change to other purines and pyrimidines to pyrimidines
- Transition: change within group $A \Leftrightarrow G, C \Leftrightarrow T$,
Transversion: change to a different group $\{A, G\} \Leftrightarrow \{C, T\}$
- Two parameters: rate of transitions α , rate of transversions β

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

HKY model (Hasegawa, Kishino, Yano)

- Extension of Kimura model, which allows different probabilities of A, C, G, T in the equilibrium
- If we set time to infinity, original base is not important, base frequencies stabilize in an equilibrium.
- Jukes-Cantor has probability of each base in the equilibrium 1/4.
- In HKY the equilibrium frequencies $\pi_A, \pi_C, \pi_G, \pi_T$ are parameters (summing to 1)
- Parameter κ : transition / transversion ratio (α/β)
- Rate matrix:
$$\mu_{x,y} = \begin{cases} \kappa\pi_y & \text{if mutation from } x \text{ to } y \text{ is transition} \\ \pi_y & \text{if mutation from } x \text{ to } y \text{ is transversion} \end{cases}$$

From rate matrix R to transition probabilities $S(t)$

- J-C and some other models have explicit formulas for $S(t)$
- For more complex models, such formulas are not available
- In general, $S(t) = e^{Rt}$
- Exponential of a matrix A is defined as $e^A = \sum_{k=0}^{\infty} \frac{1}{k!} A^k$
- If R is diagonalized $R = UDU^{-1}$, where D is a diagonal matrix, then $e^{Rt} = Ue^{Dt}U^{-1}$ and the exponential function is applied to the diagonal elements of D
- Diagonalization always exists for symmetric matrices R (the diagonal contains eigenvalues)