

ÁRBOLES FILOGENÉTICOS

SEMANA 04

27 de Abril del 2020



BIOINFORMÁTICA Y BIOESTADÍSTICA 2020-1

LUIS JARAMILLO MSc

Índice

- Modelos de Árboles Filogenéticos
- Comparación de Modelos



BIOINFORMÁTICA Y BIOESTADÍSTICA 2020-1

LUIS JARAMILLO MSc

1

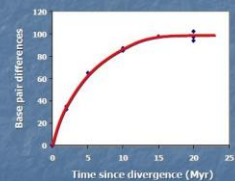
ÁRBOLES FILOGENÉTICOS

BIOINFORMÁTICA Y BIOESTADÍSTICA 2020-1

LUIS JARAMILLO MSc

Midiendo el cambio evolutivo

- Medida simple: Contar el número de sitios diferentes.
- Estimador muy inexacto:
 - Sitios pueden tener sustituciones repetidas.
 - Divergencia de secuencias llega a ser menos exacta en su estimación



BIOINFORMÁTICA Y BIOESTADÍSTICA 2020-1

LUIS JARAMILLO MSc

Modelo Jukes-Cantor (JC)

Asume que las cuatro bases tienen igual frecuencia y que las sustituciones son igualmente probables.

$$P_{ij} = \begin{bmatrix} 1 & \alpha & \alpha & \alpha \\ \alpha & 1 & \alpha & \alpha \\ \alpha & \alpha & 1 & \alpha \\ \alpha & \alpha & \alpha & 1 \end{bmatrix} \quad f = [1/4, 1/4, 1/4, 1/4]$$

Modelo Kimura de 2 parámetros (K2P)

Toma en cuenta diferencias entre transiciones vs. transversiones.

$$P_{ij} = \begin{bmatrix} 1 & \beta & \alpha & \beta \\ \beta & 1 & \alpha & \beta \\ \alpha & \alpha & 1 & \beta \\ \beta & \alpha & \beta & 1 \end{bmatrix} \quad f = [1/4, 1/4, 1/4, 1/4]$$

Existen varios modelos que estiman la tasa de cambio entre las secuencias de nucleótidos



BIOINFORMÁTICA Y BIOESTADÍSTICA 2020-1

LUIS JARAMILLO MSc

Felsenstein (1981) (F81)

- Toma en cuenta diferencias en la composición de las bases.
- Porcentaje (G + C) puede variar entre 25% - 75%.
- F81 permite que la frecuencia de cada nucleótido sea diferente. No permite variación en las frecuencias entre genes y especies.

$$P_{ij} = \begin{bmatrix} 1 & \pi_{AT} & \pi_{GT} & \pi_{CT} \\ \pi_{AT} & 1 & \pi_{GT} & \pi_{CT} \\ \pi_{GT} & \pi_{CT} & 1 & \pi_{AT} \\ \pi_{CT} & \pi_{AT} & \pi_{GT} & 1 \end{bmatrix} \quad f = [\pi_A, \pi_C, \pi_G, \pi_T]$$

Hasegawa, Kishino y Yano (1985) (HKY85)

- Esencialmente mezcla modelos K2P and F81, permitiendo la ocurrencia de transiciones y transversiones a distintas tasas y a su vez permitiendo que la frecuencia de bases varie.

$$P_{ij} = \begin{bmatrix} 1 & \pi_{AT} & \pi_{GT} & \pi_{CT} \\ \pi_{AT} & 1 & \pi_{GT} & \pi_{CT} \\ \pi_{GT} & \pi_{CT} & 1 & \pi_{AT} \\ \pi_{CT} & \pi_{AT} & \pi_{GT} & 1 \end{bmatrix} \quad f = [\pi_A, \pi_C, \pi_G, \pi_T]$$

Los modelos también pueden estimar la frecuencia de los nucleótidos



BIOINFORMÁTICA Y BIOESTADÍSTICA 2020-1 LUIS JARAMILLO MSC

modelo General reversible (REV)

Modelo más general – cada sustitución tiene su propia probabilidad.

$$P = \begin{bmatrix} \pi_A & \pi_C & \pi_G & \pi_T \\ \pi_C & \pi_A & \pi_G & \pi_T \\ \pi_G & \pi_C & \pi_A & \pi_T \\ \pi_T & \pi_G & \pi_C & \pi_A \end{bmatrix} \quad f = [\pi_A \pi_C \pi_G \pi_T]$$

Comparando los modelos

Se debe escoger el modelo que mejor se ajuste a los datos

UTEC

BIOINFORMÁTICA Y BIOESTADÍSTICA 2020-1 LUIS JARAMILLO MSC

Comparando los modelos

Permite sesgos transición/transversión: JC, K2P, HKY85, REV

Permite que la frecuencia de bases varíe: JC, K2P, HKY85, REV

Permite sesgos transición/transversión: JC, K2P, HKY85, REV

Existen en la actualidad modelos de sustitución nucleotídica

UTEC

BIOINFORMÁTICA Y BIOESTADÍSTICA 2020-1 LUIS JARAMILLO MSC

1. Calcule la probabilidad para cada sitio.
2. Sume los valores de L para todos los sitios en el árbol.
3. Compare los valores de L para todos los árboles posibles.
4. Elija el árbol con el valor más alto de L.

UTEC

BIOINFORMÁTICA Y BIOESTADÍSTICA 2020-1 LUIS JARAMILLO MSC

The theory says a lot, but does not really bring us any closer to the secret of the 'old one'. I, at any rate, am convinced that **He does not throw dice**.

UTEC

BIOINFORMÁTICA Y BIOESTADÍSTICA 2020-1 LUIS JARAMILLO MSC

Jugando a los dados

- 90 dados normales + 10 dados trucados

Observation Fair Biased

1	1/6	1/21
2	1/6	2/21
3	1/6	3/21
4	1/6	4/21
5	1/6	5/21
6	1/6	6/21

Probabilidad de sacar un dado trucado? hipótesis=trucado

$$\Pr(\text{trucado}) = \frac{10}{90 + 10} = 0.1$$

Sacamos y lanzamos 2 dados: 4 y 6

UTEC

BIOINFORMÁTICA Y BIOESTADÍSTICA 2020-1 LUIS JARAMILLO MSC

- Probabilidad de que no estén trucados:

$$\Pr(4y6|\text{buenos}) = \frac{1}{6} \times \frac{1}{6} = \frac{1}{36} \quad \mathbf{0,0278}$$

- Probabilidad de que estén trucados:

$$\Pr(4y6|\text{trucados}) = \frac{4}{21} \times \frac{6}{21} = \frac{24}{441} \quad \mathbf{0,0544}$$

∴ Dados trucados

Observation Fair Biased

1	1/6	1/21
2	1/6	2/21
3	1/6	3/21
4	1/6	4/21
5	1/6	5/21
6	1/6	6/21

Probabilidad de los datos dado el modelo, i.e. el LIKELIHOOD

UTEC

Hallando la probabilidad posterior: Probabilidad del modelo dados los datos

$$\Pr(\text{trucado}) = \frac{10}{90 + 10} = 0.1$$

$$\Pr(4y6|\text{trucados}) = \frac{4}{21} \times \frac{6}{21} = \frac{24}{441}$$

$$\Pr(\text{trucados}|4y6) = \frac{\Pr(\text{trucados}) \Pr(4y6|\text{trucados})}{\Pr(4y6)}$$

$$\Pr(4y6|\text{trucados}) \times 0.1 + \Pr(4y6|\text{buenos}) \times 0.9$$

$$\Pr(4y6|\text{trucados}) = \frac{4}{21} \times \frac{6}{21} = \frac{24}{441}$$

$$\Pr(4y6|\text{buenos}) = \frac{1}{6} \times \frac{1}{6} = \frac{1}{36}$$

$$\Pr(\text{trucados}|4y6) = 0.17877 > P(\text{trucados}) = 0.10$$

BIOINFORMÁTICA Y BIOESTADÍSTICA 2020-1

LUIS JARAMILLO MSc

Conclusiones

- ✓ Modelos que estiman la tasa de cambio entre las secuencias de nucleótidos (JP y K2P).
- ✓ Los modelos también pueden estimar la frecuencia de los nucleótidos (F81 y HKY85).
- ✓ Probabilidad de modelos de sustitución nucleotídica.



Gracias

ERES MÁS QUE CAPAZ
DE ENCONTRAR EL
ÉXITO, PERO SOLO
OCURRIRÁ SI TE
PONES A ELLO !!

