

### Introducción a la Inferencia Filogenética y Evolución Molecular

23-26 Junio 2008, Fac. C. Biológicas - UANL



Pablo Vinuesa ([vinuesa@ccg.unam.mx](mailto:vinuesa@ccg.unam.mx))



Centro de Ciencias Genómicas-UNAM, México  
<http://www.ccg.unam.mx/~vinuesa/>

Todo el material del curso lo puedes descargar desde:

<http://www.ccg.unam.mx/~vinuesa/UANL08>

#### • Tema 6: Modelos de evolución de secuencias y prueba de bootstrap

1. El uso de modelos en ciencia y en filogenética
2. Modelos empíricos vs. paramétricos
3. Derivación de matrices de sustitución empíricas a partir de alineamientos múltiples de proteínas
4. Modelos paramétricos de evolución de secuencias de DNA - la familia GTR
5. Modelos y corrección de distancias genéticas
6. La prueba de bootstrap para determinar la confiabilidad estadística de biparticiones

### Modelos de evolución de secuencias - introducción

#### • Modelos de evolución del proceso de sustitución y métodos de reconstrucción filogenética: consideraciones generales

- 1.- La reconstrucción o estima filogenética es un **problema de inferencia estadística**, y como tal **requiere un modelo** de sustitución de residuos (aa o nt), es decir, un modelo de evolución molecular de las secuencias. Todos los modelos, por no ser más que aproximaciones de los procesos naturales, hacen una serie de **suposiciones** (simplificaciones)
- 2.- Los **modelos de evolución de secs.** son usados en filogenética **para describir las probabilidades con las que se dan los distintos eventos de sustitución entre aa o nt**, con el fin de **corregir o compensar las sustituciones no observadas a lo largo de la filogenia**
- 3.- Mientras que los métodos de MP asumen un **modelo implícito** de evolución (número mínimo de sustituciones a lo largo de la filogenia), los métodos de distancia (UPGMA, NJ), los de ML y Bayesianos requieren de un **modelo explícito** de evolución
- 4.- Los **métodos de distancia** estiman finalmente un sólo parámetro (no. sust./sitio) dado el modelo y el valor de los parámetros del mismo; en cambio, los **métodos de ML y Bayesianos pueden estimar el valor de cada uno de los parámetros del modelo explicitado**, dada una topología y la matriz de datos (alineamiento)

### Modelos de evolución de secuencias - introducción

• Para el análisis filogenético de secuencias alineadas virtualmente todos los métodos describen la evolución de las secuencias usando un modelo que consta de dos componentes:

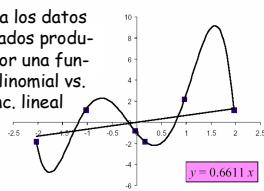
1. un árbol filogenético
2. una descripción de las probabilidades con las que se dan las sustituciones de aa o nts a lo largo de las ramas del árbol

• ¿Porqué necesitamos modelos y para qué sirven?

- Los modelos nos sirven para interpolar adecuadamente entre nuestras observaciones con el fin de poder hacer predicciones inteligentes sobre observaciones futuras

$$y = -1.5972x^5 + 23.167x^4 - 126.18x^3 + 319.17x^2 - 369.22x + 155.67$$

ajuste a los datos observados producidos por una función polinomial vs. una func. lineal



- **añadir parámetros** a un modelo generalmente mejora su ajuste a los datos observados
- **modelos infra-parametrizados** conducen a un pobre ajuste a los datos observados
- **modelos supra-parametrizados** conducen a una pobre predicción de eventos futuros
- existen métodos estadísticos para **seleccionar modelos ajustados** a cada set de datos

### Modelos de evolución de secuencias - introducción

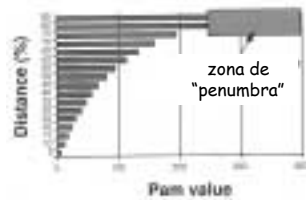
#### • Modelos de evolución del proceso de sustitución y métodos de reconstrucción filogenética: consideraciones generales

#### Corolario:

1. El grado de confianza que tengamos en una filogenia particular realmente depende de la que tengamos en el modelo subyacente
2. Por lo tanto, siempre que usemos un método basado en un modelo explícito de evolución (NJ, ML, By) es necesario usar rigurosas pruebas estadísticas para seleccionar el modelo y el valor de sus parámetros que mejor se ajusten a la matriz de datos a analizar



**Alineamiento de proteínas: selección de matrices de ponderación - consejos prácticos para la identificación de homólogos**



Distancias observadas vs. evolutivas (PAM) entre prots.

Diferencia % obs.	Dist. evol. PAM
1	1
5	5
10	11
15	17
20	23
30	38
40	56
50	80
60	112
70	159
80	246
85	328

← z. penumbra

A medida que el nivel de divergencia evolutiva entre pares de proteínas incrementa (distancias PAM) disminuye el número de diferencias observadas, debido a fenómenos de **reversión (homoplasia)**. Por tanto, si no se cuenta con evidencia estructural, el análisis filogenético de proteínas debe restringirse a aquellas con  $\geq 20\%$  de identidad. **Los alns. tampoco son confiables**

- A medida que el nivel de divergencia entre pares de proteínas alcanza el valor de PAM250 (~ 20% identidad), comienza a ser dudosa su relación de homología, pudiendo tratarse de secuencias que presentan cierto grado de similitud por azar, en base a composiciones de AAs similares en ambas secuencias !!!
- Al entrar en esta **zona de penumbra**, es esencial considerar información adicional, particularmente motivos estructurales, para validar o descartar una posible relación de homología

**Modelos de evolución de secuencias -DNA**

**Modelos de sustitución de nucleótidos**

El modelaje de la evolución a nivel del DNA se ha concentrado en la aproximación paramétrica. Se manejan **tres tipos principales de parámetros** en estos modelos:

- parámetros de **frecuencia**
- parámetros de **tasas de intercambio**
- parámetros de **heterogeneidad de tasas de sustitución** entre sitios

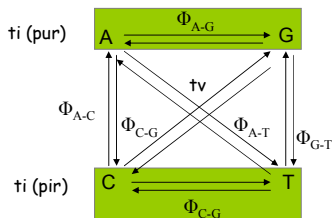
**Modelos de evolución de sustitución de nucleótidos -modelos paramétricos**

los diversos modelos evolutivos se distinguen por su grado de parametrización

**I. Frecuencias de nt** :  $\pi_A = \pi_C = \pi_G = \pi_T = 0.25$  ó  $\pi_A \neq \pi_C \neq \pi_G \neq \pi_T$

- modelos de = frecuencia: JC69; K2P, K3P ...
- modelos de  $\neq$  frecuencia: F81, HKY85, TrN93, GTR ...

**II. Tasas de sustitución transicionales/transversionales**

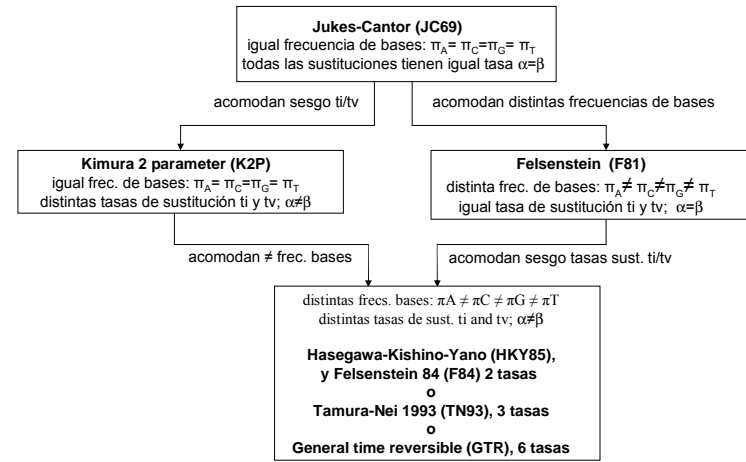


Existen 4 tipos de sustituciones  $t_i$  y 8  $t_v$ ; cuando  $t_i/t_v \neq 0.5$  existe un sesgo en sustituciones  $t_i$  (o  $t_v$ ) en el set de datos.  $t_i$  generalmente  $\gg 1$

los modelos evolutivos se diferencian también en la cantidad de parámetros que utilizan para acomodar diversas tasas de sustitución:

tasas	modelo
1	JC69 ( $t_i=t_v$ )
2	K2P ( $t_i \neq t_v$ )
3	TrN ó K3P (2 $t_i$ , 1 $t_v$ )
6	GTR (cada sust. su tasa)

**Modelos básicos de evolución de DNA: la familia de modelos anidados GTR o REV**



### Modelos básicos de evolución de DNA: la familia de modelos anidados GTR o REV

• Matriz de tasas de sustitución instantáneas del modelo GTR

	A	C	G	T
A	-	$\pi_C a \mu$	$\pi_G b \mu$	$\pi_T c \mu$
C	$\pi_A a \mu$	-	$\pi_G d \mu$	$\pi_T e \mu$
G	$\pi_A b \mu$	$\pi_C d \mu$	-	$\pi_T f \mu$
T	$\pi_A c \mu$	$\pi_C e \mu$	$\pi_G f \mu$	-

9 parámetros

 $\pi_A$   
 $\pi_C$   
 $\pi_G$   
 $a$   
 $b$   
 $c$   
 $d$   
 $e$   
 $f$   
 $\mu$

←  $-\mu (\pi_A c + \pi_C e + \pi_G f)$

El modelo GTR es idéntico al de JC69 si  $a = b = c = d = e = f = 1$  y todas las bases se asumen que tienen igual frecuencia ( $\frac{1}{4}$ )

$\mu$  = tasa del proceso generador de todos los tipos de sustituciones  
 $a, \dots, e$  = modificadores de tasa relativa de cada tipo particular de sustitución  
 $\pi$  = frecuencia de cada nt

### Modelos básicos de evolución de DNA: la familia de modelos anidados GTR o REV

1 parámetro ( $\alpha$ )

Incremento en el número de parámetros  
modelos más generales

11 parámetros libres a estimar

$\pi_A, \pi_C, \pi_G$   
 $a, b, c, d, e$   
 $\mu,$   
 $I, T$

→

En total existen 203 modelos posibles en la familia GTR al combinar params. de frec., tasa, G e I. La mayoría de ellos carecen de nombre.

### Comparación empírica de modelos sust. de DNA

• Comparación de los modelos de **JC69** y **K2P** en su capacidad de corregir distancias observadas ( $p$ ) entre pares de secuencias según su grado de divergencia

$$d_{JC69} = -\frac{3}{4} \ln \left( 1 - \frac{4}{3}p \right) \quad \text{vs.} \quad d_{K2P} = \frac{1}{2} \ln \left( \frac{1}{1-2P-Q} \right) + \frac{1}{4} \ln \left( \frac{1}{1-2Q} \right)$$

• Escenario I:

- sean 2 secs. de long. = 200 nt, que difieren en 20 *ti* y 4 *tv*

por lo tanto  $L = 200, P = 20/200 = 0.1$  y  $Q = 4/200 = 0.02$

$p = 24/200 = 0.12$   
 $d_{JC69} \approx 0.13$  (sust./sitio)  
 no. de sust. esperadas =  $0.13 \times 200 \approx 26$

$d_{K2P} \approx 0.13$  (sust./sitio)  
 no. de sust. esperadas =  $0.13 \times 200 \approx 26$

### Comparación empírica de modelos sust. de DNA

• Comparación de los modelos de **JC69** y **K2P** en su capacidad de corregir distancias observadas ( $p$ ) entre pares de secuencias según su grado de divergencia

$$d_{JC69} = -\frac{3}{4} \ln \left( 1 - \frac{4}{3}p \right) \quad \text{vs.} \quad d_{K2P} = \frac{1}{2} \ln \left( \frac{1}{1-2P-Q} \right) + \frac{1}{4} \ln \left( \frac{1}{1-2Q} \right)$$

• Escenario II:

- sean 2 secs. de long. = 200 nt, que difieren en 50 *ti* y 16 *tv*

por lo tanto  $L = 200, P = 50/200 = 0.25$  y  $Q = 16/200 = 0.08$

$p = 66/200 = 0.33$   
 $d_{JC69} \approx 0.43$  (sust./sitio)  
 no. de sust. esperadas =  $0.43 \times 200 \approx 86$

$d_{K2P} \approx 0.48$  (sust./sitio)  
 no. de sust. esperadas =  $0.48 \times 200 \approx 96$

