

Introducción a la Filoinformática:  
LEG-UM5, Rabat. 10-14 Junio 2019.

Pablo Vinuesa (vinuesa@ccg.unam.mx)

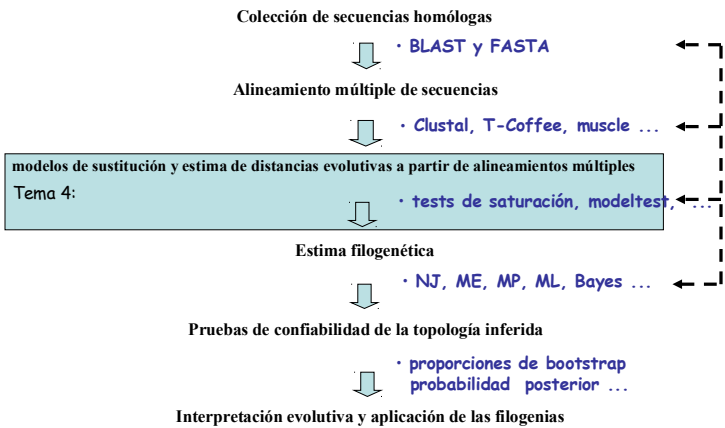
Programa de Ingeniería Genómica, CCG-UNAM, México  
<http://www.ccg.unam.mx/~vinuesa/>

Todo el material del curso (presentaciones, tutoriales y datos de secuencias) lo encontrarás en:  
<https://github.com/vinuesa/intro2phyloinfo>

• Tema 4: Intro a la filogenética y  
modelos de evolución de secuencias


- 1. Para qué sirven los modelos en ciencia
- 2. Modelos paramétricos vs. empíricos de evolución de secuencias
- 3. Parametrización de los modelos de sustitución de DNA
- 4. Condiciones (supuestos) de aplicabilidad de los modelos
- 5. Modelos de sustitución y distancias evolutivas
- 6. La familia GTR(+I+G) de modelos de sustitución para secuencias nucleotídicas

Protocolo básico para un análisis filogenético de  
secuencias moleculares



Inferencia filogenética molecular –  
clasificación de métodos

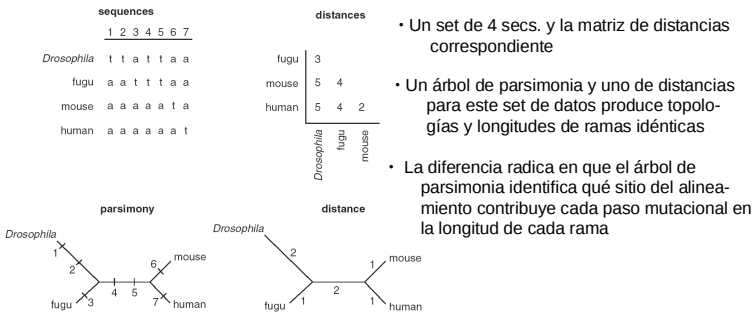
- Podemos clasificar a los métodos de reconstrucción filogenética en base al tipo de datos que emplean (**caracteres discretos vs. distancias**) y si usan un **método algorítmico** o un **criterio de optimización** para encontrar la topología

		Tipo de datos	
Método de reconstrucción	algoritmo de agrupamiento	distancias	caracteres discretos
		UPGMA  Neighbour joining	
	criterio de optimización	Evolución mínima	Máxima parsimonia
			Máxima verosimilitud

Métodos de reconstrucción filogenética – una clasificación

I.- Tipos de datos: distancias vs. caracteres discretos

- Los **métodos de distancia** primero convierten los alineamientos de secuencias en una matriz de distancias genéticas en base al modelo evolutivo seleccionado, la cual es usada por el método algorítmico de reconstrucción para calcular el árbol (**UPGMA y NJ**)
- Los **métodos discretos (MP, ML, Bayesianos)** consideran cada sitio del alineamiento (o una función probabilística para cada sitio) directamente



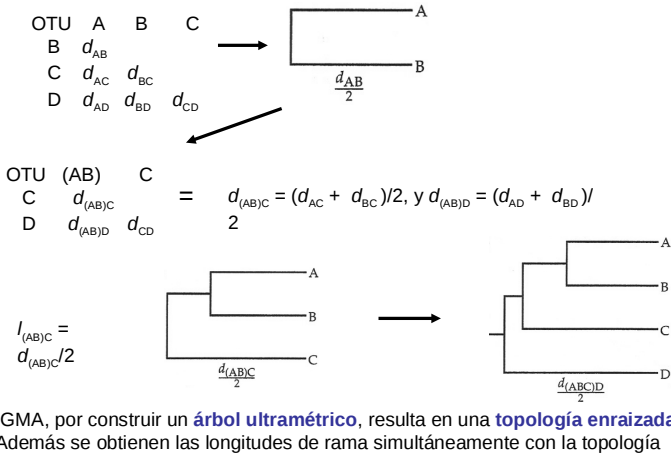
Inferencia filogenética molecular – métodos basados en matrices de distancias

• Unweighted pair group method with arithmetic means (UPGMA)

- este es uno de los pocos métodos que construye **árboles ultramétricos** (todas las hojas equidistantes de la raíz), es decir **asume un reloj molecular** perfecto a lo largo de toda la topología, lo que resulta en una **topología enraizada**. Además se obtienen las longitudes de rama simultáneamente con la topología
- se puede concebir como un método heurístico para encontrar la topología ultramétrica de mínimos cuadrados para una matriz de distancias pareadas

Inferencia filogenética molecular – métodos basados en matrices de distancias

• Unweighted pair group method with arithmetic means (UPGMA)



Ejercicio:

Calcula una matriz de distancias pareadas en base al número observado de diferencias entre OTUs, y en base a ella dibuja un árbol de UPGMA, indicando las longitudes de cada rama

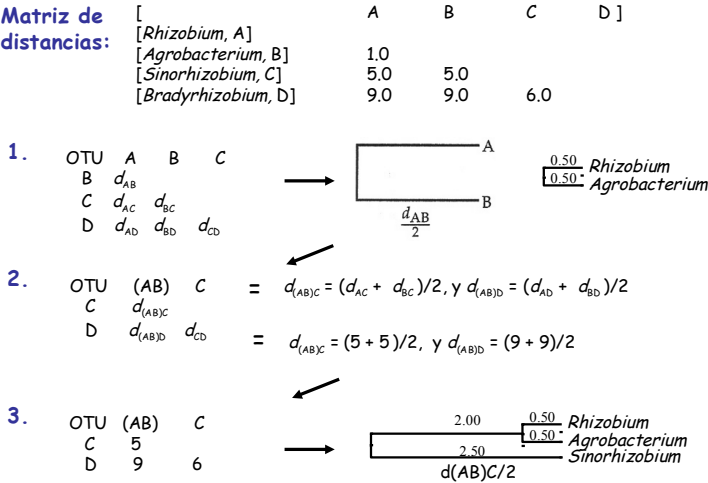
1. Alineamiento: No. sitios : 15; OTUs (taxa) = 4

<i>Rhizobium</i>	GGA	GGG	AGG	AGG	CCT
<i>Agrobacterium</i>	GGC	GGG	AGG	AGG	CCT
<i>Sinorhizobium</i>	GGG	GGA	AGG	TGT	CCG
<i>Bradyrhizobium</i>	GGT	CGT	AGC	TGT	GTG

2. Matriz de distancias:  $d$  : distancia (no. de diferencias observadas)

	A	B	C	D
[ <i>Rhizobium</i> , A]				
[ <i>Agrobacterium</i> , B]	1.0			
[ <i>Sinorhizobium</i> , C]	5.0	5.0		
[ <i>Bradyrhizobium</i> , D]	9.0	9.0	6.0	

Inferencia de un árbol UPGMA usando el no. de dif. obs. como medida de la distancia genética entre OTUs



**Inferencia de un árbol UPGMA usando el no. de dif. obs. como medida de la distancia genética entre OTUs**

**Matriz de distancias:**

	A	B	C	D
[Rhizobium, A]				
[Agrobacterium, B]	1.0			
[Sinorhizobium, C]	5.0	5.0		
[Bradyrhizobium, D]	9.0	9.0	6.0	

4.  $d_{(ABC)D} = (d_{AB} + d_{BC} + d_{CD}) / 3$   
 $d_{(ABC)D} = (9 + 9 + 6) / 3 = 8$

5.

**Inferencia de un árbol UPGMA usando el no. de dif. obs. como medida de la distancia genética entre OTUs**

**Matriz de distancias:**

	A	B	C	D
[Rhizobium, A]				
[Agrobacterium, B]	1.0			
[Sinorhizobium, C]	5.0	5.0		
[Bradyrhizobium, D]	9.0	9.0	6.0	

• ¿Notan alguna inconsistencia entre las distancias topológicas y observadas?

- La distancia entre C y D no es aditiva y no queda adecuadamente reflejada en la correspondiente longitud de rama

**Métodos de reconstrucción filogenética - una clasificación**

**III. máxima parsimonia: dados dos árboles, se prefiere el que requiere menos cambios en estados de carácter**

• El **método de máxima parsimonia (MP)** considera cada sitio filogenéticamente informativo (**Pi**) el alineamiento (al menos 2 pares de secuencias que compartan un polimorfismo distinto). Los sitios constantes (**C**) no son considerados y los singletons (**S**) no son Pars. informativos

• El supuesto teórico (modelo de evolución) implícito al método es que el árbol más verosímil es aquel que requiere el mínimo número de sustituciones para explicar los datos del alineamiento. **El criterio de optimización de la MP es el de cambio o evolución mínima.**

• Para cada sitio del alineamiento el objetivo es reconstruir su evolución bajo la constrictión de **invocar el número mínimo de pasos evolutivos**. El número total de cambios evolutivos sobre un árbol de MP (longitud en pasos evolutivos del árbol) es simplemente la suma de cambios de estados de carácter (p. ej. sustituciones) de cada sitio variable

**Clases de sitios:**  
Pi= Pars. inform.  
C= Constant  
S= Singleton

**reconstrucciones para el sitio 2**

**tree 1**

**tree 2**

**tree 3**

**events per tree**

1	0	2	0	1	2	2	0	2	1	2	12
2	0	2	0	1	2	1	0	2	2	2	12
3	0	1	0	1	1	2	0	1	2	2	10

**Total**

**Métodos de búsqueda de árboles**

**I. - el problema del número de topologías**

El número de topologías posibles incrementa factorialmente con cada nuevo taxon o secuencia que se añade al análisis

No. de árboles no enraizados =  $(2n-5)!!/2^{n-3}(n-3)$

No. de árboles enraizados =  $(2n-3)!!/2^{n-2}(n-2)$

Taxa	árboles no enraiz*	árb. enraiz.
4	3	15
8	10,395	135,135
10	2,027,025	34,459,425
22	$3 \times 10^{23}$	...
50	$3 \times 10^{74}$	...

\*por ej. para sólo 15 OTUs tenemos 213,458,046,676,875 topologías

- i si pudiésemos evaluar  $1 \times 10^6$  topol./seg. necesitaríamos 6 años y 9 meses para completar la búsqueda! El no. de Avogadro es  $\sim 6 \times 10^{23}$  (átomos/mol). Según la teor. de la relatividad de la estructura del universo de Einstein, existen  $10^{80}$  átomos de  $H_2$  en el universo ...

[http://en.wikipedia.org/wiki/Observable\\_universe](http://en.wikipedia.org/wiki/Observable_universe)

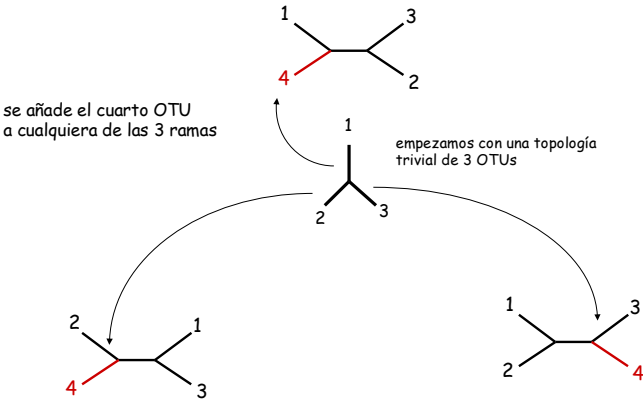
Por tanto se requieren **estrategias heurísticas de búsqueda árboles** cuando se emplean métodos basados en criterios de optimización y  $n > \sim 25$

Métodos de búsqueda de árboles

- Pasos lógicos de los métodos filogenéticos basados en criterios de optimización (MP, ML ...)
- 1. definir el criterio de optimización (descrito formalmente en una **función objetiva**)
- 2. Construir un árbol de partida que contenga todos los OTUs
- 3. Emplar **algoritmos de búsqueda** que tratan de encontrar árboles mejores bajo el criterio de optimización escogido que el árbol actual o de partida.

1. Criterios de optimización	2. Estrategias de búsqueda	
Parsimonia	Enumeración exhaustiva ( $n \leq 12$ ) (exhaustive enumeration)	<b>Métodos exactos:</b> garantizan encontrar la topología óptima
Máxima verosimilitud	Ramificación y límite ( $n \leq 25$ ) (branch-and-bound)	
Evolución Mínima	Decomposición en estrella (star decomposition)	<b>Métodos heurísticos:</b> no garantizan encontrar la topología óptima
Mínimos cuadrados	Adición secuencial (stepwise addition)	
	(Inter-)cambio de rama (branch swapping)	

Métodos de búsqueda de árboles  
-enumeración exhaustiva ( $n \leq 12$ )

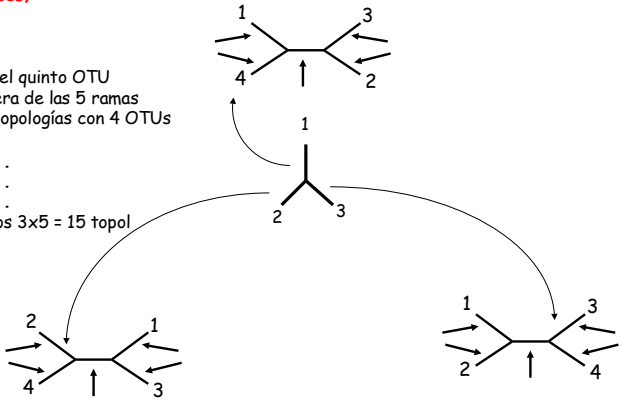


Métodos exactos de búsqueda de árboles  
-enumeración exhaustiva ( $n \leq 12$ )

PAUP\* command:  
**alltrees;**

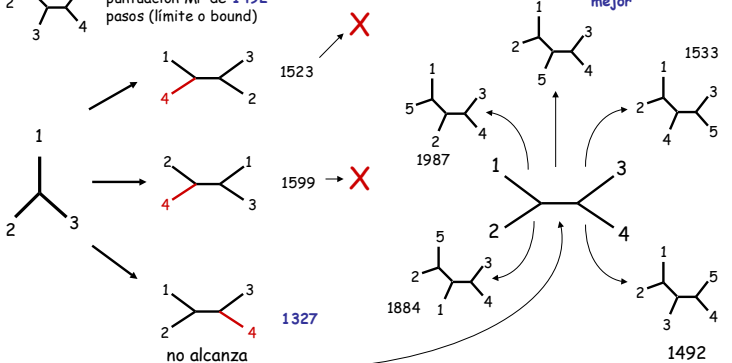
se añade el quinto OTU a cualquiera de las 5 ramas de las 3 topologías con 4 OTUs

obtenemos  $3 \times 5 = 15$  topol



Métodos exactos de búsqueda de árboles  
- "branch and bound" ( $n \leq 25$ )

árbol obtenido por un método heurístico ó NJ con puntuación MP de **1492** pasos (límite o bound)

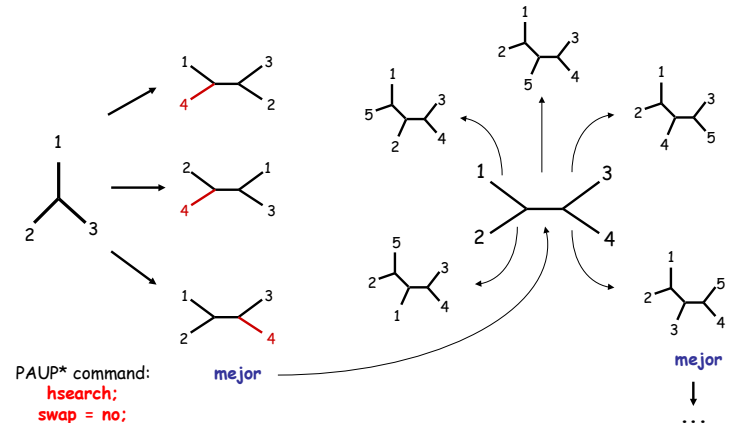
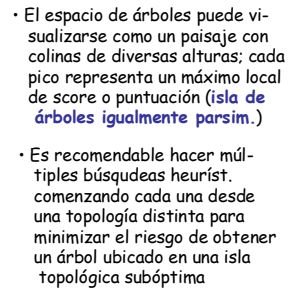


• PAUP\* command:  
**bandb;**



• Al igual que la búsqueda exhaustiva, **garantiza encontrar el árbol óptimo**

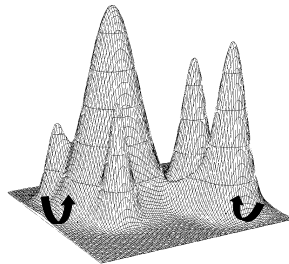
## Métodos heurísticos de búsqueda de árboles - adición secuencial (aleatorizada)

- Este método se usa con frecuencia para generar distintos "árboles semilla" a partir de los cuales comenzar búsquedas heurísticas, partiendo de "distintos puntos del espacio de árboles"



## Métodos heurísticos de búsqueda de árboles

- (a)  árbol estrella para NOTUS
- (b) 
- PAUP\* command:  
**stardecomp;**
- mejor  $\rightarrow$  hasta unir las (N-3) posibles ramas internas  
puntuación
- $N(N-1)/2$   
· modos de  
· buscar pares

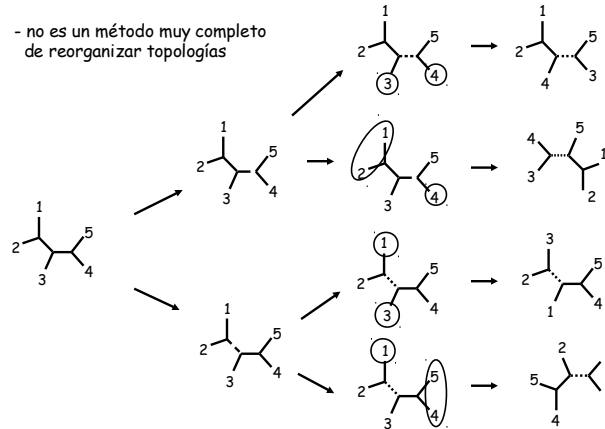


- NJ usa este método junto al criterio de evolución mínima
- una vez que 2 OTUs han sido unidos ya no pueden ser desacoplados más adelante; en esto difiere del algoritmo de adición secuencial
- sensible al orden en que se van uniendo los OTUs; problema incrementa con el no. de OTUs
- no debe ser por tanto usado como método de búsqueda definitivo
- buena estrategia para producir árboles iniciales que sean mejorados mediante otras estrategias heurísticas

### Métodos heurísticos de búsqueda de árboles - intercambio de ramas (branch swapping)

- Intercambio entre vecinos más próximos (Nearest Neighbor Interchange, NNI)

- no es un método muy completo de reorganizar topologías



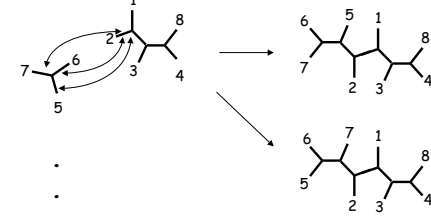
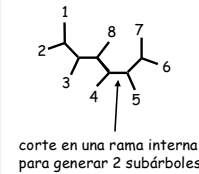
PALP\* cmd: hsearch swap=nni start=stepwise addseq=random;

### Métodos heurísticos de búsqueda de árboles - intercambio de ramas (branch swapping)

- Bisección-reconexión de árboles (Tree Bisection-Reconnection, TBR)

-Este método evalúa muchas más topols. que el NNI

se **reconectan** los dos subárboles en **todas las posiciones posibles** (ej:  $3 \times 5 = 15$  subarreglos en nuestro ejemplo)



se repite esta operación para reconectar el subárbol chico en las ramas terminales 1, 8, 4 y 3 del subárbol grande

PALP\* cmd: hsearch swap=tbr start=stepwise addseq=random;

### 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:

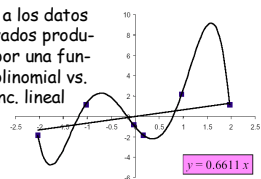
- un árbol filogenético
- una descripción de las probabilidades con las que se dan las sustituciones de aa o nts a lo largo de las ramas del árbol

- ¿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



$$y = 0.6611x$$

- añadir parámetros** a un modelo generalmente mejora su ajuste a los datos observados

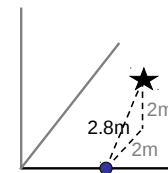
- modelos infra-parametrizados** conducen a un pobre ajuste a los datos observados

- modelos supra-parametrizados** conducen a una pobre predicción de eventos futuros

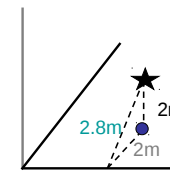
- existen métodos estadísticos para **seleccionar modelos ajustados** a cada set de datos

### Modelos de evolución de secuencias -introducción

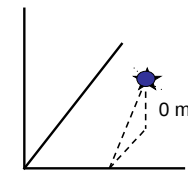
**dimensiones de un modelo:** cada parámetro en un modelo estadístico puede ser concebido como la adición de una nueva dimensión, tal y como se ilustra en el ejemplo siguiente:



- En este **modelo 1D** lo más cerca que podemos llegar al pto. marcado en un espacio 3D es 2.8 m



- En este **modelo 2D** lo más cerca que podemos llegar al pto. marcado en un espacio 3D es 2 m

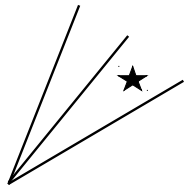


- En este **modelo 3D** podemos aproximar el punto exactamente. El modelo 3D se ajusta 100% a la realidad espacial.

**En el caso de modelos de sustitución nunca obtendremos un ajuste del 100% entre el modelo y la realidad.** Todos los modelos son sólo aproximaciones de la realidad, pero algunos modelos son útiles para describir el proceso de sustitución (y otros mucho menos).

### Modelos de evolución de secuencias -introducción

**Dimensiones de un modelo:** en realidad, los **parámetros** de un modelo complejo **no** son siempre **independientes**, existiendo diversos grados de **colinealidad**. En el peor de los casos, los parámetros pueden ser totalmente colineales, en cuyo caso uno de ellos es 100% redundante, por lo que no aporta nada a la fuerza del modelo para explicar los datos observados



- uno de los objetivos primordiales de los modelos de sustitución de nt y aa es el de **incorporar los parámetros más relevantes, que expliquen características fundamentales de las secuencias** cuya evolución tratan de modelar de la manera más realista posible

- En este **modelo 3D** existe un nivel significativo de **colinealidad** entre sus dimensiones (o parámetros)

### 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

### Modelos de evolución de secuencias -introducción

- **Modelos de evolución del proceso de sustitución y métodos de reconstrucción filogenética: consideraciones generales**

- Existen dos aproximaciones para construir modelos de evolución de secuencias.

1. construcción de **modelos empíricos** basados en propiedades del proceso de sustitución calculadas a partir de comparaciones de un gran número de secuencias. Los modelos empíricos **resultan en valores fijos de los parámetros**, los cuales son estimados sólo una vez, suponiéndose que son adecuados para el análisis de otros sets de datos. Esto los hace fácil de usar e implementar en términos computacionales, pero su utilidad real para cada caso particular ha de ser evaluada críticamente
2. construcción de **modelos paramétricos** basado en el modelaje de propiedades químicas o genéticas de los aas y nts. Los modelos paramétricos tienen la ventaja de que los **valores de los parámetros pueden ser derivados de cada set de datos** al hacer un análisis de los mismos usando métodos de ML o By, por tanto ajustándolos a cada matriz de datos particular

### 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:

1. parámetros de **frecuencia**
2. parámetros de **tasas de intercambio**
3. 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

ti (pur)

A

G

$\Phi_{A-G}$

$\Phi_{G-A}$

C

T

$\Phi_{C-G}$

$\Phi_{G-C}$

A

C

$\Phi_{A-C}$

$\Phi_{C-A}$

G

T

$\Phi_{G-T}$

$\Phi_{T-G}$

A

T

$\Phi_{A-T}$

$\Phi_{T-A}$

C

T

$\Phi_{C-T}$

$\Phi_{T-C}$

ti (pir)

Existen 4 tipos de sustituciones ti y 8 tv; cuando  $ti/tv \neq 0.5$  existe un sesgo en sustituciones ti (o tv) en el set de datos. ti 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 (ti=tv)
2	K2P (ti $\neq$ tv)
3	TrN ó K3P (2 ti, 1 tv)
6	GTR (cada sust. su tasa)

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

Jukes-Cantor (JC69)  
igual frecuencia de bases:  $\pi_A = \pi_C = \pi_G = \pi_T$   
todas las sustituciones tienen igual tasa  $\alpha = \beta$

acomodan sesgo ti/tv

acomodan distintas frecuencias de bases

Kimura 2 parameter (K2P)  
igual frec. de bases:  $\pi_A = \pi_C = \pi_G = \pi_T$   
distintas tasas de sustitución ti y tv;  $\alpha \neq \beta$

Felsenstein (F81)  
distinta frec. de bases:  $\pi_A \neq \pi_C \neq \pi_G \neq \pi_T$   
igual tasa de sustitución ti y tv;  $\alpha = \beta$

acomodan  $\neq$  frec. bases

acomodan sesgo tasas sust. ti/tv

distintas frecs. bases:  $\pi_A \neq \pi_C \neq \pi_G \neq \pi_T$   
distintas tasas de sust. ti and tv;  $\alpha \neq \beta$

Hasegawa-Kishino-Yano (HKY85),  
y Felsenstein 84 (F84) 2 tasas

Tamura-Nei 1993 (TN93), 3 tasas

General time reversible (GTR), 6 tasas

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

JC

JC+I

JC+ $\Gamma$

F81

K80

JC+I+ $\Gamma$

F81+I

K80+I

F81+ $\Gamma$

K80+ $\Gamma$

HKY

SYM

F81+I+ $\Gamma$

K80+I+ $\Gamma$

HKY+I

SYM+I

HKY+ $\Gamma$

SYM+ $\Gamma$

GTR

HKY+I+ $\Gamma$

SYM+I+ $\Gamma$

GTR+I

GTR+ $\Gamma$

GTR+I+ $\Gamma$

11 parámetros libres a estimar

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

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

Modelos de evolución de sustitución de nucleótidos  
-modelos paramétricos

Condiciones de aplicabilidad de los modelos (supuestos)

1.- Supuesto de independencia: las mutaciones en un sitio no afectan a otros en la secuencia. Violado por ej. en el caso de rRNAs suelen seleccionarse mutaciones compensatorias (evolución covariada entre sitios)

2.- Supuesto de homogeneidad de tasas de sustitución a lo largo del tiempo y entre linajes: en este supuesto se basa el reloj molecular y de su cumplimiento depende la posibilidad de poder utilizar un "reloj molecular" para datar clados

3.- Las frecuencias de nucleótidos son homogéneas entre linajes: este supuesto es frecuentemente violado cuando usamos secuencias de linajes muy distantes, particularmente en procariontes, ya que los contenidos de G+C de distintos grupos microbianos varía mucho, del 22 % (*Wigglesworthia, gamma-Proteobacteria*) al 75 % (*Anaeromyxobacter, delta-Proteobacteria*)

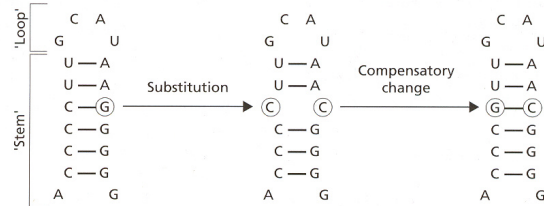
4.- Las probabilidades de sustitución son las mismas para cada sitio: este supuesto es violado casi sin excepción. Así por ejemplo, las 3as. pos. de los codones acumulan mutaciones mucho más rápidamente que la 2a y 1a. Los distintos dominios de una proteína o rRNA también evolucionan con tasas distintas. Distribución Gamma ( $\Gamma$ )

© Pablo Vinuesa 2019,  
vinuesATccg[dot]unam[dot]mx,  
<http://www.ccg.unam.mx/~vinuesa/>



### Condiciones de aplicabilidad de los modelos (supuestos)

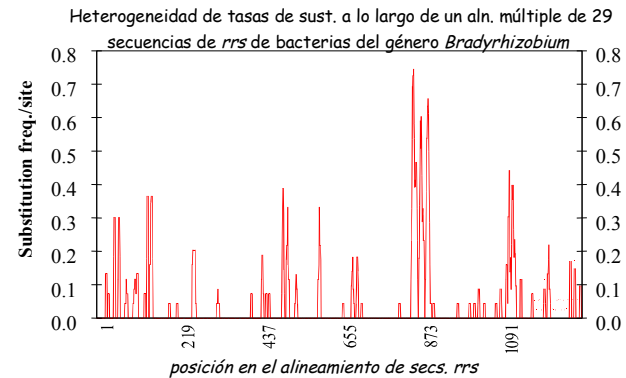
1.- **Supuesto de independencia y modelos de covariación:** las mutaciones en un sitio no afectan a otros en la secuencia. En el caso de rRNAs suelen seleccionarse mutaciones compensatorias (**covariación de sitios**). Existen modelos que acomodan este hecho.



### Condiciones de aplicabilidad de los modelos (supuestos)

Acomodo de la heterogeneidad de tasas de sustitución entre sitios

(**I**) acomoda las posiciones invariables (**proporción de sitios invariables**)  
(**Γ**) acomoda la **heterogeneidad de tasas de sust.** entre las posiciones variables



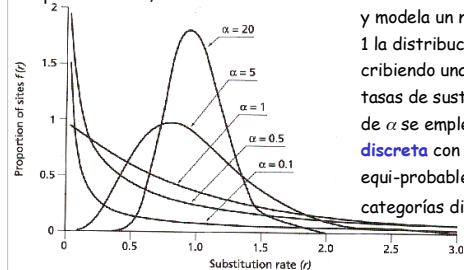
### Condiciones de aplicabilidad de los modelos (supuestos)

2.- **Distribución gamma y heterogeneidad de tasas de sust. entre sitios:** Para modelar con cierto realismo el proceso de sustitución es esencial acomodar adecuadamente la heterogeneidad de tasas de sustitución entre sitios de un alineamiento

$$Pdf(r) = \alpha^\alpha r^{\alpha-1} / \exp(\alpha r) \Gamma(\alpha)$$

Diversas formas de la distribución gamma ( $\Gamma$ ) para un rango de valores del parámetro  $\alpha = 1/CV^2$ , donde CV = coef. de var. de las tasas.

Así para CV = 0.3,  $\alpha = 1/0.09 = 11.1111$



Para ello se asume generalmente una **distribución gamma ( $\Gamma$ )** de las tasas y que cada sitio tiene una tasa tomada aleatoriamente de dicha distribución, e independientemente de los demás sitios. El parámetro  $\alpha$  controla la forma de la distribución. Para  $\alpha > 1$  la distribución tiene forma de campana y modela un nivel bajo de heterogeneidad. Para  $\alpha < 1$  la distribución toma forma de L invertida, describiendo una situación de fuerte heterog. de tasas de sust. entre sitios. Para calcular el valor de  $\alpha$  se emplea generalmente una **distribución  $\Gamma$  discreta** con un número  $c$  finito de tasas des sust. equi-probables ( $q_1, q_2, \dots, q_c$ ). El uso de 4 a 8 categorías discretas permite obtener una buena aprox. de la distrib. continua.

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

• El método de momentos es de utilidad limitada en estadística (y filogenética) ya que no permite obtener una fórmula explícita para calcular la distancia entre secuencias usando modelos más complejos como el HKY85 o GTR

• La fórmula explícita de distancia para el modelo K2P es:

$$d = \frac{1}{2} \ln \left( \frac{1}{1 - 2P - Q} \right) + \frac{1}{4} \ln \left( \frac{1}{1 - 2Q} \right)$$

-este modelo tiene 2 parámetros,  $P$  y  $Q$  (proporción de  $ti$  y  $tv$  en que difieren 2 secuencias, donde  $p = P + Q$ )

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

- 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  $Q = 4/200 = 0.02$

$$p = 24/200 = 0.12$$

$$d_{JC69} \approx 0.13 \text{ (sust./sitio)}$$

$$d_{K2P} \approx 0.13 \text{ (sust./sitio)}$$

$$\text{no. de sust. esperadas} = 0.13 \times 200 \approx 26$$

$$\text{no. de sust. esperadas} = 0.13 \times 200 \approx 26$$

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

- 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  $Q = 16/200 = 0.08$

$$p = 66/200 = 0.33$$

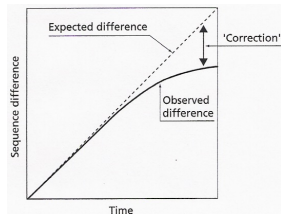
$$d_{JC69} \approx 0.43 \text{ (sust./sitio)}$$

$$d_{K2P} \approx 0.48 \text{ (sust./sitio)}$$

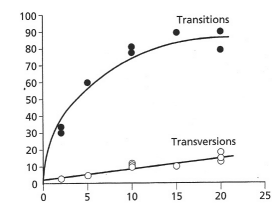
$$\text{no. de sust. esperadas} = 0.43 \times 200 \approx 86$$

$$\text{no. de sust. esperadas} = 0.48 \times 200 \approx 96$$

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



- El objetivo de los modelos de sustitución es el de **compensar para los eventos homoplásicos de múltiples sustituciones**, y así obtener estimas de distancias evolutivas corregidas



- El número de  $ti$ s generalmente > que el de  $tv$ , fenómeno que se acentúa cuanto mayor es la divergencia entre las secuencias a comparar. De ahí que en nuestro ejemplo las diferencias entre los escenarios I y II sólo se hicieron notar en el caso en el que la divergencia entre las secuencias era mayor (escenario II)

### Estima de la confianza que podemos tener en distintas partes de una filogenia: el método de bootstrap

"Filogenias bien soportadas

vs. pobremente apoyadas

por los datos"

