

Introducción a la Filoinformática:
Pan-genómica y filogenómica

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/TIB-filoinfo>

Tema 5: 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

Inferencia filogenética molecular –
clasificación de métodos

Podemos clasificar a los métodos de reconstrucción filogenética en base a:

1. el tipo de datos que emplean (**caracteres discretos vs. distancias**)
2. uso de un **método algorítmico** o un **criterio de optimización** para encontrar la topología

Tipo de datos

| | distancias | caracteres discretos |
|--------------------------|----------------------------|---|
| Método de reconstrucción | UPGMA Neighbour joining | |
| criterio de optimización | Evolución mínima | Máxima parsimonia Máxima verosimilitud |

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

Protocolo básico para un análisis filogenético de
secuencias moleculares

Colección de secuencias homólogas

• **BLAST y FASTA**

Alineamiento múltiple de secuencias

• **Clustal, T-Coffee, muscle ...**

modelos de sustitución y estima de distancias evolutivas a partir de alineamientos múltiples

Tema 5:

• **Selección de modelos...**

Estima filogenética

• **NJ, ME, MP, ML, Bayes ...**

Pruebas de confiabilidad de la topología inferida

• **proporciones de bootstrap probabilidad posterior ...**

Interpretación evolutiva y aplicación de las filogenias

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

sequences

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|------------|---|---|---|---|---|---|---|
| Drosophila | t | t | a | t | t | a | a |
| fugu | a | a | t | t | t | a | a |
| mouse | a | a | a | a | a | t | a |
| human | a | a | a | a | a | a | t |

distances

| | | | |
|-------|---|---|---|
| fugu | 3 | | |
| mouse | 5 | 4 | |
| human | 5 | 4 | 2 |

Drosophila

fugu

mouse

parsimony

distance

- Un set de 4 secs. y la matriz de distancias correspondiente
- Un árbol de parsimonia y uno de distancias para este set de datos produce topologías y longitudes de ramas idénticas
- La diferencia radica en que el árbol de parsimonia identifica qué sitio del alineamiento contribuye cada paso mutacional en la longitud de cada rama

Inferencia filogenética molecular –
métodos basados en matrices de distancias

Unweighted pair group method with arithmetic means (**UPGMA**)

OTU A B C

B d_{AB}

C d_{AC} d_{BC}

D d_{AD} d_{BD} d_{CD}

$\frac{d_{AB}}{2}$

OTU (AB) C

C $d_{(AB)C}$

D $d_{(AB)D}$ d_{CD}

$d_{(AB)C} = (d_{AC} + d_{BC})/2$, y $d_{(AB)D} = (d_{AD} + d_{BD})/2$

$I_{(AB)C} = d_{(AB)C}/2$

$d_{(AB)CD}$

UPGMA, por construir un **árbol ultramétrico**, resulta en una **topología enraizada**.
Además se obtienen las longitudes de rama simultáneamente con la topología

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

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

Matriz de distancias:

| | | | | | |
|------------------------------|--|-----|-----|-----|----|
| [| | 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 | |

4. $OTU \begin{matrix} (ABC) & D \\ D & d_{AB|C|D} \end{matrix} = \begin{matrix} d_{AB|C|D} = (d_{AD} + d_{BD} + d_{CD}) / 3 \\ d_{AB|C|D} = (9 + 9 + 6) / 3 = 8 \end{matrix}$

5.

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 construcción 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

$L = \sum_{i=1}^k l_i$

reconstrucciones para el sitio 2

tree 1:

tree 2:

tree 3:

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] |
| [<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 | |

1. $OTU \begin{matrix} A & B & C \\ B & d_{AB} & d_{BC} \\ C & d_{AC} & d_{BC} \\ D & d_{AD} & d_{BD} & d_{CD} \end{matrix} \rightarrow \begin{matrix} A \\ B \\ C \\ D \end{matrix} \begin{matrix} d_{AB} \\ d_{BC} \\ d_{AC} \\ d_{AD} & d_{BD} & d_{CD} \end{matrix}$

2. $OTU \begin{matrix} (AB) & C \\ C & d_{AB|C} \\ D & d_{AB|D} & d_{CD} \end{matrix} = \begin{matrix} d_{AB|C} = (d_{AC} + d_{BC})/2, \text{ y } d_{AB|D} = (d_{AD} + d_{BD})/2 \\ d_{AB|C} = (5 + 5)/2, \text{ y } d_{AB|D} = (9 + 9)/2 \end{matrix}$

3. $OTU \begin{matrix} (AB) & C \\ C & 5 \\ D & 9 \end{matrix} \begin{matrix} 6 \\ 6 \end{matrix} \rightarrow \begin{matrix} 2.00 \\ 2.50 \\ d(AB)C/2 \end{matrix} \begin{matrix} 0.50 \\ 0.50 \end{matrix} \begin{matrix} Rhizobium \\ Agrobacterium \\ Sinorhizobium \end{matrix}$

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] |
| [<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 | |

1. $OTU \begin{matrix} A & B & C \\ B & d_{AB} & d_{BC} \\ C & d_{AC} & d_{BC} \\ D & d_{AD} & d_{BD} & d_{CD} \end{matrix} \rightarrow \begin{matrix} A \\ B \\ C \\ D \end{matrix} \begin{matrix} d_{AB} \\ d_{BC} \\ d_{AC} \\ d_{AD} & d_{BD} & d_{CD} \end{matrix}$

2. $OTU \begin{matrix} (AB) & C \\ C & d_{AB|C} \\ D & d_{AB|D} & d_{CD} \end{matrix} = \begin{matrix} d_{AB|C} = (d_{AC} + d_{BC})/2, \text{ y } d_{AB|D} = (d_{AD} + d_{BD})/2 \\ d_{AB|C} = (5 + 5)/2, \text{ y } d_{AB|D} = (9 + 9)/2 \end{matrix}$

3. $OTU \begin{matrix} (AB) & C \\ C & 5 \\ D & 9 \end{matrix} \begin{matrix} 6 \\ 6 \end{matrix} \rightarrow \begin{matrix} 2.00 \\ 2.50 \\ d(AB)C/2 \end{matrix} \begin{matrix} 0.50 \\ 0.50 \end{matrix} \begin{matrix} Rhizobium \\ Agrobacterium \\ Sinorhizobium \end{matrix}$

• ¿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 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

| Taxa | árboles no enraiz* | árb. enraiz. |
|------|--------------------|--------------|
| 4 | 3 | 15 |
| 8 | 10,395 | 135,135 |
| 10 | 2,027,025 | 34,459,425 |
| 22 | 3×10^{23} | ... |
| 50 | 3×10^{24} | ... |

*por ej. para sólo 15 OTUs tenemos 213,458,046,676,875 topologías - ¡ si pudiésemos evaluar 1×10^6 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

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 y By)
- 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) |
| Máxima verosimilitud | Ramificación y límite ($n \leq 25$) (branch-and-bound) |
| Bayesiana | Decomposición en estrella (star decomposition) |
| | Adición secuencial (stepwise addition) |
| | (Inter-)cambio de rama (branch swapping) |

Métodos exactos:
garantizan encontrar la topología óptima

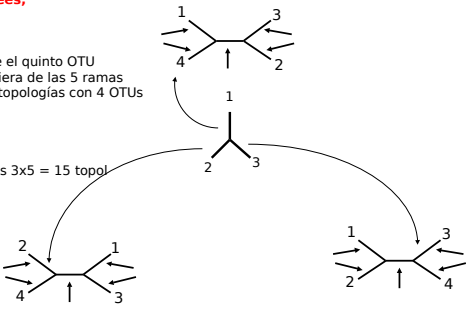
Métodos heurísticos:
no garantizan encontrar la topología óptima

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

PAUP* command:
alltrees;

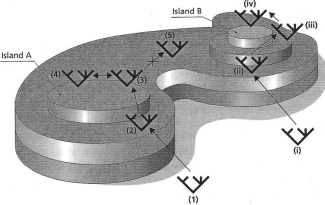
se añade el quinto OTU a cualquiera de las 5 ramas de las 3 topologías con 4 OTUs

obtenemos $3 \times 5 = 15$ topol



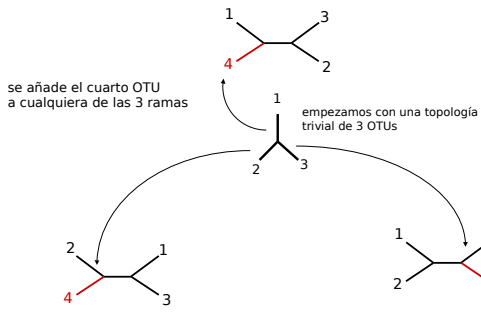
Métodos heurísticos de búsqueda de árboles - islas de árboles

- En la mayor parte de los análisis emplearán métodos heurísticos;
- éstos comienzan con un árbol (aleatorio, NJ o de adición secuencial) para realizar intercambios de ramas (**branch swapping**) sobre esta topología inicial con el propósito de encontrar topologías de mejor puntuación (según la func. de objetividad) que la de partida
- estos métodos heurísticos no garantizan encontrar la topología óptima pero trabajan muy bien cuando se comparan con sets de datos de ≤ 25 secs. analizados mediante B&B

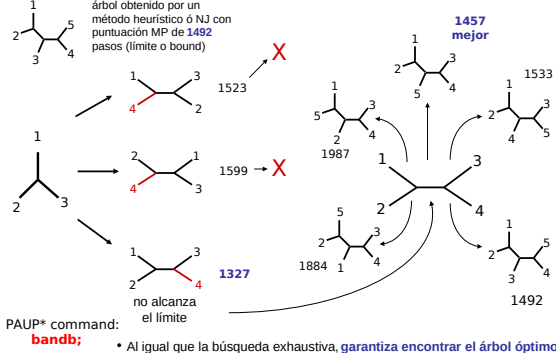


- El espacio de árboles puede visualizarse como un paisaje con colinas de diversas alturas; cada pico representa un máximo local de score o puntuación (**isla de árboles igualmente parsim.**)
- Es recomendable hacer múltiples búsquedas heuríst. comenzando cada una desde una topología distinta para minimizar el riesgo de obtener un árbol ubicado en una isla topológica subóptima

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

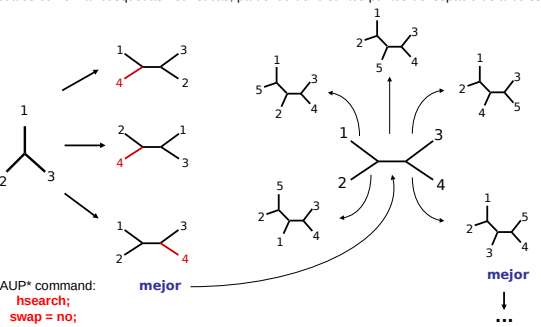


Métodos exactos de búsqueda de árboles - "branch and bound" ($n \leq 25$)



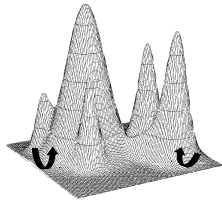
Métodos heurísticos de búsqueda de árboles - adición secuencial (aleatorizada)

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 - adición secuencial (aleatorizada)

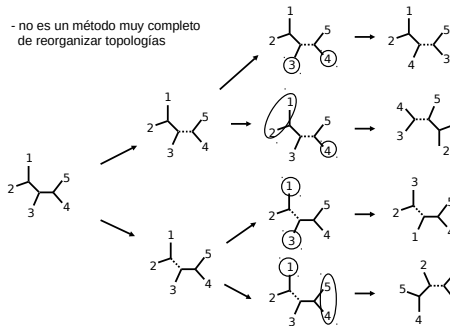
- El orden en el que se añaden los OTUs puede cambiar los resultados
- Por ello suele repetirse varias veces, añadiendo OTUs en cada ciclo de manera aleatorizada
- Sirven por lo tanto como **árboles semilla** para iniciar distintas búsquedas heurísticas partiendo de topologías potencialmente diferentes para eficientizar la exploración del espacio de topologías (pero **no adecuados como hipótesis filogenética en sí mismos**)



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



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:

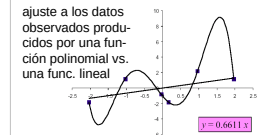
- un árbol filogenético
- una descripción de las probabilidades con las que se dan las sustituciones de aa o nts a lo largo de las ramas del árbol

- ¿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

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

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



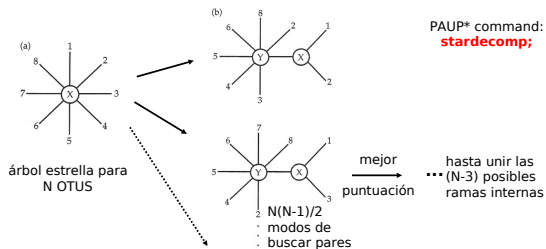
- añadir parámetros** a un modelo generalmente mejora su ajuste a los datos observados

- 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

Métodos heurísticos de búsqueda de árboles - decomposición de estrella



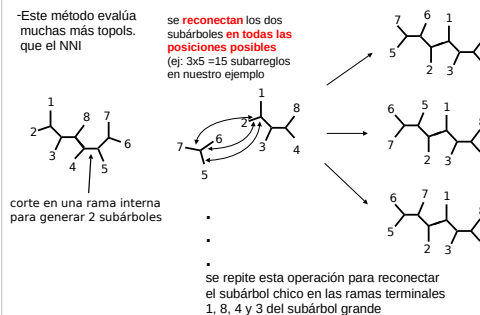
- NJ usa este método junto al criterio de evolución mínima
- una vez que 2 OTUs han sido unidos ya no pueden ser desacoplados más adelante; en esto difiere del algoritmo de adición secuencial
- sensible al orden en que se van uniendo los OTUs; problema incrementa con el no. de OTUs
- no debe ser por 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)

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

-Este método evalúa muchas más topol. que el NNI

se **reconectan** los dos subárboles en todas las posiciones posibles (ej. 3x5=15 subarreglos en nuestro ejemplo)

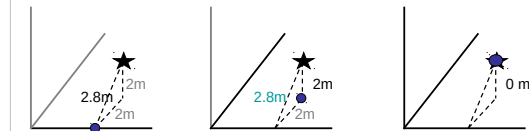


corte en una rama interna para generar 2 subárboles

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

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, dos 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
- Existen dos aproximaciones para construir modelos de evolución de secuencias.
 - 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
 - construcción de modelos paramétricos basado en el modelaje de propiedades químicas o genéticas del 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 ajustándose a cada matriz de datos particular

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 ≠ frecuencia: F81, HKY85, TrN93, GTR ...

II. Tasas de sustitución transicionales/transversales

- Existen 4 tipos de sustituciones ti y 8 tv; cuando ti/tv ≠ 0.5 existe un sesgo en sustituciones ti (o tv) en el set de datos. ti generalmente >> 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 ≠ tv) |
| 3 | TrN ó K3P (2 ti, 1 tv) |
| 6 | GTR (cada sust. su tasa) |

Modelos de evolución de secuencias -introducción

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

Corolario:

- El grado de confianza que tengamos en una filogenia particular realmente depende de la que tengamos en el modelo subyacente
- 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 -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 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, $\alpha = \beta$

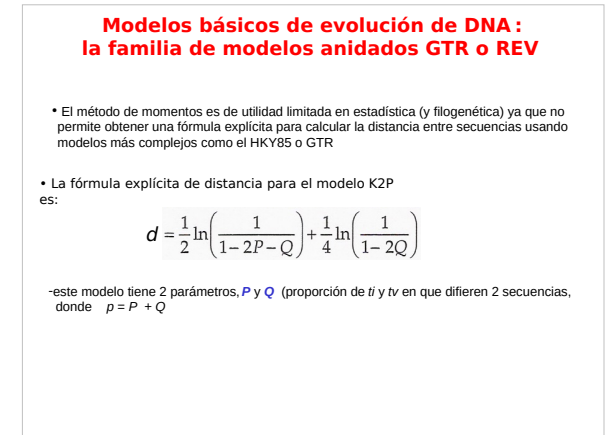
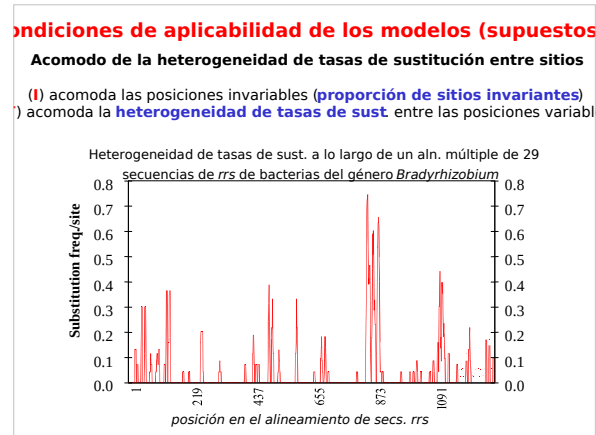
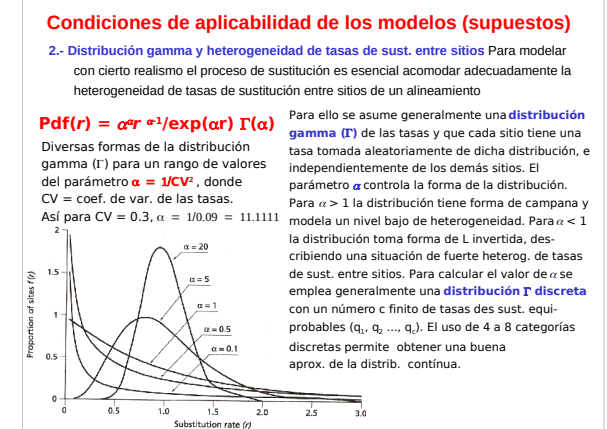
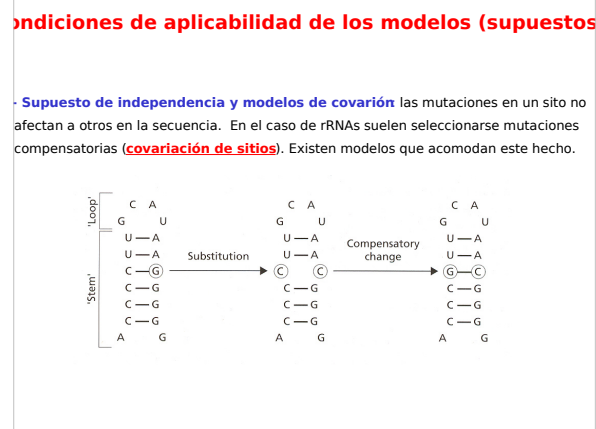
acomodan ≠ frec. bases

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
o
Tamura-Nei 1993 (TN93), 3 tasas
o
General time reversible (GTR), 6 tasas

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



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 $Q = 4/200 = 0.02$

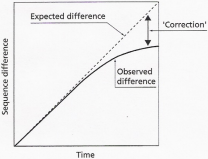
$p = 24/200 = 0.12$
 $d_{JC69} \approx 0.13$ (sust./sitio)

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

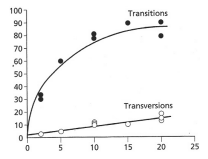
no. de sust. esperadas = 0.13 X 200 \approx 26

no. de sust. esperadas = 0.13 X 200 \approx 26

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

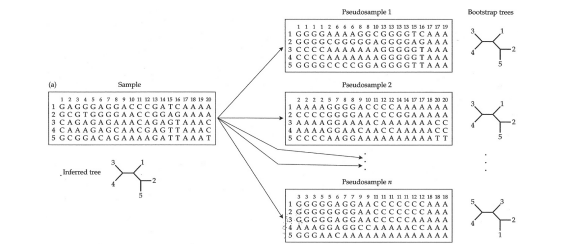


- El objetivo de los modelos de sustitución es el de **compensar para los eventos homoplásicos de múltiples sustituciones**, y así obtener estimas de distancias evolutivas corregidas



- El número de ti es 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 del error de muestreo mediante el método de bootstrap



BOOTSTRAPPING

- generación de n (100-1000) **pseudoréplicas** (muestreo aleatorio con reemplazo)
- estíma de la filogenia para cada pseudoréplica
- cálculo de un **árbol consenso**
- mapeo de las proporciones de bootstrap sobre la topología inferida de los datos originales

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 $Q = 16/200 = 0.08$

$p = 66/200 = 0.33$
 $d_{JC69} \approx 0.43$ (sust./sitio)

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

no. de sust. esperadas = 0.43 X 200 \approx 86

no. de sust. esperadas = 0.48 X 200 \approx 96

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"

