

Introducción a la Inferencia Filogenética y Evolución Molecular

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

Pablo Vinuesa (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 8: Criterios de optimización II: máxima verosimilitud

1. El criterio de máxima verosimilitud y su aplicación en filogenética
2. Estima de máxima verosimilitud de parámetros del modelo de sustitución
3. Selección de modelos de sustitución usando pruebas de razones de verosimilitud (LRTs)
4. Modelos paramétricos de evolución de secuencias de DNA - la familia GTR
5. Prácticas con PAUP*, ModelTest y proml

Métodos de reconstrucción filogenética - Máxima Verosimilitud

Máxima verosimilitud: dadas dos topologías, la que hace los datos observados más probables ("menos sorprendentes") es la preferida

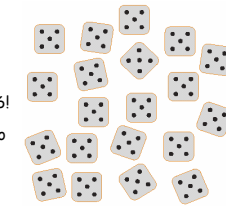
El **método de máxima verosimilitud (ML)** considera cada sitio variable del alineamiento (incluidos singletons). Bajo el criterio de ML se busca la topología que hace más verosímil el patrón de sustituciones de un alineamiento dado un modelo evolutivo explícito!

Así, para un set de datos D y una hipótesis evolutiva (topología) H , la verosimilitud de dichos datos viene dado por la expresión:

$L_D = \Pr(D|H)$ que es la probabilidad de obtener D dada H (una **probabilidad condicional**)!

Por tanto la topología que hace nuestros datos el resultado evolutivo más probable corresponde a la estima de máxima verosimilitud de la filogenia (likelihood score ó valor de verosimilitud).

- La probabilidad está relacionada con la "sorpresividad" de los datos
- Estaríamos sorprendidos de obtener este resultado, dada su bajísima probabilidad $(1/6)^{20}$ ó 1 en 3,656,158, 440,062,976!
- Pero la probabilidad depende del modelo probabilístico asumido
- En filogenética, las distintas topologías representan a los distintos modelos, y se selecciona aquel modelo que nos hace sorprendernos menos de los datos que hemos coleccionado



Máxima verosimilitud y estima de parámetros de modelos de sustitución

- La inferencia filogenética bajo el criterio de máxima verosimilitud se basa en el uso de una cantidad llamada **log-likelihood** para evaluar topologías alternativas con el fin de encontrar aquella que **maximiza este valor**.
- El **log-likelihood** es el ln de la verosimilitud, que es igual a la probabilidad de los datos observados dadas una topología particular (τ), set de longitudes de rama (ν) y modelo de sustitución (ϕ).
- Nótese que **la verosimilitud no representa la probabilidad de que un árbol sea correcto**; ésta viene determinada por la **probabilidad posterior** de la estadística bayesiana.
- Hablar de la "verosimilitud de un conjunto de datos" no es correcto ya que la verosimilitud es un función de los parámetros de un modelo estadístico, y no de los datos (D). **Los datos son constantes siendo el modelo lo que es variable al calcular verosimilitudes**. Se puede por lo tanto hablar de verosimilitudes como funciones de modelos o hipótesis (H). La verosimilitud de una hipótesis dado un set de datos es igual a la probabilidad condicional de los datos dada una hipótesis.

Formalmente: $L(H|D) = \Pr(D|H) = \Pr(D|\tau\nu\phi)$

Máxima verosimilitud y estima de parámetros de modelos de sustitución

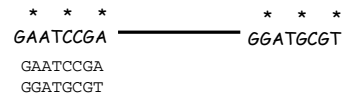
$$L(H|D) = \Pr(D|H) = \Pr(D|\tau\nu\phi)$$

- Lo mejor es pensar en los **árboles como modelos**. La verosimilitud de una topología particular (τ) será la probabilidad de los datos dada esa topología. Cada topología tiene como parámetros las longitudes de rama (ν), y la verosimilitud de un modelo (ϕ) cambia según varíen los valores de los parámetros de longitud de rama
- Por lo tanto se puede concebir la filogenética bajo el criterio de máxima verosimilitud como un **problema de selección de modelos**. Se trata de encontrar las estimas de los valores de cada parámetro del modelo y luego comparar las verosimilitudes de los distintos modelos, escogiendo el mejor (topología) en base a su verosimilitud
- La topología que hace de nuestros datos el resultado evolutivo más probable (dado un modelo de sust.) es la estima de máxima verosimilitud de nuestra filogenia. Por tanto, al contrario que bajo los criterios de optimización de MP, LS o ME, **bajo ML se trata de seleccionar modelos y parámetros que maximicen la función de optimización**.

Tema 8: máxima verosimilitud, estima de parámetros y selección de modelos

Máxima verosimilitud y estima de parámetros de modelos de sustitución

- Esquema del procedimiento del cálculo del valor de verosimilitud de un árbol con 4 OTUs

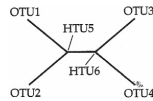


En un "árbol" con sólo 2 OTUs no tenemos ningún nodo interior o ancestral. El cómputo lo realizamos directamente sobre los datos observados

$$L = L_1 L_2 \dots L_8 = [1/16 (1 + 3e^{-4\alpha t})]^5 [1/16 (1 - e^{-4\alpha t})]^3$$

- La complicación adicional que encontramos para el cálculo de verosimilitudes de árboles con > 3 OTUs radica esencialmente en que tenemos ahora nodos interiores para los que carecemos de observaciones. Se trata de unidades taxonómicas hipotéticas HTUs. En este caso, para calcular la verosimilitud del árbol **tenemos que considerar cada posible estado de carácter para cada nodo interior y para cada topología !!!**.

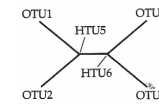
	1	2	3	4	5	6	7	8	9	...	n
OTU1	A	A	G	A	C	T	T	C	A	...	N
OTU2	A	G	C	C	C	T	T	C	T	...	N
OTU3	A	G	A	T	A	T	C	C	A	...	N
OTU4	A	G	A	G	G	T	C	C	T	...	N



Máxima verosimilitud y estima de parámetros de modelos de sustitución

- Esquema del procedimiento del cálculo del valor de verosimilitud de un árbol con 4 OTUs

	1	2	3	4	5	6	7	8	9	...	n
OTU1	A	A	G	A	C	T	T	C	A	...	N
OTU2	A	G	C	C	C	T	T	C	T	...	N
OTU3	A	G	A	T	A	T	C	C	A	...	N
OTU4	A	G	A	G	G	T	C	C	T	...	N



Para 4 OTUs existen 3 topologías posibles. Por ello hemos de repetir este cálculo para cada una de ellas con el fin de encontrar la topol. más verosímil

$L_{(5)} = \text{Prob} \left(\begin{matrix} C \\ \swarrow \quad \searrow \\ A-A \quad G \end{matrix} \right) + \text{Prob} \left(\begin{matrix} C \\ \swarrow \quad \searrow \\ A-C \quad G \end{matrix} \right) + \text{Prob} \left(\begin{matrix} C \\ \swarrow \quad \searrow \\ A-T \quad G \end{matrix} \right) + \text{Prob} \left(\begin{matrix} C \\ \swarrow \quad \searrow \\ A-G \quad G \end{matrix} \right)$
 $+ \text{Prob} \left(\begin{matrix} C \\ \swarrow \quad \searrow \\ C-A \quad G \end{matrix} \right) + \text{Prob} \left(\begin{matrix} C \\ \swarrow \quad \searrow \\ C-C \quad G \end{matrix} \right) + \text{Prob} \left(\begin{matrix} C \\ \swarrow \quad \searrow \\ C-T \quad G \end{matrix} \right) + \text{Prob} \left(\begin{matrix} C \\ \swarrow \quad \searrow \\ C-G \quad G \end{matrix} \right)$
 $+ \text{Prob} \left(\begin{matrix} C \\ \swarrow \quad \searrow \\ T-A \quad G \end{matrix} \right) + \text{Prob} \left(\begin{matrix} C \\ \swarrow \quad \searrow \\ T-C \quad G \end{matrix} \right) + \text{Prob} \left(\begin{matrix} C \\ \swarrow \quad \searrow \\ T-T \quad G \end{matrix} \right) + \text{Prob} \left(\begin{matrix} C \\ \swarrow \quad \searrow \\ T-G \quad G \end{matrix} \right)$
 $+ \text{Prob} \left(\begin{matrix} C \\ \swarrow \quad \searrow \\ G-A \quad G \end{matrix} \right) + \text{Prob} \left(\begin{matrix} C \\ \swarrow \quad \searrow \\ G-C \quad G \end{matrix} \right) + \text{Prob} \left(\begin{matrix} C \\ \swarrow \quad \searrow \\ G-T \quad G \end{matrix} \right) + \text{Prob} \left(\begin{matrix} C \\ \swarrow \quad \searrow \\ G-G \quad G \end{matrix} \right)$

La verosimilitud para cada sitio representa la suma sobre todas las posibles asignaciones de estados de carácter en todas las ramas interiores de un árbol. La verosimilitud total es el producto de las veros. por sitio.

$$L = L_{(1)} \times L_{(2)} \times L_{(3)} \times \dots \times L_{(n)} = \prod_{i=1}^n L_{(i)}$$

$$\ln L = \ln L_{(1)} + \ln L_{(2)} + \ln L_{(3)} + \dots + \ln L_{(n)} = \sum_{i=1}^n \ln L_{(i)}$$

Máxima verosimilitud y estima de parámetros de modelos de sustitución

- La inferencia filogenética bajo el criterio de máxima verosimilitud implica **MUCHISIMO TRABAJO COMPUTACIONAL** (= > mucho tiempo de trabajo de procesador)
- Las verosimilitudes globales han de ser maximizadas para cada topol. Para ello necesitamos:
 - encontrar EMV para cada long. de rama y cada parámetro del modelo de sust.
 - ello implica calcular la verosimilitud global muchas, pero que muchas veces
- En la práctica los árboles de ML se estiman en múltiples ciclos, en los que se van optimizando secuencialmente los diversos parámetros del modelo de sustitución y longitudes de rama. La estima conjunta de todos los parámetros se hace computacionalmente prohibitiva
- Por lo general se comienzan estos ciclos partiendo de una topología obtenida por un método rápido, tal como NJ o MP. Sobre esta topología se ajustan los valores de los parámetros del modelo. A continuación se emplea algún método de reajuste de topología (branch swapping) y se ajustan las longitudes de rama, cerrando un ciclo. En múltiples ciclos consecutivos se va optimizando la topología y long. de rama, hasta que convergen en la estima de máxima verosimilitud global

Máxima verosimilitud y estima de parámetros de modelos de sustitución

1. La relevancia e impacto de los modelos de evolución en filogenética y evol. molecular

- Los modelos no son sólo importantes por sus consecuencias en la estima filogenética, sino que además lo son porque la **caracterización del proceso evolutivo** a nivel molecular es **objeto de estudio en sí mismo**, en el ámbito de la evolución molecular.
- Los modelos de evolución **son herramientas poderosas** siempre y cuando, a pesar de las simplificaciones que hacen, **describan adecuadamente las características más salientes de los datos** y permitan **hacer predicciones precisas** sobre el problema bajo estudio.
- El rendimiento de un método se maximiza cuando se satisfacen los supuestos subyacentes.**
- Es conveniente por tanto seleccionar el modelo más adecuado para cada set particular de datos y cuantificar el ajuste de los datos al modelo seleccionado.

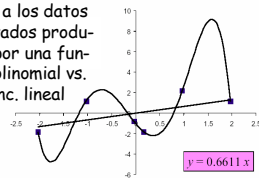
Máxima verosimilitud y estima de parámetros de modelos de sustitución

2. Selección de modelos de sustitución de secuencias de DNA

- En términos generales modelos complejos se ajustan a los datos mejor que los simples. Idealmente **se ha de seleccionar un modelo lo suficientemente complejo** (rico en parámetros) como **para describir adecuadamente las características más notables del patrón de sust.** del set de datos, pero no sobreparametrizado para **evitar colinearidad de parámetros (redundancia)**, tiempos excesivamente largos de cómputo y estimas poco precisas de los parámetros por excesiva varianza.

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

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

Máxima verosimilitud y estima de parámetros de modelos de sustitución

2. Selección de modelos de sustitución de secuencias de DNA y Proteína

- Se deben de usar **pruebas estadísticas para seleccionar el modelo que mejor se ajusta a los datos de entre los disponibles**. Este ajuste de los modelos a los datos puede ser evaluado usando pruebas de razones de verosimilitud (likelihood ratio tests, **LRTs**) o usando criterios de información de Akaike o bayesianos (**AIC** y **BIC**, respectivamente). Se puede usar una prueba de LRT para evaluar la capacidad que tiene un modelo particular en ajustar los datos.

- Idealmente debemos de seleccionar el mejor modelo para cada gen o región genómica que queramos analizar. No conviene hacerlo para una supermatriz de alineamientos concatenados. El uso de **modelos particionados** en los que se ajusta el modelo para cada posición de los codones, por cada gen a analizar, resultan generalmente en ajustes globales significativamente mejores que **modelos promediados** para cada gen.

Máxima verosimilitud y estima de parámetros de modelos de sustitución

3. Prueba de razón de verosimilitudes (LRT)

• Una manera natural y muy usada de comparar el ajuste relativo de dos modelos alternativos a una matriz de datos es contrastar las verosimilitudes resultantes mediante la prueba de razones de verosimilitud (RV) ó likelihood ratio test (LRT):

$$\Delta = 2(\log_e L_1 - \log_e L_0)$$

donde L_1 es el valor de ML global para la hipótesis alternativa (modelo más rico en parámetros) y L_0 es el valor de ML global para la hipótesis nula (el modelo más simple).

$\Delta \geq 0$ siempre, ya que los parámetros adicionales van a dar una mejor explicación de la variación estocástica en los datos que el modelo más sencillo.

• Cuando los modelos a comparar están anidados (L_0 es un caso especial de L_1) el estadístico Δ sigue aproximadamente una **distribución χ^2 con q grados de libertad**, donde q = diferencia entre el no. de parámetros libres entre L_1 y L_0 .

Máxima verosimilitud y estima de parámetros de modelos de sustitución

3. Prueba de razón de verosimilitudes (LRT)

- El LRT es por tanto una prueba estadística para cuantificar la bondad relativa de ajuste entre dos modelos anidados. Veamos un ejemplo. Vamos seleccionar entre los modelos JC69, F81, HKY85 y TrN93 para el set de datos de mtDNA-primates.nex, considerando sólo las regiones codificadoras y eliminando Lemur_catta, Tarsius_syrichtha y Saimiri_scireus y usando un árbol NJ sobre el cual estimar parámetros

Modelo	$-\ln L$	• ¿ Qué podemos concluir de estos valores de
JC69	3585.54820	$-\ln L$ en cuanto a la importancia relativa de
F81	3508.04085	los parámetros considerados por estos
HKY85	3233.34395	modelos en cuanto al nivel de ajuste a los datos
TrN93	3232.29439	que alcanzan ?

Máxima verosimilitud y estima de parámetros de modelos de sustitución

3. Prueba de razón de verosimilitudes (LRT)

Modelo	-lnL	H_0 a rechazar (o hipótesis anidadas a evaluar)
JC69	3585.54820	1. igual frec. de bases
F81	3508.04085	
HKY85	3233.34395	
TrN93	3232.29439	2. $T_i = T_v$
		3. tasas de T_i iguales
		...

modelos	diff. GL = q	χ^2	P
JC-F81	3 - 0 = 3	155	0
JC-HKY85	4 - 0 = 4	704.4	0
JC-TrN	5 - 0 = 5	706.4	0
F81-HKY85	4 - 3 = 1	549.4	0
F81-TrN	5 - 3 = 2	551.4	0
KHY-TrN	5 - 4 = 1	2.1	0.15

Por lo tanto el modelo seleccionado es el HKY

<http://www.fourmilab.ch/rpkp/experiments/analysis/chiCalc.html>

Máxima verosimilitud y estima de parámetros de modelos de sustitución

3. Prueba de razón de verosimilitudes (LRT)

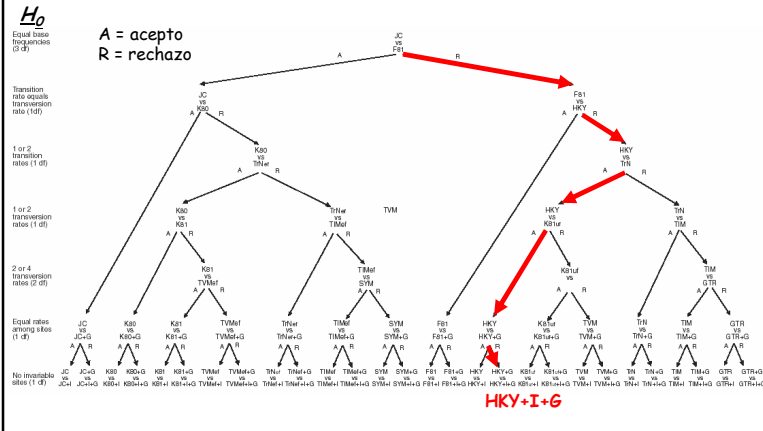
Modelo	-lnL	H_0 a rechazar (o hipótesis anidadas a evaluar)
HKY85	3233.34395	1. tasa homogénea de sust entre sitios
HKY85 +G	3145.29031	
HKY85 +I+G	3142.36439	2. no existe proporción de sitios invariantes

modelos	diff. GL = q	χ^2	P
HKY85-vs. +G	1	176	0
HKY85+G vs. I+G	1	5.85	0.015

Por lo tanto el modelo seleccionado es el HKY+G si tomamos 0.01 como punto de corte, o HKY+I+G si usamos $\alpha = 0.05$.

Máxima verosimilitud y estima de parámetros de modelos de sustitución

3. Esquema jerárquico de efectuar LRTs partiendo desde el modelo más sencillo (JC69)



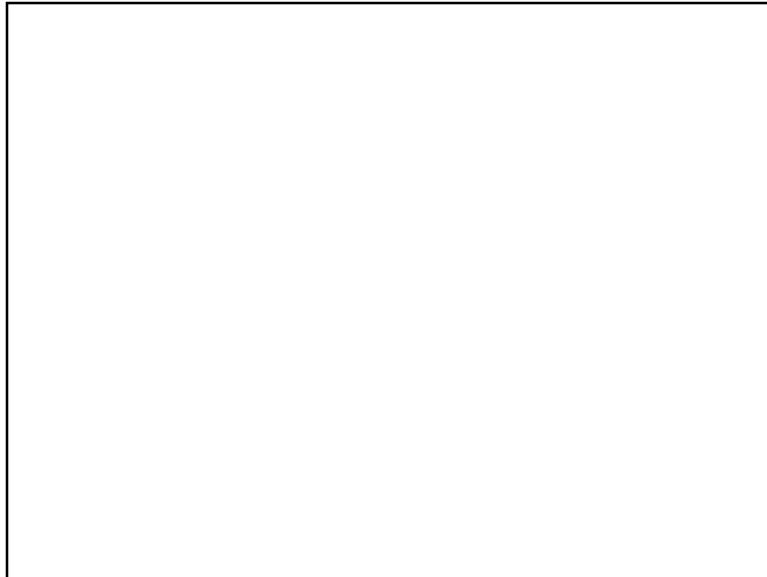
Máxima verosimilitud y estima de parámetros de modelos de sustitución

3. Resumen de algunos modelos y sus parámetros libres

- Dado que en los modelos de sust. de DNA la tasa promedio de sustitución se considera = 1 y los parámetros de tasa relativa se escalan de tal manera que la tasa promedio de sust. en equilibrio = 1, el modelo más sencillo (JC69) no tiene ningún parámetro libre, dado que el único parámetro (α) a estimar valdrá $\frac{1}{4}$ en este contexto.

Modelo	características	no. de parámetros libres
JC	nst=1 basefreq= equal	0
F81	nst=1 basefreq=uneq	3 para las <i>frec. de bases</i>
K2P	nst=2 basefreq=eq	1 para el <i>tratio</i> (ti/tv)
HKY85	nst=2 basefreq=uneq	4 (1 para <i>tratio</i> y 3 para <i>frec. de bases</i>)
TrN93	nst=3 basefreq=uneq	5 (2 tasas de ti y 3 para <i>frec de bases</i>)
GTR	nst=6 basefreq=uneq	8 (5 para tasas de subst y 3 para <i>frec. de bases</i>)

proporción de sitios invariantes (I) 1 parámetro libre adicional para pinv
 distribución gamma (G) 1 parámetro libre adicional para G
 ambos combinados (I+G) 2 parámetros libres adicionales



Máxima verosimilitud y estima de parámetros de modelos de sustitución

3. Selección de modelos usando criterios de información

- LRT compara pares de modelos anidados. Los criterios de información como el **Akaike information criterion (AIC)** y **Bayesian information criterion (BIC)** comparan simultáneamente todos los modelos en competición y permiten seleccionar modelos aunque no sean anidados.
- Se trata nuevamente de incorporar tanta complejidad (parámetros) al modelo como requieran los datos. La verosimilitud para cada modelo es penalizada en función del número de parámetros: **a mayor cantidad de parámetros mayor penalización.**

Máxima verosimilitud y estima de parámetros de modelos de sustitución

3. Selección de modelos usando criterios de información

- **AIC.** Es un estimador no sesgado del parámetro de contenido de información de Kullback-Leibler, el cual es una **medida de la información perdida al usar un modelo para aproximar la realidad.** Por tanto, **a menor valor de AIC mejor ajuste** del modelo a los datos. Al penalizar por cada parámetro adicional, **considera tanto la bondad de ajuste como la varianza asociada a la estima de los parámetros.**

$$AIC_i = -2 \ln L_i + 2 N_i$$

N_i = no. de parámetros libres en el modelo i
 L_i = verosimilitud bajo el modelo i

Máxima verosimilitud y estima de parámetros de modelos de sustitución

3. Selección de modelos usando criterios de información: **AIC**

- Se pueden usar los **estadísticos de diferencias en AIC (Δ_i)** y **ponderaciones de Akaike para cuantificar el nivel de incertidumbre en la selección del modelo.** Las Δ_i son AICs re-escalados con respecto al modelo con el AIC más bajo (minAIC), de modo que $\Delta_i = AIC_i - \text{minAIC}$. Las Δ_i son fáciles de interpretar y permiten ordenar los modelos candidatos. Así, **modelos con Δ_i en un rango de 1-2 con respecto al modelo ganador tienen un soporte sustancial** y deben de ser considerados como modelos alternativos. Modelos con Δ_i **en un rango de 3-7 con respecto al modelo ganador tienen un soporte significativamente inferior**, y modelos con $\Delta_i > 10$ **carecen de soporte.**

Máxima verosimilitud y estima de parámetros de modelos de sustitución

3. Selección de modelos usando criterios de información: AIC

- Las ponderaciones o pesos de Akaike (w_i) son los AIC relativos normalizados para cada modelo en competición y pueden ser interpretados como la probabilidad de que un modelo es la mejor abstracción de la realidad dados los datos. Para R modelos candidatos a evaluar:

$$w_i = \frac{\exp\left(-\frac{1}{2}\Delta_i\right)}{\sum_{r=1}^R \exp\left(-\frac{1}{2}\Delta_r\right)}$$

- Una aplicación muy útil de los w_i es que la inferencia se puede promediar a partir de los modelos que muestran valores de no w_i triviales. Así, una estima del valor del parámetro α de la distribución gamma promediada a partir de varios modelos se calcularía así:

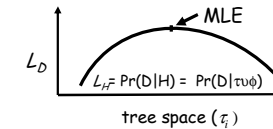
$$\hat{\alpha} = \sum_{i=1}^R w_i \hat{\alpha}_i$$

También podemos **reconstruir filogenias bajo los distintos modelos con peso significativo y combinar los árboles resultantes acorde a sus pesos de Akaike**. Esta estrategia es particularmente útil en un contexto bayesiano.

Criterios de optimización: la alternativa Bayesiana

- Aproximaciones tradicionales (matrices de distancia, ME, ML, MP)

- la búsqueda tiene por objetivo encontrar la topología óptima (**estima puntual**)
- no pueden establecer el soporte relativo de las biparticiones a partir de una única búsqueda



- Aproximación Bayesiana

- no busca una sola topología óptima sino una **población de árboles muestreados en función de su probabilidad posterior** (algoritmos MCMC)
- la muestra de árboles obtenidos en una sola sesión de "búsqueda" es usada para valorar el soporte de cada split en términos de probabilidad posterior

