
TEMA 3.5. Obtención de árboles mediante el método de inferencia bayesiana

Contacto: Isabel Draper (Isabel.draper@uam.es)

PROGRAMAS NECESARIOS

- (1) Para la obtención de árboles mediante Inferencia Bayesiana vamos a utilizar el programa MrBayes (Huelsenbeck & Ronquist, 2001), que puede ser descargado gratuitamente desde:

<http://mrbayes.sourceforge.net/>

En esa misma página se explica cómo instalar el programa y se puede acceder al manual de instrucciones. Para realizar la siguiente práctica ten en cuenta que el formato de archivo necesario para este programa es el de tipo nexus. Además, no olvides que este archivo debe estar en la misma carpeta en la que está el ejecutable de MrBayes.

Programas alternativos:

BEAST (Drummond & Rambaut 2007) es otro programa gratuito que utiliza cadenas de Markov para la aplicación de inferencia bayesiana en el análisis de secuencias moleculares. Este programa está orientado a filogenias con tiempos de divergencia.

En el siguiente enlace puedes encontrar el link a las descargas además de un manual para su uso y algunos otros enlaces de interés.

http://beast.bio.ed.ac.uk/Main_Page

METODOLOGÍA Y PRÁCTICA**I. Preparación del archivo nexus para MrBayes**

- Paso 1. MrBayes trabaja con formato nexus. Para preparar tu matriz de datos, comprueba que el archivo tiene la siguiente estructura inicial, teniendo especial cuidado de respetar los puntos y coma finales [OJO, no olvidar ningún “,”]:

```
# nexus
begin data;
dimensions ntax=N nchar=M
format datatype=X
interleave=Y
gap= — missing=?;
matrix
```

Abre la matriz "Glaucorea_seda ITS_LF" y escribe estos comandos al inicio de tu matriz sustituyendo:

N por el número de muestras

M por el número de caracteres de cada muestra

X por el tipo de datos:

dna (como en nuestro caso)

standard (si son datos morfológicos)

restriction (si son gaps codificados)

mixed (si se trata de una matriz mixta en la que se combinan diferentes tipos de datos). En este caso habría que indicar a continuación y entre paréntesis qué caracteres son de cada tipo, por ejemplo:

format datatype = mixed (dna:1-25, restriction:26-30)

Y por NO (ya que las secuencias en nuestra matriz no están divididas en bloques, sino que se encuentran como una cadena seguida).

Recuerda también que las secuencias no pueden empezar por el símbolo > y que después de la matriz se debe indicar que acaba el archivo mediante (véase tema 3.2):

;end; [end of file]

Con este tipo de información en el archivo, MrBayes va a tratar por defecto todos los datos de ADN como una única partición. Sin embargo, si hemos combinado varios genes suele ser conveniente dividir los datos en particiones (una para cada región genética). De este modo, podemos analizarlas por separado (por ejemplo si cada región presenta un modelo evolutivo diferente). Para separar los datos en particiones, lo más conveniente es añadir, utilizando un editor de texto (WordPad para PC, o TextWrangler para MAC; OJO es muy importante que no lo abráis con Word), un bloque informativo más en el archivo que va a analizar MrBayes (a continuación de la matriz de datos). En concreto se tiene que añadir la siguiente información (cada comando seguido de un signo de igualdad y la información que se quiera asociar al comando, y cada línea siempre terminada en punto y coma):

begin mrbayes

[para indicar que se trata de un bloque informativo para ese programa]

charset

[para asociar un nombre con un conjunto de caracteres]

partition favored

[para definir una partición asociada a cada nombre dado]

set partition

[para indicar a MrBayes que queremos trabajar con nuestras particiones en lugar de con la partición por defecto; las particiones a las que hace referencia pueden ser las nombradas con *charset* o indicarse con los números de caracteres: por ejemplo *partition favored = 4: 1-275, 276-427, 428-639, 640-1375;*]

end

[para cerrar el bloque]

En nuestra matriz combinada de ITS y trnLF el bloque debería ser así:

```
begin mrbayes;  
charset ITS1 = 1-275;  
charset 5.8S = 276-427;  
charset ITS2 = 428-639;  
charset trnLF = 640-1375;  
partition favored = 4: ITS1, 5.8S, ITS2, trnLF;  
set partition = favored;  
end;
```

Ejercicio 3.5.1. Abre la matriz “Glaucoseda_ITS_LF” y prepárala para su análisis en MrBayes.

Ejercicio 3.5.2. De acuerdo a estas directrices, abre la matriz combinada de ITS y LF (Glaucoseda_ITS_LF.nex) con un editor de texto (WordPad para PC, o TextWrangler para MAC; OJO es muy importante que no lo abráis con Word), establece las particiones y guarda el archivo bajo el nombre

“Glaucoseda_ITS_LF_particiones.nex”.

II. Análisis de Inferencia Bayesiana con MrBayes

A continuación hacemos un breve resumen de los comandos básicos del programa. Te proponemos que realices una reconstrucción filogenética utilizando para ello A) la matriz combinada de ITS y LF con las particiones especificadas (Glaucoseda_ITS_LF_particiones.nex), y B) la misma matriz combinada pero incluyendo al final los gaps codificados según Simmons y Ochoterena (2000) (Glaucoseda_ITS_LF_gaps_particiones.nex). De esta manera podrás comprobar la influencia que tiene considerar o no la información de los gaps a la hora de reconstruir la filogenia.

Guarda el archivo que vayas a utilizar (Glaucoseda_ITS_LF_particiones.nex y Glaucoseda_ITS_LF_gaps_particiones.nex) en la carpeta raíz del programa de MrBayes. Ten en cuenta que el nombre del archivo no puede contener espacios.

Paso 1. Abre el programa MrBayes y ejecuta el archivo con el comando *execute Glaucoseda_ITS_LF_particiones.nex*



Ten en cuenta que MrBayes por defecto considera como outgroup el taxon que aparece en primer lugar en la matriz. Si quieres que se considere otro se puede indicar utilizando el comando *outgroup N* [donde N es el número de orden en la matriz del taxon que quieres que sea el *outgroup*]

Paso 2. Especifica el modelo evolutivo que mejor se ajusta a tus datos (véase tema 3.4):

El comando *showmodel* te permite ver el tipo de modelo que MrBayes aplica por defecto.

```

C:\Archivos de programa\MrBayes\mrbayes-3.1.2\mrbayes.exe
Exiting MrBayes block
Reached end of file
MrBayes > showmodel

Model settings:

Settings for partition 1 --
Datatype = DNA
Nucmodel = 4by4
Nst = 1
Covarion = No
# States = 4
State frequencies have a Dirichlet prior
<1.00,1.00,1.00,1.00>
Rates = Equal

Settings for partition 2 --
Datatype = DNA
Nucmodel = 4by4
Nst = 1
Covarion = No
# States = 4
State frequencies have a Dirichlet prior
<1.00,1.00,1.00,1.00>
Rates = Equal

Settings for partition 3 --
Datatype = DNA
Nucmodel = 4by4
Nst = 1
Covarion = No
# States = 4
State frequencies have a Dirichlet prior
<1.00,1.00,1.00,1.00>
Rates = Equal

Settings for partition 4 --
Datatype = DNA
Nucmodel = 4by4
Nst = 1
Covarion = No
# States = 4
State frequencies have a Dirichlet prior
<1.00,1.00,1.00,1.00>
Rates = Equal

Active parameters:

Parameters      Partition(s)
1 2 3 4
-----
Statefreq       1 1 1 1
Topology        2 2 2 2
BrLens          3 3 3 3

All parameters can be linked or unlinked across partitions

1 -- Parameter = Statefreq
   Prior = Dirichlet
   Partitions = 1, 2, 3, and 4
2 -- Parameter = Topology
   Prior = All topologies equally probable a priori
   Partitions = 1, 2, 3, and 4
3 -- Parameter = BrLens
   Prior = Branch lengths are Unconstrained:Exponential(10.0)
   Partitions = 1, 2, 3, and 4
MrBayes >

```

Por defecto, se aplica el modelo F81 a todas las particiones. Por lo que, aquellas particiones para las que el modelo que mejor se haya ajustado a tus datos no sea este (véase tema 3.4), debes cambiarlo. El tipo de modelo se puede cambiar utilizando los siguientes comandos:

lset Nst=X Rates=Y

Donde *lset* sirve para cambiar el modelo

Nst sirve para indicar qué modelo se quiere seleccionar (sustituir X por 1 si el modelo es JC o F81, 2 si es HKY o K80, y 6 si es GTR o SYM)

rates indica la tasa de sustitución (sustituir Y por equal, gamma [+G], propinv [+I] o invgamma [+I+G])

Además hay que especificar a qué particiones queremos aplicar el cambio utilizando el comando *applyto*. Por ejemplo, para indicar que ITS1 (nuestra primera partición) se ajusta a un modelo evolutivo K80, teclea *lset applyto=(1) nst=2 statefreqpr=fixed(equal)*

```

MrBayes > lset applyto=(1) nst=2
Setting Nst to 2 for partition 1
Successfully set likelihood model parameters to
partition 1 (if applicable)
MrBayes > lset applyto=(3) rates=propinv
Setting Rates to Propinv for partition 3
Successfully set likelihood model parameters to
partition 3 (if applicable)
MrBayes >

```

Repíte este comando para ajustar a cada partición el modelo que hemos seleccionado como en el tema 3.4. Con ello conseguirás que a cada región genética se le aplique su modelo evolutivo.

Hace falta además que cada partición tenga independientes el resto de los parámetros (los *priors*: frecuencia estacionaria de los nucleótidos o *statefreq*, tasa de sustitución de los nucleótidos o *revmat*, proporción de sitios invariables o *pinvar* y forma de la distribución gamma o *shape*). Para ello, hay dos maneras: (1) especificar para cada partición los valores utilizando el comando *statefreqpr=* [en nuestra partición 1 sería *=fixed(equal)*] o (2) que lo haga automáticamente el programa para cada partición tecleando el comando *unlink statefreq=(all) revmat=(all) shape=(all) pinvar=(all)*.

```
MrBayes > lset applyto=(1) nst=2
Setting Nst to 2 for partition 1
Successfully set likelihood model parameters to
partition 1 (if applicable)
MrBayes > lset applyto=(3) rates=propinv
Setting Rates to Propinv for partition 3
Successfully set likelihood model parameters to
partition 3 (if applicable)
MrBayes > unlink statefreq=(all) revmat=(all) shape=(all) pinvar=(all)_
```

Finalmente, hay que indicar que la tasa de variación general puede ser variable entre las particiones. Esto se consigue con el parámetro *ratepr* del comando *prset*. Teclea *prset applyto=(all) ratepr=variable*.

```
MrBayes > unlink statefreq=(all) revmat=(all) shape=(all) pinvar=(all)
Unlinking
MrBayes > prset applyto=(all) ratepr=variable
Setting Ratepr to Variable [Dirichlet(...,1,...)] for partition 1
Setting Ratepr to Variable [Dirichlet(...,1,...)] for partition 2
Setting Ratepr to Variable [Dirichlet(...,1,...)] for partition 3
Setting Ratepr to Variable [Dirichlet(...,1,...)] for partition 4
Successfully set prior model parameters to all
applicable data partitions
MrBayes >
```

Paso 3. Selecciona el número de generaciones que quieres correr utilizando el comando *Ngen=* y empieza el análisis:

El comando *mcmc* te permite cambiar los parámetros de análisis sin empezarlo. Teclea *mcmc Ngen=1000000*, con ello estarás especificando que cuando se haga el análisis se realicen 1000000 de generaciones.

```
MrBayes > mcmc Ngen=1000000
Setting number of generations to 1000000
Successfully set chain parameters
MrBayes >
```

Paso 4. Ya hemos establecido todas las especificaciones necesarias y podemos proceder al análisis. El comando *mcmc* comienza el análisis con los parámetros previamente seleccionados o, en caso de no haberlo establecido, con los que se indiquen en el momento. Por ejemplo si tecleas *mcmc Ngen=1000000* se inicia el análisis con 1000000 de generaciones [Ten en cuenta que para una publicación deberías utilizar al menos 30.000.000 de generaciones].

El comando *mcmc* incrementa la velocidad a la que se estabiliza la varianza porque implica que se salte de un pico a otro.

Paso 5. El análisis debe continuar hasta que la varianza se estabilice por debajo de 0.01. En la pantalla se va mostrando el número de generaciones que

lleva junto con los valores del *likelihood* obtenidos en cada generación para cada una de las cuatro cadenas. Además, se indica la varianza alcanzada y el tiempo que falta para que se termine el análisis. Cuando se alcanza el número de generaciones solicitado el programa pregunta si se quiere continuar con el análisis (yes) o no (no). Si el número de generaciones solicitado no ha sido suficiente como para estabilizar la varianza por debajo de 0.01 contestaremos que sí queremos continuar con el análisis, indicando cuántas generaciones adicionales queremos hacer hasta obtener la varianza requerida.

Ejercicio 3.5.3. ¿Debes realizar más generaciones o con 1.000.000 es suficiente?

```

998000 -- (-2675.895) (-2685.192) (-2675.542) [-2666.282] * [-2674.977] (-2698.562) (-2675.946) (-2688.298) -- 0:00:00
Average standard deviation of split frequencies: 0.003633
998100 -- (-2675.588) (-2708.690) (-2679.198) [-2662.911] * (-2676.326) (-2691.116) [-2670.112] (-2683.457) -- 0:00:00
998200 -- (-2671.975) (-2694.964) (-2675.874) [-2660.959] * (-2681.878) (-2691.213) [-2662.256] (-2684.246) -- 0:00:00
998300 -- (-2674.489) (-2691.445) (-2675.514) [-2674.840] * (-2686.254) (-2702.080) [-2661.361] (-2698.086) -- 0:00:00
998400 -- (-2681.567) (-2689.400) (-2674.770) [-2672.517] * (-2681.759) (-2699.987) [-2671.099] (-2688.506) -- 0:00:00
998500 -- (-2677.823) (-2701.476) [-2674.743] [-2667.424] * (-2688.178) (-2700.878) [-2666.029] (-2691.885) -- 0:00:00
998600 -- (-2672.328) (-2714.794) (-2672.932) [-2667.511] * (-2686.609) (-2697.405) [-2666.468] (-2681.939) -- 0:00:00
998700 -- (-2683.475) (-2716.428) (-2676.194) [-2664.824] * (-2686.995) (-2696.856) [-2663.468] (-2678.562) -- 0:00:00
998800 -- (-2693.579) (-2724.861) [-2674.737] [-2672.136] * (-2682.043) (-2695.392) [-2663.754] (-2679.737) -- 0:00:00
998900 -- (-2692.782) (-2713.125) (-2681.756) [-2667.349] * (-2676.787) (-2703.474) [-2668.129] (-2676.938) -- 0:00:00
999000 -- (-2687.066) (-2702.809) (-2683.842) [-2670.616] * (-2677.298) (-2701.769) [-2664.702] [-2677.738] -- 0:00:00
Average standard deviation of split frequencies: 0.003542
999100 -- (-2681.620) (-2699.882) (-2677.420) [-2664.897] * (-2671.533) (-2708.172) [-2664.814] (-2681.737) -- 0:00:00
999200 -- (-2676.063) (-2715.438) (-2679.623) [-2659.561] * (-2667.838) (-2707.812) [-2664.261] (-2675.555) -- 0:00:00
999300 -- (-2692.626) (-2716.712) (-2675.824) [-2658.511] * (-2665.971) (-2703.980) [-2665.387] (-2674.326) -- 0:00:00
999400 -- (-2679.282) (-2715.184) (-2676.330) [-2653.746] * (-2668.105) (-2714.611) [-2669.443] (-2676.661) -- 0:00:00
999500 -- (-2679.076) (-2709.655) (-2680.346) [-2654.180] * (-2668.075) (-2702.218) (-2675.123) [-2673.563] -- 0:00:00
999600 -- (-2670.376) (-2701.977) (-2685.324) [-2657.172] * (-2672.376) (-2716.745) (-2675.777) [-2675.171] -- 0:00:00
999700 -- (-2682.781) (-2690.757) (-2681.371) [-2657.808] * (-2670.387) (-2721.298) (-2670.794) [-2669.836] -- 0:00:00
999800 -- (-2681.941) (-2695.989) (-2684.354) [-2660.504] * (-2674.327) (-2716.708) (-2671.249) [-2671.203] -- 0:00:00
999900 -- (-2687.586) (-2688.693) (-2691.983) [-2658.759] * (-2669.516) (-2706.411) [-2675.263] (-2671.985) -- 0:00:00
1000000 -- (-2677.353) (-2686.726) (-2688.604) [-2662.431] * (-2679.579) (-2714.765) [-2667.803] (-2677.377) -- 0:00:00
Average standard deviation of split frequencies: 0.003615
Continue with analysis? (yes/no): no

C:\Archivos de programa\MrBayes\mrBayes-3.1.2\mrBayes.exe
999700 -- (-2682.781) (-2690.757) (-2681.371) [-2657.808] * (-2670.387) (-2721.298) (-2670.794) [-2669.836] -- 0:00:00
999800 -- (-2681.941) (-2695.989) (-2684.354) [-2660.504] * (-2674.327) (-2716.708) (-2671.249) [-2671.203] -- 0:00:00
999900 -- (-2687.586) (-2688.693) (-2691.983) [-2658.759] * (-2669.516) (-2706.411) [-2675.263] (-2671.985) -- 0:00:00
1000000 -- (-2677.353) (-2686.726) (-2688.604) [-2662.431] * (-2679.579) (-2714.765) [-2667.803] (-2677.377) -- 0:00:00
Average standard deviation of split frequencies: 0.003615
Continue with analysis? (yes/no): no

Analysis completed in 493 seconds
Analysis used 492.58 seconds of CPU time
Likelihood of best state for "cold" chain of run 1 was -2647.98
Likelihood of best state for "cold" chain of run 2 was -2649.35
Acceptance rates for the moves in the "cold" chain of run 1:
With prob. Chain accepted changes to
62.37 % param. 1 (ratio) with Dirichlet proposal
42.23 % param. 2 (state frequencies) with Dirichlet proposal
56.52 % param. 3 (state frequencies) with Dirichlet proposal
47.83 % param. 4 (state frequencies) with Dirichlet proposal
25.29 % param. 5 (state frequencies) with Dirichlet proposal
94.98 % param. 6 (prop. invar. sites) with sliding window
65.15 % param. 7 (rate multiplier) with Dirichlet proposal
24.29 % param. 8 (topology and branch lengths) with extending TBR
50.92 % param. 8 (topology and branch lengths) with LOCAL
Acceptance rates for the moves in the "cold" chain of run 2:
With prob. Chain accepted changes to
62.40 % param. 1 (ratio) with Dirichlet proposal
42.21 % param. 2 (state frequencies) with Dirichlet proposal
57.22 % param. 3 (state frequencies) with Dirichlet proposal
47.02 % param. 4 (state frequencies) with Dirichlet proposal
25.11 % param. 5 (state frequencies) with Dirichlet proposal
94.75 % param. 6 (prop. invar. sites) with sliding window
65.08 % param. 7 (rate multiplier) with Dirichlet proposal
24.34 % param. 8 (topology and branch lengths) with extending TBR
50.84 % param. 8 (topology and branch lengths) with LOCAL

Chain swap information for run 1:
      1      2      3      4
1 :      0.43      0.15      0.04
2 : 166544      0.51      0.20
3 : 166369 167476      0.54
4 : 166653 165998 166960

Chain swap information for run 2:
      1      2      3      4
1 : 166884      0.44      0.16      0.04
2 : 166139 166603      0.51      0.20
3 : 166567 166912 166895

Upper diagonal: Proportion of successful state exchanges between chains
Lower diagonal: Number of attempted state exchanges between chains

Chain information:
ID -- Heat
1 -- 1.00 (cold chain)
2 -- 0.83
3 -- 0.71
4 -- 0.65

Heat = 1 / (1 + T * (ID - 1))
where T = 0.20 is the temperature and ID is the chain number

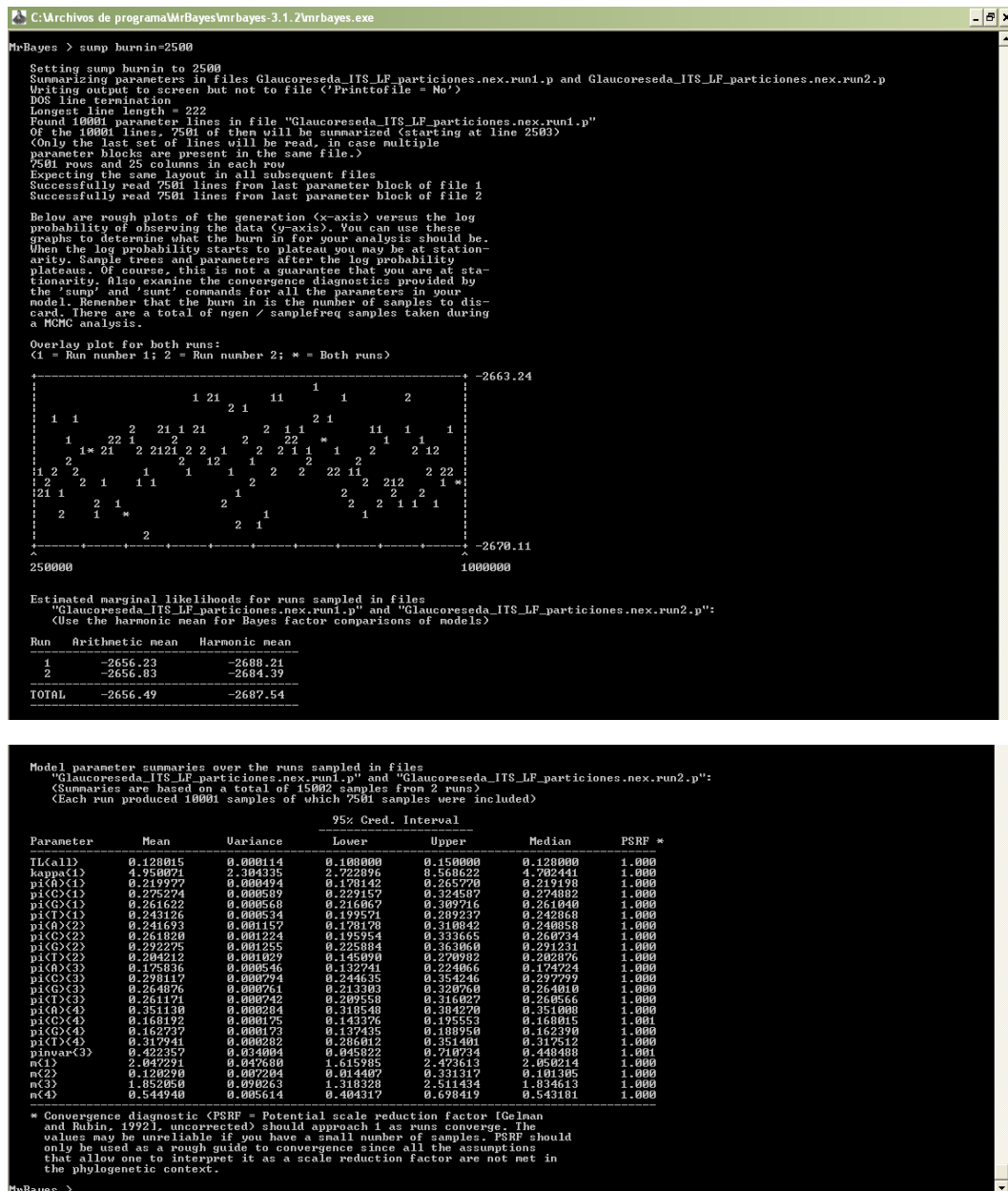
MrBayes > sump burnin=2500

```

Paso 6. Una vez estabilizada la varianza y finalizado el análisis debemos visualizar los resultados. Para ello, teclea el comando *sump*. Lo que nos proporciona este comando es una tabla (en la que debemos comprobar que el parámetro PSRF es próximo a 1), y un gráfico de dispersión (que no debe mostrar tendencias si la varianza estaba estabilizada). Sin

embargo, estos resultados no nos interesan ya que incluyen también resultados previos a la estabilización. Para matrices no muy complejas, como la nuestra, se puede asumir que eliminando el 25% inicial de los árboles construidos estaremos eliminando aquellos árboles obtenidos antes de la estabilización. Para ello necesitamos el comando *burnin*. Por defecto, MrBayes guarda los resultados de la búsqueda cada 100 generaciones (un árbol y los valores de los parámetros asociados), por lo que si hemos realizado un análisis de 1000000 de generaciones, debemos eliminar el burnin tecleando la línea de comando:

```
sump burnin=2500
```



- Paso 7. Para obtener un resumen de los árboles el comando que se utiliza es *sumt*, también descartando el 25% inicial de muestras. En el ejemplo anterior, la línea de comando sería:

```
sumt burnin=2500
```

El comando *sumt* crea tres archivos en la carpeta raíz: uno con extensión .parts, en el que se incluyen las particiones (clados), su probabilidad a posteriori (pp) y la longitud de rama asociada; uno con extensión .con, en el que se incluyen dos árboles consenso, uno con la longitud de las ramas y otro con la probabilidad; y uno con extensión .trprobs, en el que se guardan los árboles obtenidos ordenados por su probabilidad posterior.

Por defecto, para la construcción del árbol de consenso, MrBayes condensa los clados con menos de 50 de probabilidad posterior (*halfcompat*). Esto se puede cambiar con el comando *contype=* [Por ejemplo, *contype=allcompat* no condensa ningún clado]

- Paso 8. Una vez obtenidos los árboles salimos del programa con el comando *quit*.

Otra opción alternativa consiste en incluir todos los comandos del análisis en un bloque de comandos a continuación del que hemos añadido especificando las particiones al final del archivo de la matriz. Para ello, puedes abrir la matriz con un editor de texto (WordPad para PC, o TextWrangler para MAC; OJO es muy importante que no lo abráis con Word) y a continuación del bloque donde has definido las particiones escribir lo siguiente:

```
lset applyto=(1) nst=2 rates=equal;
prset applyto=(1) statefreqpr=fixed(equal);
lset applyto=(2) nst=1 rates=equal;
prset applyto=(2) statefreqpr=fixed(equal);
lset applyto=(3) nst=1 rates=propinv;
prset applyto=(3) statefreqpr=fixed(equal);
lset applyto=(4) nst=1 rates=equal;
prset applyto=(4) statefreqpr=dirichlet(1,1,1,1);
```

[hasta aquí le hemos indicado a MrBayes que al ejecutar esta matriz aplique el modelo evolutivo K80 a la primera partición (ITS1), JC a la 2ª (5.8S), JC+I a la 3ª (ITS2) y F81 a la 4ª (*trnL-F*)]

```
set autoclose=yes;
```

[Con este comando le indicamos al programa que se cierre una vez terminado el análisis]

```
mcmc ngen=1000000 printfreq=100 samplefreq=100 nruns=2
nchains=4 savebrlens=yes burninfrac=0.25;
```

[Aquí le estamos indicando los parámetros para realizar la búsqueda: mcmc -que inicie las cadenas-, ngen=1000000 -que realice un millón de generaciones-, printfreq=100 -que nos muestre en la pantalla los resultados cada 100 generaciones-, samplefreq=100 -que guarde los

resultados de la búsqueda (árbol, más parámetros) una vez cada 100 generaciones-, nruns=2 –que inicie dos análisis al mismo tiempo-, nchains=4 –que corra cuatro cadenas de Markov por cada análisis-, savebrlens=yes –que guarde la longitud de las ramas-, burninfrac=0.25 –que aplique un burnin del 25%-]

```
sump;  
sumt contype=halfcompat;  
end;
```

[Estos últimos comandos le indican que realizado el análisis: sump; -nos compile los resultados-, sumt contype=halfcompat; -que compile los árboles y compute el árbol de consenso-, end; -que termine el análisis-]

Ejercicio 3.5.4. Escribe los comandos para el análisis dentro del archivo de la matriz Glaucoreseda_ITS_LF_particiones.nex.