

MCB 372

Positive, and purifying selection.
Neutral theory

Peter Gogarten
Office: BSP 404
phone: 860 486-4061.
Email: gogarten@uconn.edu

the gradualist point of view

Evolution occurs within populations where the fittest organisms have a selective advantage. Over time the advantageous genes become fixed in a population and the population gradually changes.

Note: this is not in contradiction to the theory of neutral evolution. (which says what?)

Processes that MIGHT go beyond inheritance with variation and selection?

- Horizontal gene transfer and recombination
- Polyploidization (botany, vertebrate evolution) see [here](#)
- Fusion and cooperation of organisms (Kefir, lichen, also the eukaryotic cell)
- Targeted mutations (?), genetic memory (?) (see [Foster's](#) and [Hall's](#) reviews on directed/adaptive mutations; see [here](#) for a counterpoint)
- Random genetic drift (i.e. traits are fixed even though they do not provide an advantage)
- [Gratuitous complexity](#) (introns, split intein)
- Selfish genes (who/what is the subject of evolution??)
- Parasitism, altruism, Morons (Gene Transfer Agents)

Assignments:

- Read through chapter 9
- Work on your student project
- Analyze one dataset of your choice in MrBayes.

Alternative Approaches to Estimate Posterior Probabilities

Bayesian Posterior Probability Mapping with MrBayes
(Huelsenbeck and Ronquist, 2001)

Problem:

Strimmer's formula $p_i = \frac{L_i}{L_1 + L_2 + L_3}$ only considers 3 trees (those that maximize the likelihood for the three topologies)

Solution:

Exploration of the tree space by sampling trees using a biased random walk (Implemented in MrBayes program)

Trees with higher likelihoods will be sampled more often

$$p_i = \frac{N_i}{N_{\text{total}}}$$

where N_i - number of sampled trees of topology i , $i=1,2,3$
 N_{total} - total number of sampled trees (has to be large)

Illustration of a biased random walk

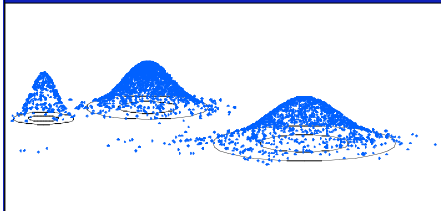


Figure generated using MCRobot program (Paul Lewis, 2001)

selection versus drift

see Kent Holsinger's java simulations at <http://darwin.ecb.uconn.edu/simulations/simulations.html>

The law of the gutter.

compare **drift** versus **select + drift**

The larger the population the longer it takes for an allele to become fixed.

Note: Even though an allele conveys a strong selective advantage of 10%, the allele has a rather large chance to go extinct.

Note#2: Fixation is faster under selection than under drift.

BUT

s=0

Probability of fixation, P, is equal to frequency of allele in population.
Mutation rate (per gene/per unit of time) = u
freq. with which allele is generated in diploid population size N = u*2N
Probability of fixation for each allele = 1/(2N)

Substitution rate =

frequency with which new alleles are generated * Probability of fixation = u*2N * 1/(2N) = u

Therefore:

If $s=0$, the substitution rate is independent of population size, and equal to the mutation rate !!!! (NOTE: Mutation unequal Substitution!)

This is the reason that there is hope that the molecular clock might sometimes work.

Fixation time due to drift alone:

$t_{av} = 4 * N_e$ generations

(N_e = effective population size; For n discrete generations

$$N_e = n(1/N_1 + 1/N_2 + \dots + 1/N_n)$$

s>0

Time till fixation on average:

$t_{av} = (2/s) \ln(2N)$ generations
(also true for mutations with negative "s" ! discuss among yourselves)

E.g.: $N=10^6$,

$s=0$: average time to fixation: $4 * 10^6$ generations

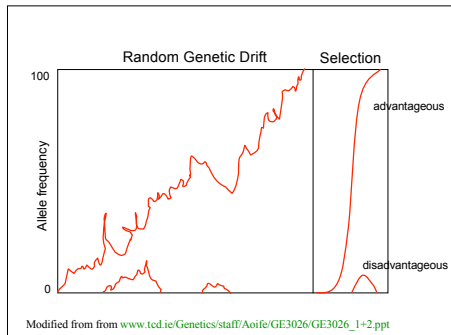
$s=0.01$: average time to fixation: 2900 generations

$N=10^4$,

$s=0$: average time to fixation: 40.000 generations

$s=0.01$: average time to fixation: 1.900 generations

=> substitution rate of mutation under positive selection is larger than the rate with which neutral mutations are fixed.



Modified from www.tcd.ie/Genetics/staff/Aoife/GE3026/GE3026_1+2.ppt

Measuring Selection on Genes

- Null hypothesis = neutral evolution
- Under neutral evolution, synonymous substitutions should accumulate at a rate equal to mutation rate
- Under neutral evolution, amino acid substitutions should also accumulate at a rate equal to the mutation rate

From: mentor.lscf.ucsb.edu/course/spring/eemb102/lecture/Lecture7.ppt

Counting #s/#a

	Ser	Ser	Ser	Ser	Ser
Species1	TGA	TGC	TGT	TGT	TGT
	Ser	Ser	Ser	Ser	Aia
Species2	TGT	TGT	TGT	TGT	GGT

#s = 2 sites
 #a = 1 site
 #a/#s=0.5

To assess selection pressures one needs to calculate the rates (Ka, Ks), i.e. the occurring substitutions as a fraction of the possible syn. and nonsyn. substitutions.

Things get more complicated, if one wants to take transition transversion ratios and codon bias into account. See chapter 4 in Nei and Kumar, Molecular Evolution and Phylogenetics.

Modified from: mentor.lscf.ucsb.edu/course/spring/eemb102/lecture/Lecture7.ppt

dambe

Two programs worked well for me to align nucleotide sequences based on the amino acid alignment,

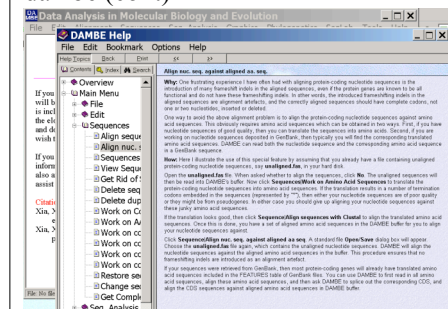
One is **DAMBE** (only for windows). This is a handy program for a lot of things, including reading a lot of different formats, calculating phylogenies, it even runs codeml (from PAML) for you.

The procedure is not straight forward, but is well described on the help pages. After installing DAMBE go to HELP -> general HELP -> sequences -> align nucleotide sequences based on ...->

If you follow the instructions to the letter, it works fine.

DAMBE also calculates Ka and Ks distances from codon based aligned sequences.

dambe (cont)



aa based nucleotide alignments (cont)

An alternative is the tralign program that is part of the emboss package. On bbxsrvi you can invoke the program by typing tralign.

Instructions and program description are [here](#).

If you want to use your own dataset in the lab on Wednesday, generate a codon based alignment with either *dambe* (on PCs only) or *tralign* (Emboss, installed on cluster) and save it as a nexus file and as a phylog formatted multiple sequence file (using either clustalw, PAUP (export or tonexus), dambe, or [readseq](#) on the web)

PAML (codeml) the basic model

$$q_{ij} = \begin{cases} 0, & \text{if the two codons differ at more than one position,} \\ \pi_j, & \text{for synonymous transversion,} \\ \kappa\pi_j, & \text{for synonymous transition,} \\ \omega\pi_j, & \text{for nonsynonymous transversion,} \\ \omega\kappa\pi_j, & \text{for nonsynonymous transition,} \end{cases}$$

The equilibrium frequency of codon j (π_j) can be considered a free parameter, but can also be calculated from the nucleotide frequencies at the three codon positions (control variable CodonFreq). Under this model, the relationship holds that $\omega = d_s/d_n$, the ratio of nonsynonymous/synonymous substitution rates. This basic model is fitted by specifying model = 0 NSsites = 0, in the control file codeml.ctl. It forms the basis for more sophisticated models implemented in codeml.

sites versus branches

You can determine omega for the whole dataset; however, usually not all sites in a sequence are under selection all the time.

PAML (and other programs) allow to either determine omega for each site over the whole tree, Branch Models, or determine omega for each branch for the whole sequence, Site Models.

It would be great to do both, i.e., conclude codon 176 in the vacuolar ATPases was under positive selection during the evolution of modern humans – alas, a single site does not provide any statistics

Sites model(s)

work great have been shown to work great in few instances. The most celebrated case is the influenza virus HA gene.

A talk by Walter Fitch (slides and sound) on the evolution of this molecule is [here](#). This [article by Yang et al. 2000](#) gives more background on ml approaches to measure omega. The dataset used by Yang et al is here: [flu_data.paup](#).

Nexus files:

This is the file format used by many popular programs like MacClade, Mesquite, ModelTest, MrBayes and PAUP*. Nexus file names often have a .nxs or .nex extension.

A formal description of the NEXUS format can be found in Maddison et al. (1997).

Conversion of an interleaved NEXUS file to a non-interleaved NEXUS file: execute the file in PAUP*, and export the file as non-interleaved NEXUS file. You can also type the commands:

```
export file=yourfile.nex format=nexus interleaved=no;
clustalw saves and reads Nexus sequence and tree files
(check on gap treatment and label as DNA or aa)
```

sample DNA file

```
#nexus
begin data;
dimensions ntax=10 nchar=705;
format datatyp=dnal interleave=yes gap=- missing=?;
matrix
Cow      ATGGCATATCCCATACAGACATAGAGTTCCTCAGATGCGACATCAACATCATAGAGAGACTA
Carp     ATGGCAGACCCAGCCAGCAGTAGGTGTTTCAGAGAGCCAGCCATACCCCTTATAGAGAGACTT
Chickens ATGGCCAGACCCCTCCAGCAGTAGGTGTTTCAGAGAGCCAGCCATACCCCTTATAGAGAGACTT
Human    ATGGCAGACCCAGCCAGCAGTAGGTGTTTCAGAGAGCCAGCCATACCCCTTATAGAGAGACTT
Leach    ATGGCAGACCCAGCCAGCAGTAGGTGTTTCAGAGAGCCAGCCATACCCCTTATAGAGAGACTT
Mouse    ATGGCAGACCCAGCCAGCAGTAGGTGTTTCAGAGAGCCAGCCATACCCCTTATAGAGAGACTT
Rat      ATGGCAGACCCAGCCAGCAGTAGGTGTTTCAGAGAGCCAGCCATACCCCTTATAGAGAGACTT
Seal     ATGGCAGACCCAGCCAGCAGTAGGTGTTTCAGAGAGCCAGCCATACCCCTTATAGAGAGACTT
Whale    ATGGCAGACCCAGCCAGCAGTAGGTGTTTCAGAGAGCCAGCCATACCCCTTATAGAGAGACTT
Frog     ATGGCAGACCCAGCCAGCAGTAGGTGTTTCAGAGAGCCAGCCATACCCCTTATAGAGAGACTT
//
Cow      CTCACCTTCACGACACAGCCATTAATATATTCCTTCCTTAATAGCTCATATAGACTTTCAC
Carp     CTCACCTTCACGACACAGCCATTAATATATTCCTTCCTTAATAGCTCATATAGACTTTCAC
Chickens CTCACCTTCACGACACAGCCATTAATATATTCCTTCCTTAATAGCTCATATAGACTTTCAC
Human    CTCACCTTCACGACACAGCCATTAATATATTCCTTCCTTAATAGCTCATATAGACTTTCAC
Leach    CTCACCTTCACGACACAGCCATTAATATATTCCTTCCTTAATAGCTCATATAGACTTTCAC
Mouse    CTCACCTTCACGACACAGCCATTAATATATTCCTTCCTTAATAGCTCATATAGACTTTCAC
Rat      CTCACCTTCACGACACAGCCATTAATATATTCCTTCCTTAATAGCTCATATAGACTTTCAC
Seal     CTCACCTTCACGACACAGCCATTAATATATTCCTTCCTTAATAGCTCATATAGACTTTCAC
Whale    CTCACCTTCACGACACAGCCATTAATATATTCCTTCCTTAATAGCTCATATAGACTTTCAC
Frog     CTCACCTTCACGACACAGCCATTAATATATTCCTTCCTTAATAGCTCATATAGACTTTCAC
//
Frog     AACTGCATCTTCATACATACTA-----GAAGCATCACTA-----AGA
end;
```

sample aa file

```
#NEXUS
Begin data;
Dimensions ntax=10 nchar=234;
Format datatyp=proal dgap=- list=leave;
Matrix
Cow      NAYVYVQLGQFQDAMPDREELLEFPRFHTLMLVFLISESLVLYLLEMLTFLTWTFPHWDGE
Carp     NAYVYVQLGQFQDAMPDREELLEFPRFHTLMLVFLISESLVLYLLEMLTFLTWTFPHWDGE
Chickens NAYVYVQLGQFQDAMPDREELLEFPRFHTLMLVFLISESLVLYLLEMLTFLTWTFPHWDGE
Human    NAYVYVQLGQFQDAMPDREELLEFPRFHTLMLVFLISESLVLYLLEMLTFLTWTFPHWDGE
Leach    NAYVYVQLGQFQDAMPDREELLEFPRFHTLMLVFLISESLVLYLLEMLTFLTWTFPHWDGE
Mouse    NAYVYVQLGQFQDAMPDREELLEFPRFHTLMLVFLISESLVLYLLEMLTFLTWTFPHWDGE
Rat      NAYVYVQLGQFQDAMPDREELLEFPRFHTLMLVFLISESLVLYLLEMLTFLTWTFPHWDGE
Seal     NAYVYVQLGQFQDAMPDREELLEFPRFHTLMLVFLISESLVLYLLEMLTFLTWTFPHWDGE
Whale    NAYVYVQLGQFQDAMPDREELLEFPRFHTLMLVFLISESLVLYLLEMLTFLTWTFPHWDGE
Frog     NAYVYVQLGQFQDAMPDREELLEFPRFHTLMLVFLISESLVLYLLEMLTFLTWTFPHWDGE
//
Loach    QTAFIARRPVQVQCSIECCANBSFDPVYVAVVPLRFENWSSMLKQASLQSLG
Mouse    QATVYVNRDGLFVQCSIECCANBSFDPVYVAVVPLRFENWSSMLKQASLQSLG
Rat      QATVYVNRDGLFVQCSIECCANBSFDPVYVAVVPLRFENWSSMLKQASLQSLG
Seal     QTFLYVNRDGLFVQCSIECCANBSFDPVYVAVVPLRFENWSSMLKQASLQSLG
Whale    QTFLYVNRDGLFVQCSIECCANBSFDPVYVAVVPLRFENWSSMLKQASLQSLG
Frog     QTFLYVNRDGLFVQCSIECCANBSFDPVYVAVVPLRFENWSSMLKQASLQSLG
//
End;
```

Another example is [here](#)

More information on Nexus files and PAUP and MrBayes commands are in the respective manuals:

<http://paup.csjt.fsu.edu/>, manual [here](#), tutorial [here](#).

<http://mrbayes.csjt.fsu.edu/>, manual [here](#), [Wikki](#)

sites model in MrBayes

The MrBayes block in a nexus file might look something like this:

```
begin mrbayes;
set autoclose=yes;
lset nst=2 rates=gamma nuemodel=codon omegavar=Ny98;
mcmc samplefreq=500 printfreq=500;
mcmc ngen=500000;
sump burnin=50;
sumt burnin=50;
end;
```

Vincent Daubin and Howard Ochman: Bacterial Genomes as New Gene Homes: The Genealogy of ORFans in *E. coli*. *Genome Research* 14:1036-1042, 2004

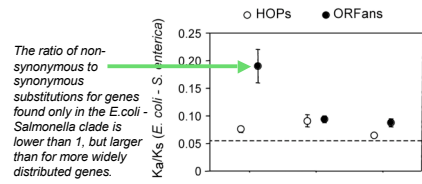


Fig. 3 from Vincent Daubin and Howard Ochman, *Genome Research* 14:1036-1042, 2004

Trunk-of-my-car analogy: Hardly anything in there is the result of providing a selective advantage. Some items are removed quickly (purifying selection), some are useful under some conditions, but most things do not alter the fitness.



Could some of the inferred purifying selection be due to the acquisition of novel detrimental characteristics (e.g., protein toxicity)?

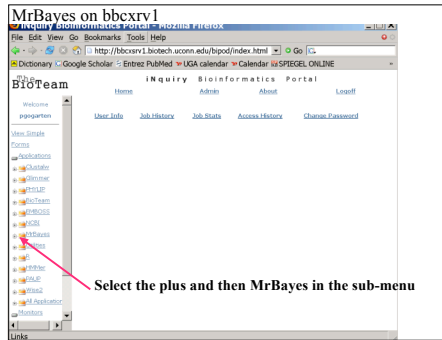
MrBayes on bbcxrv1

Create the nexus file on your computer.

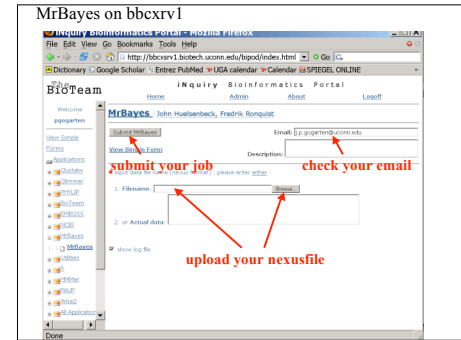
It will help to have MrBayes installed locally, this way you can check that you don't have any typos in the MrBayes block.

Direct your browser to

<http://bbcxrv1.biotech.uconn.edu/bipod/index.html>



Select the plus and then MrBayes in the sub-menu

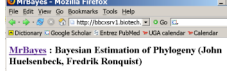


MrBayes on bbcbxrv1

You will receive the results per email, and you will receive the link of a web page that lists all the output files. In this case:

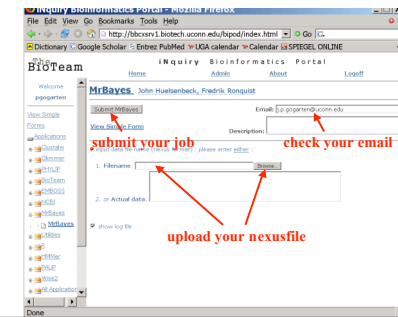
<http://bbcbxrv1.biotech.uconn.edu/nise/tmp/A1070011431640/results.html>

You can save the files from your browser, or open the email attachments.



The files we are particularly interested in are the parameter file and the MrBayes output (to check for potential problems).

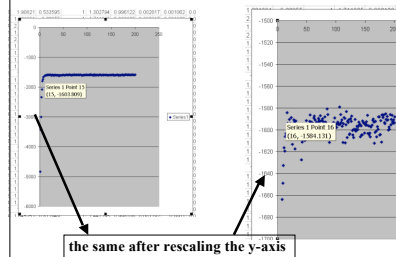
MrBayes on bbcbxrv1



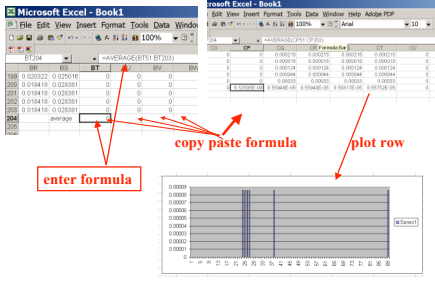
MrBayes analyzing the *.nex.p file

1. The easiest is to load the file into excel (if your alignment is too long, you need to load the data into separate spreadsheets – see [here](#) exercise 2 item 2 for more info)
2. plot LogL to determine which samples to ignore
3. for each codon calculate the the average probability (from the samples you do not ignore) that the codon belongs to the group of codons with omega>1.
4. plot this quantity using a bar graph.

plot LogL to determine which samples to ignore



for each codon calculate the the average probability



MrBayes on bbcbxrv1

- If you do this for your own data,
 - run the procedure first for only 50000 generations (takes about 30 minutes) to check that everything works as expected,
 - then run the program overnight for at least 500 000 generations.
 - Especially, if you have a large dataset, do the latter twice and compare the results for consistency. (I prefer two runs over 500000 generations each over one run over a million generations.)

The preferred wa to run mrbayes is to use the command line:
>mb
Do example on threonlyRS

PAML – codeml – sites model

the paml package contains several distinct programs for nucleotides (baseml) protein coding sequences and amino acid sequences (codeml) and to simulate sequences evolution.

The input file needs to be in phyml format.

By default it assumes a sequential format (e.g. [here](#)).

If the sequences are interleaved, you need to add an "1" to the first line, as in these example headers:

```

# 467
# 1
# 2
# 3
# 4
# 5
# 6
# 7
# 8
# 9
# 10
# 11
# 12
# 13
# 14
# 15
# 16
# 17
# 18
# 19
# 20
# 21
# 22
# 23
# 24
# 25
# 26
# 27
# 28
# 29
# 30
# 31
# 32
# 33
# 34
# 35
# 36
# 37
# 38
# 39
# 40
# 41
# 42
# 43
# 44
# 45
# 46
# 47
# 48
# 49
# 50
# 51
# 52
# 53
# 54
# 55
# 56
# 57
# 58
# 59
# 60
# 61
# 62
# 63
# 64
# 65
# 66
# 67
# 68
# 69
# 70
# 71
# 72
# 73
# 74
# 75
# 76
# 77
# 78
# 79
# 80
# 81
# 82
# 83
# 84
# 85
# 86
# 87
# 88
# 89
# 90
# 91
# 92
# 93
# 94
# 95
# 96
# 97
# 98
# 99
# 100
# 101
# 102
# 103
# 104
# 105
# 106
# 107
# 108
# 109
# 110
# 111
# 112
# 113
# 114
# 115
# 116
# 117
# 118
# 119
# 120
# 121
# 122
# 123
# 124
# 125
# 126
# 127
# 128
# 129
# 130
# 131
# 132
# 133
# 134
# 135
# 136
# 137
# 138
# 139
# 140
# 141
# 142
# 143
# 144
# 145
# 146
# 147
# 148
# 149
# 150
# 151
# 152
# 153
# 154
# 155
# 156
# 157
# 158
# 159
# 160
# 161
# 162
# 163
# 164
# 165
# 166
# 167
# 168
# 169
# 170
# 171
# 172
# 173
# 174
# 175
# 176
# 177
# 178
# 179
# 180
# 181
# 182
# 183
# 184
# 185
# 186
# 187
# 188
# 189
# 190
# 191
# 192
# 193
# 194
# 195
# 196
# 197
# 198
# 199
# 200
# 201
# 202
# 203
# 204
# 205
# 206
# 207
# 208
# 209
# 210
# 211
# 212
# 213
# 214
# 215
# 216
# 217
# 218
# 219
# 220
# 221
# 222
# 223
# 224
# 225
# 226
# 227
# 228
# 229
# 230
# 231
# 232
# 233
# 234
# 235
# 236
# 237
# 238
# 239
# 240
# 241
# 242
# 243
# 244
# 245
# 246
# 247
# 248
# 249
# 250
# 251
# 252
# 253
# 254
# 255
# 256
# 257
# 258
# 259
# 260
# 261
# 262
# 263
# 264
# 265
# 266
# 267
# 268
# 269
# 270
# 271
# 272
# 273
# 274
# 275
# 276
# 277
# 278
# 279
# 280
# 281
# 282
# 283
# 284
# 285
# 286
# 287
# 288
# 289
# 290
# 291
# 292
# 293
# 294
# 295
# 296
# 297
# 298
# 299
# 300
# 301
# 302
# 303
# 304
# 305
# 306
# 307
# 308
# 309
# 310
# 311
# 312
# 313
# 314
# 315
# 316
# 317
# 318
# 319
# 320
# 321
# 322
# 323
# 324
# 325
# 326
# 327
# 328
# 329
# 330
# 331
# 332
# 333
# 334
# 335
# 336
# 337
# 338
# 339
# 340
# 341
# 342
# 343
# 344
# 345
# 346
# 347
# 348
# 349
# 350
# 351
# 352
# 353
# 354
# 355
# 356
# 357
# 358
# 359
# 360
# 361
# 362
# 363
# 364
# 365
# 366
# 367
# 368
# 369
# 370
# 371
# 372
# 373
# 374
# 375
# 376
# 377
# 378
# 379
# 380
# 381
# 382
# 383
# 384
# 385
# 386
# 387
# 388
# 389
# 390
# 391
# 392
# 393
# 394
# 395
# 396
# 397
# 398
# 399
# 400
# 401
# 402
# 403
# 404
# 405
# 406
# 407
# 408
# 409
# 410
# 411
# 412
# 413
# 414
# 415
# 416
# 417
# 418
# 419
# 420
# 421
# 422
# 423
# 424
# 425
# 426
# 427
# 428
# 429
# 430
# 431
# 432
# 433
# 434
# 435
# 436
# 437
# 438
# 439
# 440
# 441
# 442
# 443
# 444
# 445
# 446
# 447
# 448
# 449
# 450
# 451
# 452
# 453
# 454
# 455
# 456
# 457
# 458
# 459
# 460
# 461
# 462
# 463
# 464
# 465
# 466
# 467

```

PAML – codeml – sites model (cont.)

the program is invoked by typing codeml followed by the name of a control file that tells the program what to do.

paml can be used to find the maximum likelihood tree, however, the program is rather slow. Phyml is a better choice to find the tree, which then can be used as a user tree.

An example for a codeml.ctf file is [codeml_hv1_sites.ctf](#)
This file directs codeml to run three different models: one with an omega fixed at 1, a second where each site can be either have an omega between 0 and 1, or an omega of 1, and third a model that uses three omegas as described before for MrBayes.
The output is written into a file called [Hv1_sites.codeml_out](#) (as directed by the control file).

Point out log likelihoods and estimated parameter line (kappa and omegas)

Additional useful information is in the [rst](#) file generated by the codeml

Discuss overall result.

PAML – codeml – branch model

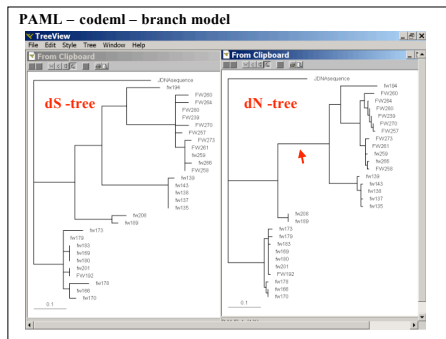
For the same dataset to estimate the dN/dS ratios for individual branches, you could use this file [codeml_hv1_branches.ctf](#) as control file.

The output is written, as directed by the control file, into a file called [Hv1_branch.codeml_out](#)

A good way to check for episodes with plenty of non-synonymous substitutions is to compare the dn and ds trees.

Also, it might be a good idea to repeat the analyses on parts of the sequence (using the same tree). In this case the sequences encode a family of spider toxins that include the mature toxin, a propeptide and a signal sequence (see [here](#) for more information).

Bottom line: one needs plenty of sequences to detect positive selection.



where to get help

read the manuals and help files
check out the discussion boards

else

hy-phy (hypothesis testing using phylogenetics) does very well in analyzing selection pressures.

The easiest is probably to run the analyses on the authors [datamonkey](http://datamonkey.com).

→

