# modelestimator v2: inferring amino acid replacement rates from multi-sequence alignments

**Ruben Ridderström**[1] **and Lars Arvestad** [1,2]

**1** Department of Mathematics, Stockholm University, Sweden **2** Swedish e-science Research Centre

## Summary

Phylogenetic inference is often based on generative probabilistic models, under which one seeks to maximize the likelihood of an estimated tree or estimate a posterior distribution on trees for some given DNA or protein sequence data. The core of most evolution models describe how symbols in the input sequences are replaced. For models of DNA substitution, the few involved parameters can be inferred during tree estimation, but that approach is difficult for protein models due to the many free parameters. However, there is a range of standard models that provide good approximations in many cases, for example Dayhoff (Dayhoff et al., 1978), WAG (Whelan & Goldman, 2001), JTT (Jones et al., 1992), LG (Le & Gascuel, 2008), VT (Müller & Vingron, 2000), and more.

Methods have been developed to estimate "empirical models" – replacement rate matrices directly from peptide multiple sequence alignments (Arvestad, 2006; Holmes & Rubin, 2002; Minh et al., 2021). This paper reports a re-implementation of the modelestimator method (Arvestad, 2006).

## Improvements

For the re-implementation, we sought improvements in some basic aspects:

1. User friendliness
2. Interoperability
3. Computational efficiency
4. Maintainability

We consider *User friendliness* to include easy deployment and following Unix practices for command-line tools. The new version is written in Python and we therefore use PyPI, the Python Package Index (*PyPI*, n.d.), to provide straightforward software installation. Proper care has been taken to provide a clean command-line user interface, which improves on the previous version of modelestimator.

*Interoperability* has been improved in two ways. Firstly, we get easy parsing of several standard input formats using Biopython (Cock et al., 2009). Secondly, there are command-line options that ensure that the output is suitable for popular phylogenetic inference tools, including MrBayes (Ronquist et al., 2012), IQTREE2 (Nguyen et al., 2015), PhyML (Guindon et al., 2010), RAxML (Stamatakis, 2014), and PAML (Yang, 2007).

Choosing Python for the implementation may not be optimal for *Computational efficiency*, but since NumPy (Oliphant, 2006) is used for core computations, modelestimator is fast enough for problem-free interactive usage even on large input alignments.

To ensure *Maintainability*, we have implemented a set of unit tests, put the source code at GitHub, where development and community infrastructure is provided, and implemented

40 continuous integration tests. Furthermore, leaving the old code base for `modelestimator` (a
41 combination of Perl and Octave) is a step towards maintainability.

42 We recommend installing `modelestimator` using PyPI: "`pip install modelestimator-v2`".

## Statement of need

44 As shown in (Arvestad, 2006), inferring an amino-acid replacement-model for an alignment
45 often improves the likelihood of a phylogeny inference. Consider estimating a phylogeny for an
46 alignment of some Cres/Testatin homologs [CTES; Frygelius et al. (2010)], with 65 sequences
47 and 441 columns. In the following, we use `modelestimator` to estimate replacement models
48 for this specific alignment and compare it to standard models and methods using PhyML
49 (Guindon et al., 2010) and IQTREE2 (Nguyen et al., 2015).

### PhyML with modelestimator

51 Running PhyML to infer a phylogenetic tree using default options means using the LG model
52 and empirical equilibrium frequencies:

```
$ phyml -d aa -i CTES.phy
```

53 This estimation took about 2m40s on a MacBook Pro with an M1 Pro CPU and returns a
54 tree with a log-likelihood of -12418.6. LG might not be the right model, so it is a good idea
55 to consider alternative models. However, we can also try `modelestimator` and get a model
56 tailored for the present alignment:

```
$ modelestimator -a phyml -f phylip CTES.phy > ratematrix.txt
```

57 Option `-a phyml` makes sure the output works for PhyML, but the format is the same for
58 IQTREE2. This estimation took about a second on the same computer.

59 The resulting model can be used by PhyML in a new attempt. Using "`--model custom`"
60 instructs PhyML to use a user-provided rate matrix and "`--aa ratematrix.txt`" locates the
61 file containing it.

```
$ phyml -d aa -i CTES.phy  --model custom --aa ratematrix.txt
```

62 The resulting tree is slightly different and has a log-likelihood of -12206.7, an improvement by
63 about 200. That is a significant improvement for a long-running ML computation given that
64 finding the rate matrix only took a second.

### Using IQTREE2, ModelFinder, and modelestimator

66 IQTREE2 provides a convenient method [ModelFinder; Kalyaanamoorthy et al. (2017)] for
67 automating model search among a large number of published models, including combinations
68 with models for rate heterogeneity. In addition, support for making a statistically sensible
69 model choice (with respect to the number of free parameters in the model), makes IQTREE2
70 helpful.

71 Applying IQTREE2 to our alignment, using default options,

```
$ iqtree2 -s CTES.phy
```

72 compares 1225 model combinations and takes almost 8 minutes using a single CPU core. The
73 JTT+R3 model is chosen based on the Bayesian Information Criterion (BIC), and it allows a
74 tree with a log-likelihood of -12289.5 to be inferred.

75 Using the model estimated by `modelestimator` is done like this:

```
$ iqtree2 -s CTES.phy -m ratematrix.txt
```

<sub>76</sub> IQTREE2 now finds a slightly different tree with this specific rate matrix and the log-likelihood
<sub>77</sub> increases to -12278.4. The computation takes about 40 seconds, by avoiding ModelFinder, so
<sub>78</sub> we have a 12x speedup with a better log-likelihood.

<sub>79</sub> One can also request that IQTREE adds a rate heterogeneity model by specifying the model
<sub>80</sub> "ratematrix.txt+I+G4". The same tree is inferred, but the log-likelihood improves to -12208.3.

### Using IQTREE2, QMaker, and modelestimator

<sub>82</sub> IQTREE has support for using ML for estimating an empirical model with the QMaker method
<sub>83</sub> (Minh et al., 2021). We can apply QMaker with the already inferred tree and note JTT as a
<sub>84</sub> good starting point for model search:

```
$ iqtree2 -seed 1 -T 1 -s CTES.phy -te treefile --init-model JTT \
        --model-joint GTR20+FO --prefix GTR20_FO
```

<sub>85</sub> This inference has a better log-likelihood, -12136.2, than what was achieved with
<sub>86</sub> modelestimator. However, it requires 7 minutes to compute (single core) even though JTT is
<sub>87</sub> given as a suitable starting point and the previously inferred tree is used to initialize the search.

<sub>88</sub> Now try the same search with the rate matrix from modelestimator as a starting point.

```
$ iqtree2 -seed 1 -T 1 -s CTES.phy -te treefile --init-model ratematrix.txt \
        --model-joint GTR20+FO --prefix GTR20_FO_with_modelestimator
```

<sub>89</sub> This computation takes about a minute and finds a model that calculates the likelihood of the
<sub>90</sub> tree to -12135.9. That is, we get a seven-fold speedup due to a better starting point for the
<sub>91</sub> parameter search, and the cost is a computation that took us a second.

## Conclusion

<sub>93</sub> The reimplemented modelestimator is easy to install and run, and speeds up model choice
<sub>94</sub> or inference. With enough data, it estimates rate matrices that yield higher likelihoods than
<sub>95</sub> standard models.

## References

<sub>97</sub> Arvestad, L. (2006). Efficient methods for estimating amino acid replacement rates. *Journal*
<sub>98</sub> *of Molecular Evolution*, *62*(6), 663–673. https://doi.org/10.1007/s00239-004-0113-9

<sub>99</sub> Cock, P. J. A., Antao, T., Chang, J. T., Chapman, B. A., Cox, C. J., Dalke, A., Friedberg,
<sub>100</sub> I., Hamelryck, T., Kauff, F., Wilczynski, B., & Hoon, M. J. L. de. (2009). BioPython:
<sub>101</sub> Freely available Python tools for computational molecular biology and bioinformatics.
<sub>102</sub> *Bioinformatics*, *25*(11), 1422–1423. https://doi.org/10.1093/bioinformatics/btp163

<sub>103</sub> Dayhoff, M., Schwartz, R., & Orcutt, B. (1978). A model of evolutionary change in proteins.
<sub>104</sub> *Atlas of Protein Sequences and Structures*, 345–352.

<sub>105</sub> Frygelius, J., Arvestad, L., Wedell, A., & Töhönen, V. (2010). Evolution and human tissue
<sub>106</sub> expression of the cres/testatin subgroup genes, a reproductive tissue specific subgroup of
<sub>107</sub> the type 2 cystatins. *Evolution & Development*, *12*(3), 329–342. https://doi.org/10.1111/
<sub>108</sub> j.1525-142X.2010.00418.x

<sub>109</sub> Guindon, S., Dufayard, J.-F., Lefort, V., Anisimova, M., Hordijk, W., & Gascuel, O. (2010).
<sub>110</sub> New Algorithms and Methods to Estimate Maximum-Likelihood Phylogenies: Assessing
<sub>111</sub> the Performance of PhyML 3.0. *Systematic Biology*, *59*(3), 307–321. https://doi.org/10.
<sub>112</sub> 1093/sysbio/syq010

113 Holmes, I., & Rubin, G. M. (2002). An expectation maximization algorithm for training
114 hidden substitution models. *Journal of Molecular Biology*, *317*(5), 753–764. https:
115 //doi.org/10.1006/jmbi.2002.5405

116 Jones, D. T., Taylor, W. R., & Thornton, J. M. (1992). The rapid generation of mutation data
117 matrices from protein sequences. *Computer Applications in the Biosciences : CABIOS*,
118 *8*(3), 275–282. https://doi.org/10.1093/bioinformatics/8.3.275

119 Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K., Von Haeseler, A., & Jermiin, L. S. (2017).
120 ModelFinder: Fast model selection for accurate phylogenetic estimates. *Nature Methods*,
121 *14*(6), 587–589. https://doi.org/10.1038/nmeth.4285

122 Le, S. Q., & Gascuel, O. (2008). An improved general amino acid replacement matrix. *Molecular*
123 *Biology and Evolution*, *25*(7), 1307–1320. https://doi.org/10.1093/molbev/msn067

124 Minh, B. Q., Dang, C. C., Vinh, L. S., & Lanfear, R. (2021). QMaker: Fast and Accurate
125 Method to Estimate Empirical Models of Protein Evolution. *Systematic Biology*, *70*(5),
126 1046–1060. https://doi.org/10.1093/sysbio/syab010

127 Müller, T., & Vingron, M. (2000). Modeling amino acid replacement. *Journal of Computational*
128 *Biology*, *7*(6), 761–776. https://doi.org/10.1089/10665270050514918

129 Nguyen, L.-T., Schmidt, H. A., Haeseler, A. von, & Minh, B. Q. (2015). IQ-TREE: A fast and
130 effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular*
131 *Biology and Evolution*, *32*(1), 268–274. https://doi.org/10.1093/molbev/msu300

132 Oliphant, T. E. (2006). *Guide to NumPy*. Trelgol Publishing USA.

133 *PyPI: The python package index*. (n.d.). pypi.org.

134 Ronquist, F., Teslenko, M., Mark, P. van der, Ayres, D. L., Darling, A., Höhna, S., Larget, B.,
135 Liu, L., Suchard, M. A., & Huelsenbeck, J. P. (2012). MrBayes 3.2: Efficient Bayesian
136 Phylogenetic Inference and Model Choice Across a Large Model Space. *Systematic Biology*,
137 *61*(3), 539–542. https://doi.org/10.1093/sysbio/sys029

138 Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-
139 analysis of large phylogenies. *Bioinformatics*, *30*(9), 1312–1313. https://doi.org/10.
140 1093/bioinformatics/btu033

141 Whelan, S., & Goldman, N. (2001). A general empirical model of protein evolution derived
142 from multiple protein families using a maximum-likelihood approach. *Molecular Biology*
143 *and Evolution*, *18*(5), 691–699. https://doi.org/10.1093/oxfordjournals.molbev.a003851

144 Yang, Z. (2007). PAML 4: Phylogenetic Analysis by Maximum Likelihood. *Molecular Biology*
145 *and Evolution*, *24*(8), 1586–1591. https://doi.org/10.1093/molbev/msm088