

Available online at www.sciencedirect.com



Gene 365 (2006) 125-129



www.elsevier.com/locate/gene

Statistical properties of the methods for detecting positively selected amino acid sites

Yoshiyuki Suzuki*

Center for Information Biology and DNA Data Bank of Japan, National Institute of Genetics, 1111 Yata, Mishima-shi, Shizuoka-ken 411-8540, Japan

Received 28 April 2005; received in revised form 30 June 2005; accepted 7 September 2005 Available online 26 October 2005

Abstract

Parsimony and Bayesian methods have been developed for detecting positively selected amino acid sites. It has been reported that the parsimony method is generally conservative. In contrast, the Bayesian method is known to identify more positively selected sites than the parsimony method, especially when the number of sequences analyzed is small, although the interpretation of results obtained from the former method is controversial. Here I show that the likelihood-ratio test (LRT) of the Bayesian method corresponds to the parsimony method with window analysis, by analyzing the nucleotide sequences encoding the trans-activator (*tax*) gene of human T-cell lymphotropic virus type I (HTLV-I). It is also indicated that in the parsimony method, the test of selective neutrality using the binomial probability tends to be conservative, but the Monte Carlo simulation is useful for solving this problem. In addition, in the Bayesian method, the bootstrap method appears to produce similar results to the LRT. This information may be useful for improving the methods for detecting positively selected amino acid sites. © 2005 Elsevier B.V. All rights reserved.

Keywords: Positive selection; Parsimony; Bayesian; Human T-cell lymphotropic virus; Tax

1. Introduction

Positive selection operating at the amino acid sequence level can be inferred when the rate of nonsynonymous substitution (r_N) is higher than that of synonymous substitution (r_S) , or the nonsynonymous/synonymous rate ratio $(\omega = r_N/r_S)$ is greater than 1 for protein-coding nucleotide sequences (Hughes and Nei, 1988, 1989). In a protein molecule, different amino acid sites usually perform different biological functions, and therefore the type and strength of selection is expected to vary

* Tel.: +81 55 981 6847; fax: +81 55 981 6848.

E-mail address: yossuzuk@lab.nig.ac.jp.

with the site. Parsimony (Suzuki and Gojobori, 1999) and Bayesian (Yang et al., 2000) methods have been developed for detecting positively selected amino acid sites.

These methods utilize the phylogenetic tree of sequences analyzed. In the parsimony method, we compute the total numbers of synonymous (c_S) and nonsynonymous (c_N) substitutions as well as the average numbers of synonymous (s_S) and nonsynonymous (s_N) sites per codon over the phylogenetic tree for each codon site according to the maximum parsimony principle (Fitch, 1971; Hartigan, 1973). The null hypothesis of selective neutrality ($r_S = r_N$ or $\omega = 1$) is tested for each site by computing the probability (p) of obtaining the observed or more biased values for c_S and c_N , which are assumed to follow a binomial distribution with the probabilities of occurrence of synonymous and nonsynonymous substitutions given by $s_S/$ (s_S+s_N) and $s_N/(s_S+s_N)$, respectively. Positive selection is inferred when p < 0.05 and $c_N/s_N > c_S/s_S$.

In the Bayesian method, codon sites are usually classified into categories 0 and 1, which exist with proportions p_0 and p_1 (=1 - p_0), respectively. In category 0, ω (ω_0) follows a β distribution among codon sites with a range of $0 < \omega_0 < 1$, whereas in category 1, ω (ω_1) is constant for all sites with $\omega_1 > 1$.

Abbreviations: $r_{S(N)}$, rate of synonymous (nonsynonymous) substitution; [$c_{S(N)}$, $C_{S(N)}$, $c_{S(N)}$], number of synonymous (nonsynonymous) substitutions per codon; [$s_{S(N)}$, $S_{S(N)}$], number of synonymous (nonsynonymous) sites per codon; $p_{0(1)}$, proportion of category 0 (1); ω , nonsynonymous/synonymous rate ratio; $\omega_{0(1)}$, ω in category 0 (1); ln*L*, log-likelihood; LRT, likelihood-ratio test; BEB, Bayes empirical Bayes; NEB, naive empirical Bayes; *R*, transition/ transversion ratio; BP, bootstrap probability; HTLV-I, human T-cell lymphotropic virus type I; *tax*, trans-activator; HLA, human leukocyte antigen; INSD, International Nucleotide Sequence Database; *E*, exponential of the entropy for an amino acid site.

This model is called M8B (Massingham and Goldman, 2005). Free parameters are estimated by the maximum likelihood method. We then consider the case of $\omega_1 = 1$ in M8B, which is called M8A (Swanson et al., 2003). The likelihood-ratio test (LRT) with a 5% significance level is conducted to examine whether the log-likelihood (lnL) of M8B is greater than that of M8A under the assumption that twice the difference in $\ln L$ is asymptotically distributed as an equal mixture of a point mass on 0 and a χ^2 distribution with one degree of freedom (Wong et al., 2004). If the LRT is significant, it is concluded that positively selected sites exist in the sequence. The posterior probability of each codon site belonging to category 1 in M8B is computed by the Bayes empirical Bayes (BEB) method, where the errors in the parameter estimates are taken into account (Yang et al., 2005). The sites with the posterior probability of >95% are inferred as positively selected.

The sensitivity and selectivity of these methods have been examined using computer simulation and real data analysis. It has been reported that the parsimony method is generally conservative (Suzuki and Gojobori, 1999; Suzuki and Nei, 2001, 2002, 2004; Wong et al., 2004). In contrast, the Bayesian method is known to identify more positively selected sites than the parsimony method, especially when the number of sequences analyzed is small. However, the interpretation of results obtained from the Bayesian method is controversial. Anisimova et al. (2001, 2002, 2003) and Wong et al. (2004) have claimed that they are reliable, whereas Suzuki and Nei (2001, 2002) suggested that they sometimes contain many falsepositives and false-negatives according to the initial ω values used in the computation because of the existence of multiple local maxima of the likelihood function. The Bayesian method could also generate many false-positives when the topology of the phylogenetic tree used was unreliable (Suzuki and Nei, 2004), although in general, the effect of errors in the phylogenetic tree on the result appeared to be small for both the parsimony and Bayesian methods (Yang et al., 2000; Suzuki, 2004a).

The purpose of this paper is to further characterize the statistical properties of the parsimony and Bayesian methods. In particular, it is shown that the LRT of the Bayesian method corresponds to the parsimony method with window analysis. Additional theoretical and practical problems of these methods are discussed, and the new methods for solving them are proposed.

2. Materials and methods

2.1. Sequence data

Twenty nucleotide sequences with 181 codon sites of the trans-activator (*tax*) gene of human T-cell lymphotropic virus type I (HTLV-I), which were previously analyzed by Suzuki and Nei (2004) and Yang et al. (2005), were used to examine the statistical properties of the parsimony and Bayesian methods for detecting positively selected amino acid sites. The accession numbers for these sequences in the International Nucleotide Sequence Database (INSD) are given in Suzuki and Nei (2004).

2.2. Data analysis

The computer program CLUSTAL W (Thompson et al., 1994) was used to make a multiple alignment of tax gene sequences, which did not contain any gaps. The phylogenetic tree for these sequences was constructed by the neighborjoining method (Saitou and Nei, 1987) using the p-distance. However, since all the nucleotide substitutions observed were singletons and the sequence divergence was small for these sequences, as mentioned below, the phylogenetic tree could be obtained guite easily and was a star phylogeny (Suzuki and Nei, 2004). ADAPTSITE (version 1.3) (Suzuki et al., 2001) was used to conduct the parsimony method for detecting positively selected amino acid sites. The transition/transversion ratio (R)was estimated as the ratio of the transitional/transversional nucleotide diversity, which was 10.528. A one-tailed test was adopted for the test of selective neutrality. PAML (version 3.14) (Yang, 1997) was used to conduct the Bayesian method for detecting positively selected amino acid sites. Since multiple local maxima often exist for the likelihood function, the initial ω values of 0.4 and 3.14 were used in the computation. However, the results that were obtained were essentially the same for both cases. Therefore, only the results obtained using the initial ω value of 0.4 were presented.

3. Results and discussion

3.1. Statistical properties of the parsimony and Bayesian methods

When the parsimony method was used for analyzing the 20 *tax* gene sequences of HTLV-I, the $c_{\rm N}/c_{\rm S}$ values were estimated as 1/0 and 0/1 for the total of 21 and 2 codon sites, respectively. The remaining 158 codon sites were invariable. In the parsimony method, the test of selective neutrality is conducted for a given codon site by using the c_N and c_S values of that site only. However, since the sample size $(c_N + c_S)$ for the test was too small (0 or 1) for obtaining statistical significance, no site was inferred as positively selected. In contrast, in the Bayesian method, ω_1 =4.655 in M8B and the LRT was significant (lnL was -857.101 in M8B and -860.348 in M8A) (Table 1), suggesting that positively selected sites existed in the sequence. By computing the posterior probability, positive selection was inferred for 21 codon sites, where the $c_{\rm N}/c_{\rm S}$ values were estimated as 1/0 in the parsimony method. These results were consistent with Suzuki and Nei (2004) and Yang et al. (2005).

It should be noted, however, that each of the positively selected sites identified by the Bayesian method did not appear to contain enough information for inferring positive selection by itself, as mentioned above. It is therefore likely that positive selection was inferred for each of these sites using information from other sites. To show that this is indeed the case, *tax* gene sequences were analyzed using the Bayesian method by eliminating the positively selected sites cumulatively from the 3'-end of the sequence. When one site (position 181) was eliminated, two of the remaining 20 sites were not identified as

1	2	-
1	4	1

Results obtained non- the analysis of <i>tax</i> gene sequences using the Bayesian method and the parsimony method with window and	119515
Results obtained from the analysis of tar gene sequences using the Bayesian method and the parsimony method with window ana	lycie
Table 1	

Models Excluded ^a	Bayesian					Parsimony			
	M8B	M8B				M8A		Binomial ^e	Monte
	p_1	ω_1	Category 1 ^b	lnL	BP ^c	lnL	LRT ^d		Carlo ^f
None	1.000	4.655	2, 4, 39, 43, 53, 60, 62, 69, 81, 85, 92, 101, 108, 115, 146, 152, 154, 157, 161, 166, 181	-857.101	1.00	-860.348	S	0.008	0.006
181	1.000	4.443	2, 39, 43, 53, 60, 62, 69, 81, 85, 92, 101, 108, 115, 146, 152, 154, 157, 166	-845.960	1.00	-848.946	S	0.011	0.008
166, 181	1.000	4.232	2, 39, 43, 53, 60, 81, 85, 92, 108, 115, 152, 154, 157	-833.618	1.00	-836.346	S	0.014	0.011
161, 166, 181	1.000	4.050	2, 60, 81, 85, 92, 115, 154, 157	-823.992	1.00	-826.497	S	0.020	0.015
157, 161, 166, 181	1.000	3.809	None	-812.078	1.00	-814.305	S	0.027	0.021
154, 157, 161, 166, 181	1.000	3.598	None	-800.261	0.98	-802.246	S	0.036	0.028
152, 154, 157, 161, 166, 181	1.000	3.386	None	-790.229	0.99	-791.978	S	0.048	0.037
146, 152, 154, 157, 161, 166, 181	1.000	3.137	None	-779.526	0.97	-781.012	S	0.064	0.049
115, 146, 152, 154, 157, 161, 166, 181	1.000	2.885	None	-764.342	0.95	-765.572	NS	0.085	0.066

^a Only the results obtained by eliminating none to eight of 21 positively selected sites are indicated to save space. Similar results as observed for the case of eliminating eight sites were obtained when 9 to 21 sites were eliminated.

^b Positions inferred to be included in category 1 with the posterior probability of >95% by the BEB analysis of the Bayesian method.

^c Bootstrap probability that $p_1 > 0.000$ in M8B of the Bayesian method.

^d "S" indicates that the LRT was significant and "NS" indicates that the LRT was not significant in the Bayesian method.

^e The p values computed by the parsimony method with window analysis assuming the binomial distribution.

^f The *p* values computed by the parsimony method with window analysis using the Monte Carlo simulation.

positively selected (Table 1). The number of remaining sites, which were not identified as positively selected, increased to four, six, and ten when two, three, and four sites were eliminated, respectively. Eventually, none of the remaining sites were identified as positively selected when five or more sites were eliminated. These results were obtained because the posterior probability of each of the remaining sites belonging to category 1 became smaller than 95%, suggesting that positive selection was inferred for these sites using information from other sites (Suzuki and Nei, 2004). This apparently resulted from the assumption that the ω values followed a particular model (M8B) among codon sites, and also resulted in the observation that the Bayesian method generated more positively selected sites than the parsimony method. However, since the model assumed is arbitrary and thus unlikely to be realistic in the real data analysis, the inference of positively selected sites may be unreliable (Suzuki, 2004a; Kosakovsky Pond and Frost, 2005; Massingham and Goldman, 2005).

Similarly, the LRT appeared to be affected by the elimination of positively selected sites in the analysis of tax gene sequences (Table 1). The difference in lnL between M8B and M8A decreased as more sites were eliminated, and eventually the LRT became non-significant when eight or more sites were eliminated. These results indicate that each of the 21 positively selected sites contributed to increasing the difference in lnLbetween M8B and M8A, and positive selection was detected for these sites as a whole by the LRT. These observations suggest that the LRT in the Bayesian method corresponds to the window analysis, where a certain number of sites are grouped together to detect positive selection as a whole. In the ordinary window analysis, the codon sites included in a window are pre-

determined based on the linear sequence (Clark and Kao, 1991) or three-dimensional structure (Suzuki, 2004b) of proteins. In contrast, in the LRT, a virtual window (category 1) is constructed by grouping the sites that do not follow the β distribution of category 0. The exact sites included in category 1 usually remain unknown. However, in the case of the analysis of tax gene sequences, they could be identified because $p_1 = 1.000$, which suggested that all sites in the sequence were included in category 1 (Table 1). To show that the LRT in the Bayesian method was indeed similar to the window analysis, the test of selective neutrality was conducted using the parsimony method regarding category 1 in M8B of the Bayesian method as a window. $c_{\rm S}$ as well as $c_{\rm N}$ were summed over all sites in the sequence, and the results were denoted by $C_{\rm S}$ and $C_{\rm N}$, respectively. In addition, $s_{\rm S}$ as well as $s_{\rm N}$ were averaged over all sites, and the results were denoted by $S_{\rm S}$, and $S_{\rm N}$, respectively. The p value was computed as indicated above, by replacing $c_{\rm S}$, $c_{\rm N}$, $s_{\rm S}$, and $s_{\rm N}$ with $C_{\rm S}$, $C_{\rm N}$, $S_{\rm S}$, and $S_{\rm N}$, respectively. In eight cases where the LRT in the Bayesian method was significant, positive selection was detected for seven cases in the parsimony method with window analysis, suggesting that these methods roughly correspond to each other (Table 1). Note, however, that the former method generated one more positive result than the latter method. This is probably because the Bayesian method corrects for the multiple substitutions by assuming the codon substitution model (Goldman and Yang, 1994; Muse and Gaut, 1994) or because it tends to generate many false-positives (Massingham and Goldman, 2005). It is also possible that the parsimony method is conservative, because $C_{\rm S}$ and $C_{\rm N}$ may not always follow a binomial distribution, as discussed below.

3.2. Additional problems of the parsimony and Bayesian methods

There appeared to be additional theoretical and practical problems in the parsimony and Bayesian methods. In the parsimony method, $c_{\rm S}$ and $c_{\rm N}$ are assumed to follow a binomial distribution with the probabilities of occurrence of synonymous and nonsynonymous substitutions given by $s_S/(s_S+s_N)$ and s_N/s_N $(s_{\rm S}+s_{\rm N})$, respectively. However, this assumption may not always hold, especially when the numbers of synonymous and nonsynonymous sites are different among codons at the interior and exterior nodes of the phylogenetic tree. In this case, the results obtained may be biased. In fact, it has been observed that the false-positive rate of the parsimony method was sometimes smaller than expected, suggesting that this method was conservative (Suzuki and Gojobori, 1999; Wong et al., 2004). The extended binomial distribution may be useful for obtaining unbiased results (Kosakovsky Pond and Frost, 2005). However, the Monte Carlo simulation may also be useful for directly obtaining the *p* value for the test of selective neutrality, without assuming a binomial distribution. In the parsimony method, it is possible to identify each nucleotide substitution as well as a pair of codons which is associated with it over the phylogenetic tree for a given codon site. In the simulation, we generate a nucleotide substitution corresponding to each of the inferred substitutions. The probability that a generated substitution is synonymous and nonsynonymous is proportional to the numbers of synonymous and nonsynonymous sites, respectively, of the ancestral codon among the pair of codons which is associated with the inferred substitution. In this way, we obtain the total numbers of synonymous (c'_{s}) and nonsynonymous $(c'_{\rm N})$ substitutions over the phylogenetic tree, which are generated by a simulation (note that $c'_{\rm S} + c'_{\rm N} = c_{\rm S} + c_{\rm N}$). By replicating this process a sufficiently large number of times, we obtain the p value as a proportion of replications where $c'_{\rm S} \leq c_{\rm S}$ and $c'_{\rm N} \geq c_{\rm N}$. This method was applied to the parsimony method with window analysis for tax gene sequences, by regarding the nucleotide substitutions occurring at different codon sites as those occurring at a single codon site. It was observed that the p value obtained by the Monte Carlo simulation with 1000 replications was always smaller than that obtained by assuming a binomial distribution (Table 1). Interestingly, positive selection was inferred even when seven sites were eliminated from the sequence, supporting the idea that the LRT in the Bayesian method corresponds to the parsimony method with window analysis. These results also suggest that the test of selective neutrality using the binomial distribution tends to be conservative, but the Monte Carlo simulation is useful for solving this problem in the parsimony method as well as its derivatives, e.g., the three-dimensional window analysis (Suzuki, 2004b).

As for the Bayesian method, the computational algorithm in PAML as well as the Bayesian method itself has changed several times, and thus the earlier studies using this method should be re-examined. For example, Suzuki and Nei (2001) reported that many false-positives and false-negatives were obtained according to the initial ω values used in the

computation because of the existence of multiple local maxima of the likelihood function, by analyzing human leukocyte antigen (HLA) with version 3.0a of PAML. Wong et al. (2004) indicated that this version contained serious problems in the computational algorithms. They analyzed the same data set with version 3.13 and obtained different results, although they still observed multiple local maxima and thus the conclusions obtained by Suzuki and Nei (2001) remained essentially unchanged. In addition, it was originally recommended that the LRT was conducted between M3, where each codon site had one of three ω values, and M0, where all sites had the same ω value (Yang et al., 2000). The LRT was also conducted between M8, which was the same as M8B except that $\omega_1 > 1$ was not required, and M7, where the ω values of all sites followed a β distribution. The posterior probability of each codon site belonging to the positively selected category was computed by the naive empirical Bayes (NEB) method, where the parameter estimates were assumed to be correct. These approaches, however, were found to generate many false-positives (Anisimova et al., 2002; Suzuki and Nei, 2002, 2004; Swanson et al., 2003; Zhang, 2004; Massingham and Goldman, 2005; Yang et al., 2005). It is currently recommended that the LRT is conducted between M8B and M8A, and the posterior probability is computed by the BEB method. However, even these models have been chosen empirically and their theoretical basis has not been wellestablished. Note that Masatoshi Nei (personal communication) has indicated that the LRT is not really necessary if the bootstrap method is used for the test of the reliability of ω_1 values. Positive selection may be inferred if the bootstrap probability of observing $p_1 > 0.000$ in M8B is greater than 95%. This method was applied to the analysis of tax gene sequences with 100 resamplings. The results obtained from the bootstrap method were similar to those obtained from the LRT between M8B and M8A, suggesting that the former method is useful (Table 1).

Finally, it is assumed, in the parsimony method, that all amino acids are equally acceptable for each amino acid site of the protein (Suzuki and Gojobori, 1999). In addition, in the Bayesian method, the equilibrium frequencies of amino acids are assumed to be the same for all sites (Yang et al., 2000). In reality, however, different sites are surrounded by different environments in the three-dimensional structure of the protein. Therefore, under the structural constraint, the acceptable amino acids are expected to vary with the site (Bastolla et al., 1999). Furthermore, they are also expected to vary with the time for each site, according to the change in the interacting amino acids in the three-dimensional structure. In this case, it is likely that positive selection can operate only on the acceptable amino acids. Then, the number of nonsynonymous sites may be overestimated in the parsimony method and biased in the Bayesian method. To solve this problem, the exponential (E) of the entropy for each amino acid site, which indicates the average number of acceptable amino acids during evolution under the structurally constrained neutral (SCN) model (Porto et al., 2004), may be useful for correcting the number of nonsynonymous sites, by multiplying with E/20, in the parsimony method. In addition, the equilibrium frequencies of

amino acids for each site, which can also be estimated under this model, may be useful for obtaining unbiased results in the Bayesian method. Note that this approach is applicable to essentially all methods for examining $r_{\rm S}$ and $r_{\rm N}$.

Acknowledgements

The author thanks Masatoshi Nei for valuable suggestions and comments. I am also grateful to two anonymous reviewers for valuable comments. This work was supported by KAKENHI 17770007.

References

- Anisimova, M., Bielawski, J.P., Yang, Z., 2001. Accuracy and power of the likelihood ratio test in detecting adaptive molecular evolution. Mol. Biol. Evol. 18, 1585–1592.
- Anisimova, M., Bielawski, J.P., Yang, Z., 2002. Accuracy and power of Bayes prediction of amino acid sites under positive selection. Mol. Biol. Evol. 19, 950–958.
- Anisimova, M., Nielsen, R., Yang, Z., 2003. Effect of recombination on the accuracy of the likelihood method for detecting positive selection at amino acid sites. Genetics 164, 1229–1236.
- Bastolla, U., Vendruscolo, M., Roman, H.E., 1999. Neutral evolution of model proteins: diffusion in sequence space and overdispersion. J. Theor. Biol. 200, 49–64.
- Clark, A.G., Kao, T.-H., 1991. Excess nonsynonymous substitution at shared polymorphic sites among self-incompatibility alleles of Solanaceae. Proc. Natl. Acad. Sci. U. S. A. 88, 9823–9827.
- Fitch, W.M., 1971. Toward defining the course of evolution: minimum change for a specific tree topology. Syst. Zool. 20, 406–416.
- Goldman, N., Yang, Z., 1994. A codon-based model of nucleotide substitution for protein-coding DNA sequences. Mol. Biol. Evol. 11, 725–736.
- Hartigan, J.A., 1973. Minimum mutation fits to a given tree. Biometrics 29, 53–65.
- Hughes, A.L., Nei, M., 1988. Pattern of nucleotide substitution at major histocompatibility complex class I loci reveals overdominant selection. Nature 335, 167–170.
- Hughes, A.L., Nei, M., 1989. Nucleotide substitution at major histocompatibility complex class II loci: evidence for overdominant selection. Proc. Natl. Acad. Sci. U. S. A. 86, 958–962.
- Kosakovsky Pond, S.L., Frost, S.D.W., 2005. Not so different after all: a comparison of methods for detecting amino acid sites under selection. Mol. Biol. Evol. 22, 1208–1222.

- Massingham, T., Goldman, N., 2005. Detecting amino acid sites under positive selection and purifying selection. Genetics 169, 1753–1762.
- Muse, S.V., Gaut, S.B., 1994. A likelihood approach for comparing synonymous and nonsynonymous nucleotide substitution rates, with application to the chloroplast genome. Mol. Biol. Evol. 11, 715–724.
- Porto, M., Roman, H.E., Vendruscolo, M., Bastolla, U., 2004. Prediction of sitespecific amino acid distributions and limits of divergent evolutionary changes in protein sequences. Mol. Biol. Evol. 22, 630–638.
- Saitou, N., Nei, M., 1987. The neighbor-joining method: a new method for reconstructing phylogenetic trees. Mol. Biol. Evol. 4, 406–425.
- Suzuki, Y., 2004a. New methods for detecting positive selection at single amino acid sites. J. Mol. Evol. 59, 11–19.
- Suzuki, Y., 2004b. Three-dimensional window analysis for detecting positive selection at structural regions of proteins. Mol. Biol. Evol. 21, 2352–2359.
- Suzuki, Y., Gojobori, T., 1999. A method for detecting positive selection at single amino acid sites. Mol. Biol. Evol. 16, 1315–1328.
- Suzuki, Y., Nei, M., 2001. Reliabilities of parsimony-based and likelihoodbased methods for detecting positive selection at single amino acis sites. Mol. Biol. Evol. 18, 2179–2185.
- Suzuki, Y., Nei, M., 2002. Simulation study of the reliability and robustness of the statistical methods for detecting positive selection at single amino acid sites. Mol. Biol. Evol. 19, 1865–1869.
- Suzuki, Y., Nei, M., 2004. False-positive selection identified by ML-based methods: examples from the *Sig1* gene of the diatom *Thalassiosira weissflogii* and the *tax* gene of a human T-cell lymphotropic virus. Mol. Biol. Evol. 21, 914–921.
- Suzuki, Y., Gojobori, T., Nei, M., 2001. ADAPTSITE: detecting natural selection at single amino acid sites. Bioinformatics 17, 660–661.
- Swanson, W.J., Nielsen, R., Yang, Q., 2003. Pervasive adaptive evolution in mammalian fertilization proteins. Mol. Biol. Evol. 20, 18–20.
- Thompson, J.D., Higgins, D.G., Gibson, T.J., 1994. CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. Nucleic Acids Res. 22, 4673–4680.
- Wong, W.S.W., Yang, Z., Goldman, N., Nielsen, R., 2004. Accuracy and power of statistical methods for detecting adaptive evolution in protein coding sequences and for identifying positively selected sites. Genetics 168, 1041–1051.
- Yang, Z., 1997. PAML: a program package for phylogenetic analysis by maximum likelihood. Comput. Appl. Biosci. 13, 555–556.
- Yang, Z., Nielsen, R., Goldman, N., Pedersen, A.M., 2000. Codon-substitution models for heterogeneous selection pressure at amino acid sites. Genetics 155, 431–449.
- Yang, Z., Wong, W.S.W., Nielsen, R., 2005. Bayes empirical Bayes inference of amino acid sites under positive selection. Mol. Biol. Evol. 22, 1107–1118.
- Zhang, J., 2004. Frequent false detection of positive selection by the likelihood method with branch-site models. Mol. Biol. Evol. 21, 1332–1339.