

A simple test to uncover signals of CpG hypermutability using posterior predictive simulations

Simon Laurin-Lemay¹ and Nicolas Rodrigue^{1,2,3}

¹Department of Biology, Carleton University, Ottawa, Canada

²Institute of Biochemistry, Carleton University, Ottawa, Canada

³School of Mathematics and Statistics, Carleton University,
Ottawa, Canada

Keywords

DNA sequence evolution, Markov chain Monte Carlo, phylogenetics, misspecification, confounding factors

Declarations

Conflict of interest: The authors declare that they have no competing interests.

Correspondence

Simon Laurin-Lemay
209 Nesbitt Biology Building,
1125 Colonel By Drive Ottawa,
Ontario, CANADA
K1A 0C6
evol.simon@gmail.com

Abstract

Objective

CpG hypermutability is a process known to impact vertebrate evolution and is caused by spontaneous deamination of methylated cytosines within CpG contexts. CpG hypermutability has been shown to confound the detection of negative and positive selection on codon usage and amino acids. In this work, we propose a simple test based on the use of posterior predictive sampling to detect the presence of CpG hypermutability using the GTR+G substitution model.

Results

First, we validated the CpG test using simulations. We recovered between 0-8% of false positives and no false negatives. The simulations were performed using the jump-chain algorithm with a simple codon substitution model for more realistic simulation conditions and to allow for different degree of CpG hypermutability. All mammalian genes tested were positive with the posterior predictive test we developed to detect CpG hypermutability, which is consistent with previous results obtained using complex mechanistic phylogenetic codon substitution models designed to measure the level of CpG transition rate. This test could easily accompany any study aimed at detecting negative or positive selection, even for large-scale analyses, knowing that methylation of cytosines exists in all domains of life (i.e., bacteria, archaeobacteria, and eukaryota).

1 Introduction

2 The phylogenetic method involves first the generation of data sets and then the
3 analysis of those data sets with various substitution models. This is followed
4 by a cycle of work in which molecular evolutionists compare measurements of
5 population processes (i.e. mutation, selection, even demographic processes)
6 included in the definitions of their models (e.g. ??). One of the things most
7 feared by molecular evolutionists is the discovery of confounding effects (e.g., ??)
8 which could invalidate the proposed mechanistic hypothesis testing approach
9 to detect and measure levels of specific evolutionary processes. Confounders
10 are evolutionary processes not accounted for by the substitution model. In
11 the Bayesian framework, the presence of confounders will incorrectly affect the
12 posterior values of one or more parameters of the substitution model, which was
13 designed to parameterize another population process. Finally, false positives can
14 result from the presence of confounders. For example, the hypermutability of
15 CpGs due to the process of spontaneous deamination of methylated cytosines
16 in CpG contexts may confound the detection of negative or positive selection
17 on synonymous and non-synonymous mutations (e.g., ??).

18 Once a confounding process has been identified, researchers should go back
19 to the original data sets and try to control for that aspect, if possible, by re-
20 moving the sites affected by the confounding process: (e.g., ?). On the other
21 hand, managing confounders by selecting sites for analysis may be particularly
22 difficult or even impossible when dealing with confounders that affect substi-
23 tion processes with dependencies between sites, such as the hypermutability
24 of CpGs. Ultimately, the researcher should propose ways to incorporate the
25 confounding process into a new substitution model, so that the process can be
26 conditioned jointly with other evolutionary processes, or in a two-step procedure
27 (e.g., topology is usually not co-conditioned with other mutation and selection

28 parameters in the field of molecular evolution).

29 Most substitution models assume that sites evolve independently, which is
30 particularly convenient for calculating the likelihood of an alignment; we only
31 need to multiply all the site likelihoods calculated from each site of the align-
32 ment (?). On the other hand, taking into account substitution processes in-
33 volving dependencies between sites (e.g., ???), as required to parameterize the
34 hypermutability of CpGs, is not a trivial task. Dealing with such substitution
35 processes makes the development of mathematics and code more complex and
36 limited to a small circle of researchers. This has led researchers to develop hy-
37 brid methods that integrate standard posterior sampling methods (e.g., Markov
38 Chain Monte Carlo) with so-called simulation-based methods (e.g., ?) or to
39 develop strategies that rely on data augmentation (e.g., ??).

40 Posterior predictive sampling (??) is a central tool in the Bayesian toolbox.
41 This tool is of primary use when validating new phylogenetic models (e.g., ??).
42 For example, knowing that CpGs are depleted by the high rate of spontaneous
43 deamination due to methylation of most vertebrate cytosines within CpG con-
44 texts (??), we are interested in the ability of substitution models to predict
45 dinucleotide (e.g. CpG) frequencies of mammalian protein-coding genes.

46 Following the textbook case of mammalian CpG hypermutability, we will
47 test our ability to detect the presence of CpG hypermutability using one of the
48 most widely used substitution models, the GTR+G substitution model, using
49 Bayesian posterior predictive sampling. By definition, we expect the GTR+G
50 substitution model to predict more CpG dinucleotides than would be observed
51 in real alignments because CpG hypermutability is not accounted for in the
52 model definition. Similarly, the GTR+G substitution model will allow stop
53 codons to be generated within predictive alignments. We first validated the
54 posterior predictive test using simulations generated without and with CpG hy-

55 permutability, using the GTR+G substitution model and a codon substitution
 56 model implemented in the mutation selection framework. The codon substi-
 57 tution model used uses a GTR parameterization for the background mutation
 58 process plus a λ parameter for CpG hypermutability and a parameter to cap-
 59 ture global negative selection on amino acids. We then apply the test to 137
 60 mammalian protein-coding genes for which we have already measured the de-
 61 gree of CpG hypermutability using a complex codon substitution model (?).
 62 Finally, we discuss the ranking of other dinucleotides in terms of their ability
 63 to be predicted by the GTR+G substitution model using the same approach
 64 developed for testing for the presence of CpG hypermutability.

65 **Materials and methods**

66 **Detecting CpG hypermutability using posterior predictive** 67 **sampling**

68 The CpG test consists of calculating the proportion of times, i.e., *p-value*, that
 69 the frequency of CpGs calculated from the real alignments is lower than the
 70 CpG frequency recovered from the predictive alignments, where the predictive
 71 alignments are generated from a set of parameter values, i , sampled from the
 72 posterior. The following equation details the CpG test:

$$p\text{-value} = \frac{1}{N} \sum_{i=1}^N 1(freq_{CpG}^{real} < freq_{CpG}^{pred_i}), \quad (1)$$

73 where N is the predictive sample size and 1 is an indicator function that
 74 returns 1 if the CpG frequency calculated from the true alignment is less than
 75 the predictive CpG frequency and 0 otherwise.

76 Real data and tree topology

77 We retrieved the 137 mammalian codon alignments from ? along with the tree
78 topology.

79 Simulation study

80 We first analyzed 10 mammalian protein-coding genes selected for their wide
81 range of GC content as used to validate the new approach developed in ? us-
82 ing the GTR+G substitution model as well as the M0GTR codon substitution
83 model, but with ω fixed at 1 (?) using Phylobayes-MPI (?) ** j'ai fait une pe-
84 tite bourde ici, en laissant omega fixed, mais c'est juste pour avoir des valeurs
85 de paramètres, donc ça ne dérange pas **. We generated 100 predictive align-
86 ments (10 genes \times 10 samples) using AliSim (?) for the GTR+G substitution
87 model, which will later be analyzed with the same model to generate a new set
88 of predictive alignments on which to apply the CpG test. We also generated 600
89 predictive alignments (10 genes \times 3 values of $\lambda \times$ 2 values of $\omega \times$ 10 posterior
90 predictive samples) under the M0GTR substitution model using the jump chain
91 simulation algorithm described in ? to account for CpG hypermutability. All
92 simulations were performed with a single ω value (i.e., 0.2 or 1). To evaluate
93 the false positive and false negative rates of the test, we generated predictive
94 alignments using $\lambda = 1$ and $\lambda = \{4.8\}$, respectively.

95 Each simulated predictive alignment was then analyzed with the GTR+G
96 substitution model using Phylobayes-MPI (?) software, we generate posterior
97 predictive alignments, 50 for each analysis, using the same software. The predic-
98 tive alignments were generated by sampling 50 sets of parameter values from the
99 posterior of each of the GTR+G analyses performed on synthetic alignments.
100 Each of the experimental conditions (3 values of $\lambda \times$ 2 values of ω) is replicated
101 100 times. The proportion of false positives and false negatives in each set of ex-

102 perimental conditions was then calculated by applying the posterior predictive
103 tests.

104 **Empirical study**

105 We analyzed the 137 mammalian codon sequence alignments with the GTR+G
106 substitution model implemented in Phylobayes-MPI: (?). We sample 100 pos-
107 terior predictive alignments, also using Phylobayes-MPI, for each mammalian
108 gene of interest. We then apply the posterior predictive test by calculating the
109 frequency of CpGs from the true and predictive alignments to detect the pres-
110 ence of CpG hypermutability. Similarly, for each dinucleotide context, we could
111 apply the specific predictive test as designed for the CpG context and calculate
112 the proportion of positive tests to identify dinucleotide contexts missed by the
113 GTR+G substitution model.

114 **Results**

115 **Simulation study**

116 Less than 5% of the CpG tests were significant when performed on synthetic
117 sequence alignments generated under the GTR+G substitution model, i.e., 2%
118 (Table 1). Similarly, none of the CpG tests were significant when performed
119 on synthetic sequence alignments generated without CpG hypermutability and
120 with $\omega = 0.2$ using the M0GTR codon substitution model (Table 1). On the
121 other hand, under less realistic conditions, without negative selection on amino
122 acids, $\omega = 1$, we recovered more false positives, i.e., 8% (Table 1), but still close
123 to the alpha significance threshold of 5%. No false negatives were generated in
124 the presence of CpG hypermutability, (i.e., Table 1: $\lambda = \{4, 8\}$), all CpG tests
125 were positive, for both of the omega values used, i.e., $\omega = \{0.2, 1\}$.

126 Empirical study

127 Therefore, after validating the CpG test, we analyzed 137 mammalian protein-
128 coding genes with the GTR+G substitution model and found 100% of the tested
129 genes to be significant ($p\text{-value}=0$, $\alpha=5\%$), rejecting the null hypothesis that
130 the GTR+G substitution model accounts for CpG hypermutability. In other
131 words, CpG frequencies are significantly reduced in all real alignments compared
132 to frequencies recovered from predicted alignments generated by the GTR+G
133 substitution model. As an example, in Figure 1, panel A, we show the dis-
134 crepancy between the CpG frequencies predicted by the GTR+G substitution
135 model and the observed CpG frequency in one of the examined coding sequence
136 alignments. In panel B of the same figure, we see that the predicted values of
137 ApT are distributed on either side of the actual value calculated from the same
138 gene examined in panel A.

139 Next, we applied the predictive test to all other dinucleotide contexts using
140 the same substitution model, i.e., GTR+G. After CpG, TpA, and GpT contexts,
141 it was found that the GTR+G substitution model significantly missed most of
142 the tested genes (94.9% and 81%, respectively; Table S1), considering an alpha
143 significance threshold of 5%.

144 Discussion

145 In this work, we designed a simple test, based on posterior predictive sampling, a
146 tool from the Bayesian toolbox, to detect the presence of CpG hypermutability.
147 We validated the test using a simulation study: with approximately the alpha
148 significance threshold of false positives and no false negatives. We then applied
149 the test to a set of mammalian protein-coding genes, for which we had already
150 obtained measurements of CpG hypermutability using a complex codon sub-

stitution model implemented in the mechanistic mutation-selection framework
(?).

We also showed that contexts other than CpG were systematically missed by the predictive samples generated by the GTR+G substitution model, namely TpA with 94.9% of genes detected as positive (Table S1) and GpT with 81% of positive tests (Table S1). The TpA context may be artifactually overpredicted by the GTR+G substitution model because it does not take into account the structure of the genetic code and allows the presence of stop codons in predictive protein-coding sequence alignments. Two out of three stop codons contain a TpA dinucleotide (i.e., TAA and TAG). The use of a codon substitution model implemented in the mechanistic mutation-selection framework, such as the one used here to validate the CpG test, i.e., M0GTR, should at least help to reduce the over-prediction of TpA contexts by disallowing the presence of stop codons within the predictive alignments generated when running the test, as the substitution model imposes infinite negative selection against mutations that land on stop codons. On the other hand, the etiology of TpA context depletion is highly controversial and may not be exclusively due to mutation or selection processes (e.g., ??). Interestingly, GpT contexts have been shown to be hypermutable on shorter evolutionary timescales, i.e. within human populations (?), suggesting that the evolutionary process behind the depletion of GpT context is conserved over the mammalian tree. On the other hand, considering contexts other than CpG should also require extending the scope of the simulation study, since predicting the effect of one mutation process on another mutation process is particularly complicated. For example, CpG and TpA contexts share the same transition products (TpG and CpA).

However, knowing that CpG hypermutability could affect the measurement of negative and positive selection on synonymous and nonsynonymous mutations

178 (??), the test could be systematically applied when conducting large-scale evo-
179 lutionary studies (e.g., ???) or in future projects such as the Zoonomia Project
180 (?) to promote the development of new mechanistic tests using codon substitu-
181 tion models that take into account a wider range of population processes (e.g.,
182 ??????). We could also investigate the presence of hypermutability patterns
183 from substitution mappings, similar to what has been developed by (??).

184 Limitations

185 The main limitation of this study is that we have not yet investigated how
186 a non-mechanistic test, such as the one developed here based on comparison
187 of posterior predictive alignments, might be affected by the presence of other
188 population processes (i.e. mutation, selection, and demographic processes) not
189 accounted for by the GTR+G substitution model, which could be addressed by
190 simulation studies. Because the CpG test we developed here is non-mechanistic,
191 it also does not allow quantification of the degree of CpG hypermutability, as a
192 proper mechanistic test could do (e.g., ?).

193 Availability of data and materials

194 Workflow and data is available.

195 Abbreviations

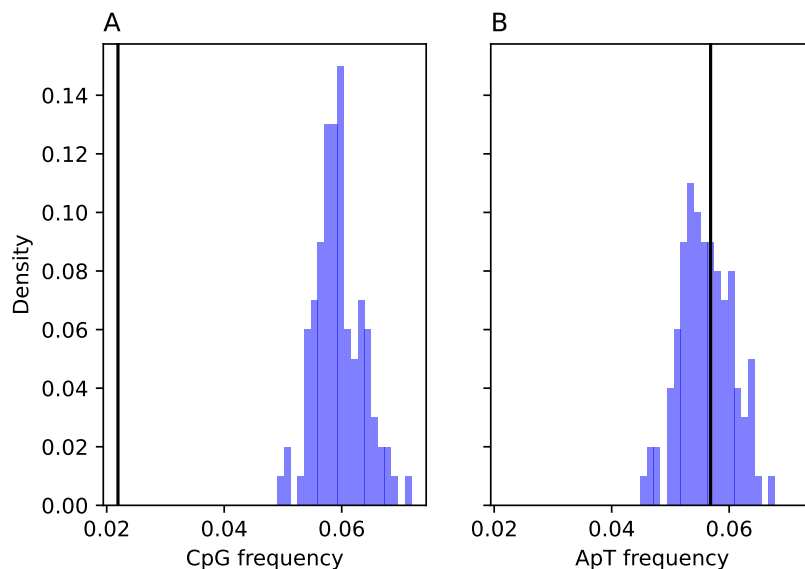


Figure 1: Comparison of CpG (A) and ApT (B) frequencies computed from the *MEP1A* mammalian protein-coding gene alignment, black vertical lines, and corresponding frequencies computed from posterior predictive alignments generated using the GTR+G substitution model, blue histograms. Note that the predicted CpG frequencies (A) overestimate the true value 100% of the time, while the true ApT frequency (B) falls within the distribution of predicted frequencies.

Tables

type of controls	models used to generate the synthetic alignments	positive tests (%)
negative	GTR+G	0
negative	M0GTR+ $(\lambda = 1)+(\omega = 0.2)$	0
negative	M0GTR+ $(\lambda = 1)+(\omega = 1.0)$	8
positive	M0GTR+ $(\lambda = 4)+(\omega = 0.2)$	100
positive	M0GTR+ $(\lambda = 4)+(\omega = 1.0)$	100
positive	M0GTR+ $(\lambda = 8)+(\omega = 0.2)$	100
positive	M0GTR+ $(\lambda = 8)+(\omega = 1.0)$	100

Table 1: Validation of CpG test using posterior predictive sampling under GTR+G substitution model with an α threshold of significance 5%. Synthetic alignments were generated using GTR+G and M0GTR codon substitution model using three CpG hypermutabilities values, $\lambda = \{1, 4, 8\}$ and a global selection on amino acids using two ω values (i.e., 0.2 and 1).

198 **References**

199 **Acknowledgments**

200 This work was funded by the Natural Sciences and Engineering Research Council
201 of Canada.

202 **Funding**

203 **Author information**

204 **Authors and Affiliations**

205 **Contributions**

206 Conceptualization: SLL, NR. Formal analysis: SLL Investigation: SLL, NR.
207 Methodology: SLL, NR. Resources: SLL, NR. Supervision: NR Writing original
208 draft: SLL, NR. Review and editing: SLL, NR.
209 All authors read and approved the final manuscript.

210 **Corresponding author**

211 Correspondance to Simon Laurin-Lemay

212 **Ethics declarations**

213 **Ethics approval and consent to participate**

214 Not applicable.

215 **Consent for publication**

216 Not applicable.

217 **Competing interests**

218 The authors declare that they have no competing interests.

219 **Additional information**

220 **Supplementary information**

221 Table S1: Proportion of positive tests calculated for each dinucleotide context
222 obtained using GTR+G substitution model over the 137 mammalian protein-
223 coding genes studied.