



Coalescent-based DNA barcoding: multilocus analysis and robustness

Olivier David, Catherine Larédo, Raphael Leblois, Brigitte Schaeffer, Nicolas Vergne

► To cite this version:

Olivier David, Catherine Larédo, Raphael Leblois, Brigitte Schaeffer, Nicolas Vergne. Coalescent-based DNA barcoding: multilocus analysis and robustness. *Journal of Computational Biology*, Mary Ann Liebert, 2012, 10.1089/cmb.2011.0122 . hal-02337205

HAL Id: hal-02337205

<https://hal-normandie-univ.archives-ouvertes.fr/hal-02337205>

Submitted on 29 Oct 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Coalescent-based DNA barcoding: multilocus analysis and robustness

Olivier David^{1,*} Catherine Larédo^{1,3} Raphael Leblois^{2,4}

Brigitte Schaeffer¹ Nicolas Vergne^{1,5}

¹UR341, Mathématiques et informatique appliquées, INRA, F-78350 Jouy-en-Josas,
France

²Muséum National d'Histoire Naturelle, UMR 5202 MNHN/CNRS, Laboratoire Origine
Structure Evolution de la Biodiversité (OSEB), Paris, France

³Laboratoire de Probabilités et Modèles Aléatoires, Universités Paris 6 et 7, UMR CNRS
7599, Paris, France

⁴Centre de Biologie et de Gestion des Populations (CBGP), UMR INRA-IRD-CIRAD
1062, Montferrier-sur-Lez, France

⁵Laboratoire de Mathématiques Raphaël Salem, UMR 6085 CNRS-Université de Rouen,
76801 Saint-Etienne-du-Rouvray, France

*Corresponding author

Abstract

DNA barcoding is the assignment of individuals to species using standardized mitochondrial sequences. Nuclear data are sometimes added to the mitochondrial data to increase the power of barcoding. A barcoding method for combining mitochondrial and nuclear informations that is based on the coalescent model is developed. This method is assessed using simulated and real data. It is found that adding nuclear data can reduce the number of ambiguous assignments. Finally, the robustness of coalescent-based barcoding to departures from model assumptions is studied using simulations. This method is found to be robust to past population size variations, to within-species population structures and to designs that poorly sample populations within species.

Key words: Bayesian inference, classification, coalescent, DNA barcoding, species assignment.

1 Introduction

DNA barcoding is the assignment of individuals to species or higher taxonomic levels using standardized genetic data observed on the target individuals and samples from each species (Frézal and Leblois, 2008; Valentini et al., 2009). The DNA barcode project is conceived as a standard system for fast and accurate identification of all eukaryotic species (Hebert et al., 2003; Miller, 2007). The DNA barcode itself consists of a 648 bp region of the cytochrome c oxidase 1 (COI) gene. Additionally to the mitochondrial COI gene, nuclear loci are sometimes also considered to improve assignment performance (Austerlitz et al., 2009; Elias et al., 2007).

DNA barcoding is a classification problem rather than a clustering one since the classes (species) are predefined and do not have to be inferred from the data (but see Pons et al. (2006) for an application of clustering to barcoding). Barcoding assignment methods can be divided into similarity methods based on the match between the query sequence and the reference sequences such as BLAST search, phylogenetic approaches (Hebert et al., 2003; Elias et al., 2007), classification algorithms with no underlying biological models such as the nearest-neighbour method and methods based on population genetics (Matz and Nielsen, 2005). Two Bayesian methods based on models have recently been developed. In the method of Munch et al. (2008a), species are assumed to evolve according to a phylogenetic model while the within-species variation is not modelled. Conversely, in TheAssigner, the method of Abdo and Golding (2007), species are assumed to evolve independently and the dependence between sequences within species is modelled using a classical population genetics model called the coalescent. The latter is a model for the genealogical tree of a random sample of genes drawn from a large panmictic population (Chapter 10, Ewens, 2004; Kingman, 1982a,b; Tajima, 1983). Model-based barcoding methods raise various issues. Current methods assume that the data are mitochondrial and cannot cope with nuclear data. Moreover, their robustness to departures from model assumptions has been

little studied.

The main objective of the present paper is to study how to take account of nuclear data in coalescent-based classification and to study the robustness of this type of classification to departures from model assumptions. First a coalescent-based classification for assigning individuals to species using mitochondrial data is developed (Section 3). Then this method is extended to take account of nuclear data (Section 4). Finally the performance and robustness of coalescent-based classification are studied using simulated and real data sets (Sections 5 and 6).

2 Bayesian classification

First we briefly review some basic material on Bayesian classification. In this method, individuals are assumed to belong to c classes. A data set y is available that includes measurements observed on reference individuals whose class is known. The objective is to predict the class $z \in \{1, \dots, c\}$ of a test individual given its data x and the reference data y .

In Bayesian classification, a test individual is assigned to the class with the largest posterior probability of membership (Abdo and Golding, 2007; Munch et al., 2008a; Ripley, 1996). The assignment may be considered as ambiguous if the latter probability does not exceed some specified threshold. According to Bayes theorem, the posterior probability that a test individual belongs to class i is equal to $P(z = i|y, x) = P(z = i, x|y)/P(x|y) = r_i / \sum_k r_k$, where:

$$r_i = P(z = i|y)P(x|y, z = i). \quad (1)$$

In this equation, $P(z = i|y)$ is the probability that the test individual belongs to class i given the reference data y prior to the knowledge of x and plays the role of a prior probability of membership. The probability $P(x|y, z = i)$ is the conditional probability that an individual sampled in class i has data x . Bayesian classification is optimal for the 0–1 loss function (Chapter 2, Ripley, 1996) and provides a measure of assignment confidence.

3 Species assignment with mitochondrial data

We now apply Bayesian classification to DNA barcoding. In this section, the data consist of mitochondrial DNA sequences. The assumed demographic model is a set of isolated and panmictic species with a common ancestry at a given time in the past (i.e., the divergence time). This demographic model is the same as the one of Abdo and Golding (2007). The mitochondrial locus is assumed to evolve according to the coalescent model within each species independently. Following the standard coalescent, it is assumed that species sizes do not vary over time, that there is no migration between species and that all alleles are neutral. All individuals are assumed to be sampled at the same time and the species of any test individual is assumed to be represented in the reference data y . In this model, mutations occur on each ancestral lineage of species i according to a Poisson process with parameter $\theta_i/2$. The assumed mutation model is the infinitely many-sites model (ISM), in which a gene is considered as an infinitely long DNA sequence and each new mutant site is sampled uniformly and independently along the sequence (Chapter 9, Ewens, 2004). Finally, it is assumed that at each site it is known which base is the mutant base or the ancestral base (Section 7) and that there are no missing data or errors in the data.

The mutation parameters θ_i are first assumed to be known. Then, under the assumption that species evolve independently, the probability $P(x|y, z = i)$ in (1) is equal to:

$$P(x|y, z = i) = P(x|y_i, z = i),$$

where y_i denotes the data of species i in the reference data base. This probability will be written for simplicity as $P(x|y_i)$ in what follows. Generally it cannot be calculated explicitly under the ISM but it can be estimated as follows. It is equal to (p. 420, De Iorio

and Griffiths, 2004, supplementary materials A):

$$P(x|y_i) = \frac{n_i(x) + 1}{n_i + 1} \frac{P^0(x, y_i)}{P^0(y_i)}, \quad (2)$$

where P^0 is the probability of an unordered sample, n_i is the number of genes in the sample of species i and $n_i(x)$ is the number of genes with sequence x in the sample of species i . The probabilities $P^0(x, y_i)$ and $P^0(y_i)$ can be estimated using importance sampling (IS) (De Iorio and Griffiths, 2004, supplementary materials A). Note that the probability $P^0(y_i)$ needs to be estimated only once if there are several individuals to assign.

Mutation processes are generally unknown for most species and the vector θ of mutation parameters is thus usually not known. In this case, the posterior probabilities of membership can be estimated by plug-in, that is by assuming that θ is known and equal to an estimate $\hat{\theta}$, computed from the reference data set as in Abdo and Golding (2007). The vector θ may be estimated, for example, using the method of Watterson (1975), by coalescent-based maximum likelihood or Bayesian methods (Bahlo and Griffiths, 2000; Kuhner et al., 1995). Alternatively, a predictive approach can be used in which the dependence of probabilities on θ is removed by integration (Chapter 2, Ripley, 1996, supplementary materials A).

4 Species assignment with mitochondrial and nuclear data

Individuals are now assumed to be genotyped at l diploid nuclear loci in addition to the mitochondrial locus. The two sequences of an individual at a nuclear locus are assumed to be known (Section 7). The genetic data of a test individual are denoted by $x = (x_0, \dots, x_l)$, where x_0 is the mitochondrial sequence and x_j ($j \geq 1$) is the pair of sequences at nuclear locus j . Each locus is assumed, as in the previous section, to evolve according to the coalescent model within each species independently. All the loci are assumed to evolve independently (Hudson, 1991; Nordborg, 2001) and there is no recombination within a locus. Mutations are assumed to occur according to the ISM with parameter $\theta_{ij}/2$ for species i and locus j . Finally, for simplicity, all parameters θ_{ij} are assumed to be known in this section.

With independent loci, the quantity r_i in (1) is equal to (Chapter 8, Ripley, 1996):

$$r_i = P(z = i|y) \prod_{j=0}^l P(x_j|y_{ij}, z = i), \quad (3)$$

where y_{ij} is the reference data for species i and locus j . This equation allows us to easily combine the mitochondrial and the nuclear informations. For a nuclear locus in a diploid species, (2) becomes:

$$P(x_j|y_{ij}, z = i) = \frac{(n_{ij}(x_{j1}) + 1)(n_{ij}(x_{j2}) + 1 + \delta_j)}{(n_{ij} + 2)(n_{ij} + 1)} \frac{P^0(x_j, y_{ij})}{P^0(y_{ij})},$$

where $\delta_j = 0$ if the test individual is heterozygote at locus j and $\delta_j = 1$ if the test individual is homozygote at locus j . In this equation, x_{j1} and x_{j2} denote the two test sequences at locus j , n_{ij} denotes the number of genes sampled for species i at locus j , $n_{ij}(x_{j1})$ denotes the multiplicity of allele x_{j1} in the sample of species i and locus j . The probabilities $P^0(x_j, y_{ij})$

and $P^0(y_{ij})$ can then be estimated using IS on coalescent histories as before.

5 Simulation study

Simulations were carried out to assess the methods described above. In these simulations, one ancestral species split T generations ago into two new species with effective size N_e and mutation parameter θ . There were n reference individuals in each species. First sequences were simulated for a mitochondrial locus and diploid nuclear locus to study the effect of adding nuclear data. Then to test the robustness of the methods developed, mitochondrial sequences were simulated assuming that species size varied over time or that each species was divided into several populations exchanging migrants. To mimic extreme sampling strategies that can be done in structured populations, we considered an “extended” sampling, in which the reference individuals were sampled in all populations for each new species, and a “clustered” sampling, in which all reference individuals were sampled from a single population in each new species. Details on these simulations are presented in supplementary materials B.

The simulated data were analysed with the nearest-neighbour classification (1NN) and the developed Bayesian assigner (BA) (supplementary materials B). The 1NN method was used because it had been found to be efficient compared with other barcoding methods (Austerlitz et al., 2009) and it was expected to be robust since it was not based on a specific biological model. This method was implemented with bagging in order to obtain a measure of confidence for an assignment (Hastie et al., 2001, supplementary materials B). Assignment performance was quantified using sensitivity and specificity (Munch et al., 2008a,b). Specificity is the fraction of non-ambiguous assignments (Section 2) that are correct. Sensitivity is the fraction of all the assignments that are correct.

The simulations with nuclear data first showed that performances were the best for the combination of the mitochondrial and the nuclear data, intermediate for the mitochondrial data (Fig. 1), and the least good for the nuclear data alone (Fig. S2 and S3). The poor results for the nuclear data alone were probably due to the larger effective size we used for

the nuclear locus, leading to smaller scaled divergence times T/N_e and thus lower levels of differentiation between the two new species. Nevertheless adding nuclear data clearly increased sensitivity (Fig. 1). This was mainly due to a reduction of the number of ambiguous assignments since specificity did not increase much (Fig. 1). Our simulations also showed that 1NN and BA had similar performances, except for the nuclear data alone for which 1NN had a low sensitivity (Fig. S3). However, we can note that BA had more ambiguous assignments than 1NN but made fewer errors among the non-ambiguous assignments (Fig. 1). Another important result was that the estimation of mutation parameters did not change the BA performance much (Fig. S1). Finally and as it was expected, increasing the values of θ , T or n improved the performance of both methods as in Austerlitz et al. (2009).

For past population size variations, the main results were that past expansions strongly increased specificity, sensitivity and the rate of non-ambiguous assignments, whereas past contractions had the opposite effect of decreasing specificity and sensitivity (Fig. 2). Our simulations also showed that past expansions affected both methods similarly, but 1NN always showed a slightly better performance than BA. On the contrary, it is interesting to note that the effect of past contractions was more pronounced for 1NN than for BA, resulting in much better performances for BA. Finally, the effect of past population size variations was found to be important for all the growth rate values we used and to be stronger for expansions than for contractions.

The effect of population structure was more complex because it depended on the sampling strategy. Compared with the unstructured species results, a population structure with a weak migration mainly affected sensitivity and the rate of non-ambiguous assignments, that both increased for the “clustered” samples and decreased for the “extended” samples (Fig. 2). This result was unexpected as the population of origin of a test individual was represented by two individuals in the reference samples for the “extended” samples but not for the “clustered” samples. Finally, we note that population structure affected both

methods similarly and that the effect of population structure became noticeable only when migration was weak enough.

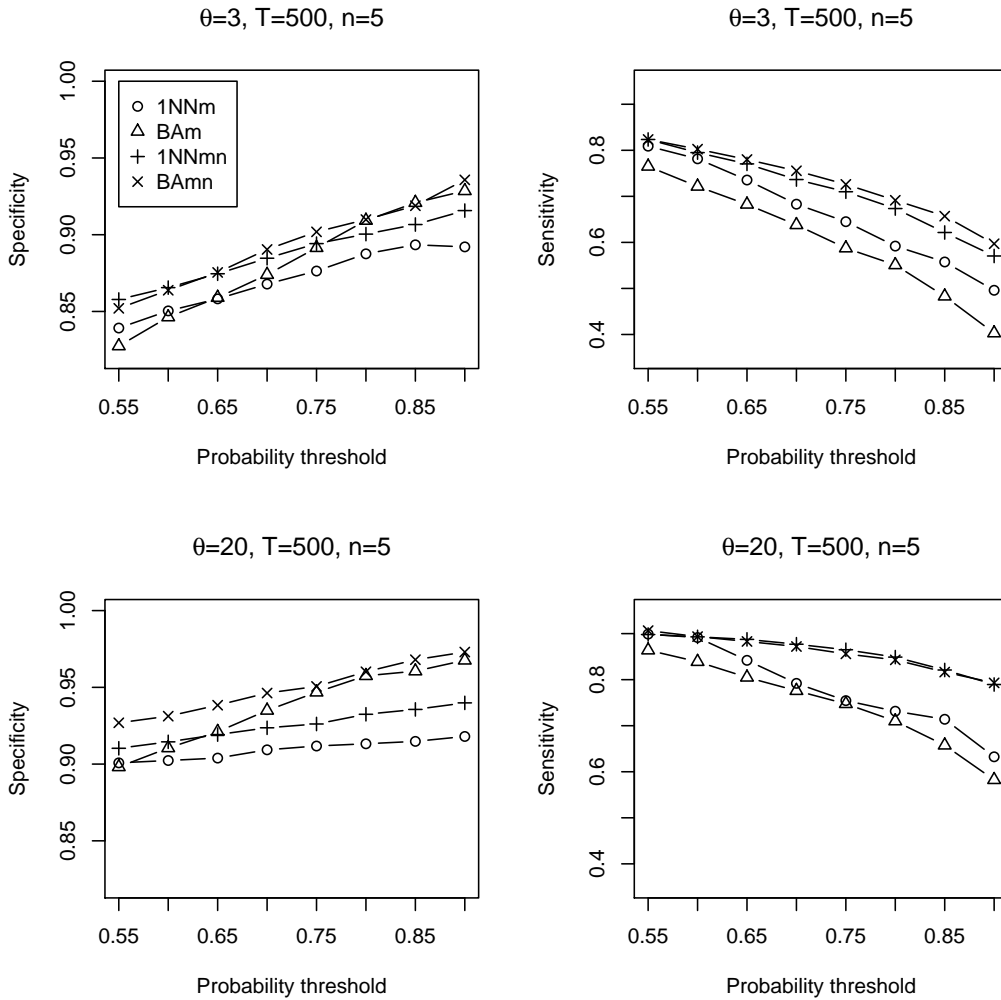


Figure 1: *Effect of adding nuclear data on the performance of coalescent-based barcoding.* Specificity is the fraction of non-ambiguous assignments that are correct. Sensitivity is the fraction of all the assignments that are correct. The probability threshold is the threshold used to decide if an assignment is ambiguous. 1NN and BA are the nearest-neighbour classification and the developed Bayesian assigner with a known value of θ . The subscripts m and mn denote the mitochondrial data and the combination of mitochondrial and nuclear data, respectively. Adding nuclear data increases sensitivity and reduces the ambiguity of assignments.

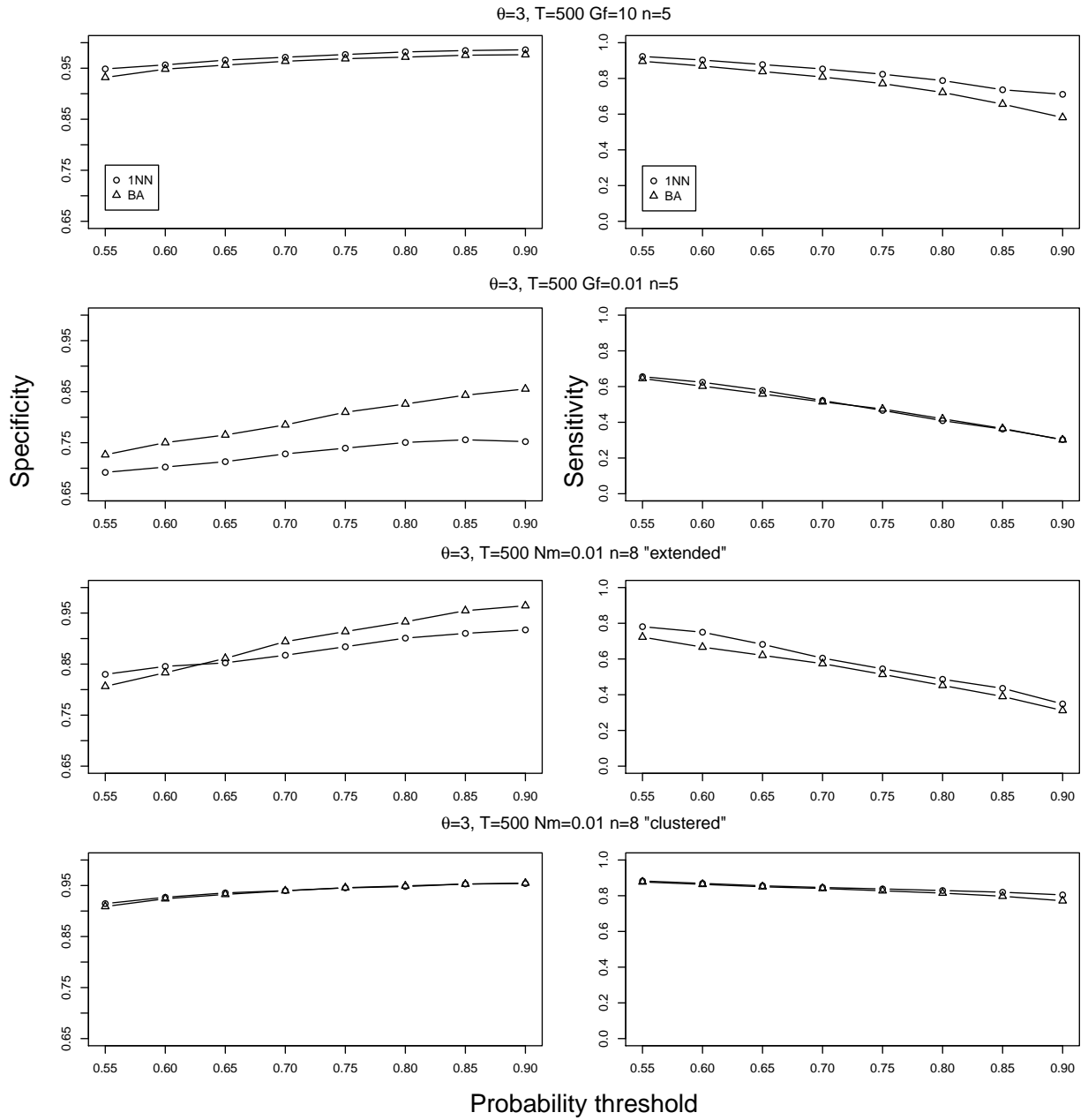


Figure 2: *Robustness of coalescent-based barcoding to past population size changes and population structures.* 1NN and BA are the nearest-neighbour classification and the developed Bayesian assigner with estimated mutation parameters. Results for past population size changes are presented on the first two lines, with Gf being the growth factor. A growth factor larger than one indicates a population expansion from divergence to present, whereas a growth factor smaller than one indicates a population decline. Results for population structures are presented on the last two lines, with Nm being the number of migrants exchanged between adjacent populations in one generation. BA appears robust since its performance is similar to the one of 1NN that is model-free.

6 Analysis of real data sets

We chose to test our method on two different data sets that contained both differentiated and undifferentiated species. The first data set used came from the study of Hebert et al. (2004) on *Astraptes* species and consisted of mitochondrial sequences (CO1 locus). The second data set used came from the study of Elias et al. (2007) on *Ithomiinae* species and consisted of mitochondrial (CO1 locus) and nuclear data (EF1 α locus). The data were analysed with 1NN, BA and TheAssigner (Abdo and Golding, 2007). The performance of each method was quantified using a leave-one-out analysis in which each haplotype was used as a test sequence after reducing its multiplicity by one in the reference data. Details on these data sets and their analyses are given in supplementary materials C.

The results first showed that adding nuclear data reduced the ambiguity of the BA assignments (Fig. 3). The analyses also showed that no method had the highest specificity in all cases (Fig. 3). Moreover BA had a lower sensitivity than the other methods and thus assigned fewer individuals (Fig. 3), except for the nuclear *Ithomiinae* data alone (Fig. S4). Another result of our analyses was that some posterior probabilities of membership were sensitive to the choice of the ancestral bases (supplementary materials C). Finally a few conditional probabilities were estimated with the predictive method (supplementary materials A) and the corresponding estimates were close to the plug-in estimates.

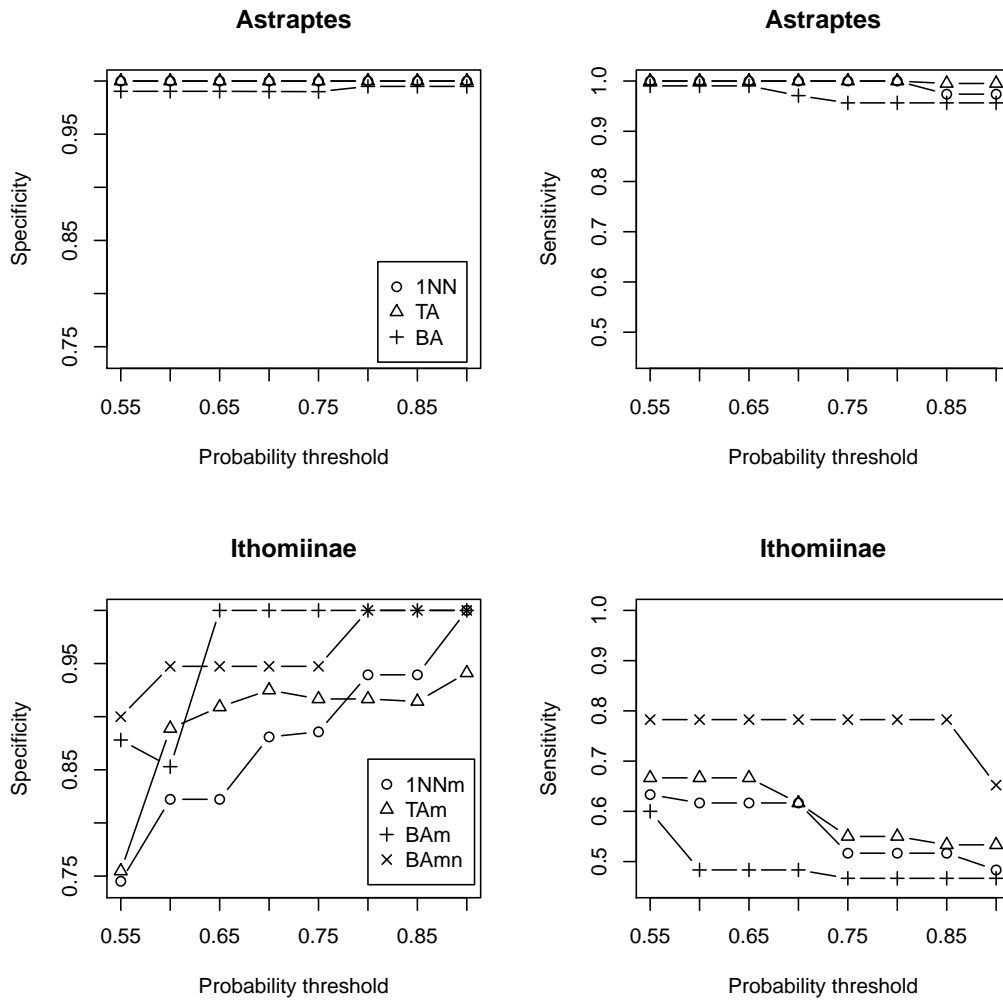


Figure 3: *Performance of coalescent-based barcoding with real data.* 1NN, TA and BA are the nearest-neighbour classification, TheAssigner and the developed Bayesian assigner. The subscripts m and mn denote the mitochondrial data and the combination of mitochondrial and nuclear data. Adding nuclear data increases sensitivity and reduces the ambiguity of BA assignments. No method has the best specificity for both data sets.

7 Discussion

Classification inputs. Bayesian classification requires prior probabilities of membership. When these probabilities are not known, they may be estimated from the reference data provided that these data can be considered as a random sample among all the species considered (page 53, Ripley, 1996) or they may be fixed to $1/c$.

The developed methods require the ancestral sequence of each locus. If this sequence is not known, it can be inferred from the data (Bahlo and Griffiths, 2000; Gascuel and Steel, 2010) or posterior probabilities of membership may be estimated using unrooted trees (Section 5, Tavaré and Zeitouni, 2004; Bahlo and Griffiths, 2000). Moreover, many sequences from the barcoding reference database could be used as outgroups and thus greatly facilitate the inference of the ancestral sequence.

Finally, both alleles of an individual at a nuclear locus were assumed to be known. Current genotyping technologies are able to determine which two bases are present at each site of a nuclear locus but not the two sequences of the locus. It is a general problem for most nuclear sequence analysis methods, and statistical methods, known as phasing methods, can infer these two sequences from unphased data together with missing data (Scheet and Stephens, 2006).

Classification assumptions. The mutation model considered in this paper was the ISM, a model that requires fewer computations than models with a finite number of sites. However it assumes that a particular mutation can only occur once so that in particular there is no homoplasy. It is more adapted to situations where species are closely related since the assumption of absence of homoplasy is more likely to be satisfied in this case. This does not seem to be a problem for DNA barcoding since species that are distantly related to a test individual can be discarded using simpler methods (Austerlitz et al., 2009; Munch et al., 2008b). In our study, classification methods were compared using data sets that were compatible with the ISM so that all the methods had the same amount of information.

Species classification based on the ISM could be extended to account for different mutation rates for transitions and transversions.

The species of a test individual was assumed to be represented in the reference data. The conditional probabilities of an allele $P(x|y_i)$ can be used to check if this assumption is satisfied: low probabilities are an indication that this assumption may not be satisfied.

The developed methods are based on various simplifying assumptions. It would be interesting to relax some of these assumptions to improve classification performance. The program genetree can perform likelihood estimations with varying population size and population structures under the ISM (Bahlo and Griffiths, 2000). Divergence models and models that combine phylogenetic and population genetics models do not assume that species are independent (Matz and Nielsen, 2005; Pons et al., 2006).

Performance of the developed methods. The method developed to combine mitochondrial and nuclear informations appeared satisfactory. Adding nuclear data reduced the ambiguity of assignments in our analyses.

We showed that coalescent-based classification was robust to departures from demographic stability and panmixia and to designs that did not sample the within-species variation efficiently. It performed similarly to a model-free method (1NN) in the robustness study. Demographic expansion was found to increase the power of barcoding. This is an expected result, however, considering that speciation events are probably often associated with founder events followed by demographic expansions or selective sweeps on the mitochondria, it may highlight the reasons why DNA barcoding works so well with a limited sequence information.

Finally, no assignment method was found to be always the best in our analyses. Similar results were obtained by Austerlitz et al. (2009) when comparing phylogenetic and statistical methods. However the developed Bayesian assigner generally appeared more cautious than the other methods in the sense that it assigned fewer individuals but made fewer errors

among the assigned individuals.

The supplementary materials referenced in Sections 3, 5 and 6 are available at arxiv.org.

Acknowledgements

This study was funded by the Agence Nationale de la Recherche (IFORA ANR-06-BDIV-014 and EMILE NT09-611697 projects). We thank F. Austerlitz for helpful comments.

Disclosure statement

No competing financial interests exist.

References

- Abdo, Z., Golding, G.B., 2007. A step toward barcoding life: A model-based, decision-theoretic method to assign genes to preexisting species groups. *Systematic Biology* 56, 44–56.
- Austerlitz, F., David, O., Schaeffer, B., Bleakley, K., Olteanu, M., Leblois, R., Veuille, M., Larédo, C., 2009. DNA barcode analysis: a comparison of phylogenetic and statistical classification methods. *BMC Bioinformatics, Special Issue Biodiversity Informatics* .
- Bahlo, M., Griffiths, R.C., 2000. Inference from gene trees in a subdivided population. *Theor. Popul. Biol.* 57, 79–95.
- De Iorio, M., Griffiths, R.C., 2004. Importance sampling on coalescent histories. I. *Adv. Appl. Prob.* 36, 417–433.
- Elias, M., Hill, R.I., Willmott, K.R., Dasmahapatra, K.K., Brower, A.V., Mallet, J., Jiggins, C.D., 2007. Limited performance of DNA barcoding in a diverse community of tropical butterflies. *Proc R. Soc. B.* 274, 2881–9.
- Ewens, W.J., 2004. *Mathematical population genetics. volume 27 of Interdisciplinary Applied Mathematics.* Springer. second edition.
- Frézal, L., Leblois, R., 2008. Four years of DNA barcoding: Current advances and prospects. *Infection, Genetics and Evolution* 8, 727 – 736.
- Gascuel, O., Steel, M., 2010. Inferring ancestral sequences in taxon-rich phylogenies. *Mathematical Biosciences* 227, 125 – 135.
- Hastie, T., Tibshirani, R., Friedman, J., 2001. *The elements of statistical learning: data mining, inference, and prediction.* Springer Series in Statistics, Springer.

- Hebert, P.D.N., Penton, E.H., Burns, J.M., Janzen, D.H., Hallwachs, W., 2004. Ten species in one: Dna barcoding reveals cryptic species in the neotropical skipper butterfly *Asdrastes fulgurator*. *Proceedings of the National Academy of Sciences of the United States of America* 101, 14812–14817.
- Hebert, P.D.N., Ratnasingham, S., deWaard, J.R., 2003. Barcoding animal life: cytochrome c oxidase subunit 1 divergences among closely related species. *Proc. R. Soc. B* 270, S96–S99.
- Hudson, R.R., 1991. Gene genealogies and the coalescent process. *Oxford Surveys in Evolutionary Biology* 7, 1–44.
- Kingman, J.F.C., 1982a. The coalescent. *Stochastic Processes and their Applications* 13, 235 – 248.
- Kingman, J.F.C., 1982b. On the genealogy of large populations. *Journal of Applied Probability* 19, 27–43.
- Kuhner, M.K., Yamato, J., Felsenstein, J., 1995. Estimating effective population size and mutation rate from sequence data using Metropolis-Hastings sampling. *Genetics* 140, 1421–1430.
- Matz, M.V., Nielsen, R., 2005. A likelihood ratio test for species membership based on DNA sequence data. *Philosophical Transactions of the Royal Society B - Biological Sciences* 360, 1969–1974.
- Miller, S.E., 2007. DNA barcoding and the renaissance of taxonomy. *Proceedings of the National Academy of Sciences* 104, 4775–4776.
- Munch, K., Boomsma, W., Huelsenbeck, J.P., Willerslev, E., Nielsen, R., 2008a. Statistical

- assignment of DNA sequences using bayesian phylogenetics. *Systematic Biology* 57, 750–757.
- Munch, K., Boomsma, W., Willerslev, E., Nielsen, R., 2008b. Fast phylogenetic DNA barcoding. *Philosophical Transactions of the Royal Society B* 363, 3997 – 4002.
- Nordborg, M., 2001. Coalescent theory, in: Balding, D.J., Bishop, M.J., Cannings, C. (Eds.), *Handbook of Statistical Genetics*, John Wiley & Sons, Inc., Chichester, U.K.. pp. 179–212.
- Pons, J., Barraclough, T., Gomez-Zurita, J., Cardoso, A., Duran, D., Hazell, S., Kamoun, S., Sumlin, W., Vogler, A., 2006. Sequence-based species delimitation for the DNA taxonomy of undescribed insects. *Systematic Biology* 55, 595–609.
- Ripley, B.D., 1996. *Pattern recognition and neural networks*. Cambridge University Press, Cambridge, UK.
- Scheet, P., Stephens, M., 2006. A fast and flexible statistical model for large-scale population genotype data: applications to inferring missing genotypes and haplotypic phase. *American Journal of Human Genetics* 78, 629–44.
- Tajima, F., 1983. Evolutionary relationship of DNA sequences in finite populations. *Genetics* 105, 437–460.
- Tavaré, S., Zeitouni, O., 2004. *Lectures on probability theory and statistics : Ecole d'été de probabilités de Saint-Flour XXXI - 2001*. Lecture notes in mathematics, Springer.
- Valentini, A., Pompanon, F., Taberlet, P., 2009. DNA barcoding for ecologists. *Trends in ecology & evolution* 24, 110–7.
- Watterson, G.A., 1975. On the number of segregating sites in genetical models without recombination. *Theoretical Population Biology* 7, 256–276.