

Dear reviewers,

we would like to thank you for your careful reading, and for the many constructive suggestions that have helped us progress in our understanding of the subject. We submit here a corrected version of our article, and we thank you in advance for your feedback, while remaining open to any new suggestion. The modifications we propose are numerous without being major, here is a detailed and argued presentation.

First of all, after much reflection, we have preferred to persist in our presentation, perhaps more appropriate for a mathematically trained readership, on the one hand by keeping the formulas in the text, and on the other hand by keeping the first part essentially useless for population geneticists. We explain this explicitly in a "note to the reader" in the preamble. But this persistence is not a *sine qua non* condition, if you think that it could be prohibitive for a recommendation.

One of your main requests was to better link our work to the existing bibliography. It took us a lot of time to read and reread, but we gladly agreed to this global request which we found quite legitimate, and it was very profitable for us. A new sub-section of section 2 was thus created (page 7), and many new references were also added in sections 4 and 5.

It seems that you were particularly intrigued by the results of section 4, so we have completely reworked this section in an attempt to respond to your requests for clarification and explanations, which also seemed quite legitimate to us.

Finally, the topic of "effective population size" is indeed very relevant to our work and very complex. We are not arguing that the PSMC should not be said to estimate an effective size, but recall that there is not only one definition of effective size (the three most common being based on 1) the probability of identity-by-descent of two alleles chosen at random, 2) the variance in offspring allele frequency, and 3) the leading non-unit eigen-value of the allele frequency transition matrix), and we only underline, as Sjödin et al. does¹, that « one too often reads of *the* effective population size without reference to the particular notion being considered », and in the case of the PSMC, the notion considered is the IICR, which cannot either be summarized in the 3rd above, although it is the closest.

Please find below detailed answers (in red in the text) to your specific requests and suggestions.

Olivier Mazet, Toulouse, le 10 mars 2023

¹ Sjödin, P., Kaj, I., Krone, S., Lascoux, M., & Nordborg, M. (2005). On the meaning and existence of an effective population size. *Genetics*, 169(2), 1061-1070.

Reviewer 1 – Alan Rogers

Minor comments

p. 2: « Each individual independently generates a number of descendants following a Poisson distribution. » This is only approximately correct. The marginal distribution of an individual is binomial with parameters $1=2N$ and $2N$. The joint distribution of all individuals is multinomial. **The whole sentence is « each individual independently generates a number of descendants following a Poisson distribution, conditioned by the fact that the size of the population must remain constant ». This is furthermore mathematically equivalent to the fact that backward in time, each individual choose its parent uniformly in the previous generation.**

Section « Structure and IICR: sampling strategy »: This one-paragraph section is cryptic. It needs an example, or a fuller explanation, to make clear what is going on.

Three lines of explanation were added, we hope that it is clearer.

p. 13: It took me a few minutes to figure out what was going on here. I came up with the following prose to explain this to myself. Here it is, in case the authors want to paraphrase it. If the genome is divided into several zones, each with a different effective population size, then the coalescence rate in the recent past will be high because of all the coalescent events happening in zones of small N_E . Consequently, N_E will be small in the recent past. In the distant past, all lineages within zones of small N_E will have coalesced into a single ancestral lineage, and multiple lineages will remain only in zones of large N_E . Consequently N_E will be large in the distant past. This will look like a decline in population size. As explained above, however, it seems strange that this effect would be large, in view of the small fraction of the genome that is likely to be affected by selection.

Thanks you for your suggestion, we tried to explain better what is going on.

p. 14: I was unable to understand the paragraph after Fig. 10.

Idem, we added explicitations, in order to make it clearer.

Reviewer 2 – Anonymous – 17 Aug 2022

My other general comment would be about the importance of the problem in practice. Ignoring population structure is only a problem if populations are structured, and sufficiently so that the IICR is substantially affected. So a naive question would be: are typical natural populations sufficiently structured for the problem to be serious? I am pretty sure the answer is yes, and indeed some of the authors' papers address this, but I feel like this review could be an opportunity to make kind of a general argument here. Could the authors give an idea, maybe based on reviews of the empirical literature, of what proportion of analyses are expected to be seriously flawed vs. more or less robust to this problem? Even a broad picture would be a plus in my opinion.

The answer to this question is not obvious. This requires a quantification of the structuring of the population which is not easy to define. On the other hand, typical natural populations, in addition to being structured, have also undergone changes in size in their evolutionary history which make the analysis more complex. A sentence has been added to the conclusion. Finally I have a number of point-by-point comments/suggestions which I hope might help improve the clarity of this important, well-done manuscript.

- [p1; section 1; paragraph 2] "until successive common ancestors are found" -> sounds like an awkward process-ending condition to me; maybe "until a common ancestor is found" ?
Agreed, and corrected

- [p1; section 1; paragraph 3] "The mathematical objects of interest..."
-> I think I would use singular instead of plural ("The mathematical object of interest is the joint distribution...")
Agreed, and corrected

-> "express [...] as a function of [...]" : maybe just "predict"?
The expression of the function of parameters is a way to predict their values with statistical tools, the prediction does not replace the expression, they are successive.

-> this sentence more or less implies that the whole coalescence process is entirely described by coalescence times, whereas, one could say, tree topology also matters (e.g. you don't expect the same SFS depending on whether tree is symmetrical vs. pectinated, as soon as $n > 3$); "coalescence times" and "tree" are often taken as synonyms in the ms; maybe clarify by adding a section/sentence about tree topology, its distribution and independence with respect to demography and selection?

We have added a sentence stating that topology is not taken into account (anyway in a large part of the paper we consider trees with two branches).

-> this sentence mentions "a family tree", suggesting that we're here considering a single locus, whereas the previous paragraph mentions "loci", and the next sentence mentions recombination as a relevant parameter, implying several loci - could the text be more consistent with this respect?

We have put « tree » in the plural, and explicitated that each locus has its own tree.

- [fig1 legend]
-> "in the past" -> "ago"? (twice)

Corrected

- the first three equations of p3 are not numbered

They do not need, since we never refer to them...

- [p3, second equation]

-> I am not sure every reader will know the "floor" symbol so would suggest defining it (or dropping it)

Agreed, dropped.

- [p3, third equation]

-> Tk and t in the third equation are not expressed in the same unit as Tk and t in the second equation (2N generations vs. generations), but the same symbols were kept

It's now made explicit

- overall I fell like the treatment of time scale normalization by Ne could be improved; it is introduced very lightly in the first equations, then recalled in several sections, sometimes in a lengthy way

We are afraid we did not understand what you meant...

- [p3, last section] "The absence of the panmictic assumption" -> "To relax the assumption of panmixy" ?

OK, corrected.

- [p4, first section] "assuming that the sizes of each population are sufficiently large" -> "assuming that populations are sufficiently large" ?

OK, corrected.

- equation 3: I do not understand why the left-hand term is not $Q(\alpha, \beta)$ instead of $Q(n_\alpha, n_\beta)$; I do not understand what these n_α and n_β terms represent; I am apparently missing a level of complexity here

It is just a way to have integers as indices of the matrix.

- [p4, "Mutation and genetic diversity" section] Reference to Tajima's D probably misplaced here. Tajima's D is a statistic that combines two estimators of $4N_e\mu$ (topic of the section) to learn about deviations from the standard coalescent. Maybe the intended reference is Tajima 1983 Genetics 105:437 ?

You are completely right, there was a mistake now corrected.

- [p6, first section] the key sentence starting "Considering that..." seems to lack a verb

OK, corrected.

- [p6, middle, first consequence ("The sole data...")] I would suggest being explicit and replacing "demographic model" with "population structure". At first reading I mentally interpreted "demographic model" as "model of Ne change in a single pop", thus missing the point. I know this is because I'm biased in a priori considering a single pop when thinking coalescent. Still I might not be alone, and given the importance of the sentence I would suggest avoiding any ambiguity.

OK, precision added.

- [p8] "so the matrices can be time dependent as piecewise constant functions" -> "so the matrices can be piecewise constant functions of time" ?

OK, corrected.

- [p10] "it is natural to want to increase" -> I would suggest rephrasing as "A natural way to increase the precision of the estimation of demographic parameters is to increase the sample size."

OK, rephrased.

- [p13, middle]: "this rate being linked to the reproductive capacity"; this seems confusing, and the nature of the link is unclear; there is no such a thing as distinct portions of a genome differing in their "reproductive capacities"; instead a genomic portion hosts genes at which distinct alleles conferring distinct reproductive capacities to their carriers used to appear and segregate, affecting the IICR; the sentence seems to entail a simple link between coalescence rate and the strength of selection, whereas in reality things are more complicated; for instance both selective sweeps (positive s) and background selection (negative s) tend to reduced coalescence times. Maybe instead "this rate being linked to the selective regime at work"?

This part has been substantially rewritten.

- Section 4 could discuss the appropriateness of modeling linked selection via a variable N_e across the genome. Selection has some intrinsic property of being variable in time, whereas the model presented here assumes a constant in time coalescence rate in any portion of the genome. For instance selective sweeps are expected to induce brief periods of very high coalescence rate (eg see papers by Barton, Hermisson, Petrov, Jensen and many others), i.e., affect the shape of genealogies and the IICR in a way not easily captured by the model used here. This model however is probably great in capturing the across-genome variations in constant-in-time selective regime, e.g. regions under balancing selection vs regions under recurrent background selection. Overall I feel like connecting section 4 a bit more firmly to the (heavy) body of literature on selection detection in pop genomics would be a great addition. Note that these considerations on tree shape essentially disappear when sample size=2, which is still often the case in PSMC-related literature, and this could also be mentioned as a justification for the approach presented here.

Thanks to your comments, we have added precisions and references.

- [last section] in addition to Charlesworth 2009 the authors might like to cite the recent Waples 2022 J Hered 113:371

Agreed, reference added.

Reviewer 3 – Anonymous – 15 Sep 2022

A central point of this manuscript is that changes of IIRC inferred by PSMC should only be cautiously interpreted as actual changes in effective population sizes (N_e) if one is confident that the panmictic model holds. However, $N_e \neq N$. This distinction between N_e & N should be made explicit throughout the manuscript, as sometimes the term “population size” is used in a blurry ways. The utility of this paper would be significantly increased if the author(s) were to further discuss how different evolutionary and demographic processes can shape N_e /IIRC and PSMC plots.

We are sorry, we don't see very well were exactly the term « population size » is used in a blurry way... We have checked, and each time that the term « effective » does not appear, the meaning of « (sub)population size » is truly the real size, the census.

Please, ensure to cite previous relevant work at all appropriate positions throughout the text. Sometimes it is not clear which study is discussed (e.g., discussion of Figure 7 on Page 9).

Right, we tried to fight against LaTeX in order to make the reading simpler...

Also other relevant computational tools for inferring population structure and changes in IICR should be cited (e.g., SMC++ and diCAL).

A large number of tools aim to infer the history of a population from whole sequence data, but to our knowledge only two of them infer the IICR : PSMC and MSMC (IICRk for the latter if there are $k > 2$ lineages). For instance, SMC++ involves the SFS into its analysis, and diCAL mixes the coalescence rates of all the T_k of the genealogical tree. So none of them infers the IICR, nor even directly analyses changes in it.

Additional suggestions and queries are listed below:

Page 1: I appreciate that the authors are up-front about being an anonymous collective (“Camille Noûs”). However, I think this can be shortened or even mentioned in a footnote. If the goal of give credit to other members of a research team, why not name them explicitly in an acknowledgements section at the end of the manuscript?

Agreed, acknowledgement section added.

Page 4: The interpretations of the conditions in Eq. 3 should be provided.

OK, conditions provided.

Page 5: Recent years have seen advances in the application of ARGs. It would be a good idea to add some citations in this regard, and perhaps go into additional detail about the strength and weaknesses of ARGs.

ARGs are not used in this work, since the SMC theory is intended to avoid them and their complexity. We think that it would be out of the scope to focus on them...

Page 7: I understand that Figure 3 was taken from another publication but simulation parameters in the figure legend would be helpful. Also please add a citation to the original paper in the Figure 3 legend so that one knows it is not an original result. Similar suggestions hold for subsequent figures.

OK, we added precisions and citations in the legends.

Pages 7 and 9 (related to Figure 5): What exact models were simulated? I get that they are

described in detail in Chikhi et al. 2018 but the reader shouldn't need to go to the original publication to understand what is simulated.

OK, precisions added.

Page 9 (related of Figure 7). This is an interesting point. It might be worth going into more detail, and providing the relevant citation.

OK, we went a bit into more detail.

Page 10: Computation of the IICR_k: This is interesting theory but why is introduced? Was it applied to data yet? If yes, what were the results?

The computation of the IICR₃ has been done for the n-island in [Grusea et al. 2018], which shows that this new statistics with IICR₂ is sufficient to distinguish structure from panmixia. But it has not been applied to real data yet, hopefully soon.

Page 11: In Figure 8, please define what a component is in the figure legend. What are the units of t and M ?

OK, precisions added.

Page 13: The discussion of the SFS comes out of the blue and it is not entirely clear how it fits into the broader context. As such, it could be deleted without harming the paper.

We removed this part, and kept only the reference to his perspective in the introduction of the « increase of the sample size » part.