

SOFTWARE

Open Access



MSA: reproducible mutational signature attribution with confidence based on simulations

Sergey Senkin*

*Correspondence:
senkins@iarc.fr
Genomic Epidemiology
Branch, International Agency
for Research on Cancer
(IARC/WHO), Lyon, France

Abstract

Background: Mutational signatures proved to be a useful tool for identifying patterns of mutations in genomes, often providing valuable insights about mutagenic processes or normal DNA damage. De novo extraction of signatures is commonly performed using Non-Negative Matrix Factorisation methods, however, accurate attribution of these signatures to individual samples is a distinct problem requiring uncertainty estimation, particularly in noisy scenarios or when the acting signatures have similar shapes. Whilst many packages for signature attribution exist, a few provide accuracy measures, and most are not easily reproducible nor scalable in high-performance computing environments.

Results: We present Mutational Signature Attribution (MSA), a reproducible pipeline designed to assign signatures of different mutation types on a single-sample basis, using Non-Negative Least Squares method with optimisation based on configurable simulations. Parametric bootstrap is proposed as a way to measure statistical uncertainties of signature attribution. Supported mutation types include single and doublet base substitutions, indels and structural variants. Results are validated using simulations with reference COSMIC signatures, as well as randomly generated signatures.

Conclusions: MSA is a tool for optimised mutational signature attribution based on simulations, providing confidence intervals using parametric bootstrap. It comprises a set of Python scripts unified in a single Nextflow pipeline with containerisation for cross-platform reproducibility and scalability in high-performance computing environments. The tool is publicly available from <https://gitlab.com/s.senkin/MSA>.

Keywords: MSA, Mutational signatures, NNLS, Parametric bootstrap, Nextflow

Background

Mutational signatures are distinctive combinations of somatic mutations which can be of various origin, such as exogenous or endogenous exposures, defective DNA repair pathways, DNA replication infidelity or DNA enzymatic editing [1, 2]. Currently, on the order of 100 signatures have been discovered in human cancer, but for the majority of them the aetiology remains unknown. This has given rise to a rich new field of



© The Author(s), 2021. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

mutational signature discovery, as well as linking signatures to various risk factors, in projects like Mutographs [3].

De novo extraction of signatures is commonly performed using Non-Negative Matrix Factorisation (NMF) [4] for somatic mutations under various mutational classifications [5], with tools such as SigProfilerExtractor [6]. Such signature extraction has been extremely informative in the analysis of many cancer types and shed light into mutagenesis of endogenous and exogenous risk factors [2].

It has also become apparent that certain signatures show a dose-response relationship with risk factors, for example COSMIC signature SBS4 with tobacco smoking [7]. In order to characterise such relationships, it is increasingly important to attribute mutational signatures, i.e. estimate their activities in any given sample, with confidence intervals. Some signatures with similar shapes are difficult to differentiate between each other, for instance COSMIC signatures SBS5 and SBS40 that both have a relatively flat profile. For such signatures, point estimates of attribution can often be inaccurate, leading to false positive or false negative findings.

Ideally, statistical uncertainty of signature attribution would be best estimated by performing repeated measurements. Given the high cost and complexity of such measurements (especially sample preparation, DNA extraction and sequencing), one needs to look for other alternatives. Bootstrapping has been proposed as a practical method to estimate uncertainty of signature attribution [8, 9], however, not all investigators are explicit in the precise version of bootstrapping used. Whereas some suggest that simple resampling with replacement of a mutational catalogue can give a meaningful result, we argue that classic bootstrap is not applicable in mutational signature attribution, and propose the parametric bootstrap approach under assumption that mutations accumulate according to Poisson processes for each given mutation class, such as a trinucleotide context.

Another common difficulty of signature attribution packages is the low reproducibility as well as scalability of computations in high-performance computing environments. This is a particular concern when analysis is performed on large datasets, and when the number of bootstrap variations is considerable. We resolve this issue by utilising Nextflow [10], a domain-specific language (DSL) designed to primarily address computational irreproducibility and efficient parallel execution in a large number of computing environments, from individual workstations to server clusters and cloud computing services.

Implementation

Consider a mutational catalogue (such as a set of somatic mutations in cancer genomes) as a matrix \mathbf{M} ($m \times n$). Here, m is the number of mutation types, such as 96 single base substitution types in the trinucleotide context, and n is the number of samples. Let \mathbf{S} ($m \times k$) be the non-negative matrix of k mutational signatures (e.g. COSMIC catalogue [2]). The goal is to find the non-negative matrix of activities (or exposures) \mathbf{A} ($k \times n$):

$$\mathbf{M} = \mathbf{S} \times \mathbf{A} \tag{1}$$

In a more general version of the problem where the matrix \mathbf{S} is unknown, de novo extraction of signatures can be performed using algorithms such as Non-negative Matrix Factorisation (NMF). Here, we assume that the signature matrix \mathbf{S} is known, thereby scrutinising a more specific problem of finding activities of predefined signatures. This is particularly relevant when estimating the uncertainty of each signature activity is important, or simply when the number of samples n is low and one needs to quickly examine how well these samples can be explained by known mutational signatures.

The signature attribution problem can therefore be studied independently for each given sample \mathbf{y} from the mutational catalogue \mathbf{M} , as it can be seen as a linear combination of signatures \mathbf{S} and their per-sample activities \mathbf{x} :

$$\mathbf{y} = \mathbf{S} \times \mathbf{x} \quad (2)$$

Several algorithms can be used to solve this kind of problem, such as quadratic programming (QP) or simulated annealing (SA). Here, we utilise the Lawson-Hanson algorithm for non-negative least squares (NNLS) [11], which is essentially a constrained version of the ordinary least squares problem:

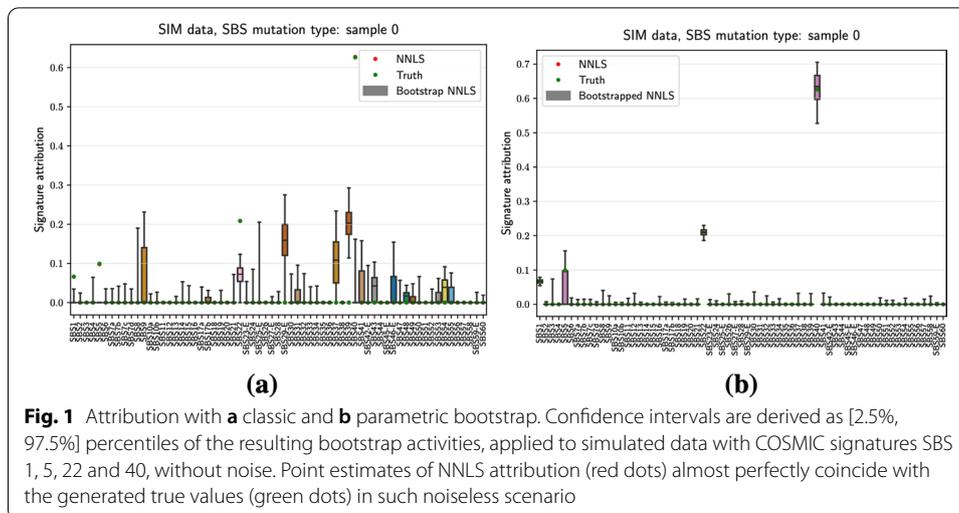
$$\arg \min_{\mathbf{x}} \|\mathbf{S}\mathbf{x} - \mathbf{y}\|_2, \mathbf{x} \geq 0 \quad (3)$$

Application of NNLS out of the box can be performed using one of the popular packages (*npls* in R or *optimize.npls* function in *scipy* library in Python), however, such approach is known to over-fit the data and lead to a large number of false-positive findings, as shown using simulations below. To mitigate this, ad-hoc approaches can be utilised, such as the optimisation of signature attribution using additional penalties on each signature's contribution, or by using pre-existing biological knowledge such as strand bias rules [2, 12]. Although such optimisation is implemented in MSA using the penalty loops described below, the main advantage of the tool is the ability to quantify the confidence of attribution for each signature.

Classic bootstrap

To estimate the confidence of signature attribution, several bootstrap approaches have been explored when developing this tool. In principle, bootstrapping can be applied to the full genome sequence data, but in the context of mutational signatures we apply bootstrap to collapsed mutational catalogues [matrix \mathbf{M} in Eq. (1)] for the sake of computational facility. Furthermore, it appears problematic to validate bootstrapped full sequences with simulations of predefined signatures.

Initially, a simple bootstrap, meaning resampling with replacement, was attempted on both simulated and real data. In this approach, N observations are randomly drawn from an initial sample with N mutation counts with replacement. In any given resample, each observation occurs 0, 1 or more times according to the binomial distribution $Binomial(N, 1/N)$. Since the total number of observations is N , all counts jointly form a multinomial distribution $Multinomial(N, 1/N, \dots, 1/N)$. Importantly, this approach gives meaningful results only under the i.i.d. assumption, i.e. when underlying random variables are independent and identically distributed [13]. This is not the case in mutation profiles: existing mutational signatures suggest that accumulation of mutations in particular trinucleotide contexts are inter-dependent and not identically distributed.



For example, samples with high levels of APOBEC activity are characterised by mutations enriched within TCA and TCT contexts, clearly not identically distributed to all other trinucleotide contexts. Not surprisingly, using simple bootstrap to derive confidence intervals of signature attribution leads to rather poor results (Fig. 1a). In a simulated example, four COSMIC signatures (SBS 1, 5, 22 and 40) were used to mimic real kidney RCC (Renal Cell Carcinoma) data without noise. Whilst running NNLS on the input sample yielded attributions nearly identical to the simulated true values, results on bootstrapped samples were vastly different, producing grossly inaccurate confidence intervals.

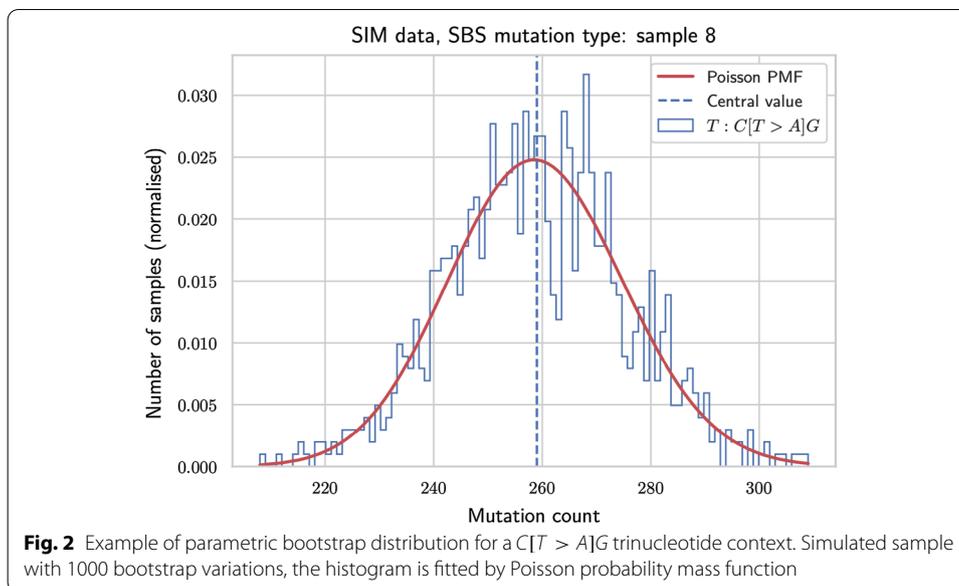
Some investigators argue that classic bootstrap approach can be regarded as conservative for non-i.i.d. data [14]. Other investigators argue that classic bootstrap is inadequate for high-dimensional non-i.i.d. data [15]. In the context of high-dimensional mutational profiles, we choose to derive confidence intervals on signature attributions using parametric bootstrap.

Parametric bootstrap using Multinomial distribution

Parametric bootstrap assumes that the data follow a known underlying distribution. This implies making certain assumptions on the original dataset, so that the bootstrap samples can be drawn from the estimated parametric model.

Here, we make an assumption that mutations are accumulated following Poisson distributions for each mutation class, such as a specific trinucleotide context, i.e. that in each class, mutations accumulate randomly, independently and at a constant rate.

The idea of parametric bootstrap applied to mutational processes was inspired by *mut-Signatures* package [9] (which is itself based on the original MATLAB framework for deciphering signatures [16]), where the input mutation matrix is bootstrapped according to the multinomial distribution $Multinomial(M, p_1, \dots, p_m)$, where M is the total mutational burden in a given sample, and probabilities p_i are normalised mutation counts for each mutation type. This distribution is chosen since the conditional distribution of a vector of independent Poisson variables is equivalent to multinomial distribution [17].



Since the mutational burden of each bootstrap sample is fixed and equal to M , we slightly modify the method by drawing counts from independent binomial distributions, so that the total mutational burden is no longer fixed. Nevertheless, for any given mutation category (e.g. C[T > A]G trinucleotide context on Fig. 2) the distribution of bootstrapped mutation counts follows a Poisson distribution.

For each bootstrap sample, NNLS attribution is applied to derive the vector of signature activities. 95% confidence intervals are then derived for each signature attribution by taking [2.5%, 97.5%] percentiles of the resulting bootstrap activities. Direct comparison with classic bootstrap on a simple simulated case without noise shows a clear advantage of the parametric method (Fig. 1b).

Optimisation of signature attribution

Since pure NNLS is based on a simple fitting approach, it is generally prone to overfitting, particularly in noisy environments [18]. To mitigate this, a form of regularisation can be applied using the penalisation loops to add or remove signatures based on their contribution to the fit. This approach is based on the penalised attribution used in Pan-Cancer Analysis of Whole Genomes (PCAWG) consortium signature assignment [2].

The default optimisation strategy, called removal strategy, starts with a full set of available signatures. Base L_2 similarity of the reconstructed profile (linear combination of input signatures) to the input mutational profile is calculated, normalised by its L_2 norm. Afterwards, a removal loop is executed, where all least contributing signatures increasing the L_2 similarity by less than a given penalty (called “weak” threshold) are sequentially removed. The resulting set of remaining signatures is used to describe the input sample by applying the final NNLS fit.

On the other hand, high penalties used in optimisation can lead to under-fitting of data, meaning that optimal penalties need to be derived. We perform this by running

simulations and measuring sensitivities, specificities as well as other metrics for all acting signatures.

Automation and integration

The MSA tool is implemented as a set of scripts in Python language, with simulation, optimisation and final attribution steps fully automatised using Nextflow workflow management system [10]. All dependencies, including packages such as *pandas*, *numpy*, *scipy*, *matplotlib* and *seaborn*, are automatically handled by Nextflow via containerisation using Docker technology, as well as Singularity where Docker is not available. Users not willing to use Docker or Singularity may opt to use a Conda environment that is also automatically handled by Nextflow, yet this approach has inferior reproducibility compared to container technology.

Native support of SigProfilerMatrixGenerator [5] and SigProfilerExtractor [6] tools is implemented in MSA for convenience. Where SigProfiler outputs are available, running MSA is as simple as executing a single line containing the output paths:

```
nextflow run https://gitlab.com/s.senkin/MSA -profile  
docker --dataset test --SP_matrix_generator_output_path  
path/to/SP_ME/ --SP_extractor_output_path path/to/SP_  
extractor/
```

Mutational classifications currently supported are SBS-96, SBS-192 and SBS-288 for single base substitution signatures, DBS-78 for doublet base substitution signatures, ID-83 for small insertion and deletion signatures, SV-32 for structural variant signatures. The SV signatures are implemented as an experimental feature.

Results

Execution of the MSA tool using Nextflow automatically produces various output files, including tables with absolute and relative signature attributions, statistical information with goodness-of-fit measures, optimisation plots and tables, fitted mutation spectra, residuals and all bootstrap output including confidence intervals of attribution. All corresponding plots are produced in PDF format.

Validation with simulations

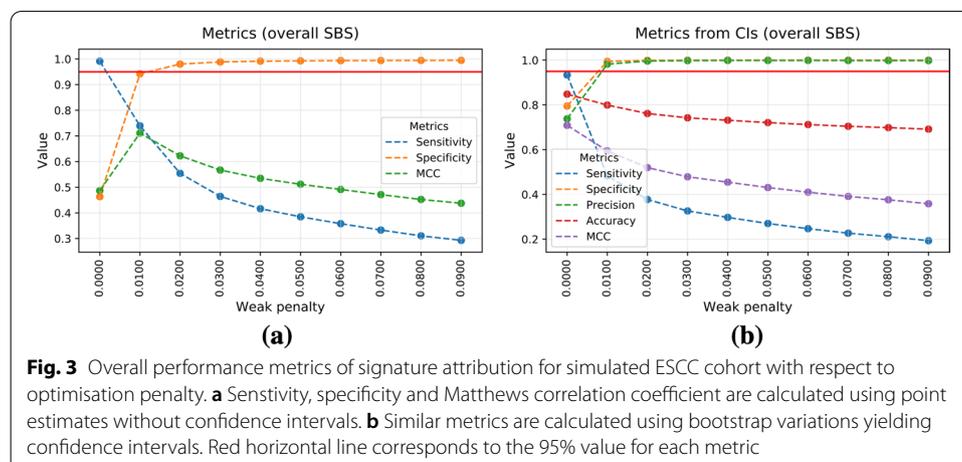
When simulating a dataset that is supposed to resemble real data, one generally has to make certain assumptions about underlying distributions. Firstly, we explored scenarios mimicking real cohorts (e.g. shown on Fig. 1), with a defined set of acting reference COSMIC signatures, where generated signature attributions follow non-negative zero-inflated Gaussian distributions. Generally, such distributions do not describe real attributions well, especially for rare signatures. Therefore, as a default approach fully automatised within the MSA Nextflow pipeline, we use data-driven simulations based on a simple bootstrap of signature activities derived without regularisation (i.e. with zero penalties). These activities are used to generate mutational profiles which are then resampled with replacement and injected with noise. Although these activities are bound to be over-fitted, they provide a way to sample signature attributions from distributions

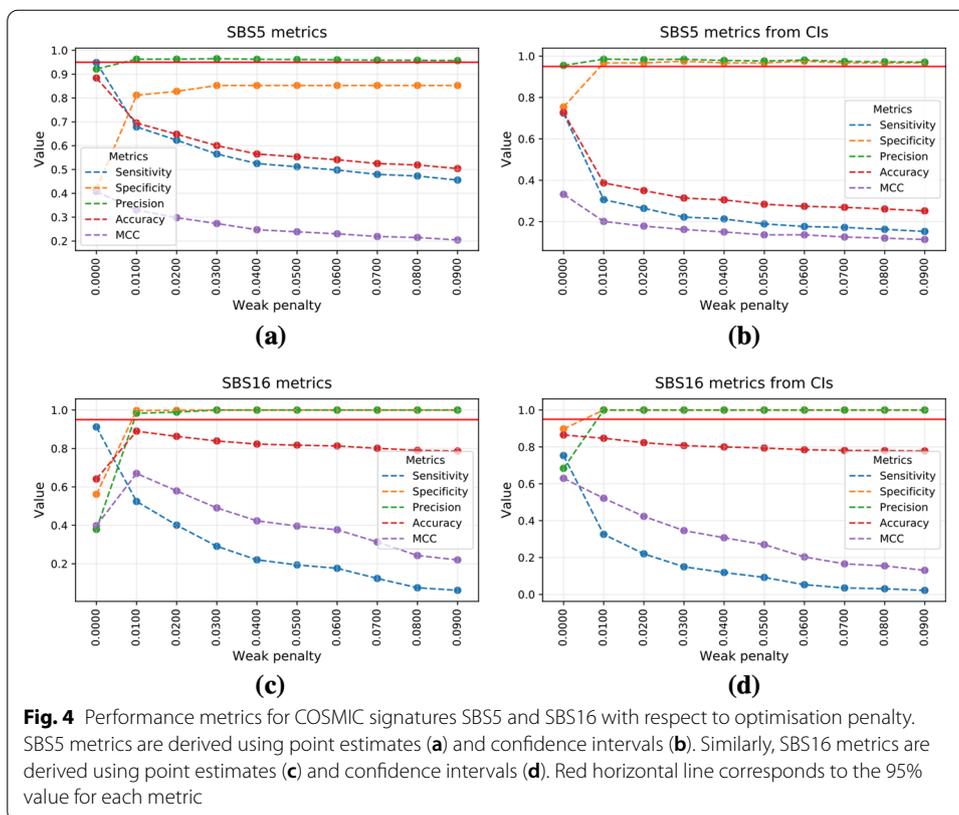
of observed real data, capturing signatures specific to any given cohort. The over-fitted generated signatures generally have low attribution levels, and can be regarded as noise in addition to the default noise model. We use the Gaussian noise model by default, with the standard deviation equal to 10% of mutational burden for each sample. Alternatively, Poisson model can be used, where the variance is equal to the mean generated mutation burden for any given mutation type.

Simulations allow to estimate the performance of signature attribution overall, as well as for each signature in different optimisation scenarios. Since the simulated truth is always known, one is able to calculate various metrics in order to estimate the accuracy of signature attribution, such as sensitivity and specificity for each signature, or for all signatures on average. As an example, Fig. 3 shows such metrics for all COSMIC SBS signatures combined, based on the simulated model of Esophageal Squamous Cell Carcinoma (ESCC) signature activities. Figure 4 shows these metrics for COSMIC signatures SBS5 and SBS16, for a range of “weak” penalties discussed previously. Signatures SBS5 and SBS16 were picked as typical examples of signatures with flat and non-flat profiles, respectively. The L2 penalty of zero corresponds to simple NNLS fit without any optimisation. Metrics are estimated with and without utilising confidence intervals, with the latter approach using the lower limit of confidence intervals when calculating metrics.

First of all, it is evident that higher penalties lead to lower sensitivity yet higher specificity of signature attribution. Secondly, confidence intervals allow to reach higher levels of specificity for signatures that are particularly difficult to attribute, such as signatures SBS5 and SBS40 due to their relatively flat profile. In a given example without confidence intervals, i.e. only using point estimates of attributions derived for simulated samples, specificity of SBS5 never reaches 90%. On the other hand, confidence intervals do drive weakly-acting signatures with low attributions to zero, hence lowering the overall sensitivity to such signatures, which can be partially recovered by applying a lower penalty.

Generally, finding an optimal penalty is always a trade-off. The advantage of the MSA tool is its flexibility with respect to optimisation. By default, the pipeline automatically runs optimisation across the default range of L_2 penalties, prioritising specificity of all signatures. However, investigators are free to choose a more or less conservative strategy by prioritising sensitivity or specificity of selected signatures.





Finally, performance of MSA was benchmarked against existing tools using both real and synthetic data published by the PCAWG consortium [19] (Additional file 1), as well as simulations based on the real PCAWG data, for SBS, DBS and Indel mutation types (Additional file 2). MSA generally achieves higher attribution performance than that of other tools due to its automatised optimisation of the penalties applied, coupled with the application of confidence intervals.

Discussion

We explored parametric bootstrap based on multinomial distribution as an effective method to derive confidence intervals of mutational signature attribution for any given mutational profile and available set of de novo extracted, or reference signatures. The multinomial distribution is chosen as the one corresponding to an assumption that mutations are accumulated according to Poisson processes for each mutation class, in line with the original signature extraction algorithm [16], but other assumptions can be potentially investigated—such as Monte Carlo simulations of mutations following distributions of increasing complexity.

The main limitation of the parametric bootstrap we use is its bias towards observed data, since empirical probabilities p_i in the multinomial distribution $Multinomial(N, p_1, \dots, p_m)$ are taken from real data which can be inaccurate. However, this method aims to estimate the statistical uncertainty of signature attribution method rather than total uncertainty of attribution, and for such purposes remains adequate.

Parametric bootstrap can be applied using both simple NNLS attribution and the one based on penalised optimisation designed to maximise sensitivity and specificity of signature attribution. We developed a set of tools assisting in the validation of optimisation parameters for any given scenario using automatised data-driven simulations of different cancer types and population cohorts. In general, it appears that such optimisation is unique for any given cohort due to the uniqueness of acting signatures. It also depends on noise, therefore investigators exploring cohorts of interest ideally need to test different noise scenarios. We opted to pick the Gaussian and Poisson noise models, yet other models such as negative binomial noise are worth consideration.

Estimating the systematic uncertainties of signature attribution, such as the ones due to sequencing artefacts, is a more challenging task that requires comparison of multiple sequencing technologies, not feasible in most settings. However, an uncertainty due to the variant calling can potentially be estimated for any given variant, and propagated to the signature attribution uncertainty. Evaluating these and other uncertainties remains a matter for further studies.

In summary, considering different sources of errors is an important exercise for any measurement, and as we show here, particularly for the attribution of mutational signatures. We hope that investigators will continue to ask such questions and strive to advance the methods shedding light on them. As a step in this direction, we present MSA—a computational tool to attribute mutational signatures with confidence intervals in an easily reproducible and scalable manner.

Conclusions

Mutational signature attribution is a problem distinct from signature extraction, requiring uncertainty estimation, particularly in noisy scenarios or when the acting signatures have similar shapes. Whilst many packages for signature attribution exist, a few provide accuracy measures, and practically none apply automatised regularisation based on simulations. Furthermore, most tools are not easily reproducible nor scalable in high-performance computing environments.

In this study, we propose MSA, a computational tool for optimised mutational signature attribution based on simulations, providing confidence intervals using parametric bootstrap. It is the first tool to perform automatic optimisation of regularisation based on data-driven simulations specific to any given input cohort. The tool comprises a set of Python scripts unified in a single Nextflow pipeline with containerisation, specifically designed for cross-platform reproducibility and scalability in high-performance computing environments.

MSA is publicly available at <https://gitlab.com/s.senkin/MSA> with an extensive documentation at <https://gitlab.com/s.senkin/MSA/-/wikis>.

Availability and requirements

- **Project name:** MSA
- **Project home page:** <https://gitlab.com/s.senkin/MSA>
- **Project wiki page:** <https://gitlab.com/s.senkin/MSA/-/wikis>
- **Operating system(s):** Platform independent
- **Programming language:** Python, Nextflow DSL

- **Other requirements:** Java 8 or later, docker, conda or singularity
- **License:** GNU GPL
- **Any restrictions to use by non-academics:** None

Abbreviations

DBS: Doublet base substitution; ESCC: Esophageal squamous cell carcinoma; ID: Small insertion and deletion; Indel: Small insertion and deletion; NMF: Non-negative matrix factorisation; NNLS: Non-negative least squares; SBS: Single base substitution; SV: Structural variant.

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12859-021-04450-8>.

Additional file 1: Comparison of MSA performance with existing tools. Performance of MSA benchmarked using real PCAWG consortium data and simulations, compared with SigProfiler and SignatureAnalyzer tools.

Additional file 2: MSA performance across different mutation types on data-driven simulations. Performance of MSA validated for SBS, DBS and ID mutation types using data-driven simulations derived from the PCAWG data.

Acknowledgements

We would like to thank Drs Ludmil Alexandrov, Graham Byrnes, Liacine Bouaoun, Vivian Viallon, Karl Smith-Byrne and Sarah Moody for immensely useful discussions and feedback.

Disclaimer

Where authors are identified as personnel of the International Agency for Research on Cancer / World Health Organization, the authors alone are responsible for the views expressed in this article and they do not necessarily represent the decisions, policy or views of the International Agency for Research on Cancer / World Health Organization.

Author contributions

SS was the sole author of the manuscript and performed all the reported work. The author read and approved the final manuscript.

Funding

This work was supported by a Cancer Grand Challenges Mutographs team award funded by Cancer Research UK [C98/A24032]. The funding bodies played no role in the design of the study and collection, analysis, and interpretation of data and in writing the manuscript.

Availability of data and materials

The datasets generated and/or analysed during the current study are available in the MSA repository, <https://gitlab.com/s.senkin/MSA>. Mutational catalogues and signatures of PCAWG consortium data, including real and synthetic data are available at <https://doi.org/10.7303/syn11726601>.

Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Received: 16 March 2021 Accepted: 13 October 2021

Published online: 04 November 2021

References

1. Alexandrov LB, et al. Signatures of mutational processes in human cancer. *Nature*. 2013;500(7463):415–21. <https://doi.org/10.1038/nature12477>.
2. Alexandrov LB, et al. The repertoire of mutational signatures in human cancer. *Nature*. 2020;578(7793):94–101. <https://doi.org/10.1038/s41586-020-1943-3>.
3. Moody S, Senkin S, et al. Mutational signatures in esophageal squamous cell carcinoma from eight countries with varying incidence. *Nat Genet*. 2021. <https://doi.org/10.1038/s41588-021-00928-6>.

4. Lee DD, Seung HS. Learning the parts of objects by non-negative matrix factorization. *Nature*. 1999;401(6755):788–91. <https://doi.org/10.1038/44565>.
5. Bergstrom EN, et al. SigProfilerMatrixGenerator: a tool for visualizing and exploring patterns of small mutational events. *BMC Genom*. 2019;20(1):685. <https://doi.org/10.1186/s12864-019-6041-2>.
6. Islam SMA, et al.: Uncovering novel mutational signatures by de novo extraction with SigProfilerExtractor. *bioRxiv* (2020). <https://doi.org/10.1101/2020.12.13.422570>.
7. Alexandrov LB et al.: Mutational signatures associated with tobacco smoking in human cancer. *Science* (New York, N.Y.) 354(6312), 618–622 (2016). <https://doi.org/10.1126/science.aag0299>.
8. Huang X, Wojtowicz D, Przytycka TM. Detecting presence of mutational signatures in cancer with confidence. *Bioinformatics*. 2017;34(2):330–7. <https://doi.org/10.1093/bioinformatics/btx604>.
9. Fantini D, et al. MutSignatures: an R package for extraction and analysis of cancer mutational signatures. *Sci Rep*. 2020;10(1):18217. <https://doi.org/10.1038/s41598-020-75062-0>.
10. Di Tommaso P, et al. Nextflow enables reproducible computational workflows. *Nat Biotechnol*. 2017;35(4):316–9. <https://doi.org/10.1038/nbt.3820>.
11. Lawson CL, Hanson RJ. Solving least squares problems. *Classics in applied mathematics*, vol. 15, p. 337. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA (1995). Revised reprint of the 1974 original
12. Li S, Crawford FW, Gerstein MB. Using sigLASSO to optimize cancer mutation signatures jointly with sampling likelihood. *Nat Commun*. 2020;11(1):3575. <https://doi.org/10.1038/s41467-020-17388-x>.
13. Efron B. Bootstrap methods: another look at the jackknife. *Ann Stat*. 1979;7(1):1–26. <https://doi.org/10.1214/aos/1176344552>.
14. Liu RY, Singh K. Using i.i.d. bootstrap inference for general non-i.i.d. models. *J Stat Plan Inference*. 1995; 43(1), 67–75. [https://doi.org/10.1016/0378-3758\(94\)00008-J](https://doi.org/10.1016/0378-3758(94)00008-J). *Statistics '91 Canada Conference Papers*
15. Owen AB, Eckles D. Bootstrapping data arrays of arbitrary order. *Ann Appl Stat*. 2012;6(3):895–927. <https://doi.org/10.1214/12-AOAS547>.
16. Alexandrov L, Nik-Zainal S, Wedge D, Campbell P, Stratton M. Deciphering signatures of mutational processes operative in human cancer. *Cell Rep*. 2013;3(1):246–59. <https://doi.org/10.1016/j.celrep.2012.12.008>.
17. Steel GD. Relation between Poisson and multinomial distributions. *Biometrics Unit Technical Reports*, BU-39-M (1953).
18. Slawski M, Hein M. Non-negative least squares for high-dimensional linear models: consistency and sparse recovery without regularization. *Electron J Stat*. 2013;7:3004–56. <https://doi.org/10.1214/13-EJS868>.
19. The ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium. Pan-cancer analysis of whole genomes. *Nature*. 2020;578(7793):82–93. <https://doi.org/10.1038/s41586-020-1969-6>.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

