

Developing Machine Learning Applications for Population Genetic Inference: Ensuring Precise Terminology and Robust Implementation

Xin Huang^{1, 2, *}

1 Department of Evolutionary Anthropology, University of Vienna, Vienna, Austria

2 Human Evolution and Archaeological Sciences (HEAS), University of Vienna, Vienna, Austria

* Corresponding author: xin.huang@univie.ac.at

Abstract

Machine learning applications for population genetic inference are emerging due to their potential to leverage large-scale genomic datasets, offering insights that traditional statistical methods may overlook. However, I have identified certain recurring issues. First, there is sometimes confusion between power and recall, and between the false discovery rate and one minus precision. These terms are specifically designed for hypothesis testing and are not appropriate for directly evaluating classification outcomes, as classification is a different task. Second, the lack of robustness in machine learning applications complicates their verification and application across different datasets, limiting their broader impact and slowing research progress. Robustness can be improved through strategies such as employing object-oriented programming for design, utilizing version control systems during development, and adopting package managers and workflow managers for distribution. I suggest by adhering to precise terminology and refining implementation practices, the impact of machine learning in population genetics can be maximized.

Keywords: Population genetics; Machine learning; Classification; Hypothesis testing; Software development.

Recent advances in machine learning, particularly in deep learning, have introduced innovative methodologies for population genetic inference, enabling the development of numerous applications for tasks such as inferring population structure, identifying genomic introgression, and investigating natural selection (Schridder and Kern 2018; Korfmann et al. 2023; Yelman and Jay 2023; Huang et al. 2024). These advancements capitalize on large-scale genomic data and sophisticated models, revealing intricate patterns that traditional statistical methods might overlook. While interdisciplinary approaches can inspire

31 novel applications, it is crucial to apply these methods with precision, particularly regarding the
32 terminology across different fields. Such precision enhances communication among disciplines and helps
33 prevent misunderstandings. As machine learning models grow increasingly complex, testing and
34 validating these approaches becomes progressively challenging (Raff 2019). Therefore, ensuring the
35 robustness of machine learning applications is essential for maximizing their impact. In this perspective, I
36 will first address the misuse of statistical terminology when applying machine learning to population
37 genetics, followed by an examination of the issues of non-robustness in certain machine learning
38 applications for population genetic inference, along with proposed strategies to enhance reproducibility
39 and reliability in the development of these applications.

40 **Utilize Precise Terminology**

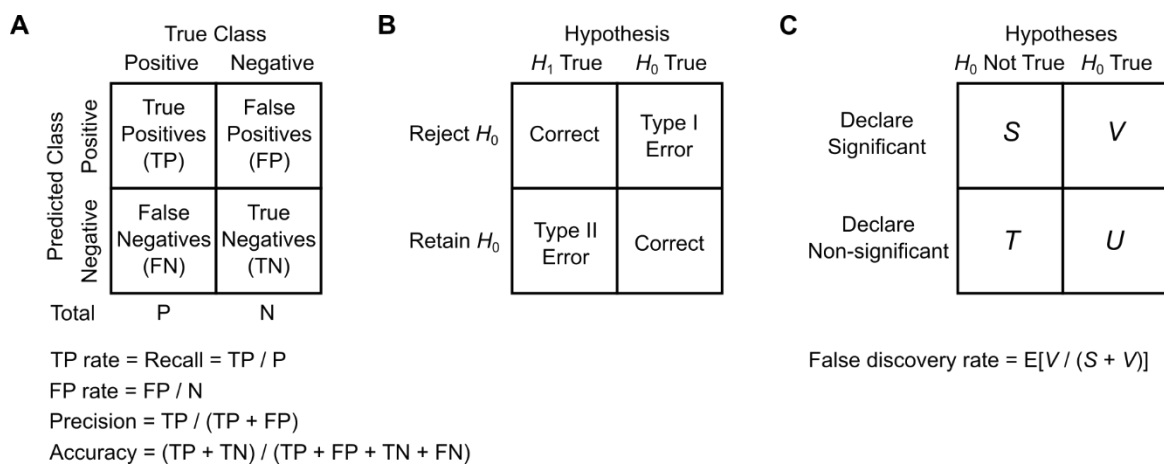
41 Several recent machine learning applications for population genetic inference (Durvasula and
42 Sankararaman 2020; Amin et al. 2023; Arnab et al. 2023; Lauterbur et al. 2023a; Zhang et al. 2023) have
43 used power and the false discovery rate (FDR) to evaluate the performance of their classification models,
44 equating power to recall (or the true positive rate) and FDR to one minus precision. However, the terms
45 ‘power’ and ‘FDR’ seem to be used colloquially to describe performance in terms of aptitude or
46 effectiveness, which can lead to confusion. These terms have specific meanings in hypothesis testing,
47 particularly in frequentist inference, and should not be conflated with machine learning metrics, as they
48 are not directly interchangeable. For a statistical test, power is defined as the probability of correctly
49 rejecting the null hypothesis when it is false, while FDR, in the context of multiple testing, is defined as
50 the expected proportion of false discoveries—not false positives—among the rejected null hypotheses
51 (Benjamini and Hochberg 1995; Wasserman 2003; Descôteaux 2007; Shreffler and Huecker 2024).

52 Hypothesis testing uses sample data to assess whether a population characteristic is more likely to be
53 consistent with the null hypothesis (H_0) or the alternative hypothesis (H_1), for the population from which
54 the sample was drawn. To conduct a hypothesis testing, an appropriate test statistic is needed (Wasserman
55 2003; Alpaydin et al. 2014). For example, suppose we sample two groups of data $\{X_1, X_2, \dots, X_n\}$ and $\{Y_1,$
56 $Y_2, \dots, Y_n\}$, a typical hypothesis task is to determine whether their means are equal. In this case, H_0 is $\mu_X =$
57 μ_Y , and H_1 is $\mu_X \neq \mu_Y$, where μ_X and μ_Y represent the true means of the populations from which the samples
58 $\{X_1, X_2, \dots, X_n\}$ and $\{Y_1, Y_2, \dots, Y_n\}$ are drawn. If we further assume that both groups are approximately
59 normally distributed, then we can calculate a t statistic based on the collected data to test our hypothesis.
60 If such a test is performed multiple times, a multiple testing correction may be applied.

61 In contrast, classification is a predictive task that focuses on assigning a label or category to each newly
62 observed data point, rather than making inferences about population characteristics. Unlike hypothesis

63 testing, classification does not rely on a test statistic to make decisions. Instead, it employs an algorithmic
 64 approach to learn patterns from the data and apply those patterns to new observations. For example,
 65 consider a dataset $\{(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)\}$, where $\{X_1, X_2, \dots, X_n\}$ are instances and $\{Y_1, Y_2, \dots,$
 66 $Y_n\}$ are their corresponding classes; classification uses these data to find a function $f: X \rightarrow Y$ that predicts
 67 the class Y_{n+1} of a new instance X_{n+1} , where $Y_{n+1} = f(X_{n+1})$ (Wasserman 2003). Since no hypothesis testing
 68 is involved in this prediction, I suggest that when applying techniques such as classification to fields like
 69 population genetics, it is important and necessary to adhere to the established terminology specific to
 70 those techniques.

71 Several metrics to assess the performance of a classifier exist, such as error, which represents the
 72 proportion of incorrectly classified instances out of the total number of instances (Alpaydin et al. 2014).
 73 Error is equal to one minus accuracy, where accuracy is the proportion of correctly classified instances
 74 (Alpaydin et al. 2014). However, accuracy is often an inadequate metric for evaluating performance,
 75 which is why receiver operating characteristic (ROC) curves have become popular in the machine
 76 learning community (Fawcett 2006). To conduct ROC analysis, various performance metrics are
 77 calculated from elements organized in a confusion matrix, which may sometimes be mistakenly equated
 78 with the outcomes of hypothesis testing or multiple testing due to their similar structure (Figure 1).



79
 80 **Figure 1 Outcomes in different scenarios.** (A) Binary classification outcomes. (B) Hypothesis testing
 81 outcomes. (C) Multiple testing outcomes.

82 One should note that ‘positive’ and ‘negative’ in machine learning are defined on a case-by-case basis.
 83 For example, when detecting introgressed fragments, ‘positive’ can refer to the detection of introgressed
 84 fragments, whereas when detecting introgression deserts, ‘positive’ can refer to the detection of these
 85 deserts. Instead, conventional statistical terminology emphasizes concepts like rejecting or failing to reject
 86 the null hypothesis, identifying significant results, or making discoveries, without framing these outcomes

87 in terms of ‘positive’ or ‘negative’ (Figure 1B). The null hypothesis typically represents a statement of no
88 effect, such as no difference between groups, which does not align naturally with the binary language of
89 ‘positive’ or ‘negative.’ This is because statistical results often involve complex interpretations that cannot
90 be easily reduced to such terms. For instance, a statistically significant result does not necessarily imply
91 scientific significance, and a non-significant result does not necessarily indicate the absence of an
92 effect—it may simply mean the effect was not detectable with the given data (Wasserman 2003). In this
93 context, a ‘discovery’ typically refers to the rejection of the null hypothesis (Soric 1989). The original
94 definition of FDR describes it as ‘the expected proportion of falsely rejected hypotheses,’ without
95 employing the term ‘positive’ (Benjamini and Hochberg 1995). This suggests that, in hypothesis testing,
96 conclusions are made carefully and conditionally based on the data, while in classification, the focus is on
97 confidently assigning labels to instances, reflecting the fundamental difference in how results are
98 interpreted.

99 Moreover, FDR is the expected value of the false discovery proportion (FDP; Wasserman 2003), a
100 random variable $Q = V / (S + V)$ (Figure 1C). Although FDP, $V / (S + V)$, resembles the formula for $FP /$
101 $(TP + FP)$, it is important to note that S and V are unobserved random variables; only the total number of
102 tests declared significant ($S + V$) is known after experimentation and analysis of the collected data
103 (Benjamini and Hochberg 1995). This differs from common metrics for evaluating classification model
104 performance, such as recall and precision, which are calculated using elements (TP, FP, TN, and FN)
105 directly observed in the data (Figure 1A). These metrics directly measure performance based on
106 prediction outcomes, without requiring a statistical hypothesis. Hence, the claim that FDR represents ‘the
107 proportion of false positives among all positives detected’ is inaccurate (Lauterbur et al. 2023a; Romieu et
108 al. 2024).

109 FDP, precision, and recall are dataset-dependent metrics, meaning they vary based on the characteristics
110 of the data and the model performance, which are often beyond direct control. In contrast, the power of a
111 statistical test remains theoretically consistent for a given null and alternative hypothesis, as long as test
112 conditions—such as sample size, significance level, and effect size—are held constant. This is because
113 power reflects the probability of detecting an effect when one exists, and these test conditions can be
114 controlled in advance (Shreffler and Huecker 2024). Although we can adjust the significance level in
115 hypothesis testing to influence the power of the test, this differs from how thresholds are used in
116 classification to determine whether an instance belongs to the positive or negative class based on the
117 probability output from a classifier. The significance level controls the probability of making a type I
118 error (Figure 1B) in deciding whether to reject the null hypothesis based on the entire dataset, and
119 adjusting it involves balancing the risk of making a type I error with the power of the test. Conversely, the

120 threshold in machine learning is a flexible parameter used to classify individual instances, and adjusting it
121 balances precision and recall, impacting how confidently we can classify each instance. For example,
122 when detecting introgressed fragments in genomes, we might choose a threshold that prioritizes high
123 precision to obtain highly confident results for candidate segments. Additionally, power analysis is
124 typically conducted during the experimental design phase, before data collection, to ensure that the study
125 is adequately powered to detect the expected effects (Descôteaux 2007). Although retrospective power
126 analysis can be conducted after experimentation, it is a controversial practice and may be fundamentally
127 flawed (Thomas 1997; Hoenig and Heisey 2001). These fundamental differences highlight why directly
128 equating FDR with one minus precision, or power with recall (or the true positive rate), is inappropriate.

129 Consequently, misusing these metrics can lead to misunderstandings about the distinct focuses,
130 workflows, and methodologies, thereby hindering effective communication and collaboration across
131 disciplines. However, these approaches can be complementary. For example, machine learning can
132 introduce novel methods for controlling FDR, offering new perspectives beyond traditional approaches
133 like the Benjamini-Hochberg procedure and independent hypothesis weighting (Xia et al. 2017). Although
134 hypothesis testing can be employed to evaluate the performance of classification algorithms or to compare
135 the efficacy of different classifiers (Dietterich 1998; Demšar 2006; Alpaydin et al. 2014), it is crucial to
136 clearly define the hypothesis, as hypothesis testing requires a well-defined hypothesis (Wasserman 2003).
137 For example, when assessing the performance of a classifier, one might test whether the probability p that
138 the classifier makes an error is less than or equal to a specified value p_0 . In this scenario, H_0 would be $p \leq$
139 p_0 , while H_1 would be $p > p_0$; statistical tests, such as the binominal test, can then be utilized to evaluate
140 this hypothesis (Alpaydin et al. 2014). In this context, we could discuss the power of these statistical tests
141 rather than the power of the classifier.

142 **Develop Robust Implementation**

143 Machine learning is fundamentally dependent on code, making robust implementation crucial for
144 ensuring the accuracy and validity of results. A robust implementation encompasses both reliable
145 programs and reproducible analyses. Such an implementation not only allows other researchers to
146 replicate experiments and verify outcomes but also facilitates the application of established machine
147 learning approaches to different datasets. This, in turn, enhances our understanding of evolution across
148 various species and populations, which is certainly a motivating intention of the researchers developing
149 such methods. More generally, there is a reproducibility crisis in life sciences (Collins and Tabak 2014;
150 Baker 2016), which can only be addressed by proper implementations. Conversely, a lack of robustness
151 forces potential users to painstakingly examine the code and analysis, possibly leading to reinventing the
152 wheel, which almost certainly slows down the pace of research and causes unnecessary duplication of

153 efforts. As noted by Nekrutenko et al. (2018), truly usable software in evolutionary biology is currently in
 154 short supply, and this issue could be exacerbated in the context of machine learning applications.

155 For example, a recent study by Romieu et al. (2024) evaluated the performance of several methods for
 156 detecting adaptive introgression, including two newly developed machine learning approaches,
 157 genomatnn (Gower et al. 2021) and MaLAdapt (Zhang et al. 2023). This study found that the Q95 statistic
 158 (Racimo et al. 2017) outperformed other methods, which differs from the results reported in the original
 159 studies for genomatnn and MaLAdapt. This discrepancy is likely due to the reliance on pre-trained
 160 models from the original studies, which were trained using simulated data based on specific human
 161 demographic models that differ from the new dataset, particularly since it originates from a different
 162 species. As a result, this mismatch may lead to model misspecification when the pre-trained models are
 163 applied to the new data. The current implementation of these methods makes training new models on
 164 novel datasets difficult, restricting their application to pre-trained models or a limited set of predefined
 165 human demographic models for simulating training data. This may lead to reduced performance when
 166 applied to diverse datasets, especially those from different species.

167 Hence, this suggests that non-robust implementation could lead to inconsistencies in model performance,
 168 particularly when new datasets are involved, ultimately hindering the broader applicability of these
 169 methods. Several recommendations to enhance the reproducibility of machine learning applications in life
 170 sciences have been proposed (Heil et al. 2021; Walsh et al. 2021). For instance, Heil et al. (2021)
 171 proposed a gold standard for reproducibility and emphasized that the entire analysis is reproducible with a
 172 single command, requiring minimal effort to replicate. Building on these recommendations and best
 173 practices from fields such as software development, I propose additional practical approaches to further
 174 promote the development of robust machine learning applications in population genetics (Table 1).

175 **Table 1 Practical approaches for implementing robust machine learning applications**

Phase	Design	Development	Distribution
Approach	<ul style="list-style-type: none"> Object-oriented programming for organizing the code Unified Modeling Language for mapping relation among different modules 	<ul style="list-style-type: none"> Version control tools for documenting development history Unit tests for ensuring code functionality Code review for improving code quality 	<ul style="list-style-type: none"> Package managers for maintaining dependencies and facilitating installation Workflow managers for automating and managing complex tasks

176 To develop reliable code, developers can start with object-oriented programming (OOP), a paradigm that
 177 efficiently organizes complex software by structuring code around objects representing real-world entities

178 or concepts. OOP encapsulates data and behavior within self-contained modules, allowing for
179 independent development, testing, and maintenance. Through inheritance, OOP enables the construction
180 of complex modules from simpler ones, increasing code reusability and reducing redundancy. These
181 principles—encapsulation and inheritance—enhance code flexibility and extensibility, thereby improving
182 reliability by making it easier to modify, extend, and maintain. Familiarity with OOP is also beneficial for
183 understanding machine learning frameworks like PyTorch (Paszke et al. 2019), which are built on OOP
184 concepts. Additionally, Unified Modeling Language can aid in complex software design by visualizing
185 the structure and behavior of software systems.

186 During code development, developers can utilize version control tools such as Git, which enable tracking
187 changes and automatically documenting the development history of the application—a practice
188 recommended by DOME (Walsh et al. 2021). In addition to writing code, developers should implement
189 unit tests to ensure that individual components function correctly and to maintain the integrity of the
190 codebase as it evolves. Developers can use Git and deposit their code on various platforms, such as
191 GitHub, Bitbucket, and Hugging Face (specialized for machine learning). These platforms provide
192 automated workflows for running unit tests, streamlining the coding and testing process, and ensuring
193 continuous integration. Furthermore, these platforms facilitate collaboration, which can improve code
194 quality through peer review, collective problem-solving, and user feedback. For instance, the development
195 guidelines of the PopSim Consortium (Adrion et al. 2020) require that code contributions undergo peer
196 review for validation in their GitHub repository, exemplifying how collaboration on platforms like
197 GitHub can enhance code quality and reliability. This collaborative approach aligns with
198 recommendations to employ reproducibility collaborators for validating machine learning studies (Heil et
199 al. 2021).

200 To further enhance reproducibility, developers can apply package managers for distributing code
201 (Nekrutenko et al. 2018) and workflow managers for streamlining analysis. Package managers—such as
202 Conda and Mamba—automate the installation of dependencies, manage different software environments,
203 and ensure that all necessary packages are available and correctly configured. This is particularly
204 important in machine learning, where packages evolve rapidly, and version conflicts are a common cause
205 of reproducibility issues. By controlling the exact versions of all dependencies (Figure 2A), these tools
206 help mitigate inconsistencies and ensure that the environment in which the code runs remains stable and
207 reproducible. Workflow managers, such as NextFlow (Tommaso et al. 2017) and Snakemake (Mölder et
208 al. 2021), allow developers to define and automate complex workflows, ensuring that each step in the
209 analysis is executed in the correct order and under consistent conditions each time. This is notably
210 valuable in machine learning, as it often involves multiple steps, such as data preprocessing, model

211 training, hyperparameter tuning, and performance evaluation (Figure 2B). Additionally, workflow
212 managers can optimize the use of computing resources by distributing tasks across different machines and
213 even integrating local and cloud servers (Huang et al. 2023). This capability is especially beneficial when
214 working with massive datasets, ensuring that machine learning applications are both scalable and
215 efficient. Together, package managers and workflow managers provide a robust framework for achieving
216 the reliability and reproducibility required in modern machine learning applications. To demonstrate, I
217 implemented a Snakemake workflow (see Data availability) to reproduce the results for detecting
218 introgressed loci using pop_gen_cnn (Flagel et al. 2019). Users can replicate the entire analysis with a
219 single command, aligning with the gold standard suggested by Heil et al. (2021).

```
220 A                                     B  
1 name: prml                           1 ##### Target rules #####  
2 channels:                             2 rule all:  
3   - conda-forge                       3   input:  
4   - bioconda                          4     "results/plots/pop_gen_cnn.original.model.performance.png",  
5 dependencies:                         5     "results/plots/pop_gen_cnn.reproduced.model.performance.png",  
6   - matplotlib=3.8.4                 6 ##### Modules #####  
7   - scikit-learn=1.5.0               7 include: "rules/download.smk"  
8   - snakemake=8.11.6                 8 include: "rules/train_pop_gen_cnn.smk"  
9   - tensorflow=2.16.1                 9 include: "rules/test_pop_gen_cnn.smk"
```

221 **Figure 2 An example illustrating how package managers and workflow managers can facilitate the**
222 **implementation of robust machine learning applications for population genetic inference.** (A) The
223 dependencies and their versions for an application can be recorded in a YAML file, which is then used by
224 the package manager Conda to create a virtual environment. This ensures a reproducible and consistent
225 environment for executing the application. (B) Different steps in the Snakemake workflow, such as data
226 downloading, model training, and testing, can be streamlined into a single file. The specific details of
227 each step can then be documented in separate modules. This approach enables users to configure the
228 appropriate environments for each step within the workflow. For instance, the training process in
229 pop_gen_cnn uses Python 3 and TensorFlow 2 (Abadi et al. 2016), while the evaluation process requires
230 Python 2 and TensorFlow 1 to maintain compatibility with the original model.

231 Besides source code, it is essential to make training and test datasets, model checkpoints, and trained
232 models publicly available. This is particularly important because machine learning often relies on
233 stochastic algorithms, where sharing model checkpoints and trained models helps in examining how
234 randomness affects model performance. For supervised learning-based applications that utilize simulated
235 data as training data—a new paradigm in population genetics (Schrider and Kern 2018)—sharing the
236 simulation scripts and simulated datasets is crucial. These resources not only help users verify findings
237 but also enable them to generate their own training data for similar research questions, for example, when
238 working on different species. Moreover, it is vital to use established simulators, such as ms, msprime, and
239 SLiM (Hudson 2002; Baumdicker et al. 2022; Haller and Messer 2023), and curated demographic

240 models, such as those from the PopSim Consortium (Lauterbur et al. 2023b), for generating simulated
241 data. These tools have been extensively tested by the community, reducing the likelihood of producing
242 unreliable datasets. As examples of such problems, customized versions of ms were used in two recent
243 machine learning applications—ArchIE and IntroUNET—to simulate training data and detect
244 introgressed fragments or alleles in genomes (Durvasula and Sankararaman 2019; Ray et al. 2024). In
245 ArchIE, a variable representing the introgression proportion in the simulated data was introduced into ms,
246 which causes this customized version of ms to fail when no introgression occurs between populations.
247 Additionally, the recipient population—the population receiving genetic material during introgression—
248 must be defined as the first subpopulation in the demographic model; otherwise, the variable cannot
249 accurately report the introgression proportion, limiting its applicability across different species or
250 population structures. Furthermore, while a four-population human demographic model was used to
251 simulate the training data, a three-population model was reported in the paper, resulting in
252 inconsistencies. In IntroUNET, modifications to the seed argument in ms introduced bugs, causing the
253 same training data to be simulated periodically, which raises concerns about the actual performance of
254 IntroUNET. These examples underscore the importance of transparency and the use of reliable, well-
255 tested tools in the simulation process.

256 In conclusion, while machine learning holds great potential for addressing complex evolutionary
257 problems, I believe that its full power can only be realized through precise terminology and robust
258 implementation.

259 **Acknowledgements**

260 X.H. thanks Josef Hackl for testing ArchIE, IntroUNET, genomatnn, and MaLAdapt; Martin Kuhlwilm
261 and Yungang He for their discussions and comments on the manuscript; and the Life Science Compute
262 Cluster at the University of Vienna for providing computing resources.

263 **Competing interests**

264 X.H. declares no conflict of interests.

265 **Data availability**

266 The Snakemake workflow to reproduce the results for detecting introgressed loci using pop_gen_cnn with
267 a single command can be found in <https://github.com/xin-huang/prml>, last accessed August 31, 2024.

268 **References**

269 Abadi M, Barham P, Chen J, Chen Z, Davis A, Dean J, Devin M, Ghemawat S, Irving G, Isard M, et al.
270 2016. TensorFlow: A system for large-scale machine learning. arXiv.
271 <https://arxiv.org/abs/1605.08695>, last accessed August 31, 2024.

272 Adrion JR, Cole CB, Dukler N, Galloway JG, Gladstein AL, Gower G, Kyriazis CC, Ragsdale AP,
273 Tsambos G, Baumdicker F, et al. 2020. A community-maintained standard library of population
274 genetic models. *eLife* **9**: e54967.

275 Alpaydin E. 2014. *Introduction to Machine Learning 3rd edition*. Cambridge: MIT Press.

276 Amin MR, Hasan M, Arnab SP, DeGiorgio M. 2023. Tensor decomposition-based feature extraction and
277 classification to detect natural selection from genomic data. *Mol Biol Evol* **40**: msad216.

278 Arnab SP, Amin MR, DeGiorgio M. 2023. Uncovering footprints of natural selection through spectral
279 analysis of genomic summary statistics. *Mol Biol Evol* **40**: msad157.

280 Baker M. 2016. 1,500 scientists lift the lid on reproducibility. *Nature* **533**: 452–454.

281 Baumdicker F, Bisschop G, Goldstein D, Gower G, Ragsdale AP, Tsambos G, Zhu S, Eldon B, Ellerman
282 EC, Galloway JG, et al. 2022. Efficient ancestry and mutation simulation with msprime 1.0. *Genetics*
283 **220**: iyab229.

284 Benjamini Y, Hochberg Y. 1995. Controlling the false discovery rate: A practical and powerful approach
285 to multiple testing. *J R Stat* **57**: 289–300.

286 Collins FS, Tabak LA. 2014. Policy: NIH plans to enhance reproducibility. *Nature* **505**: 612–613.

287 Demšar J. 2006. Statistical comparisons of classifiers over multiple data sets. *J Mach Learn Res* **7**: 1–30.

288 Descôteaux J. 2007. Statistical power: An historical introduction. *Tutor Quant Methods Psychol* **3**: 28–34.

289 Dietterich TG. 1998. Approximate statistical tests for comparing supervised classification learning
290 algorithms. *Neural Comput* **10**: 1895–1923.

291 Durvasula A, Sankararaman S. 2019. A statistical model for reference-free inference of archaic local
292 ancestry. *PLoS Genet* **15**: e1008175.

293 Durvasula A, Sankararaman S. 2020. Recovering signals of ghost archaic introgression in African
294 populations. *Sci Adv* **6**: eaax5097.

295 Fawcett T. 2006. An introduction to ROC analysis. *Pattern Recognit Lett* **27**: 861–874.

296 Flagel L, Brandvain Y, Schrider DR. 2019. The unreasonable effectiveness of convolutional neural
297 networks in population genetic inference. *Mol Biol Evol* **36**: 220–238.

298 Gower G, Picazo PI, Fumagalli M, Racimo F. 2021. Detecting adaptive introgression in human evolution
299 using convolutional neural networks. *eLife* **10**: e64669.

300 Haller BC, Messer PW. 2023. SLiM 4: Multispecies eco-evolutionary modeling. *Am Nat* **201**: E127–
301 E139.

302 Heil BJ, Hoffman MM, Markowitz F, Lee SI, Greene CS, Hicks SC. 2021. Reproducibility standards for
303 machine learning in the life sciences. *Nat Methods* **18**: 1132–1135.

304 Hoenig JM, Heisey DM. 2001. The abuse of power. *Am Stat* **55**: 19–24.

305 Huang X, Struck TJ, Davey SW, Gutenkunst RN. 2023. dadi-cli: Automated and distributed population
306 genetic model inference from allele frequency spectra. bioRxiv.
307 <https://doi.org/10.1101/2023.06.15.545182>, last accessed August 31, 2024.

308 Huang X, Rymbekova A, Dolgova O, Lao O, Kuhlwillm M. 2024. Harnessing deep learning for
309 population genetic inference. *Nat Rev Genet* **25**: 61–78.

310 Hudson RR. 2002. Generating samples under a Wright-Fisher neutral model of genetic variation.
311 *Bioinformatics* **18**: 337–338.

312 Korfmann K, Gaggiotti OE, Fumagalli M. 2023. Deep learning in population genetics. *Genome Biol Evol*
313 **15**: evad008.

314 Lauterbur ME, Munch K, Enard D. 2023a. Versatile detection of diverse selective sweeps with Flex-
315 Sweep. *Mol Biol Evol* **40**: msad139.

316 Lauterbur ME, Cavassim MIA, Gladstein AL, Gower G, Pope NS, Tsambos G, Adrion J, Belsare S,
317 Biddanda A, Caudill V, et al. 2023b. Expanding the stdpopsim species catalog, and lessons learned for
318 realistic genome simulations. *eLife* **12**: RP84874.

319 Mölder F, Jablonski KP, Letcher B, Hall MB, Tomkins-Tinch CH, Sochat V, Forster J, Lee S, Twardziok
320 SO, Kanitz A, et al. 2021. Sustainable data analysis with Snakemake. *F1000 Res* **10**: 33.

321 Nekrutenko A, Galaxy Team, Goecks J, Taylor J, Blankenberg D. 2018. Biology needs evolutionary
322 software tools: Let’s build them right. *Mol Biol Evol* **35**: 1372–1375.

323 Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, Killeen T, Lin Z, Gimelshein N, Antiga L, et
324 al. 2019. PyTorch: An imperative style, high-performance deep learning library. *Adv Neural Inf*
325 *Process Syst* **32**: 8026–8037.

326 Racimo F, Marnetto D, Huerta-Sánchez E. 2017. Signatures of archaic adaptive introgression in present-
327 day human populations. *Mol Biol Evol* **34**: 296–317.

328 Raff E. 2019. A step toward quantifying independently reproducible machine learning research. *Adv*
329 *Neural Inf Process Syst* **32**: 5485–5495.

330 Ray DD, Flagel L, Schrider DR. 2024. IntroUNET: Identifying introgressed alleles via sematic
331 segmentation. *PLoS Genet* **20**: e1010657.

332 Romieu J, Camarata G, Crochet PA, de Navascués M, Leblois R, François Rousset. 2024. Performance
333 evaluation of adaptive introgression classification methods. bioRxiv.
334 <https://doi.org/10.1101/2024.06.12.598278>, last accessed August 31, 2024.

335 Schrider DR, Kern A. 2018. Supervised machine learning for population genetics: A new paradigm.
336 *Trends Genet* **34**: 301–312.

337 Shreffler J, Huecker MR. 2024. Type I and type II errors and statistical power. In *StatPearls*. Treasure
338 Island: StatPearls Publishing.

339 Soric B. 1989. Statistical "discoveries" and effect-size estimation. *J Am Stat Assoc* **84**: 608–610.

340 Thomas L. 1997. Retrospective power analysis. *Conserv Biol* **11**: 276–280.

341 Tommaso PD, Chatzou M, Floden EW, Barja PP, Palumbo E, Notredame C. 2017. Nextflow enables
342 reproducible computational workflows. *Nat Biotechnol* **35**: 316–319.

343 Walsh I, Fishman D, Garcia-Gasulla D, Titma T, Pollastri G, ELIXIR Machine Learning Focus Group,
344 Harrow J, Psomopoulos FE, Tosatto SCE. 2021. DOME: Recommendations for supervised machine
345 learning validation in biology. *Nat Methods* **18**: 1122–1127.

346 Wasserman L. 2003. *All of Statistics: A Concise Course in Statistical Inference*. New York: Springer.

347 Xia F, Zhang MJ, Zou J, Tse D. 2017. NeuralFDR: Learning discovery thresholds from hypothesis
348 features. *Adv Neural Inf Process Syst* **30**: 1540–1549.

349 Yelmen B, Jay F. 2023. An overview of deep generative models in functional and evolutionary genomics.
350 *Annu Rev Biomed Data Sci* **6**: 173–189.

351 Zhang X, Kim B, Singh A, Sankararaman S, Durvasula A, Lohmueller KE. *MaLAdapt* reveals novel
352 targets of adaptive introgression from Neanderthals and Denisovans in worldwide human populations.
353 *Mol Biol Evol* **40**: msad001.