

SOFTWARE

Open Access



SNVstory: inferring genetic ancestry from genome sequencing data

Audrey E. Bollas^{1,2}, Andrei Rajkovic¹, Defne Ceyhan¹, Jeffrey B. Gaither¹, Elaine R. Mardis^{1,2} and Peter White^{1,2*}

*Correspondence:
peter.white@nationwidechildrens.org

¹The Steve and Cindy Rasmussen Institute for Genomic Medicine, The Abigail Wexner Research Institute, Nationwide Children's Hospital, Columbus, OH, USA

²Department of Pediatrics, The Ohio State University College of Medicine, Columbus, OH, USA

Abstract

Background: Genetic ancestry, inferred from genomic data, is a quantifiable biological parameter. While much of the human genome is identical across populations, it is estimated that as much as 0.4% of the genome can differ due to ancestry. This variation is primarily characterized by single nucleotide variants (SNVs), which are often unique to specific genetic populations. Knowledge of a patient's genetic ancestry can inform clinical decisions, from genetic testing and health screenings to medication dosages, based on ancestral disease predispositions. Nevertheless, the current reliance on self-reported ancestry can introduce subjectivity and exacerbate health disparities. While genomic sequencing data enables objective determination of a patient's genetic ancestry, existing approaches are limited to ancestry inference at the continental level.

Results: To address this challenge, and create an objective, measurable metric of genetic ancestry we present SNVstory, a method built upon three independent machine learning models for accurately inferring the sub-continental ancestry of individuals. We also introduce a novel method for simulating individual samples from aggregate allele frequencies from known populations. SNVstory includes a feature-importance scheme, unique among open-source ancestral tools, which allows the user to track the ancestral signal broadcast by a given gene or locus. We successfully evaluated SNVstory using a clinical exome sequencing dataset, comparing self-reported ethnicity and race to our inferred genetic ancestry, and demonstrate the capability of the algorithm to estimate ancestry from 36 different populations with high accuracy.

Conclusions: SNVstory represents a significant advance in methods to assign genetic ancestry, opening the door to ancestry-informed care. SNVstory, an open-source model, is packaged as a Docker container for enhanced reliability and interoperability. It can be accessed from <https://github.com/nch-igm/snvstory>.

Keywords: Genetic ancestry prediction, Machine learning, Genetic variation, Model interpretation, Personalized medicine

Introduction

Ancestry derived from genomic data, referred to as genetic ancestry, is a measurable and biologically defined parameter. Although much of the human genome is identical across all populations, it is estimated that depending on an individual's ancestry,



© The Author(s) 2024. **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.

0.1–0.4% may differ from the human reference genome. While this genetic variation includes structural variants (SVs), copy number variants (CNVs), and small insertions or deletions (indels), by far the largest and easiest to detect category occurs in the form of single nucleotide variants (SNVs), many of which are unique to genetically distinct populations [1].

Knowledge of a patient's genetic ancestry has clinical implications, ranging from genetic testing to health screening based on ancestral disease-predisposition rates, and in some cases, may inform what medicine dosage to prescribe a patient [2–4]. However, self-reported race is frequently used in the research and clinical setting and is often inconsistent with genetic ancestry, potentially driving health disparities [5–8]. Genome sequencing-based diagnostic testing in patients suspected of having a rare genetic disorder requires accurate data filtering to remove variants common to a given population. Precise identification of the patient's ancestry improves the identification of rare disease-causal variants. Therefore, developing methods to report ancestry accurately and consistently is essential.

In addition to clinical importance, knowing the ancestral composition of an individual or a population is essential in the genetic research setting. For example, signals from genome-wide association studies (GWAS) or whole genome sequencing cohorts can be reassessed based on population stratification, whereby loci associated with disease may be more accurately identified by discarding rare variants associated with an individual's ancestry rather than with the disease in question [9, 10].

Given the importance of ancestry, several ancestry inference algorithms that operate on genomic data have been developed that can be divided into two broad types: parametric and non-parametric. Parametric learning algorithms estimate a finite set of parameters from the data to establish a relationship between the independent and dependent variables. Two widely used parametric tools are STRUCTURE [11] and ADMIXTURE [12], which estimate the proportions of different ancestries (or ancestral populations) for each individual, known as admixture. Recently, Archetypal analysis was shown to be more computationally efficient and provide more interpretable results than ADMIXTURE [13]. In contrast, non-parametric methods do not have a finite set of parameters and instead rely on the intrinsic structure of the data to determine which data points best resemble each other.

The emergence of population-scale genome sequencing datasets with a form of self-reported ancestry allows models to be built with prior knowledge of represented ancestries. In place of individualized genetic data, large databases house genomic summary results, such as aggregate variant allele frequencies stratified by population. For example, the Single Nucleotide Polymorphism database (dbSNP) is the largest genomic aggregate database with 11 different populations from over one million samples [14]. However, the 11 distinct populations contain a high degree of overlap and primarily represent continental groupings [15]. The Genome Aggregation Database (gnomAD) is another aggregate database with allele frequencies from 140,000 subjects from 26 populations [16]. In addition to these large-scale repositories of aggregate allele frequencies, there exist a few datasets at the level of the individual, such as the 1000 Genomes Project (1kGP) [1] and the Simons Genome Diversity Project (SGDP) [17], which are much smaller in sample size, with 2504 and 279 samples, respectively. Nevertheless, the 1kGP and SGDP

have been critical in characterizing ancestry and human history as they contain the most granular population labels.

Taken together, these curated variant datasets enable an alternative class of models to be used to predict ancestry based on samples labeled with known ancestry [18–28]. However, many methods suffer shortcomings, including not having discrete ancestry labels beyond the main continental groups or, for those methods using the 1kGP, not considering that many subjects are within the same families and, therefore, fail to satisfy the principle of independent and identically distributed data. As such, there is a critical need for methods to accurately predict an individual's genetic ancestry from genome sequencing data by implementing supervised models.

Here, we address some limitations surrounding supervised learning of ancestry by developing three independent models from gnomAD, 1kGP, and SGDP. Our models estimate ancestry from 36 different populations with high accuracy. Furthermore, we provide software that enables users to run our models on their data, taking the widely accepted variant call format (VCF) files as input and outputting predictions and a graphical representation of the likelihood of a given genetic ancestry. Additionally, we provide feature importance by cytolocation, to help understand model prediction results. As a form of validation, we apply these models to our in-house clinical research dataset and correlate the estimates with those of self-reported ancestry.

Methods

Training datasets

Genomic datasets from gnomAD, 1kGP, and SGDP were processed separately (Fig. 1), as described below. The gnomAD variants are provided on reference genome GRCh37, and the 1kGP and SGDP were called on reference genome GRCh38. We performed a liftover between genome versions GRCh37 and GRCh38 for all variants in the training datasets. SNVstory uses this to convert the genome coordinates of the input variants to the corresponding genome version (GRCh37 for gnomAD, GRCh38 for SGDP and 1kGP) so that variants from either reference can be used.

The genome aggregation database (gnomAD)

The gnomAD v2.1 exome and genome sequencing variant dataset provides aggregated data from 17 populations, meaning allele frequencies of each population for 17 million exome variants. We reduced the number of input features for machine learning by following a similar protocol to the one described by the MacArthur lab by filtering for high call rates, biallelic-only sites, and a frequency greater than 0.1% (<https://macarthurlab.org/2018/10/17/gnomad-v2-1/>). After this filtering, 81,398 SNVs remained, formatted as a matrix of ancestries and corresponding SNV frequencies. A comparison between a model based on all genomic regions vs. one based only on exomes indicated that restriction to exonic regions resulted in no measurable loss of accuracy. This result is in line with biological and statistical intuition, which suggests that exonic variants with MAFs > 1% are unlikely to contribute to fitness, and the vastness of the genome allows us to build a saturated model even with the limited set of exonic variants. Therefore, we chose to continue using the exon filter to reduce computational resources and time for model training.

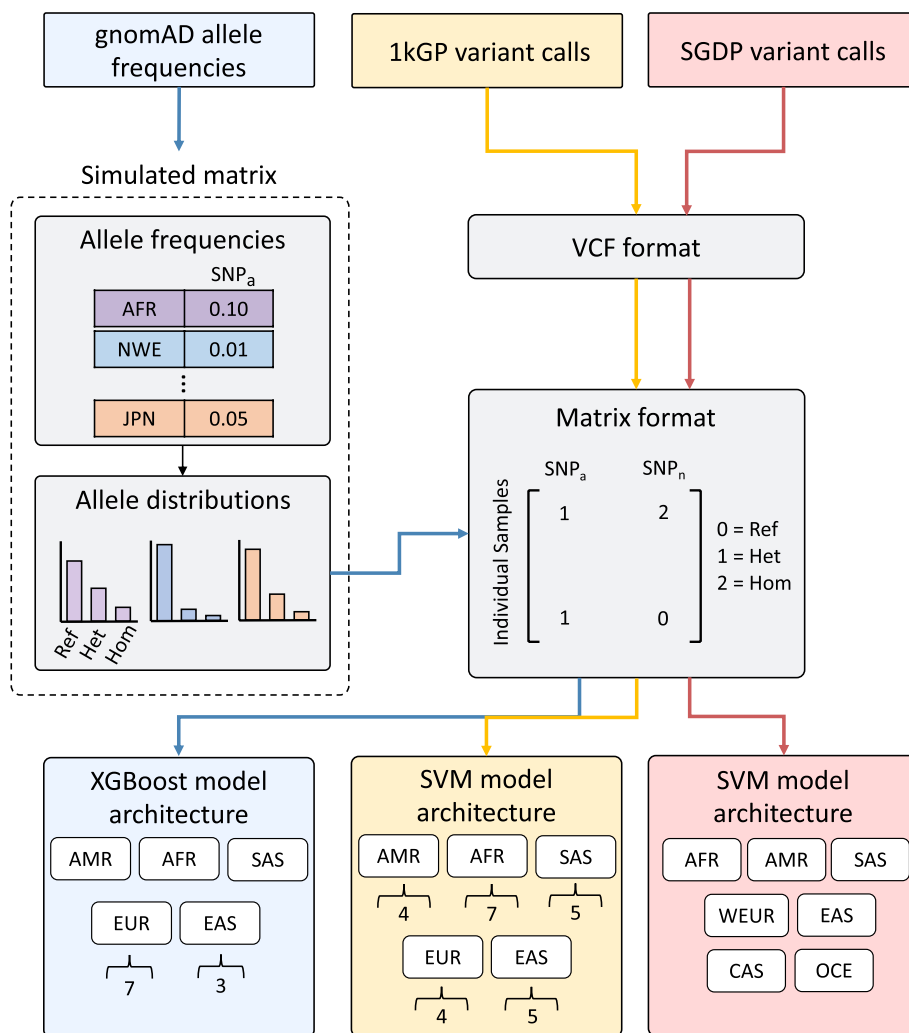


Fig. 1 Schematic of ancestry inference model strategy. The workflow visualizes each dataset separately with colored boxes and arrows: gnomAD (blue), 1kGP (yellow), and SGDP (red). For the gnomAD synthetic-based matrix, allele frequencies for each variant for each population given in gnomAD are used to create a distribution of reference, heterozygous and homozygous alleles for each population. A matrix format is created by converting the distributions into 0s, 1s, and 2s for each locus for samples in each population. For 1kGP and SGDP, a matrix format is built directly from variants in the VCF. For the model architecture, continental model labels are shown in white boxes, and the number of labels in the corresponding subcontinental models is below in brackets

To obtain SNV calls for individuals, as is provided in standard VCF format, we simulated individuals from each ancestry by generating heterozygous and homozygous SNV calls with probability weighted by their allele distributions (Fig. 1). We determined the number of samples to simulate and which SNVs to keep as features using cross-validation (see the gnomAD section in model training and cross-validation). This resulted in a matrix of simulated samples spanning the ancestry classifications in gnomAD v2.1 with SNVs, coded as reference, heterozygous, or homozygous for each SNV position. Although this approach does not capture haplotypes, the simulated samples are genetically typical examples of the chosen ancestry to a first approximation.

The 1000 Genomes Project (1kGP)

The New York Genome Center performed genome sequencing (GS) on 3202 samples, including 602 trios, from the 1kGP cohort at 30× coverage, released in 2020 [29]. The data were aligned to GRCh38 using BWA-MEM [30], and variants were called by GATK *HaplotypeCaller* (GATK version 3.5.0) using default settings. The dataset contains 126,659,422 SNVs from 26 populations spanning East and South Asia, North and South America, Africa, and Europe. Sample sizes were not uniformly represented across the different populations, i.e., the dataset was imbalanced. Due to the high genetic similarity between individuals from Utah and the United Kingdom, the Utah population was removed from the analysis.

The Simons Genome Diversity Project (SGDP)

The SGDP consists of GS of 300 individuals from seven major population groups, 75 countries, and 142 diverse populations. GS FASTQ files from 279 samples were downloaded from the European Nucleotide Archive (PRJEB9586). Sequencing reads were aligned to genome assembly GRCh38 using BWA-MEM. SNV and INDEL calling was performed with GATK version 4.1.9, described below. GATK *HaplotypeCaller* was run on each sample using the GVCF workflow to generate a per-sample intermediate GVCF. The GATK *GenotypeGVCFs* function was used to perform base calling across all samples jointly to obtain genotypes for each sample in VCF format. We then performed variant recalibration and filtering in the two-stage process using the GATK functions *VariantRecalibration* and *ApplyVQSR*. The final combined data set contained a total of 48,815,712 SNVs.

Quality control

Quality control of the gnomAD (<https://macarthurlab.org/2018/10/17/gnomad-v2-1/>) and 1kGP [29] were as previously described. For the SGDP dataset, we ran several quality-control tools to detect any issues with sequencing quality and sample contamination. We ran Picard *CollectMultipleMetrics* on the aligned bam files to collect alignment summary, quality score, and GC bias metrics (Additional file 1: Table S1). Sequencing read allocation was calculated using samtools. Coverage information was collected using mosdepth [31]. The average coverage for all realigned samples was 40× (ranging from 31 to 77×). Sample contamination level was determined by the percentage of reads inconsistent with the allele frequencies in dbSNP [14] sites, using VerifyBamID [32]. One sample was flagged for possible sample contamination due to an estimated level between 0.1 and 2%. (Additional file 1: Methods).

Removal of related samples

Related samples of the third-degree (e.g., first cousins, great-grandparents, or great-grandchildren) or closer were identified by the relationship inference tool, KING [33]. Data from the 1kGP and SDGP were preprocessed using PLINK2 with the following parameters: “*-new-id-max-allele-len 10,000 -max-alleles 2*” [34]. KING recommends performing as little filtering as possible. However, an additional filtering step was performed to prevent the computation from running out of memory. Therefore, the

analysis was restricted to variants shared by at least two individuals: “*-maf 0.0007*” in the case of the 1kGP and “*-maf 0.007*” for SDGP. After removing the variants present in only one sample, KING was executed on the resulting bed file, with the “*-kinship*” option set to report pairwise relatedness inference. Samples from the analysis were flagged that had a third-degree kinship coefficient cutoff ≥ 0.0442 , a value previously established by the authors of KING [33]. Four samples were removed from further analysis in the SGDP dataset based on the KING relatedness results (Additional file 1: Methods).

Because some samples from the 1kGP are related to more than one other individual in the cohort, the following procedure was implemented to remove the fewest number of samples. Considering only the relationships with coefficients exceeding the third-degree cutoff, a graph-based method was implemented to recursively identify nodes (samples) with the largest number of edges (relationships) and remove those nodes until all subgraphs had, at most, a single connection. For subgraphs with a single connection, one sample was randomly selected from the pair, while all singletons were included in the list of samples to keep. From 167 samples with at least one close relationship, 117 were flagged for inclusion in downstream analysis. The remaining samples were removed with PLINK2.

Variant selection and preprocessing

Variants from 1kGP and SGDP underwent a final filtering step by taking the intersection of targeted exonic regions of the exome capture reagent used routinely in our clinical lab (IDT xGen Exome Hyb Panel v2 targets hg38 BED file) with the set of genetic variants from the unrelated individuals using BEDTools *intersect* (v2.30.0) [35]. A total of 281,092 and 97,995 variants were used as features in the 1kGP and SGDP models, respectively. The resulting VCF was converted into a numerical encoding homozygous alternative = 2, heterozygous = 1, reference or missing = 0. The vectors of genotypes were combined to form a matrix of variants by genotypes. For variant selection from gnomAD, see the following gnomAD section in Model training and cross-validation below.

Model training and cross-validation

The models were trained on each dataset separately, as required by their differing labeling strategies (Fig. 1).

gnomAD

Because our gnomAD algorithm uses simulated data, we must consider two parameters: a population size that balances the model’s accuracy with training time and resources, and a p value from a Chi-Square test that removes uninformative SNVs. This was accomplished using a nested for loop to iterate over all combinations of population sizes and p values for SNV removal (Additional file 1: Fig. S1). For each combination, we generated a set of 80/20 training/validation splits of the data. A Chi-Square test was applied to each SNV (feature) in the training data to determine whether it was informative for distinguishing ancestry in the population. SNVs were removed that did not meet the p value threshold. We used a gradient-boosted decision tree from XGBoost to train the model on the training set and then test on the validation set [36]. Fold generation and

training were performed five times for each p value, and the accuracy was averaged to represent the accuracy for each p value. Once all the p values were tested, the p value with the highest accuracy was selected (Additional file 1: Fig. S2). Then, the model was retrained on all the data for that specific population size and tested on a simulated hold-out set. The accuracy for the hold-out set is representative of that population. A continental model (population size of 4084 individuals; SNV p value threshold of $7.5e-49$) was built to predict six groups: Africa, South Asia, Europe, East Asia, America, and Ashkenazi Jewish. Two sub-continental classifiers were built to predict ancestry within the East Asian (Additional file 1: Fig. S2A; population size of 13,593 individuals; SNV p value threshold of $1.78e-09$) and European groups (Additional file 1: Fig. S2B; population size of 45,243 individuals; SNV p value threshold of $1.78e-24$).

1kGP

For the 1kGP dataset, the support vector machine (SVM) library from scikit-learn [37] was used to train a classifier to predict the continental groups: Africa, Europe, South Asia, East Asia, and America. In addition, multiple classifiers were trained independently for each sub-continental group, i.e., Kenya or African Caribbean in Barbados. All SVMs were trained using the radial basis function (RBF) kernel and with the gamma parameter fixed as the default. Hyperparameter tuning of the C penalty term was accomplished by performing cross-validation using the scikit-learn stratified k-fold library. The default five splits were chosen, and the shuffle variable was set to true. The F1 macro average was selected to represent a model's performance.

SGDP

The SVM library from scikit-learn was used to train the model for the SGDP dataset. Stratified k-fold cross-validation was performed using the standard scikit-learn library. Seven continental groups were predicted from this cohort (Africa, West Eurasia, East Asia, South Asia, Oceania, Central Asia Siberia, and America), as the subcontinental groups needed more samples per group to train an accurate model. The F1 macro average was chosen as a representation of a model's performance to account for the imbalanced data.

Results

Model performance

We report the performance of the gnomAD, 1kGP, and SGDP continental models using external validation sets (Fig. 2A–F), and cross-validation results on the subcontinental models (Additional file 1: Figs. S2 and S3) were performed because additional datasets with the same subcontinental labels were not available.

Confusion matrices are shown in Fig. 2A–D, providing the ancestry prediction for each sample in the validation data. In the 1kGP and SGDP models, we see some discrepancies between the European and American groups. In the case of the 1kGP model (Fig. 2A), some SGDP samples labeled as European are predicted to be American. Similarly, in the SGDP model, some 1kGP samples labeled as American are predicted as European. This may be due to a higher similarity of the feature space between European and American samples than other groups (Additional file 1: Fig. S3). The gnomAD model is validated

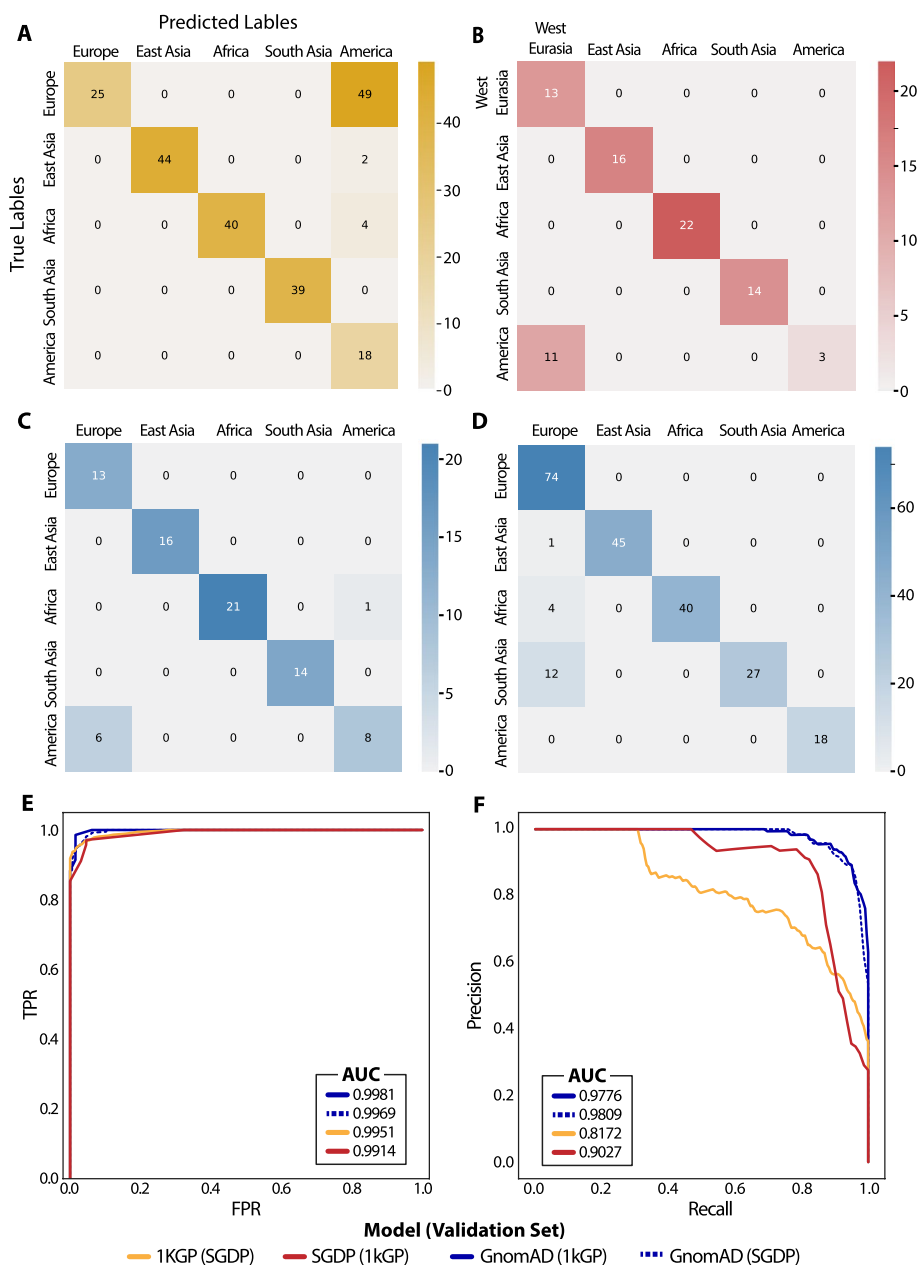


Fig. 2 Continental ancestry inference model performance. **A–D** Confusion matrices of the 1kGP model using SGDP as validation (**A**), SGDP model using 1kGP as validation (**B**), gnomAD model using 1kGP as validation (**C**), and gnomAD model using SGDP as validation (**D**). **E** Macro-averaged ROC curves. **F** Macro-averaged precision–recall curves

with 1kGP (Fig. 2C) and SGDP (Fig. 2D) samples. Overall, all continental models have a high area under the curve (AUC) in both ROC (Fig. 2E) and precision–recall (Fig. 2F) curves, described in the figure legend.

The gnomAD East Asian and European subcontinental models have accuracies of 99.90% and 80.92%, respectively (Additional file 1: Fig. S2A, B). The results for the 1kGP subcontinental model are obtained by averaging the probabilities for each sample across cross-validation folds and then computing the confusion matrix (Additional file 1: Fig.

Table 1 Model performance comparison

	Performance metric		
	ROC AUC	P-R AUC	F1 score
<i>Continental</i>			
SNVstory	0.996	0.980	0.918
RFMix	0.997	0.981	0.922
ADMIXTURE	0.946	0.477	0.537
<i>Subcontinental</i>			
SNVstory	0.986	0.965	0.885
RFMix	0.992	0.771	0.665
ADMIXTURE	0.977	0.937	0.783

A benchmarking analysis was performed using SNVstory and two alternative ancestry inference tools, ADMIXTURE and RFMix

S4). The accuracies for the 1kGP subcontinental models are as follows: Africa, 90.26%; America, 93.06%; East Asia, 87.23%; Europe, 94.29%; South Asia, 85.86%. The averaged probabilities were used to compute the AUC for the ROC and precision–recall curves.

We evaluate the performance of SNVstory compared to two other ancestry inference tools, ADMIXTURE [12] and RFMix [19]. The steps used to run and summarize results from ADMIXTURE and RFMix are briefly described (Additional file 1: Methods). Their performances are summarized using AUC of the ROC and precision–recall curves and averaged F1 scores from the model classifications (Table 1). The continental models are evaluated using a subset of samples from both the 1kGP and SGDP datasets, and the gnomAD model from SNVstory was used for the continental model comparison so that there is no overlap between samples used for training and validation. The RFMix and ADMIXTURE subcontinental results are evaluated using a two-thirds split of 1kGP samples for reference validation to provide comparable metrics to our cross-validation results.

RFMix and SNVstory perform comparably at the continental level, where both models have a slight tendency to predict some American and African samples to be European (Additional file 1: Fig. S5). ADMIXTURE has high true positive rates across all continental groups but also has an increased tendency to classify samples as American. ADMIXTURE shows an increase in model performance for the subcontinental results, while RFMix shows a decrease in subcontinental performance. SNVstory has more consistently high results across both continental and subcontinental models.

Feature interpretation

Feature importance for the gnomAD continental model was calculated using SHAP [38] values to provide insight into which SNVs and their corresponding genes have the most impact on the model predictions. SHAP values for the 1kGP and SGDP models were not calculated because the memory requirement for the kernel explainer was too high due to the number of features in the models.

Global feature importance for the gnomAD continental model is reported by aggregating SHAP values across each gene and taking the mean absolute value of each gene across 2800 of the training samples (Fig. 3). The ‘knownCanonical’ genes table was

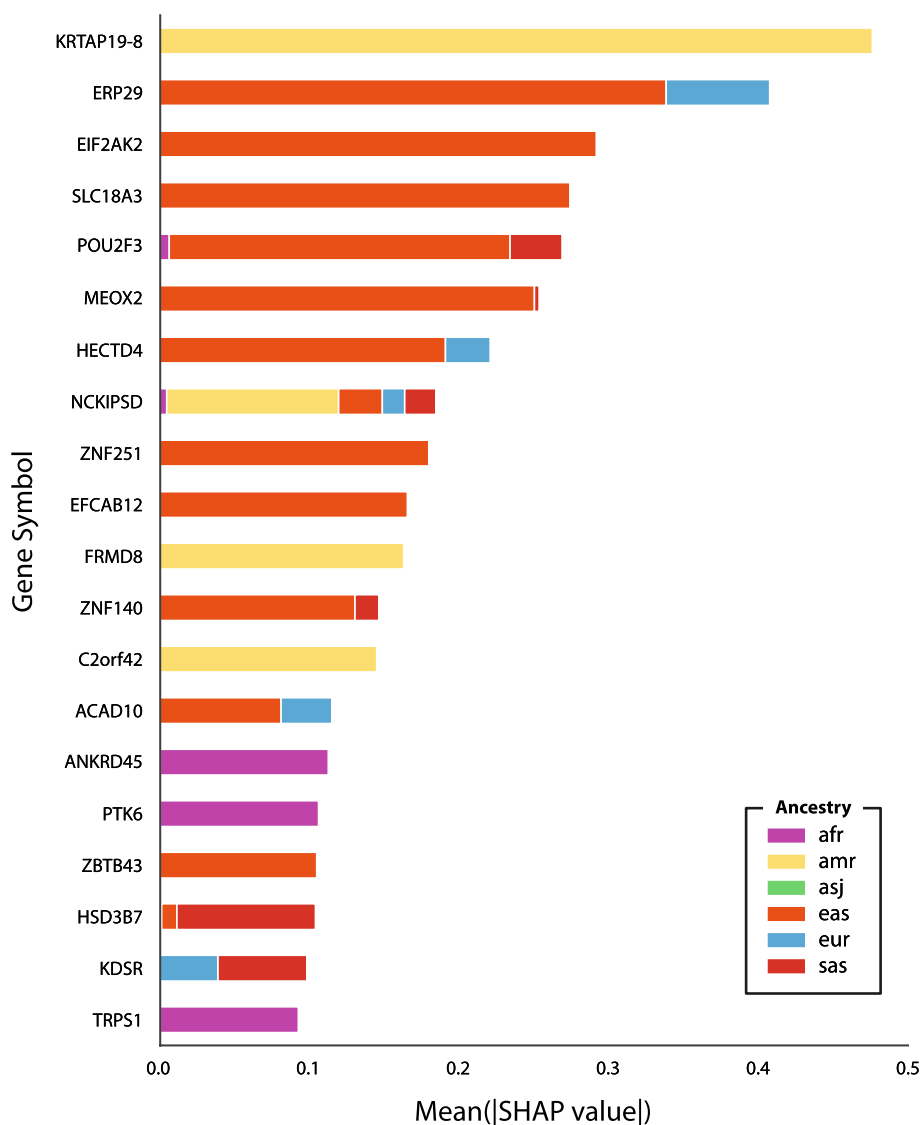


Fig. 3 Gene-level global feature importance in ancestry inference using SNVstory’s gnomAD continental model. This figure illustrates the mean absolute SHAP values aggregated for each gene, derived from 2800 training samples. The analysis highlights the top 20 genes that significantly influence ancestry inference, emphasizing the role of specific alleles in determining ancestry labels

downloaded from the UCSC Table Browser using assembly GRCh37 to get the genomic interval for each gene. If a region contains multiple genes, we combine the genes to form a non-overlapping genomic interval (e.g., ANKRD45, TEX50). Of the 77,402 SNVs used to train the model, 3231 were not located in gene regions and were removed from further analysis. The most significant gene impacting the model is KRTAP19-8 (Keratin Associated Protein 19-8). Samples with a variant in this gene are more likely to be predicted as American. However, some top-ranking genes on this list have a strong negative impact on model prediction, and variants in these locations reduce the probability of a given label. For example, ERP29 has a strong negative impact on predicting South Asian

ancestry. SNVstory provides the ability to detect variants in genes inherited from specific populations, and the highest ranking genes in terms of positive SHAP value across each continental group are provided in Additional file 1: Table S2.

We also aggregated SHAP values across larger cytoloactions to visualize which regions across the genome are most impactful in the model predictions (accessed using this file: (<https://hgdownload.soe.ucsc.edu/goldenPath/hg19/database/cytoBand.txt.gz>). Additional file 1: Fig. S6 shows the feature importance for an individual from the training data labeled as African. Regions are colored by population label with the maximum absolute SHAP value. Regions that have the most impact on predicting the sample African are 'chromosome 1: 172,900,000–176,000,000' and 'chromosome 5: 63,200,000_66,700,000'.

Comparison of genetic versus self-reported ancestry in clinical samples

SNVstory was implemented on an in-house dataset of clinical exome sequencing testing from 293 individuals generated by the Institute for Genomic Medicine Clinical Laboratory to demonstrate the application of our models. We compare the model predictions to the self-reported ancestry of the proband (Additional file 1: Table S3). Self-reported race is derived from the paternal/maternal ethnic background. Ethnicity is categorized into one of three groups: Non-Hispanic or Latino, Hispanic or Latino, and unknown/not reported ethnicity. Race is classified into one of five groups: White, Asian, bi-racial/multi-racial, Black or African American, and Unknown/Unspecified. Due to the broadness of these categories, we report the comparison between predicted genetic ancestry for the continental models only (Table 2).

Most of the individuals share agreement between genetic ancestry and ethnicity/race, e.g., for those predicted to be European, a match of White/Non-Hispanic or Latino for race/ethnicity occurs in 92.5%, 96.7%, and 89.1% of individuals by the gnomAD (Table 2A), 1kGP (Table 2B), and SGDP (Table 2C) models, respectively. However, several cases exist where individuals are self-reported as White while having a different genetic ancestry across multiple models, and vice versa. Additionally, 13 of our cases have either Unknown/Not Reported Ethnicity or Unknown/Unspecified Race. As discussed in the Introduction, the ability to refine or add genetic ancestry information in these cases is helpful for added diagnostic precision in variant filtering/prioritization.

Model interpretation for indeterminate samples

Most of our in-house dataset has agreement across all three continental models (81.9% of samples) and even more across at least two continental models (98.0%). A disproportionate number of individuals share disagreement across all three models between those that are self-reported as bi-racial or multi-racial versus those that are White, Asian, Black or African American (50% vs. 9% disagreement, respectively). Those individuals with unknown/unspecified race are not included in this calculation. These results suggest our models have worse performance on admixed samples, where two or more populations may be present. In reporting results, we use the label with the highest probability. Some discrepancies between model results may be mitigated by adding a minimum threshold on the probability required to obtain a result.

Table 2 Genetic ancestry versus self-reported ethnicity and race

Model labels	Ethnicity	Race	Counts	
<i>(A) gnomAD</i>				
afr	Non-Hispanic or Latino	Black or African American	20	
		Bi-racial/multi-racial	10	
	Unknown/not reported ethnicity	Bi-racial/multi-racial	3	
	Hispanic or Latino	Bi-racial/multi-racial	2	
amr	Non-Hispanic or Latino	White	1	
		Hispanic or Latino	White	8
	Unknown/unspecified	Black or African American	5	
		Black or African American	2	
		White	1	
asj	Non-Hispanic or Latino	Bi-racial/multi-racial	1	
	Non-Hispanic or Latino	White	1	
eas	Non-Hispanic or Latino	Asian	3	
		White	2	
eur	Hispanic or Latino	Bi-racial/multi-racial	1	
	Non-Hispanic or Latino	White	210	
		Bi-racial/multi-racial	5	
	Hispanic or Latino	Bi-racial/multi-racial	5	
		White	3	
	Unknown/not reported ethnicity	White	3	
sas	Hispanic or Latino	Bi-racial/multi-racial	1	
	Non-Hispanic or Latino	Unknown/unspecified	1	
		Asian	3	
<i>(b) 1kGP</i>	afr	Black or African American	19	
		Bi-racial/multi-racial	2	
		Bi-racial/multi-racial	1	
	amr	Unknown/not reported ethnicity	Bi-racial/multi-racial	1
		Hispanic or Latino	White	12
		Non-Hispanic or Latino	Bi-racial/multi-racial	10
		Hispanic or Latino	Bi-racial/multi-racial	8
		Non-Hispanic or Latino	White	8
		Hispanic or Latino	Unknown/unspecified	6
		Hispanic or Latino	Black or African American	2
		Unknown/not reported ethnicity	Bi-racial/multi-racial	2
	Non-Hispanic or Latino	Black or African American	1	
	eas	Non-Hispanic or Latino	Asian	3
	eur	Non-Hispanic or Latino	White	207
White			3	
Non-Hispanic or Latino		Bi-racial/multi-racial	3	
Unknown/not reported ethnicity		Bi-racial/multi-racial	1	
sas	Non-Hispanic or Latino	Asian	3	
		White	1	

Table 2 (continued)

Model labels	Ethnicity	Race	Counts
<i>(C) SGDP</i>			
Africa	Non-Hispanic or Latino	Black or African American	20
		Bi-racial/multi-racial	9
	Hispanic or Latino	Bi-racial/multi-racial	3
	Unknown/not reported ethnicity	Bi-racial/multi-racial	3
	Hispanic or Latino	Black or African American	2
CentralAsiaSiberia	Non-Hispanic or Latino	White	1
		Hispanic or Latino	Unknown/unspecified
EastAsia	Non-Hispanic or Latino	White	1
		Asian	3
SouthAsia	Hispanic or Latino	White	4
		Non-Hispanic or Latino	Asian
	Non-Hispanic or Latino	White	3
WestEurasia	Non-Hispanic or Latino	White	212
		Hispanic or Latino	White
	Non-Hispanic or Latino	Bi-racial/multi-racial	6
		Hispanic or Latino	Bi-racial/multi-racial
	Unknown/not reported ethnicity	Unknown/unspecified	3
		White	3
		Bi-racial/multi-racial	1

Value counts of genetic ancestry model predictions trained using gnomad (A), 1kGP (B), and SGDP (C) compared to self-reported ethnicity and race

Individualized ancestry report

Here, we illustrate the ability of SNVstory to provide ancestry predictions in an easily visualized format for individual samples (Fig. 4). The probabilities for the gnomAD and the 1kGP continental models were 100% European, while the SGDP continental model was 95% West Eurasia. Because the subcontinental models are trained and executed separately, SNVstory uses a weighting scheme to rate subcontinental label probabilities according to their continental output. Subcontinental model probabilities are multiplied by their corresponding continental model probability, and this result is reported. The gnomAD subcontinental model has the highest probability (48%) for North-Western European (nfe_nwe), and the 1kGP subcontinental model has the highest probability (100%) for British from England and Scotland (eur_gbr). These predictions agree with the true sample ancestry taken from the 1kGP validation set.

Discussion

We have described a method to predict ancestry from genomic data that provides multiple improvements over existing ancestry inference tools. Firstly, SNVstory incorporates samples/variants from three different curated datasets, expanding the number of labels and the granularity of the model classification beyond the main continental divisions. Secondly, drawing upon the gnomAD database produces a much larger number of variants on which our models were trained, providing the opportunity to classify ancestry on

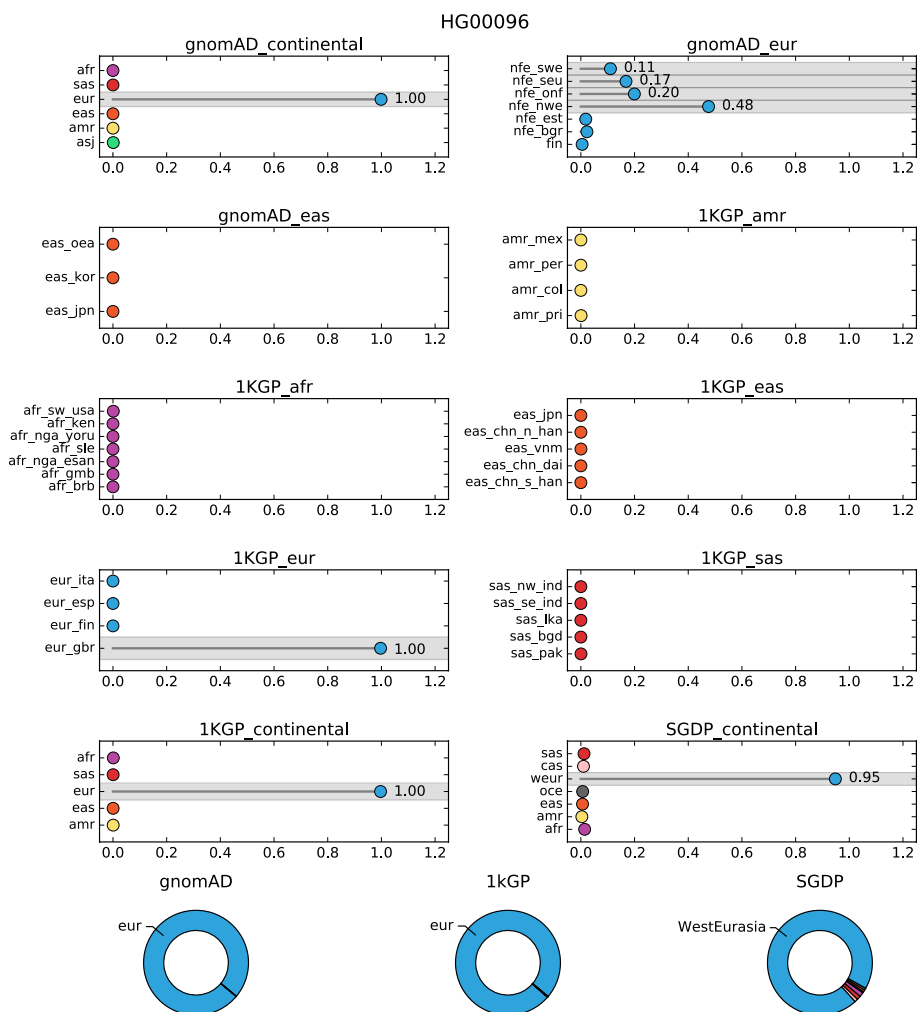


Fig. 4 SNVstory ancestry report. The representative output of model results from SNVstory for a European sample taken from the 1kGP dataset

a wider (or more diverse) range of features. Thirdly, SNVstory excludes consanguineous samples from training, ensuring that the overrepresentation of closely related individuals does not bias the model. Finally, our novel implementation is optimized for individualized results rather than clustering large cohorts of samples into shared ancestral groups.

We compared our model performance to two other popular ancestry inference tools, RFMix and ADMIXTURE. RFMix is predominately advertised as a local ancestry inference tool, but it can also supply global ancestry results by aggregating their output into one label. RFMix performed comparably to SNVstory at the continental level. However, it performed worse at the subcontinental level. This decrease may potentially be due to its output which is geared toward local ancestry. With many more labels, it becomes harder to pull out the top global ancestry from the aggregated results. SNVstory overcomes the issue of a large label set by implementing a weighting scheme to reduce the probability of picking a subcontinental label outside of the continental model results. The RFMix subcontinental results would likely improve if a similar weighting scheme was used when delivering results. ADMIXTURE showed

an opposite trend of improved results when looking at subcontinental performance, which is perhaps due to the increased similarity between reference and input samples. Using either RFMix or ADMIXTURE, the user is tasked with an additional step of supplying an appropriate reference dataset to get the most accurate results. Overall, SNVstory performs accurately at both the continental and subcontinental levels and only requires a query VCF as input.

In our gnomAD model, we introduce a method to simulate individual samples from aggregate allele frequencies of a known population. This is potentially useful for any study requiring access to reference variants from a population where data from individual samples is obfuscated. One limitation in our approach is that we did not account for linkage disequilibrium between variants when simulating individual samples. This could result in some samples with patterns of variants that do not exist in actual samples. An improvement in future models would be to remove variants with high levels of linkage disequilibrium between them. If high recognizability to actual samples is required, established metrics of linkage disequilibrium, such as the correlation coefficient r^2 , could be used to measure the 'realness' of a simulated sample based on existing variant patterns, and simulated VCFs could be validated based on this quality. However, in practice, the larger pool of variants provided by gnomAD more than compensates for the lost dependence among proximal groups of variants. We have demonstrated that the performance of the gnomAD models with simulated individuals is comparable to that of models trained with actual samples.

With the growing number of reference datasets containing individuals from diverse ancestral backgrounds, it is possible to build ancestry prediction models that reflect these populations. However, there is room for improvement, as our most diverse dataset (SGDP) includes the fewest samples. We could not build subcontinental models as granular as the labels provided because there were as few as two samples per label for many instances. Additionally, our model cannot accurately predict ancestry proportions in samples with admixed ancestry. Most admixture prediction software depends on a priori knowledge of the number of non-admixed populations and requires representation from such populations. There is limited availability of reference samples from admixed individuals, so our training data lacked representation from any admixed samples. Efforts to expand the number of reference sequences for diverse and admixed populations will provide opportunities to fill this gap.

SNVstory's feature-importance capacity is unique among ancestral tools and could have significant clinical utility. The clinical application of most ancestral prediction tools is limited to simply predicting the patient's ancestry. However, SNVstory's unique capability to describe a given locus as characteristic, or atypical, of a given ancestry could lead to improved prioritization of variants. For example, SNVstory finds the most ancestrally informative gene on average to be KRTAP19-8, which is greatly enriched for SNVs predictive of Native American/Latino ancestry (Fig. 3). This gene is a known driver of thyroid lymphoma [39], a disorder that is the second-most-common type of cancer among Hispanic women [40] but not even among the top five cancer types among women worldwide [41]. The inferred distinctiveness of Latino copies of KRTAP19-8 suggests that rare founder mutations in this gene may contribute to increased rates of thyroid cancer among women of Hispanic ancestry. The ability to target variants in genes

inherited from specific populations adds a new tool to the diagnostician's toolkit and could lead to improved patient outcomes.

Finally, our approach allows users to reliably execute our models given a single-sample or multi-sample VCF, with results tailored toward ancestry assignment for an individual sample. This provides immediately useful ancestry information in the clinical setting, where ancestry can be used to inform diagnostic or therapeutic decisions. Specifically, a subject's ancestry can be used to help prioritize variants that may be rare in one population but not another. In the clinical setting, it may be essential to recognize the difference between ethnicity, race, and genetic ancestry in determining the optimal therapy or drug dosage.

Given the widespread availability of genome sequencing data and models like SNVstory that can accurately predict ancestry, we advocate for genetic ancestry to become the standard classification reported for genetic studies and clinical applications, where appropriate. Genetic ancestry offers enormous advantages over other self-reported information, such as ethnicity or race, because it supplies biological characteristics of a population and is consistently measurable. This advantage will only increase as more populations are sequenced and ancestry prediction becomes more reliable, and we improve our ability to contextualize the impact of genetic ancestry on clinical decision-making.

Availability and requirements

Project name: SNVstory.

Project home page: <https://github.com/nch-igm/snvstory>.

Operating system(s): Platform independent.

Programming language: Python.

Other requirements: Docker.

License: 3-Clause BSD License.

Any restrictions to use by non-academics: None.

Abbreviations

SNVs	Single nucleotide variants
SVs	Structural variants
CNVs	Copy number variants
Indels	Insertions or deletions
GWAS	Genome-wide association studies
dbSNP	The single nucleotide polymorphism database
gnomAD	The genome aggregation database
1kGP	The 1000 Genomes Project
SGDP	The Simons Genome Diversity Project

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12859-024-05703-y>.

Additional file 1: Supplemental methods, figures, and tables.

Acknowledgements

We thank the Nationwide Foundation Pediatric Innovation Fund for generously supporting this project.

Author contributions

AB and AR processed data and trained models for gnomAD, 1KGP, and SGDP. AR designed the methods to simulate data from gnomAD allele frequencies and cross-validation architecture. AB and DC prepared figures and tables. AB wrote the first draft of the paper. JG, DC, AR, and PW assisted in preparing or revising the paper. AR and AB wrote the SNVstory software package. PW and EM supervised the project.

Funding

This work was supported by the Nationwide Children's Foundation and The Abigail Wexner Research Institute at Nationwide Children's. The funders had no role in study design, data collection, data analysis, the decision to publish, or manuscript preparation.

Availability of data and materials

The training data for our model are available as follows. gnomAD v2.1 data is available from <https://gnomad.broadinstitute.org/downloads/>. 1000 Genomes Project data is shared via the International Genome Sample Resource and can be accessed from <https://www.internationalgenome.org/data-portal/data-collection/30x-grch38>. Simons Genome Diversity Project data is available from the European Nucleotide Archive under project PRJEB9586. The in-house dataset of clinical exome sequencing test data used to compare self-reported ancestry to SNVstory's algorithmic prediction is not available publicly. Please contact the corresponding author with further questions about this dataset. SNVstory is published under an Open-Source Initiative approved 3-Clause BSD License to ensure that any interested academic institution can perform optimization with their own cohorts and implement their own version of the algorithm in their respective diagnostic workflows. Code for the SNVstory algorithm is available at our GitHub repository (<https://github.com/nch-igm/snvstory>).

Declarations**Ethics approval and consent to participate**

This study was reviewed and approved by the Institutional Review Board (IRB) of The Abigail Wexner Research Institute at Nationwide Children's Hospital (Office for Human Research Protections (OHRP) IORG0000326; IRB00000568) as IRB17-00206 ("Institute for Genomic Medicine Comprehensive Profiling for Cancer, Blood, and Somatic Disorders"). The participant's legal guardian/next of kin provided written informed consent to participate in this study.

Consent for publication

Not applicable.

Competing interests

No competing interests: AEB, AR, DC, JBG, and PW. ERM: Qiagen N.V., supervisory board member, honorarium, and stock-based compensation. Singular Genomics Systems, Inc., board of directors, honorarium, and stock-based compensation.

Received: 8 July 2023 Accepted: 13 February 2024

Published online: 20 February 2024

References

- Auton A, Abecasis GR, Altshuler DM, Durbin RM, Abecasis GR, Bentley DR, et al. A global reference for human genetic variation. *Nature*. 2015;526:68–74.
- Hauser D, Obeng AO, Fei K, Ramos MA, Horowitz CR. Views of primary care providers on testing patients for genetic risks for common chronic diseases. *Health Aff Proj Hope*. 2018;37:793–800.
- Jorde LB, Bamshad MJ. Genetic ancestry testing: what is it and why is it important? *JAMA*. 2020;323:1089–90.
- Ramamoorthy A, Pacanowski MA, Bull J, Zhang L. Racial/ethnic differences in drug disposition and response: review of recently approved drugs. *Clin Pharmacol Ther*. 2015;97:263–73.
- Fujimura JH, Rajagopalan R. Different differences: the use of 'genetic ancestry' versus race in biomedical human genetic research. *Soc Stud Sci*. 2011;41:5–30.
- Shraga R, Yarnall S, Elango S, Manoharan A, Rodriguez SA, Bristow SL, et al. Evaluating genetic ancestry and self-reported ethnicity in the context of carrier screening. *BMC Genet*. 2017;18:99.
- Mersha TB, Abebe T. Self-reported race/ethnicity in the age of genomic research: its potential impact on understanding health disparities. *Hum Genom*. 2015;9:1.
- Gomes MB, Gabrielli AB, Santos DC, Pizarro MH, Barros BSV, Negrato CA, et al. Self-reported color-race and genomic ancestry in an admixed population: a contribution of a nationwide survey in patients with type 1 diabetes in Brazil. *Diabetes Res Clin Pract*. 2018;140:245–52.
- Brown R, Lee H, Eskin A, Kichaev G, Lohmueller KE, Reversade B, et al. Leveraging ancestry to improve causal variant identification in exome sequencing for monogenic disorders. *Eur J Hum Genet*. 2016;24:113–9.
- Purcell SM, Wray NR, Stone JL, Visscher PM, O'Donovan MC, Sullivan PF, et al. Common polygenic variation contributes to risk of schizophrenia and bipolar disorder. *Nature*. 2009;460:748–52.
- Pritchard JK, Stephens M, Donnelly P. Inference of population structure using multilocus genotype data. *Genetics*. 2000;155:945–59.
- Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in unrelated individuals. *Genome Res*. 2009;19:1655–64.
- Gimbernat-Mayol J, Mantes AD, Bustamante CD, Montserrat DM, Ioannidis AG. Archetypal analysis for population genetics. *PLoS Comput Biol*. 2022;18: e1010301.
- Sherry ST, Ward M-H, Kholodov M, Baker J, Phan L, Smigielski EM, et al. dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res*. 2001;29:308–11.

15. Jin Y, Schaffer AA, Feolo M, Holmes JB, Kattman BL. GRAF-pop: a fast distance-based method to infer subject ancestry from multiple genotype datasets without principal components analysis. *G3 Bethesda Md.* 2019;9:2447–61.
16. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alföldi J, Wang Q, et al. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature.* 2020;581:434–43.
17. Mallick S, Li H, Lipson M, Mathieson I, Gymrek M, Racimo F, et al. The Simons Genome Diversity Project: 300 genomes from 142 diverse populations. *Nature.* 2016;538:201–6.
18. Kumar A, Montserrat DM, Bustamante C, Ioannidis A. XGMix: local-ancestry inference with stacked XGBoost. preprint. *Genomics*; 2020.
19. Maples BK, Gravel S, Kenny EE, Bustamante CD. RFMix: a discriminative modeling approach for rapid and robust local-ancestry inference. *Am J Hum Genet.* 2013;93:278–88.
20. Sheehan S, Song YS. Deep learning for population genetic inference. *PLoS Comput Biol.* 2016;12: e1004845.
21. Hwa H-L, Wu M-Y, Lin C-P, Hsieh WH, Yin H-I, Lee T-T, et al. A single nucleotide polymorphism panel for individual identification and ancestry assignment in Caucasians and four East and Southeast Asian populations using a machine learning classifier. *Forensic Sci Med Pathol.* 2019;15:67–74.
22. Durand EY, Do CB, Mountain JL, Macpherson JM. Ancestry composition: a novel, efficient pipeline for ancestry deconvolution. *Bioinformatics*; 2014.
23. Chu BB, Sobel EM, Wasiolek R, Ko S, Sinsheimer JS, Zhou H, et al. A fast data-driven method for genotype imputation, phasing, and local ancestry inference: MendelImpute.jl. *Bioinforma Oxf Engl.* 2021;37:489.
24. Shi G, Kuang Q. Ancestral spectrum analysis with population-specific variants. *Front Genet.* 2021;12: 724638.
25. Wang Y, Song S, Schraiber JG, Sedghifar A, Byrnes JK, Turissini DA, et al. Ancestry inference using reference labeled clusters of haplotypes. *BMC Bioinform.* 2021;22:459.
26. Soumare H, Rezzgui S, Gmati N, Benkahla A. New neural network classification method for individuals ancestry prediction from SNPs data. *BioData Min.* 2021;14:30.
27. Dalfovo D, Romanel A. Analysis of genetic ancestry from NGS data using EthSEQ. *Curr Protoc.* 2023;3: e663.
28. Karim MR, Cochez M, Zappa A, Sahay R, Beyan O, Schuhmann D-R, et al. Convolutional embedded networks for population scale clustering and bio-ancestry inferencing. *EEE/ACM Trans Comput Biol Bioinform.* 2020;19:369–82.
29. Byrska-Bishop M, Evani US, Zhao X, Basile AO, Regier AA, Corvelo A, et al. High coverage whole genome sequencing of the expanded 1000 Genomes Project cohort including 602 trios. *Cell.* 2022;185:3426–40.
30. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. 2013.
31. Pedersen BS, Quinlan AR. Mosdepth: quick coverage calculation for genomes and exomes. *Bioinformatics.* 2018;34:867–8.
32. Zhang F, Flickinger M, Taliun SAG, Abecasis GR, Scott LJ, McCarroll SA, et al. Ancestry-agnostic estimation of DNA sample contamination from sequence reads. *Genome Res.* 2020;30:185–94.
33. Manichaikul A, Mychaleckyj JC, Rich SS, Daly K, Sale M, Chen W-M. Robust relationship inference in genome-wide association studies. *Bioinformatics.* 2010;26:2867–73.
34. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK: rising to the challenge of larger and richer datasets. *GigaScience.* 2015. <https://doi.org/10.1186/s13742-015-0047-8>.
35. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics.* 2010;26:841–2.
36. Chen, T., & Guestrin, C. XGBoost: a scalable tree boosting system. In: Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining. 2016, pp. 785–94.
37. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: machine learning in Python. *J Mach Learn Res.* 2011;12:2825–30.
38. Lundberg S, Lee S-I. A Unified Approach to Interpreting Model Predictions. 2017.
39. Qiu K, Li K, Zeng T, Liao Y, Min J, Zhang N, et al. Integrative analyses of genes associated with Hashimoto's thyroiditis. *J Immunol Res.* 2021;2021:8263829.
40. Estrada-Florez AP, Bohórquez ME, Sahasrabudhe R, Prieto R, Lott P, Duque CS, et al. Clinical features of Hispanic thyroid cancer cases and the role of known genetic variants on disease risk. *Medicine (Baltimore).* 2016;95: e4148.
41. Ferlay J, Colombet M, Soerjomataram I, Parkin DM, Piñeros M, Znaor A, et al. Cancer statistics for the year 2020: an overview. *Int J Cancer.* 2020. <https://doi.org/10.1002/ijc.33588>.

Publisher's Note

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