

Evaluating the utility of identity-by-descent segment numbers for relatedness inference via information theory and classification

Jesse Smith ^{1,3,†} Ying Qiao ^{2,*†} Amy L. Williams ^{2,*}

¹School of Applied and Engineering Physics, Cornell University, Ithaca, NY 14853, USA,

²Department of Computational Biology, Cornell University, Ithaca, NY 14853, USA

³Present address: Applied Math Lab, Courant Institute of Mathematical Sciences, New York University, NY 10012, USA.

*Corresponding author: Department of Computational Biology, Cornell University, Ithaca, NY 14853, USA. Email: yq76@cornell.edu; * Corresponding author:

Department of Computational Biology, Cornell University, Ithaca, NY 14853, USA. Email: alw289@cornell.edu

[†]These authors contributed equally to this work.

Abstract

Despite decades of methods development for classifying relatives in genetic studies, pairwise relatedness methods' recalls are above 90% only for first through third-degree relatives. The top-performing approaches, which leverage identity-by-descent segments, often use only kinship coefficients, while others, including estimation of recent shared ancestry (ERSA), use the number of segments relatives share. To quantify the potential for using segment numbers in relatedness inference, we leveraged information theory measures to analyze exact (i.e. produced by a simulator) identity-by-descent segments from simulated relatives. Over a range of settings, we found that the mutual information between the relatives' degree of relatedness and a tuple of their kinship coefficient and segment number is on average 4.6% larger than between the degree and the kinship coefficient alone. We further evaluated identity-by-descent segment number utility by building a Bayes classifier to predict first through sixth-degree relationships using different feature sets. When trained and tested with exact segments, the inclusion of segment numbers improves the recall by between 0.28% and 3% for second through sixth-degree relatives. However, the recalls improve by less than 1.8% per degree when using inferred segments, suggesting limitations due to identity-by-descent detection accuracy. Last, we compared our Bayes classifier that includes segment numbers with both ERSA and IBIS and found comparable recalls, with the Bayes classifier and ERSA slightly outperforming each other across different degrees. Overall, this study shows that identity-by-descent segment numbers can improve relatedness inference, but errors from current SNP array-based detection methods yield dampened signals in practice.

Keywords: genetic relatedness; identity-by-descent; pedigree reconstruction; relatedness inference

Introduction

Relatedness inference in genetic data often plays a fundamental role in enabling more accurate genetic analyses—both in studies that directly leverage relatives and those that prune them to avoid modeling violations. The need and opportunity to identify genetic relatives continues to increase as the scale of genetic datasets increase (Henn *et al.* 2012; Bycroft *et al.* 2018). One notable example is the UK Biobank wherein roughly 30% of genotyped individuals have a third degree (e.g. first cousin) or closer relative in the study (Bycroft *et al.* 2018). Applications that make use of genetic relatives are numerous and varied and include pedigree reconstruction (Staples *et al.* 2014; Jewett *et al.* 2021), pedigree-based linkage analysis for disease and trait mapping (Ott *et al.* 2015), heritability estimation (Zaitlen *et al.* 2013; Young *et al.* 2018), forensic genetics (Weir *et al.* 2006), and genetic genealogy (Stallard and de Groot 2020)—a popular tool among direct-to-consumer genetic testing customers. On the other hand, traditional genome-wide association study tests and many population

genetic models assume that the study samples are unrelated, and, as such, must exclude inferred relatives to avoid spurious signals or inaccurate parameter estimates (Voight and Pritchard 2005). All these applications motivate a thorough analysis of the approaches used for relatedness inference to determine which of the various features the methods should leverage.

Many relatedness inference methods only utilize kinship coefficients (Manichaikul *et al.* 2010; Ramstetter *et al.* 2017; Dimitromanolakis *et al.* 2019; Seidman *et al.* 2020), while some such as estimation of recent shared ancestry (ERSA) leverage the number of identity-by-descent (IBD) segments between a pair (Huff *et al.* 2011). IBD segments are regions of DNA two or more relatives coinherit from a common ancestor (Thompson 2013), and kinship coefficients are a scaled measure of the total IBD length of a relative pair (Ramstetter *et al.* 2017). To date, the question of whether segment numbers provide information for relatedness inference beyond that of kinship coefficients has not been carefully explored. A recent study benchmarked 12 pairwise relatedness inference methods using thousands of real relatives

Received: September 21, 2021. Accepted: March 07, 2022

© The Author(s) 2022. Published by Oxford University Press on behalf of Genetics Society of America.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

(Ramstetter et al. 2017) and highlighted three most accurate approaches: ERSA and two IBD detection algorithms, GERMLINE (Gusev et al. 2009) and Refined IBD (Browning and Browning 2013) (i.e. using kinship coefficients derived from their output). Although ERSA models the distribution of both the number and lengths of IBD segments, that evaluation found that it does not always outperform other methods that only use kinship coefficients. One possible reason is that estimated segment numbers from most phase-based IBD detection methods are inflated due to switch errors that typically break up segments (Dimitromanolakis et al. 2019; Freyman et al. 2021; Seidman et al. 2020). Alternatively, these results may indicate that IBD segment numbers and lengths do not capture relatives' degrees of relatedness better than kinship coefficients.

To determine whether incorporating the number of IBD segments in a model with kinship coefficients (or coefficients of relatedness) improves relatedness inference, we first performed an information theory-based analysis. Feature selection based on information theory is widely used in machine learning and data mining in fields as diverse as bioinformatics and pattern recognition (Hoque et al. 2014; Lee and Kim 2015; Qian and Shu 2015). We applied a commonly used measure—mutual information (MI)—to quantify the dependency between various features and the class variable (here the degree of relatedness) and also the dependency among the features themselves. An advantage of this approach is that MI does not make an assumption of linearity between the features and can be calculated for both discrete and continuous variables (Bennasar et al. 2015). In addition, the MI analysis results do not depend on the specific classifier used downstream and can capture the relationship between variables from an information theory perspective that is distinct from classification.

We further conducted a classification-based analysis to determine the importance of IBD proportions and segment numbers for inferring degrees of relatedness. For this purpose, we developed a Bayes classifier with mathematical underpinnings that parallel those of MI. Bayes classifiers are a form of generative learning that seek to minimize the probability of misclassification by estimating the probability of a given data point being from each class (Devroye et al. 2013). In this study, we assign a pair of relatives to the maximum posterior probability degree, in contrast to approaches that map estimated kinship coefficients to degrees of relatedness using *a priori* fixed ranges of kinship (Manichaikul et al. 2010; Ramstetter et al. 2017; Seidman et al. 2020). The latter ignores the effect of population structure on IBD signals—including background IBD segments (Weir et al. 2006). These effects are important to model since they vary by population and can meaningfully influence relatedness classification. Moreover, bias in the detection of IBD segments can shift the distributions of both IBD proportions and segment numbers away from expected ranges. Such biases may especially impact classification of more distant relatives as they have smaller ranges of kinship values that correspond to a given degree. In light of these concerns, we estimate the probability of the features given the degree (i.e. the likelihood) using training data simulated using genotypes from the target population. This implicitly accounts for the influence of background IBD segments as well as any errors in IBD segment detection. Researchers with access to data from a given population can also apply this strategy by using the available samples as founders in simulated pedigrees (Caballero et al. 2019).

Finally, we benchmarked the performance of our relatedness classifier together with ERSA and identical by descent via

identical by state (IBIS) using simulated genotypes. We focus on these two methods because ERSA leverages IBD segment numbers for classification, and because IBIS (whose segments we use as input to our Bayes classifier) has comparable relatedness inference accuracy to the top-performing Refined IBD and GERMLINE detectors noted above—i.e. these are state-of-the-art methods (Ramstetter et al. 2017; Seidman et al. 2020). As our goal was to understand the impact of including segment numbers as features for classification—and given prior efforts to compare the performance of a wide range of relatedness classifiers (Ramstetter et al. 2017)—we did not include other IBD segment detectors nor allele frequency-based approaches such as KING (Manichaikul et al. 2010) and PLINK (Chang et al. 2015) in this analysis. Overall, we obtained comparable classification results for the Bayes classifier, ERSA, and IBIS, indicating that the Bayes classifier is reliable and suggesting that our approach can be used in practice given appropriate training data resources. Notably, the Bayes classifier performs similarly to IBIS (which does not use segment numbers) demonstrating that, in practice, incorporating segment numbers provides little improvement in classification rates.

All the analyses in this paper leverage IBD segments from simulated data, either exact segments produced by the simulator or segments inferred from simulated genetic data. While reliably labeled real relatives would be preferable, simulated relatives produced using both sex-specific genetic maps and crossover-interference modeling have relatedness measures that closely mirror those from real data (Caballero et al. 2019). Using the simulated relatives, we investigated (1) MI quantities based on exact segments, (2) classification rates using exact segments, and (3) classification rates from inferred segments. In this way, the MI analysis quantifies the theoretical information gain available in the limit of perfect IBD detection. In addition, the classification analysis using exact segments reveals how much improvement in relatedness inference is possible by incorporating IBD segment numbers at this same limit. Finally, comparing the classification results using exact vs inferred segments enables us to localize and quantify the influence of IBD detection errors.

Methods

We analyzed the potential for using coefficients of relatedness r (defined below) either alone or both r and n , the number of IBD segments a pair of relatives share, to infer the pair's degrees of relatedness D .

Mutual information discrete definition and binning approaches

MI is difficult to calculate for continuously valued variables without a known distribution and whose distribution must therefore be estimated from finite data. Moreover, estimating the MI between one continuous and one discrete random variable is in general nontrivial and multiple approaches exist for this estimation, such as nearest-neighbor (Ross 2014) and binning methods. To enable our MI calculations [such as $I(r; D)$], we used a procedure that bins data points of r and avoids biased MI estimates in our finite but large sample size. In computing MI, we treated the binned feature vector \vec{F} (where \vec{F} has the possibility of being 1D when representing r or n) and the degree of relatedness D as two discrete random variables with realizations f and $d \in [1, 7]$, respectively. If we know the probability mass functions (pmfs) of

the discrete random variables X and Y with realizations x and y , we can calculate MI using its definition as

$$I(X; Y) = \sum_x \sum_y p_{X,Y}(x,y) \log \frac{p_{X,Y}(x,y)}{p_X(x)p_Y(y)}, \quad (1)$$

where $p_{X,Y}$ is the joint pmf of X and Y and p_X, p_Y are the marginal pmfs of X and Y , respectively.

Binning a continuous variable in order to use Equation (1) introduces the difficulty of picking the right bin size. It has been shown (Ross 2014) that MI is sensitive to bin size and that its stability with respect to this variable is dependent on the sample distribution. Our distributions and sample sizes of r yielded a large range of bin sizes that have stable MI estimates (see the flat regions of each curve in Supplementary Fig. 1). Because the normalized MI gain G_N is a fraction composed of MI values (Results), its correct calculation relies on the unbiasedness of the various MI quantities that form it. At bin sizes smaller than 150 pairs per bin (ppb), both the means and standard deviations (Supplementary Fig. 2) of our MI quantities increase rapidly. Given this, in our calculations of MI, we binned r at 150 ppb. Our binning converts a continuous value of r to its nearest bin-value in $N/150$ evenly spaced numbers from $[\min(\vec{r}), \max(\vec{r})]$, where N represents all sampled values for the desired feature r , n , or (r, n) . Here and below \vec{r} represents all sampled training data points r .

Estimating mutual information

Calculating Equation (1) without access to the entire spaces \mathcal{X} and \mathcal{Y} —i.e. estimating MI from sampled data—is contingent on the estimation of marginal and joint probabilities p_X, p_Y , and $p_{X,Y}$. We used a simple counting approach to calculate each probability, assigning $\hat{p}_F(f) = \frac{1}{N} \sum_i \mathbb{1}(\text{bin}(F_i) = \text{bin}(f))$. Here \vec{F} is the vector of realized data points representing all N sampled values for the desired feature r , n , or (r, n) ; $\text{bin}(x)$ denotes the function that takes a continuous realization to its binned value; and $\mathbb{1}(X = Y)$ is the indicator function. By binning r to 150 ppb as noted in the previous subsection, we were able to use this discrete maximum likelihood estimator (MLE) approach for calculating every desired pmf and obtain stable results in MI.

We performed calculations of MI on the exact IBD segment data restricted to three distribution shapes: A uniform distribution, a “slow-exponential growth” distribution, and an exponential growth distribution (see Fig. 1). We accounted for different distributions of D in the calculations of $I(\vec{F}; D)$ by decomposing the joint pmf relating \vec{F} and D as $p_{\vec{F},D}(f,d) = p_{\vec{F}}(f|d)p_D(d)$, and also decomposing the marginal pmf on \vec{F} (with the law of total probability $p_{\vec{F}}(f) = \sum_{d'} p_{\vec{F}}(f|d')p_D(d')$). Equation (1) is then expressed as

$$I(\vec{F}; D) = \sum_f \sum_d p_{\vec{F}}(f|d)p_D(d) \log \frac{p_{\vec{F}}(f|d)}{\sum_{d'} p_{\vec{F}}(f|d')p_D(d')} \quad (2)$$

by canceling the $p_D(d)$ terms in the numerator and denominator. $p_{\vec{F}}(f|d)$ is the pmf of realizations of feature \vec{F} in a given degree, and $p_D(d)$ is the distribution shape (from Fig. 1). This approach removes noise associated with calculating the pmfs $\hat{p}_{\vec{F},D}$ for different distribution shapes, which stems in part from random factors in finite sample sizes (including smaller numbers of pairs in the nonuniform distributions). In particular, because the probabilities of f conditioned on any given degree of relatedness d are identical across each distribution shape, we estimated $p_{\vec{F}}(\vec{F}|d)$ once only from the uniform distribution data. Note too that the differences in MI due to the D

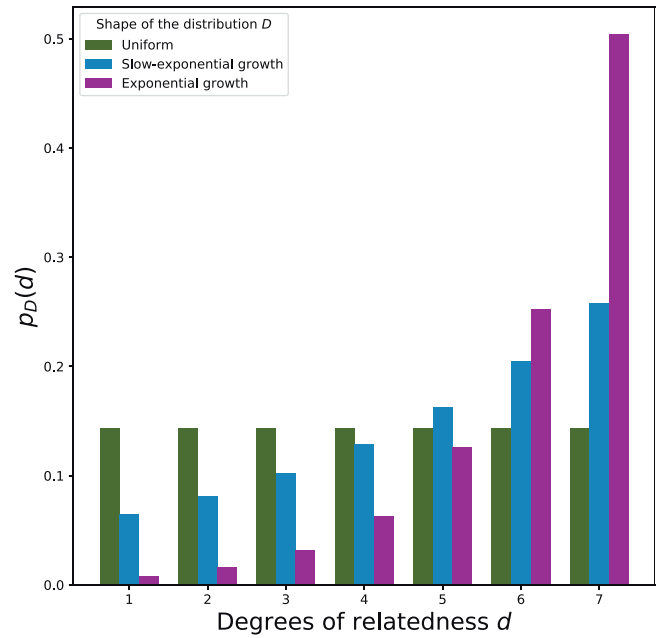


Fig. 1. Probability mass functions of different distribution shapes for D as a function of degree of relatedness d . Here uniform = $1/7$, slow-exponential growth = $(1000/15541) \times 2^{(d-1)/3}$, and exponential growth = $(160/20320) \times 2^{d-1}$. Total pair counts in the inferred segment data analysis are 21,000 for uniform, 15,541 for slow-exponential growth, and 20,320 for the exponential growth.

distribution are entirely accounted for in the probabilities $p_D(d)$, and these are exactly calculable given the equation for each distribution.

Probability density estimation of features

In the context of the Bayes classifier, we estimated the probability of a feature realization f conditioned on the training data \vec{T}^d in degree d according to the degree-wise count as

$$\hat{p}(f|d, \vec{T}^d) = \frac{1}{N_d^T} \sum_{i=1}^{N_d^T} \mathbb{1}(\text{bin}(T_i^d) = \text{bin}(f)), \quad (3)$$

T_i^d being a particular realization of the training data \vec{T}^d with total count N_d^T . However, we only had access to the frequencies of realizations f that occur at least once in the training data, so Equation (3) is only calculable for these values. The total training data \vec{T} and testing data \vec{r} are of dimension equal to their respective number of data points N^T or N^r . To generate posteriors $p(D|\tau_i)$ for realizations in \vec{r} at values where there are no training data points in $\text{bin}(\tau_i)$, we linearly interpolated the values given by Equation (3) within the convex hull (see Supplementary Fig. 3) specified by the bounds of the training data. (Strictly speaking, these posteriors are then incorrect pmfs with mass greater than 1—however, in practice this is only relevant for a vanishingly small number of points.) We used the scipy packages `interp1d` and `griddata` for the linear interpolations in 1-dimension (when \vec{F} is either r or n) and 2-dimensions (when \vec{F} is (r, n)), respectively.

In the case that \vec{F} is (r, n) , the 2D linear interpolation of $\hat{p}(f|\vec{T}^d)$ values are only well-defined inside of the 2D convex hulls of the training data. Therefore, we could not assign posteriors to realizations of the testing data that lay outside the bounds of the training data. For these data points (labeled in Supplementary

Fig. 3 as “Unscored under (r, n)”), we assigned probability values according to the 1D interpolation of $p(f|\bar{T}^d)$ values with $\bar{F} = r$. The 1D interpolations for $\bar{F} = r$ (or $\bar{F} = n$) only remained undefined when they occurred outside the interval of training values $[\min(\bar{r}), \max(\bar{r})]$ (or $[\min(\bar{n}), \max(\bar{n})]$), in which case they remained unclassified in our analysis. In the inferred segment data, there was only a maximum of one point per degree that remained unclassified.

Bayes classification

Our classifiers use the posterior probabilities $p(D|\bar{F}, \bar{T}) = \frac{p(\bar{F}|D, \bar{T})p(D)}{p(\bar{F}|\bar{T})}$ for the single and multivariate features \bar{F} to infer D in the testing data. The priors $p(D)$ are the known shape of the degree distribution (Fig. 1), and we generated the probability of our data $p(\bar{F}|\bar{T})$ as the sum across degrees according to the law of total probability $p(\bar{F}|\bar{T}) = \sum_d p(\bar{F}|d, \bar{T})p(d)$. We calculated likelihoods $p(\bar{F}|D, \bar{T})$ according to the estimator in Equation (3). To classify a testing pair τ_i to a certain degree, we calculated $\log p(D|\tau_i)$ for each degree and classified the pair as the maximum a posteriori degree:

$$D_i^p = \operatorname{argmax}_{d \in \mathcal{D}} \log p(D = d|\tau_i), \quad (4)$$

where D_i^p is the predicted degree, and \mathcal{D} is the set of possible degrees $\{1, \dots, 7\}$. The recall of a particular classifier for degree d is $\frac{1}{N_d} \sum_i 1(D_i^p = d)$.

The classifier takes input IBIS segments and calculates the IBD proportion and segment numbers for all pairs of individuals with at least one detected IBD segment. It classifies any pair with $r < 2^{-15/2}$ [a common lower bound for seventh-degree classification (Manichaikul et al. 2010; Ramstetter et al. 2017)] as unrelated, and, for all other pairs, predicts their degree using Equation 4.

Simulated data

For the exact IBD segment data, we used Ped-sim (Caballero et al. 2019) to simulate relative pairs of 13 relationship types from seven degrees of relatedness (Table 1) (replicated 80 times for the MI analysis and once for the classification-based analysis) and leveraged the IBD segments this tool prints. Thus these segments are free of error and we refer to them throughout as *exact*. We used both sex-specific genetic maps (Bhérier et al. 2017) and crossover interference modeling (Housworth and Stahl 2003) for these simulations.

For each degree, we simulated an equal number of pairs from each of two relationship types. The one exception is first-degree relatives where we only considered full sibling pairs since parent-

Table 1. Simulated relationship types for each degree of relatedness.

Degree	Relationships
1	Full siblings
2	Avuncular, Half-siblings
3	First cousins, Half avuncular
4	First cousins once removed, Half first cousins
5	Second cousins, Half first cousins once removed
6	Second cousins once removed, Half second cousins
7	Third cousins, Half second cousin once removed

Half relatives share only 1 common ancestor while other types have two common ancestors.

child pairs always have $r=0.5$ and are trivial to identify. We doubled the number of full sibling pairs (to the total number assigned from the distribution shape) so that the first-degree relatives included the full number of pairs. We calculated the IBD proportion by adding the lengths of all outputted IBD segments and dividing by the total length of the sex averaged genetic map—halving the length of IBD1 segments (see the equation for r in Results). We calculated the segment number by counting the number of outputted IBD segments from either Ped-sim (exact) or IBIS (inferred, as described next).

To simulate relatives with genetic data, we used autosomal genotypes from participants in the UK Biobank (Bycroft et al. 2018) as founders in Ped-sim runs. We used the phased data distributed by the UK Biobank (Bycroft et al. 2018) and, before simulating, filtered the samples to include the white British ancestry subjects. To filter out close relatives, we first performed SNP quality control filtering on the UK Biobank unphased genotypes [filtering SNPs with minor allele frequency less than 2%, missing data rate greater than 1%, and retaining only SNPs used for phasing in the original analysis (Bycroft et al. 2018)]. Next we ran IBIS v1.20.8 on the filtered genotypes with the `-maxDist 0.12` option and with IBD2 segment detection enabled. This provided kinship coefficients that we then input to PRIMUS (Staples et al. 2014), running it with `-no_PR` [which corresponds to not reconstructing pedigrees: executing only IMUS (Staples et al. 2013)] and `-rel_threshold 0.022` to filter out relatives with a kinship coefficient greater than 0.022 [i.e. retaining only pairs no more closely related than fifth degree (Manichaikul et al. 2010)]. We ran Ped-sim as described above (using sex-specific genetic maps and crossover interference modeling) and otherwise used default options (including genotyping error and missing data rates of 10^{-3}). Finally, we used IBIS v1.20.7 (enabling IBD2 detection with `-2` and otherwise using default options) to detect IBD segments between these simulated relatives. Note that IBIS’s default minimum segment length is 7 cM—meaningfully longer than is typical for phase-based IBD detectors (Seidman et al. 2020).

Running ERSA

To get relatedness estimates from ERSA (Huff et al. 2011), we first ran GERMLINE (Gusev et al. 2009) v1.5.1 with `-err_het=1` `-err_hom=2` `-min_m 2.5` and `-bits 64` on the simulated Ped-sim haplotypes [these are the options PADRE uses (Staples et al. 2016) and detect segments ≥ 2.5 cM long]. That is, we provided ERSA perfectly phased data output by the simulator. We then ran ERSA with default options (including a minimum IBD segment length of `-min_cm=2.5`) on the resulting GERMLINE segments.

Runtimes

To collect runtimes, we ran both ERSA and our Bayes classifier on a machine with four Intel Xeon e7 4830 v3 2.0 GHz processors and 512GB of RAM. We supplied 16 GB of RAM to ERSA and 8 GB to our Bayes classifier. Both methods are single threaded.

Results

To begin, we quantified the inherent dependency between the IBD segment features and D by analyzing MI between the features and D . MI is a quantification of the information obtained about one random variable through observing another; in this case, we analyzed the information gained about D through observing the variables r , n , or (r, n) . We compared our analysis of these MI quantities with the corresponding Bayes classification results; the conclusions we form about the classification

effectiveness of different feature sets are therefore based on both the MI and classification recall.

Throughout, we refer to IBD regions that two individuals share on only one haplotype copy as IBD1, and those the individuals share IBD on both chromosomes as IBD2.

Mutual information analysis

We used thousands of relative pairs to estimate mutual information $I(\vec{F}; D)$ between different IBD features \vec{F} and the degree of relatedness D of each pair (Methods, “Estimating mutual information”). Specifically, we compared MI values of $I(n; D)$, $I(r; D)$, and $I((r, n); D)$ calculated using units of bits. Let $k_{ij}^{(1)}$ and $k_{ij}^{(2)}$ denote the proportion of their genomes that individuals i and j share IBD1 and IBD2, respectively—i.e. the sums of genetic lengths of all IBD1 or IBD2 segments divided by the total genetic length of the genome analyzed. We calculate r as twice the kinship coefficient or $r = \frac{k_{ij}^{(1)} + k_{ij}^{(2)}}{2}$ (Ramstetter et al. 2017).

The first analysis uses exact IBD segments from pairs of individuals that each have one of 13 genetic relationships (Table 1). To reduce the influences of randomness, we replicated this analysis by performing 80 independent simulations. We also analyzed three different distributions of numbers of pairs per degree D : uniform, exponential growth, or a slow-exponential growth function where the number of pairs increases exponentially with degree for both the exponential growth and slow-exponential growth distributions (Fig. 1). The exponential growth function is potentially a more realistic distribution of relatives than the uniform, while the slow-exponential growth is intermediate between the two.

Figure 2a shows the average MI of the simulated pairs computed over all 80 runs (Methods, “Simulated data”). For each distribution shape, the MI between the multivariate feature (r, n) and univariate D is the greatest, followed by $I(r; D)$ and $I(n; D)$. To quantify the relative increase in MI when including both n and r , we used a normalized MI gain $G_N(x) \equiv \frac{I((r,n); D) - I(x; D)}{I((r,n); D)}$ where

$x \in \{r, n\}$. The normalized MI gain $G_N(r)$ (the increase in information gained from using (r, n) beyond that of only using r) is 0.030 for the uniform distribution, 0.040 for the slow-exponential growth, and 0.068 for the exponential growth. Greater MI indicates a stronger dependency between the features and D , and therefore classifying D based on features with greater MI should yield greater recall. At $G_N(r)$ of 0.068 for the exponential growth distribution, we expect that incorporating numbers of perfectly detected segments could meaningfully improve classification of degrees of relatedness compared to using r alone, especially for higher order degree pairs. In turn, the normalized gain over using segment number alone, $G_N(n)$, is 0.15 for the uniform distribution, 0.14 for the slow-exponential growth, and 0.13 for the exponential growth, consistent with the use of r dramatically improving classification rates compared to only using n , regardless of the distribution of D .

Across all three feature sets, the MI is maximal for the uniform D distribution and decreases as the distribution becomes more exponential. By construction, the exponential growth distributions have a higher proportion of high-degree relative pairs compared to the uniform distribution. Therefore, the IBD features from higher degree pairs share less information with D than lower degree pairs. This is consistent with observations from classification analyses that show that the recall of degree inference decreases as the degree increases (Ramstetter et al. 2017).

To better understand how r and n relate to each other as well as to D , we calculated MI between these two features using the exact IBD segments and conditioned on the degree of relatedness (Fig. 2b). The amount of shared information between features r and n monotonically increases with degree of relatedness, meaning that in higher degree pairs r and n have increased redundancy. Therefore, using both features has less benefit for classification in higher degrees. Nevertheless, both r and n individually become less informative about D with increasing degree, so any additional information can be useful.

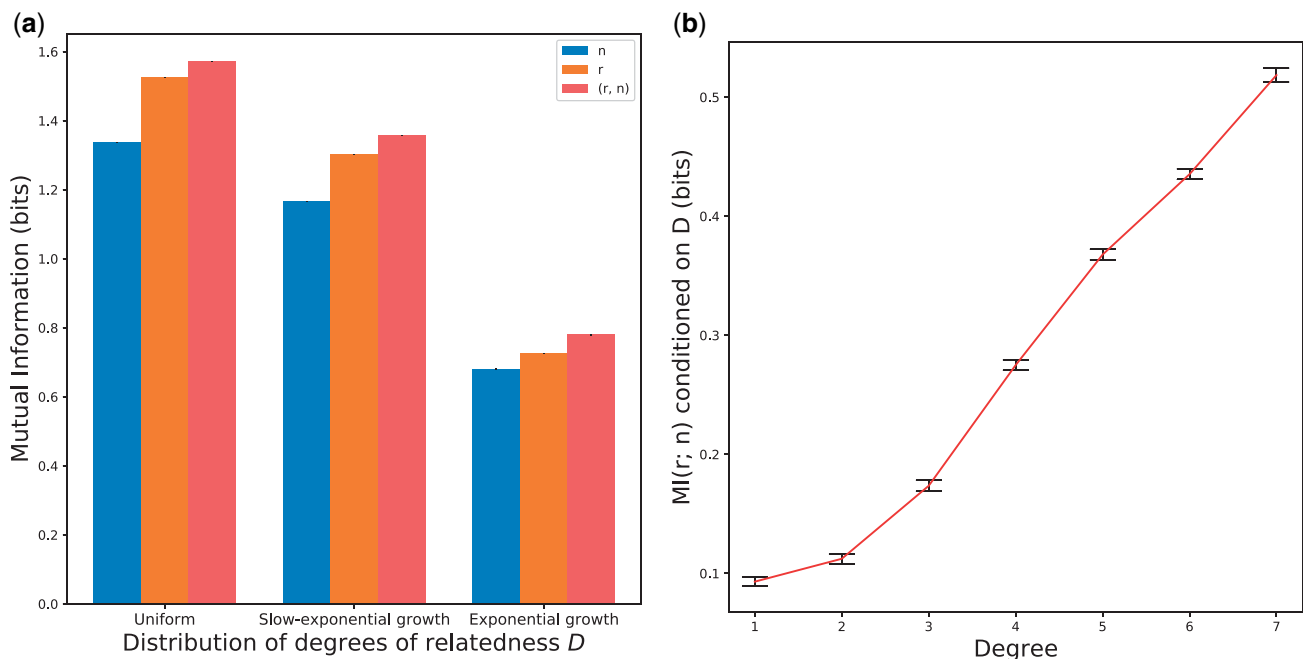


Fig. 2. Mutual information between relative pairs calculated using exact IBD segments. MI between (a) various IBD feature sets and D for each D distribution shape and (b) r and n conditioned on the relatives’ degree of relatedness. All MI quantities are averaged over 80 independent runs, and the values in (b) are calculated using the uniform distribution with 33,000 pairs per degree. Error bars indicate one standard error and are barely visible in (a) (all are of order 10^{-3}).

Bayes classification and statistical tests of exact and inferred IBD segments

As MI quantities from exact IBD segments suggest the potential for sizeable improvements by using (r, n) to determine D , we sought to understand whether parallel results arise from explicit relatedness classification. To that end, we simulated another 210,000 pairs of relatives for training, this time producing genetic data for them using genotypes from UK Biobank unrelated samples as pedigree founders (Methods, “Simulated data”). We detected IBD segments in these samples with IBIS and used the resulting r and n quantities to train Bayes classifiers. For comparison, we further trained a separate set of classifiers using the exact IBD segments from the same simulated pairs (Methods, “Bayes classification”). Using Bayes classification allowed us to incorporate our prior knowledge of the distribution of D to better determine the pairs’ degrees, and also more closely mirrors the mathematical basis of MI. For both the inferred and exact statistics, we generated a set of three classifiers, one trained only on the coefficient of relatedness r , one on the IBD segment number n , and a third on the vector (r, n) . We tested both the exact and inferred segment classifiers on 80 independent simulated datasets containing 3,000 simulated relative pairs per degree, again inferring segments with IBIS. (Genetic data for testing pairs was produced identically to the training pairs, as noted above.)

Figure 3 shows the recalls of these classifiers as a function of degree and also shows the recall differences between classifiers trained on (r, n) and r . We also show the proportions and types of misclassifications in the inferred and exact datasets in Supplementary Figs. 4 and 5. Almost all misclassified pairs are inferred as an adjacent degree of relatedness compared to the truth (i.e. one degree closer or more distant). Note that we do not report accuracy results for seventh-degree relatives as these pairs act as

an “unrelated” class that provide bounds on sixth-degree relatedness classification.

Overall, recalls for all three classifiers decrease monotonically as a function of the degree of relatedness. For first and second-degree pairs, the classifiers trained on r and (r, n) both have nearly perfect recall values of over 0.99. For higher degree pairs from third through sixth degree, the recalls of the r and (r, n) -trained classifiers fall from over 0.93 (third degree) to below 0.55. This is consistent with previous observations from real relatives (Ramstetter et al. 2017), and aligns well with our results based on MI: The features of higher degree pairs share less information with D , meaning that the IBD signals of higher degree pairs tell the classifier less about their true D (see misclassification rates in Supplementary Figs. 4 and 5). The classifier trained on n alone performs poorly in all but degree one: For second-degree relatives, the classifier trained on inferred segments has a recall of only 0.86, and in third through sixth-degree relatives its recall is 0.06 to 0.27 units lower than those of the classifier trained on r . The results for the classifier trained on exact segments are qualitatively similar to those of the inferred-segment classifier.

In general, when using either exact or inferred IBD segments, the classifiers trained on (r, n) outperform those trained on r for every degree. One exception is in the inferred IBD segments for sixth-degree pairs, where the classifier trained on r has a recall of 0.54 while the classifier trained on (r, n) has a recall of 0.53. This decrease in recall is counter-intuitive because the (r, n) classifier is trained on a strictly larger feature set and so has more information than the r classifier. In addition to general stochasticity introduced by segment detection for these distant relatives, it may be that this decrease is caused by the distributions of segment numbers inferred by IBIS

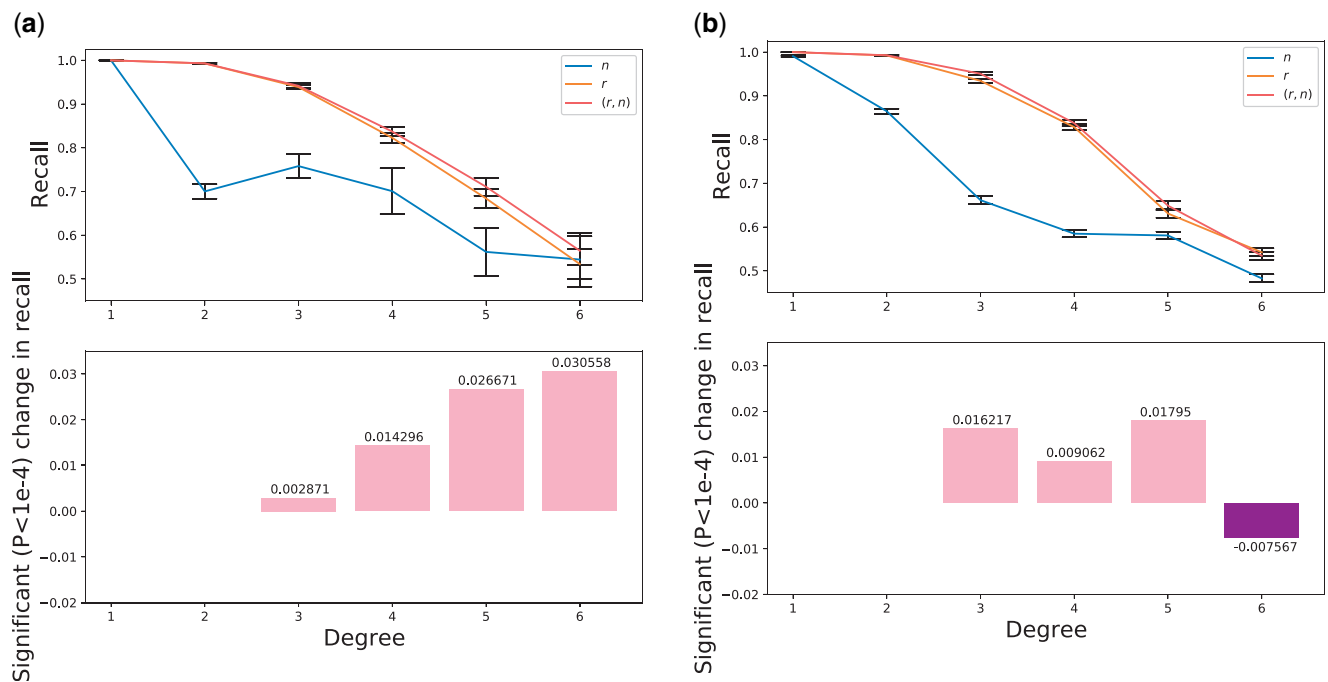


Fig. 3. Recalls of Bayes classifiers for first through sixth-degree relatives. Results are from classifiers trained on (a) exact and (b) inferred segments with features n , r , or (r, n) . The recalls for both (a) and (b) are calculated using the uniform distribution of 3,000 pairs per degree and averaged over 80 independent runs. For each degree, the lower subplot shows the corresponding significant ($P < 10^{-4}$) change in recall between classifiers (r, n) and r [positive values have greater recall in the (r, n) classifier]. Error bars indicate one standard error.

(Supplementary Figs. 6 and 7): IBIS does not detect segments smaller than 7 cM and so the distribution of numbers of detected segments for fifth and sixth-degree pairs have lower means than those of the exact segments and are more similar to each other.

We ran two-sided independent sample t-tests on the recalls from the (r, n) and r classifiers trained on the inferred IBD segments. Except for the first degree relatives, in which all three classifiers have recalls of nearly 1.0, and the second-degree pairs, in which the two classifiers containing r have above 0.99 recall, the differences in recall between the (r, n) and r classifiers are significant ($P < 10^{-7}$) but small in magnitude. These differences range from -0.00756 to 0.0179 in third through sixth-degree pairs. In turn, for the classifiers trained on exact IBD segments, the (r, n) classifier has significantly greater recall than the r classifier in third through sixth-degree relatives ($P < 10^{-4}$). In this case, the improvement in recall ranges from 0.0029 to 0.031 , suggesting that better IBD segment inference would meaningfully benefit classification with (r, n) (Fig. 3).

Comparison with IBIS and ERSA

To put these results in the context of existing methods, we compared our Bayes classifier with IBIS’s built-in relative classifier and with ERSA, another method that models relatedness using IBD segment number (as well as with segment length). This analysis uses for testing another independent set of 3,000 pairs per degree, again simulated from UK Biobank individuals. Our Bayes classifier remained trained on the same 210,000 pairs as above.

In general, all three methods performed comparably (see recalls in Fig. 4). The accuracy of the Bayes classifier closely tracks that of IBIS, which may be because the Bayes method takes IBIS segments as input. At the two extremes of relatedness we considered, all three methods have similar recalls for first and sixth degree relatives, with differences smaller than 0.01. The Bayes classifier has nearly identical recall to IBIS in second and third degree pairs (the differences are bounded above by 0.004), whereas ERSA’s recalls for these degrees are 0.06 and 0.02 units smaller, respectively. An analysis with real relatives also found

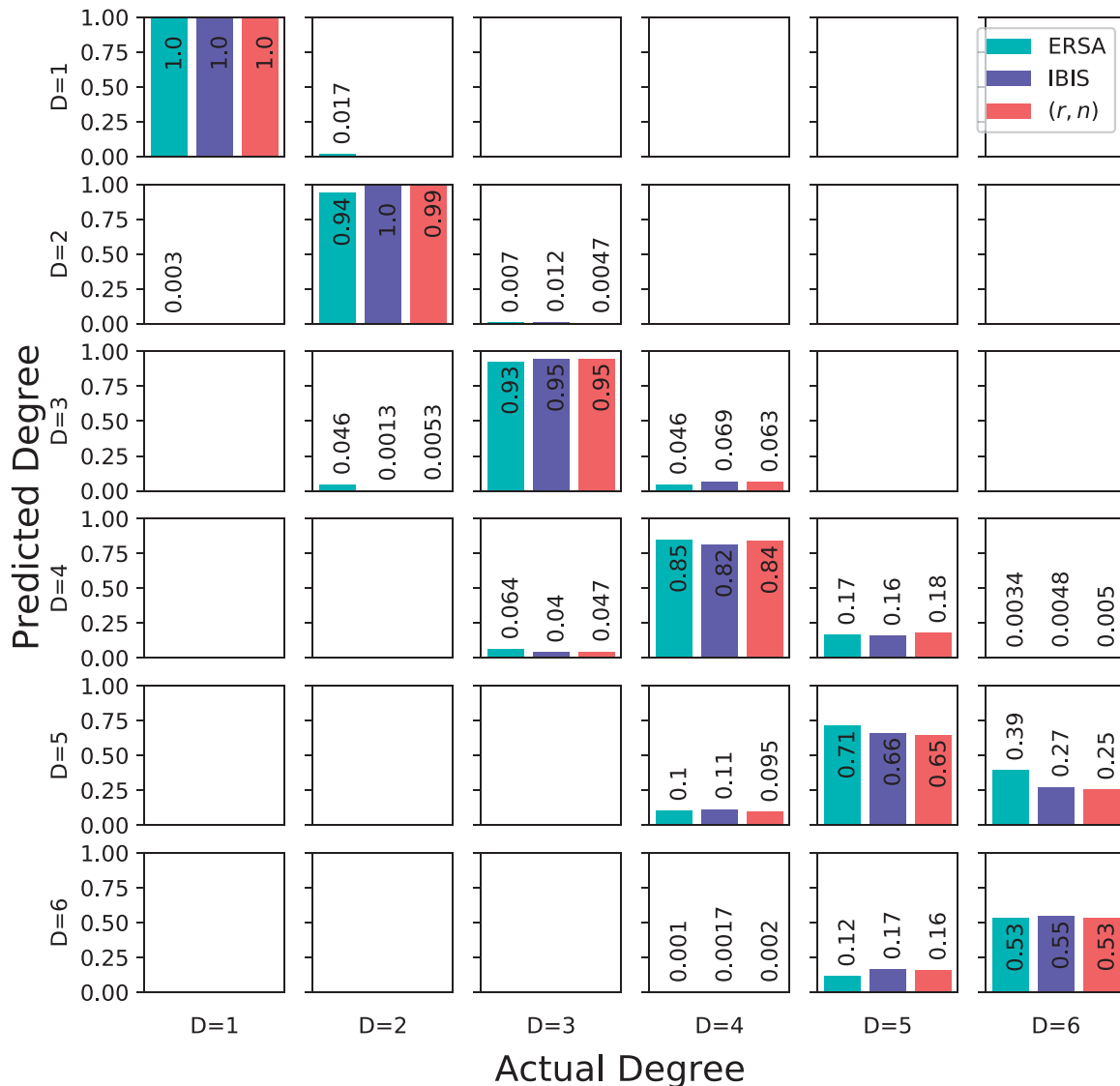


Fig. 4. Confusion matrix of recalls of ERSA, IBIS, and our Bayes classifier trained on (r, n) .

that ERSA's second-degree classification rates are reduced compared to other approaches (Ramstetter et al. 2017). For fourth-degree relatives, ERSA has a recall 0.01 units higher than the Bayes classifier, and 0.035 units higher than IBIS. ERSA also outperformed the Bayes classifier and IBIS on fifth-degree pairs by 0.067 and 0.054 units, respectively. ERSA's improved performance compared to the other two methods maybe because of its use of ≥ 2.5 cM segments (instead of ≥ 7 cM segments from IBIS). Consistent with this, simulated fourth and fifth-degree relatives have a non-trivial proportion of 3–7 cM segments (Supplementary Fig. 8)—suggesting that these undetected IBD segments may lead to more erroneous calculations by IBIS and the Bayes classifier. Another factor benefiting ERSA is its population model that accounts for background relatedness, which may help it in this and other datasets. Additionally, we used perfectly phased data as input to GERMLINE (Gusev et al. 2009), and we supplied the resulting segment calls to ERSA (Methods, “Simulated data”). Notably, ERSA's higher recalls for fourth and fifth-degree pairs are close to the range of the Bayes classifier's recalls using exact IBD segments (in fact, ERSA outperforms the exact Bayes classifier in these degrees by 0.012 and 0.0031, respectively). Finally, considering run time, the Bayes classifier is efficient, taking on average 1 min 40 s to analyze the test data and 7 s to train on the 210,000 training pairs. In contrast, ERSA takes more than 3.5 CPU days to classify the testing pairs.

Discussion

In this study, we sought to examine how much incorporating the number of IBD segments together with the coefficient of relatedness of a relative pair improves degree of relatedness inference. We thus provided both a theoretical MI analysis using simulated exact IBD segments and a machine learning-based classification analysis using exact and inferred segments. The results using exact segments show that including IBD segment numbers can nontrivially enhance relatedness inference quality, especially for distant relatives. However, the analyses using inferred segments reveal that IBD detection errors—including false negatives for segments shorter than 7 cM—meaningfully limit this improvement. Indeed, the performance of our machine learning classifier is almost indistinguishable from IBIS (Fig. 4), which does not use segment numbers. With the potential development of more accurate IBD detection tools in the future—including for whole-genome sequencing data—use of IBD segment numbers in relatedness inference models may be worth considering.

We introduced a machine learning-based classifier and demonstrated that it has comparable accuracy to two state-of-the-art methods and is computationally efficient. Because we fit the classifier to population-specific training data [instead of using fixed kinship thresholds for each degree (Ramstetter et al. 2017; Seidman et al. 2020)], it implicitly accounts for background IBD sharing and erroneous IBD signals. This approach differs from model-based methods such as ERSA in that it makes no assumptions about the distributions of IBD segment lengths or numbers with respect to relatedness degrees. Those assumptions can be violated in populations with small effective size or a historical founder effect (Huff et al. 2011). Our trials of this machine learning method suggest that even without large numbers of (labeled) real relatives, simulating relatives is a way to enable this data-driven approach to relatedness inference. In addition, both the

machine learning classifier and the MI analyses can be easily extended to include other IBD features such as the minimum or maximum IBD segment length between a pair.

An important factor in attempting to utilize IBD segment numbers is their accurate detection. Switch errors profoundly influence segment number estimates when using phase-based IBD detectors (Dimitromanolakis et al. 2019; Freyman et al. 2021; Seidman et al. 2020; Naseri et al. 2021). Our use of IBIS segments in our classifier was motivated by IBIS's ability to call IBD segments in unphased data—one of only a few methods to do so (Dimitromanolakis et al. 2019)—which is key to avoiding biased segment number estimates (Seidman et al. 2020). ERSA takes inferred IBD segments from the phase-based IBD detector GERMLINE. To exclude the possibility of phasing errors impacting ERSA's performance, the phased data we provided GERMLINE was that generated by the simulator, thus being perfect up to the limit of the haplotypes input to Ped-sim. In particular, these haplotypes do not contain switch errors in IBD segments between the simulated relatives. It is possible that ERSA's superior performance in classifying fifth-degree relatives is enhanced by its segment detection in these data.

In general, our analyses are consistent with prior work showing that relatedness inference can achieve high recall for up to third-degree relatives. However, two recent studies have focused on distinguishing relationship types of the same degree, especially three types of second-degree relatives (Williams et al. 2020; Qiao et al. 2021). In this setting, IBD segment numbers can provide useful information, such as for distinguishing avuncular from grandparent-grandchild pairs (Henn et al. 2012). Still, for degree of relatedness inference, even when using exact IBD segments, the classification recalls for distant relatives—i.e. those beyond fourth degree—are limited (Fig. 3a). This suggests that pairwise IBD information might not be sufficient to reliably infer distant relatives, regardless of segment quality. Approaches that leverage multiway IBD signals to infer more distant relatives can achieve considerably higher accuracy than those of pairwise methods (Staples et al. 2016; Ramstetter et al. 2018). Even so, these multiway methods are built on pairwise classifiers, so understanding and improving pairwise relatedness classification remains an important fundamental problem for relatedness inference.

Data availability

Data for exact segments were generated using the open source Ped-sim simulator and it is possible to generate data with the same expected summary statistics given the pedigree definition (def) files from this study. Genetic data were simulated using Ped-sim based on input UK Biobank haplotypes. The latter is available to qualified researchers from the UK Biobank. Ped-sim def files and the code we used to calculate mutual information and perform Bayes classification are available from <https://github.com/jeshaitan/mutual-information-relatedness-inference> (last accessed February 07, 2022).

Supplemental material is available at G3 online.

Acknowledgments

The authors thank Debbie Kennett for conversations about genetic genealogy and Shai Carmi and the anonymous reviewers

for helpful comments on the manuscript. Computing was performed on a cluster administered by the Biotechnology Resource Center at Cornell University. This research has been conducted using the UK Biobank Resource under Application Number 19947.

Funding

Funding for this work was provided by National Institutes of Health grant R35 GM133805.

Conflicts of interest

A.L.W. is an employee of 23andMe, has stock in 23andMe, and is the owner of HAPI-DNA LLC. The other authors declare no competing interests.

Literature cited

- Bennasar M, Hicks Y, Setchi R. Feature selection using joint mutual information maximisation. *Expert Syst Appl*. 2015;42(22):8520–8532.
- Bhéret C, Campbell CL, Auton A. Refined genetic maps reveal sexual dimorphism in human meiotic recombination at multiple scales. *Nat Commun*. 2017;8:14994.
- Browning BL, Browning SR. Improving the accuracy and efficiency of identity-by-descent detection in population data. *Genetics*. 2013;194(2):459–471.
- Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, Motyer A, Vukcevic D, Delaneau O, O'Connell J, et al. The UK Biobank resource with deep phenotyping and genomic data. *Nature*. 2018;562(7726):203–209.
- Caballero M, Seidman DN, Qiao Y, Sannerud J, Dyer TD, Lehman DM, Curran JE, Duggirala R, Blangero J, Carmi S, et al. Crossover interference and sex-specific genetic maps shape identical by descent sharing in close relatives. *PLoS Genet*. 2019;15(12):e1007979.
- 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;4:7.
- Devroye L, Györfi L, Lugosi G. *A Probabilistic Theory of Pattern Recognition*, Vol. 31. New York, NY: Springer Science & Business Media; 2013. <https://link.springer.com/book/10.1007/978-1-4612-0711-5#about>
- Dimitromanolakis A, Paterson AD, Sun L. Fast and accurate shared segment detection and relatedness estimation in un-phased genetic data via TRUFFLE. *Am J Hum Genet*. 2019;105(1):78–88.
- Freyman WA, McManus KF, Shringarpure SS, Jewett EM, Bryc K; 23 and Me Research Team, Auton A. Fast and robust identity-by-descent inference with the templated positional Burrows–Wheeler transform. *Mol Biol Evol*. 2021;38(5):2131–2151.
- Gusev A, Lowe JK, Stoffel M, Daly MJ, Altshuler D, Breslow JL, Friedman JM, Pe'er I. Whole population, genome-wide mapping of hidden relatedness. *Genome Res*. 2009;19(2):318–326.
- Henn BM, Hon L, Macpherson JM, Eriksson N, Saxonov S, Pe'er I, Mountain JL. Cryptic distant relatives are common in both isolated and cosmopolitan genetic samples. *PLoS One*. 2012;7(4):e34267.
- Hoque N, Bhattacharyya DK, Kalita JK. MIFS-ND: a mutual information-based feature selection method. *Expert Syst Appl*. 2014;41(14):6371–6385.
- Housworth E, Stahl F. Crossover interference in humans. *Am J Hum Genet*. 2003;73(1):188–197.
- Huff CD, Witherspoon DJ, Simonson TS, Xing J, Watkins WS, Zhang Y, Tuohy TM, Neklason DW, Burt RW, Guthery SL, et al. Maximum-likelihood estimation of recent shared ancestry (ERSA). *Genome Res*. 2011;21(5):768–774.
- Jewett EM, McManus KF, Freyman WA, Auton A; 23andMe Research Team. Bonsai: an efficient method for inferring large human pedigrees from genotype data. *Am J Hum Genet*. 2021;108(11):2052–2070.
- Lee J, Kim DW. Mutual information-based multi-label feature selection using interaction information. *Expert Syst Appl*. 2015;42(4):2013–2025.
- Manichaikul A, Mychaleckyj JC, Rich SS, Daly K, Sale M, Chen WM. Robust relationship inference in genome-wide association studies. *Bioinformatics*. 2010;26(22):2867–2873.
- Naseri A, Shi J, Lin X, Zhang S, Zhi D. RAFFI: accurate and fast familial relationship inference in large scale biobank studies using RaPID. *PLoS Genet*. 2021;17(1):e1009315.
- Ott J, Wang J, Leal SM. Genetic linkage analysis in the age of whole-genome sequencing. *Nat Rev Genet*. 2015;16(5):275–284.
- Qian W, Shu W. Mutual information criterion for feature selection from incomplete data. *Neurocomputing*. 2015;168:210–220.
- Qiao Y, Sannerud JG, Basu-Roy S, Hayward C, Williams AL. Distinguishing pedigree relationships via multi-way identity by descent sharing and sex-specific genetic maps. *Am J Hum Genet*. 2021;108(1):68–83.
- Ramstetter MD, Dyer TD, Lehman DM, Curran JE, Duggirala R, Blangero J, Mezey JG, Williams AL. Benchmarking relatedness inference methods with genome-wide data from thousands of relatives. *Genetics*. 2017;207(1):75–82.
- Ramstetter MD, Shenoy SA, Dyer TD, Lehman DM, Curran JE, Duggirala R, Blangero J, Mezey JG, Williams AL. Inferring identical-by-descent sharing of sample ancestors promotes high-resolution relative detection. *Am J Hum Genet*. 2018;103(1):30–44.
- Ross BC. Mutual information between discrete and continuous data sets. *PLoS One*. 2014;9(2):e87357.
- Seidman DN, Shenoy SA, Kim M, Babu R, Woods IG, Dyer TD, Lehman DM, Curran JE, Duggirala R, Blangero J, et al. Rapid, phase-free detection of long identity-by-descent segments enables effective relationship classification. *Am J Hum Genet*. 2020;106(4):453–466.
- Stallard M, de Groot J. “Things are coming out that are questionable, we never knew about”: DNA and the new family history. *Journal of Family History*. 2020;45(3):274–294.
- Staples J, Nickerson DA, Below JE. Utilizing graph theory to select the largest set of unrelated individuals for genetic analysis. *Genet Epidemiol*. 2013;37(2):136–141.
- Staples J, Qiao D, Cho MH, Silverman EK, Nickerson DA, Below JE; University of Washington Center for Mendelian Genomics. PRIMUS: rapid reconstruction of pedigrees from genome-wide estimates of identity by descent. *Am J Hum Genet*. 2014;95(5):553–564.
- Staples J, Witherspoon DJ, Jorde LB, Nickerson DA, Below JE, Huff CD, ; University of Washington Center for Mendelian Genomics. PADRE: pedigree-aware distant-relationship estimation. *Am J Hum Genet*. 2016;99(1):154–162.
- Thompson EA. Identity by descent: variation in meiosis, across genomes, and in populations. *Genetics*. 2013;194(2):301–326.

- Voight BF, Pritchard JK. Confounding from cryptic relatedness in case-control association studies. *PLoS Genet.* 2005;1(3):e32.
- Weir BS, Anderson AD, Hepler AB. Genetic relatedness analysis: modern data and new challenges. *Nat Rev Genet.* 2006;7(10):771–780.
- Williams CM, Scelza B, Gignoux CR, Henn BM. A rapid, accurate approach to inferring pedigrees in endogamous populations. *bioRxiv.* 2020. <https://doi.org/10.1101/2020.02.25.965376>
- Young AI, Frigge ML, Gudbjartsson DF, Thorleifsson G, Bjornsdottir G, Sulem P, Masson G, Thorsteinsdottir U, Stefansson K, Kong A. Relatedness disequilibrium regression estimates heritability without environmental bias. *Nat Genet.* 2018;50(9):1304–1310.
- Zaitlen N, Kraft P, Patterson N, Pasaniuc B, Bhatia G, Pollack S, Price AL. Using extended genealogy to estimate components of heritability for 23 quantitative and dichotomous traits. *PLoS Genet.* 2013;9(5):e1003520.

Communicating editor: R. Hernandez