

## DISPENSING WITH UNNECESSARY ASSUMPTIONS IN POPULATION GENETICS ANALYSIS

OLIVIER LABAYLE PABET<sup>1,2</sup>, KELSEY TETLEY-CAMPBELL<sup>1,2</sup>, MARK VAN DER LAAN<sup>4</sup>,  
CHRIS P. PONTING<sup>1</sup>, SJOERD VIKTOR BEENTJES<sup>1,3,\*</sup> AND AVA KHAMSEH<sup>1,2,\*</sup>

**ABSTRACT.** Parametric assumptions in population genetics analysis – including linearity, sources of population stratification and the gaussianity and additivity of errors – are often made, yet a principled argument for their (approximate) validity is not given. We present a unified statistical workflow, called TarGene, for targeted estimation of effect sizes, as well as two-point and higher-order epistatic interactions of genomic variants on polygenic traits, that dispenses with these unnecessary assumptions. Our approach is founded on Targeted Learning, a framework for estimation that integrates mathematical statistics, machine learning and causal inference to provide mathematical guarantees and realistic p-values. TarGene defines effect sizes of variants, as well as two-point and higher-order interactions amongst genomic variants on traits in a model-independent manner, thus avoiding all-too-common model-misspecification whilst taking advantage of a library of parametric and state-of-the-art non-parametric algorithms. TarGene data-adaptively incorporates confounders and sources of population stratification, accounts for population dependence structures and controls for multiple hypothesis testing by bounding any desired type I error rate. Extensive simulations demonstrate the necessity of this model-independent approach. We validate the effectiveness of our method by reproducing previously verified effect sizes on UK Biobank data, whilst simultaneously discovering non-linear effect sizes of additional allelic copies on trait or disease. To exemplify this, we demonstrate that for the FTO variant rs1421085 effect size on body mass index (BMI), the addition of one copy of the C allele is associated with 0.77 kg/m<sup>2</sup> (95% CI: 0.68 – 0.85) increase, while the addition of the second C copy non-linearly adds 1.31 kg/m<sup>2</sup> (95% CI: 1.19 – 1.43) to BMI. TarGene thus extends the reach of current genome-wide association studies by simultaneously (i) allowing for the classification of the types of SNPs and phenotypes for which such non-linearities occur, whilst (ii) data-adaptively incorporating complex non-linear relations between phenotype, genotype, and confounders, as well as (iii) accounting for strong population dependence such as island cohorts. The method provides a platform for comparative analyses across biobanks, or integration of multiple biobanks and heterogeneous populations to increase power, whilst controlling for population stratification and multiple hypothesis testing.

---

<sup>1</sup>MRC HUMAN GENETICS UNIT, INSTITUTE OF GENETICS & CANCER, UNIVERSITY OF EDINBURGH, EDINBURGH EH4 2XU, UNITED KINGDOM.

<sup>2</sup>SCHOOL OF INFORMATICS, UNIVERSITY OF EDINBURGH, EDINBURGH EH8 9AB, UNITED KINGDOM

<sup>3</sup>SCHOOL OF MATHEMATICS, UNIVERSITY OF EDINBURGH, EDINBURGH EH9 3FD, UNITED KINGDOM.

<sup>4</sup>DIVISION OF BIostatISTICS, UNIVERSITY OF CALIFORNIA, BERKELEY, CA

\*CORRESPONDING AUTHORS

## 1. INTRODUCTION

The principal challenge in population genetics is to pinpoint trait-causal genetic variants and the biological mechanisms through which they act [22]. Approaches to estimate phenotype-genotype relationships are well-established yet rely on strong parametric assumptions whose validity is often unsubstantiated. These assumptions include linearity of the DNA variant-trait relationship, linearity of the functional dependence of trait on covariates and confounders such as sources of population stratification and other fixed effects including batch, age and sex, as well as the gaussianity and additivity of errors, encoded as random effects, in the commonly employed linear mixed models (LMMs) [35, 23]. As these assumptions rarely hold they introduce biases that can yield false conclusions [16]. When such assumptions are invalid, both effect size and confidence interval estimates will tend to be biased, resulting in overly optimistic p-values and statistically missing the ground truth. More specifically, the desired level of statistical coverage<sup>1</sup> (*e.g.*, 95%) suffers both from bias due to model-misspecification as well as bias in the confidence interval estimates. This effect is becoming more relevant as cohort sizes grow to  $10^5 - 10^6$  or more (Fig. 1). This Curse of Big Data is especially likely because modern biobank-scale data sets, with their ever-growing size, result in very small variance, which exposes bias (Fig. 1A).

Rather than addressing the root cause of each bias, current methods often seek to ameliorate their effects post hoc [16]. An exception to this generalisation is KnockoffGWAS which makes no parametric assumptions regarding the distribution of the phenotype conditional on the genotypes [33]. This method also controls the false discovery rate (FDR) whilst accounting for population structure. Nevertheless, KnockoffGWAS does not compute effect sizes or epistatic interactions, and only reports statistical significance, and thus is unable to (i) detect genetic non-linearity, (ii) infer the strength and sign of variant effect sizes and interactions, and (iii) determine whether or not a variant of interest is stratified across the population. Estimation of these quantities is essential for explaining how variants, via biological mechanisms and regulatory functions, modify a trait or disease risk.

New Machine Learning (ML) approaches have begun to be applied to many aspects of genetics and genomics but not, until recently, to population genetics analysis. A recent ML method [25] partially addresses the issue of non-linear and interacting covariates in the phenotype-genotype relations by modeling these using deep neural networks. However, this approach relies on the independence of participants, an unrealistic assumption common to many ML methods, despite it being well understood that there may be considerable cryptic relatedness in population cohorts; see, *e.g.*, Fig. 3b in [7]. A priori, it is unknown whether and how such structure affects effect size and interaction estimates. This casts doubt on the appropriateness of ML methods ignoring dependence, and on the reliance of parametric methods such as LMMs on restrictive assumptions in order to model dependence.

---

<sup>1</sup>The statistical coverage (or *coverage probability*) is the probability that the confidence interval, as constructed in the statistical inference procedure, contains the true value. In practice, 95% coverage is often desired.

Furthermore, deep learning methods, such as [25], (i) require extensive hyper-parameter tuning that may be phenotype-genotype dependent, (ii) cannot take into account dependence structure amongst individuals, *i.e.*, require independent and identically distributed (i.i.d.) training data, and (iii) have no mathematical guarantees to adequately control for type I errors or provide coverage at the desired levels.

Others have modelled non-linearities of allelic copies on trait based on a direct association test, see *e.g.*, [18]. However, this approach still suffers from the aforementioned issues in the functional dependence of trait on population stratification and other covariates and confounders, as well as the potentially non-trivial dependence structure amongst the data of individuals as collected in various databases.

Here, we introduce Targeted Genomic Estimation (TarGene), a method that accurately estimates effect sizes, pairwise and higher-order interactions amongst variants, as well as gene-environment interactions by using flexible statistical and ML algorithms whilst equipping them with mathematical guarantees. TarGene does not suffer from any of the aforementioned shortcomings. Specifically, TarGene identifies genetic non-linearities whilst simultaneously accounting for (i) any relationship (*e.g.*, non-linear) amongst phenotype, genotype, confounders and covariates, (ii) stratification of variants across the population of interest, and (iii) population structure and dependence amongst individuals in the database [12]. It does so without the need for any parametric assumptions, whilst also controlling the desired type I error rate. Fig. 2 summarises TarGene vs LMM-based approaches currently considered as gold standard in the population genetics literature in addressing the aforementioned complexities.

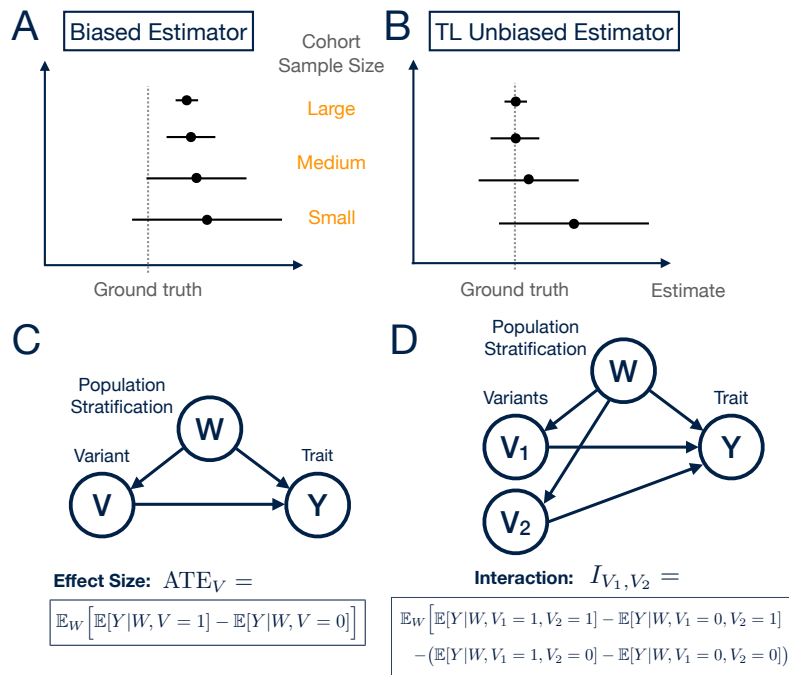
In brief, TarGene employs Super Learning, an ensemble machine learning method that estimates DNA variants' effect sizes and strengths of interactions. It does so by data-adaptively (via  $k$ -fold cross-validation) combining a library of parametric and non-parametric methods [38], with the latter making no assumptions regarding the form of the phenotype-genotype relationship [6]. The Super Learner (SL) combines methods to obtain a guaranteed optimal fit of the genotype-phenotype relationship. Any current or future, perhaps more powerful, estimation methods can be appended to its library of algorithms to improve performance. TarGene thus subsumes and supersedes any current model or ML based methods by incorporating them within the SL library. Estimates are further enhanced by Targeted Maximum Likelihood Estimation (TMLE) [37]. This step is crucial because it has the important qualities of being mathematically guaranteed to (a) reduce any residual bias due to model-misspecification, whilst (b) optimising the bias-variance trade-off, and (c) guaranteeing  $\sqrt{n}$ -convergence as the sample size  $n$  increases (Fig. 1B). Finally, TarGene updates variance estimates via a network approach [12] to accurately account for population dependence structure, thus resulting in realistic p-values.

As a proof of concept, we apply TarGene to (i) estimate the effect size of FTO intronic variant rs1421085, a candidate causal variant for obesity [11] on all 660 binary and 118

continuous traits in the UK Biobank. We demonstrate a non-linear effect size, *i.e.*, significantly differing effect sizes between the addition of the first and the second C allele, in 42 traits. This shows that TarGene can reveal non-linearity in genomic effect sizes whilst accounting for potentially complex heterogeneous population structures, as supported by its mathematical underpinnings.

To demonstrate further applications of TarGene, we investigate variants potentially interacting via the vitamin D receptor (VDR). VDR is a nuclear hormone receptor that binds to calcitriol, the active form of vitamin D, and then forms a complex with the retinoid-X receptor (RXRA). We consider three genetic variants associated with differential expression of each of these three molecules. Thus, we estimate pair-wise and 3-point epistatic interactions amongst rs7971418, rs1045570, and rs3755967, affecting 660 binary and 118 continuous traits in the UKBB. In the case of these variants, after multiple hypothesis testing, nevertheless we do not find evidence of significant interactions. This is not unexpected because, evidently, detection of epistasis is extremely challenging [43].





**FIGURE 1. A and B:** The Curse of Big Data. **A:** As sample size increases the bias of an estimator may not shrink sufficiently fast relative to the reduction in variance, leading to incorrect predictions for large sample sizes. Extensive simulations demonstrating this phenomenon are presented in Fig. 4 using different ground truth models fitted with a misspecified model, under various variant-covariate dependence structures (see Methods). **B:** Targeted Learning estimator provides correct statistical inference once it has been optimised for bias-variance trade-off in the TMLE step. **C:** Model-independent definition of effect size (Average Treatment Effect, ATE) of variant  $V$  on trait  $Y$ . We condition on sources of population stratification and other confounders,  $W$ , when estimating the ATE. We do this to get closer to a causal estimate of the variant on trait, up to linkage disequilibrium (LD). With additional molecular information to reduce LD and/or fine-mapping these estimates can approach being causal over associative. For a binary variable  $V$  the interpretation is: “Having correctly adjusted for confounders (sources of population stratification), what is the difference in the expected value of trait when a particular variant  $V$  is present ( $V = 1$ ) as compared to when it is not ( $V = 0$ )?”. Generalisation to a categorical variable  $V = 0, 1, 2$  is trivial, as all combinations can be written similarly. **D:** Model-independent definition of 2-point interaction, an extension of ATE with more than one variant, has been further generalised to higher-order interactions among  $n$  variants [4]. Interpretation: “Having adjusted for confounders, is the effect of variant  $V_1$  on trait  $Y$  modulated by the status of variant  $V_2$  and, if so, by how much and with which sign?” This defines an epistatic interaction between variants  $V_1$  and  $V_2$  with respect to a trait.

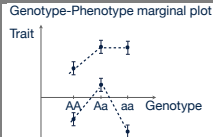
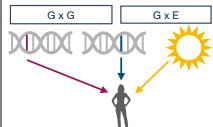

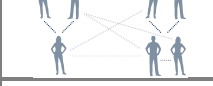

Schematic of complexities	Description and Examples	Properties	Current LMM-based methods	TarGene
	Adding another allele may increase/decrease the trait linearly or non-linearly  Detection of variants with heterozygote (dis)advantage	Can it be done in principle?	✓	✓ Built-in
		Is it done in practice?	✗	✓ Built-in
	Pair-wise and higher-order interactions: - Epistasis: $G \times G \times G \times \dots$ - Gene x Environment (categorical), e.g., Gene x Sex - Environment x Environment $\times \dots$ (categorical)	Can model (higher-order) interactions?	✓	✓
		Can explicitly target quantity of interest to maximise power?	✗	✓ TMLE
		Is (double-)robust against model-misspecification?	✗	✓ SL + TMLE
		Benefits from mathematical guarantees on inference?	✗	✓ TMLE
	Population heterogeneity and stratification, non-linear functional dependence of phenotype on genotype and covariates (e.g., PCs)	Incorporates non-linear/non-parametric flexible algorithms?	✗	✓ SL
		Data-driven cross-validation?	✗	✓ SL
		Benefits from guarantees on convergence to ground truth?	✗	✓ TMLE
	When individuals are not fully independent, this needs to be taken into account in statistical estimation. This is particularly important for inter-related populations e.g. Island communities	Can dispense with the additive noise assumption?	✗	✓ Sieve Plateau
		Can dispense with the gaussian noise assumption?	✗	✓ Sieve Plateau
	There is a trade-off between computational speed and guaranteeing ground truth coverage of estimates	Genome-wide and phenome-wide fit, at UK Biobank-scale sample size:	✓	Diverse SL library + Cross-validation ✗ Linear/LMM + TMLE have comparable speeds ✓
		Can avoid memory-intensive GRM inversion?	✗	✓

FIGURE 2. Comparison of TarGene vs current LMM-based methods considered as gold-standard in population genetics analysis. LMM-based methods use approximations that could result in reduced accuracy. In contrast, TarGene provides extensive flexibility and mathematical guarantees of ground truth coverage.

## 2. RESULTS

We introduce TarGene in five steps, first highlighting shortcomings in current methods as we proceed using simple, easily interpretable simulations before explaining how TarGene addresses these problems. The mathematical guarantees behind TarGene are presented in the Methods section. As a proof of concept, we apply TarGene to (i) estimate the effect size of FTO intronic variant rs1421085, a candidate causal variant for obesity [11], as well as (ii) pairwise interactions amongst three further variants, rs7971418, rs1045570, and rs3755967, and (iii) their 3-point interaction, on 660 binary and 118 continuous traits in the UK Biobank. These three variants are chosen due to their relevance to vitamin D receptor (VDR) function.

**Shortcomings in current methods.** We illustrate two shortcomings of current approaches to GWAS, as well as their deleterious ramifications, via simulations and examples:

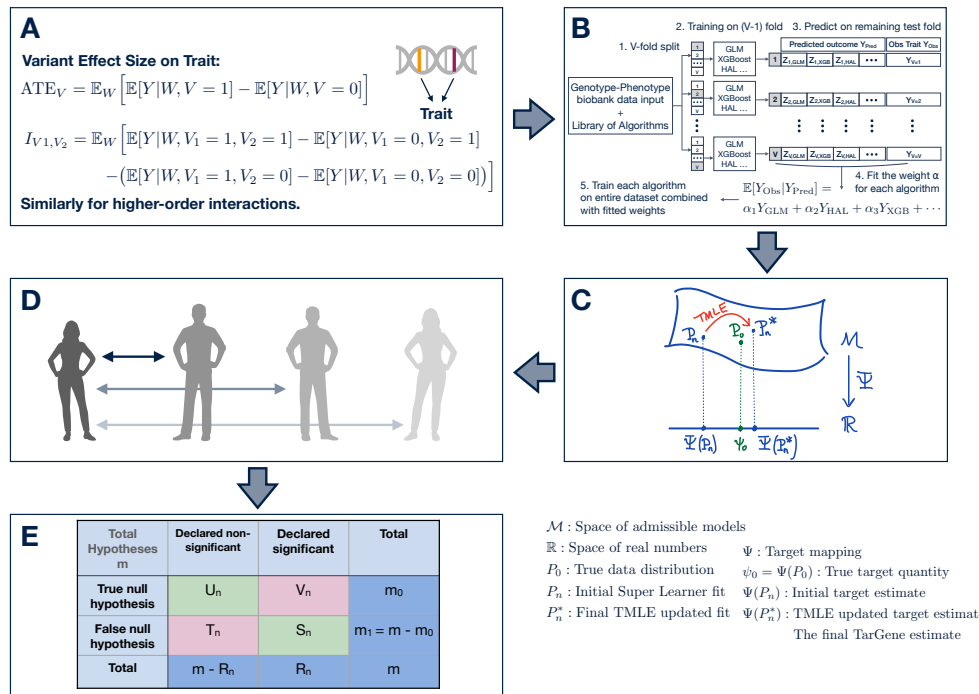


FIGURE 3. Workflow of TarGene. **A**: Model-independent definition of variant  $V$  effect size and two-point variant  $V_1, V_2$  interaction effect on trait  $Y$ , taking into account sources of population stratification  $W$ . **B**: The Super Learner fit with a library of algorithms including linear, logistic, XGBoost and Highly Adaptive Lasso (HAL) algorithms. The SL is used to obtain an initial fit of the trait as a function of variant(s) and sources of population stratification. The user is free to add any other algorithm to this library.  $k$ -fold cross-validation is performed to determine the algorithm or the combinations of algorithms with the lowest loss. **C**: Targeted Maximum Likelihood Estimation (TMLE) removes any residual model-misspecification bias for the target quantity of interest, resulting in optimal bias-variance trade-off with mathematical guarantees. **D**: Uncertainties in step C are updated to take into account dependences among individuals in the population. **E**: Multiple hypothesis testing is performed.

(i) model-misspecification, and (ii) non-principled choice of confounders and covariate relations.

*Model-misspecification in GWAS.* In genome-wide association studies (GWAS), the effect size of a variant  $V$  on trait  $Y$  in the presence of confounders (such as sources of population

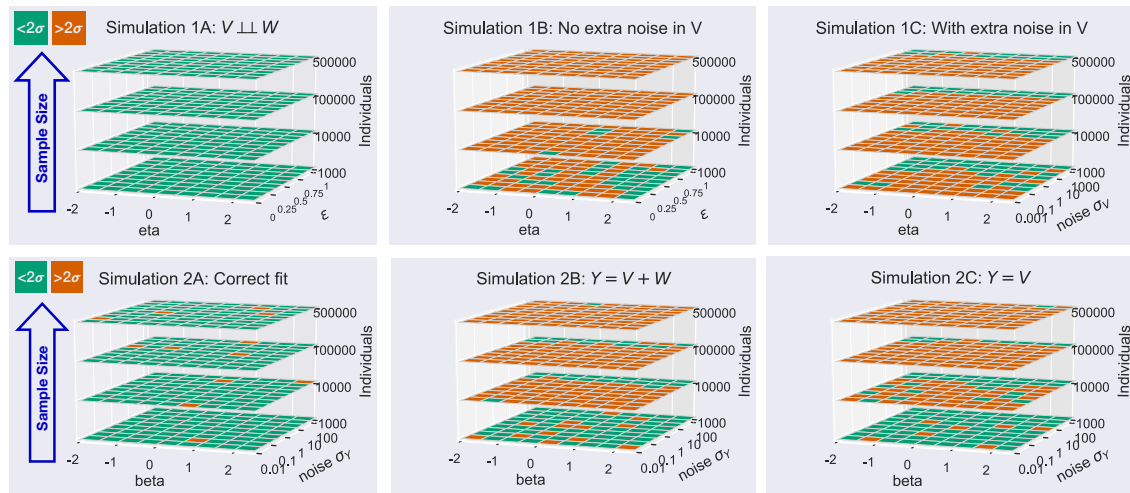


FIGURE 4. Model misspecification can yield biased estimates. In simulations, estimated effect sizes are either within (green) or outside (red) the 95% confidence interval. The striking conclusion of these simulations is that, when the model is misspecified, almost all estimates are invalid even at moderately large sample sizes. Top and bottom panels represent two distinct sets of simulations involving variables  $V$ ,  $W$  and  $Y$  representing variant, sources of population stratification and trait, respectively. **Top panels:** Trait  $Y$  is generated from  $V$  and  $W$  via an exponential model and yet is fitted with a linear model (see Methods 4 for details). Simulation **1A** (left): When  $V$  and  $W$  are *completely* independent (no population stratification present), the estimated effect size is correct (indicated in green) even when a misspecified linear model is fit to the exponential distribution. This result is irrespective of sample size. Simulation **1B** (middle): When, in a more realistic case where population stratification is present,  $V$  and  $W$  are dependent, then effect size estimates are incorrect (shown in red) when a misspecified linear fit is employed. This issue is exacerbated as the sample size grows, and manifests even with the slightest degree of Pearson or Spearman correlation (see Methods). Simulation **1C** (right): High levels of noise  $\sigma_V$  hide dependence between  $V$  and  $W$  so inference at small sample sizes may be within  $1\sigma$ . As sample size grows, model-misspecification is increasingly exposed. **Bottom panels:** The trait  $Y$  is generated from  $V$  and  $W$  via a polynomial model yet fitted with a misspecified polynomial model (see Methods). In misspecified models, effect size estimates become more incorrect as sample size increases. In all subfigures the true data distribution is generated according to the polynomial model in Eq. 17. The plots indicate  $\beta$  estimated using three different models, taking values in different parts of the parameter space. Simulation **2A** (left) indicates the estimated effect size where the true model is used. When there is no model misspecification, estimates are valid irrespective of sample size. In Simulations **2B** (middle) and **2C** (right) two misspecified models are used to estimate  $\beta$ .

stratification), is often referred to as the ‘beta’ coefficient. The following parametric form<sup>2</sup> (or its linear mixed model (LMM) equivalent) is assumed to govern the data:

$$Y = \alpha_0 + \underline{\alpha} \cdot \underline{W} + \beta V. \quad (1)$$

The effect size is then assumed to be equal to the coefficient  $\beta$ . However, in general this is not true because the data may be governed by a completely different probability distribution in which the effect size of the variant on trait is not equivalent to  $\beta$  above.

To illustrate this important point, suppose a researcher assumes that the data follows Eq. 1 with  $Y = \alpha_0 + \underline{\alpha} \cdot \underline{W} + \beta V$ . By this assumption, they thus declare the effect size of variant  $V$  on trait  $Y$  to be equal to  $\beta$  exactly. Nevertheless, they are only correct if the ground truth follows Eq. 1. If, instead, the ground truth differs, for example by addition of an interaction term between variant  $V$  and a confounder  $W_1$  (*i.e.*,  $Y = \alpha_0 + \underline{\alpha} \cdot \underline{W} + \beta V + \gamma W_1 V$ ) then the effect size of  $V$  on  $Y$  is made ambiguous: is the effect size  $\beta$  alone, or  $\gamma$ , or  $\beta + \gamma \mathbb{E}[W_1]$ ? In a second example, if the ground truth has an exponential functional form  $Y = \exp(V + (V + \epsilon)W)$  (where  $\epsilon = 0.2$  and  $W \sim \mathcal{N}(0, 1)$  is assumed standard normally distributed) then, again, the true effect size is unclear. The answer is far from  $\beta = 1$ : it equals  $\exp\left[\frac{1}{2}((1 + \epsilon)^2 + 2)\right] - \exp\left[\frac{1}{2}\epsilon^2\right] \approx 4.56$ .

We next used simulations to exemplify two ways in which model-misspecification, *i.e.*, fitting data with a model different from the one from which it was generated, results in biased and statistically incorrect estimates. In Fig. 4, we present results from these simulations describing two distinct phenotype-genotype relations in the presence of confounders such as population stratification. To demonstrate the ubiquity of incorrect conclusions produced by incorrectly specified models, we independently varied over parameter space (1) the true effect size, (2) the noise on the data, and (3) the sample size, where the latter ranges from 1000 participants to a biobank-scale of 500,000. These simulations show four features of model-misspecification (see Methods): (i) if the fitted model is far from the true data-generating distribution, then the slightest degree of measured correlation between  $V$  and  $W$  results in invalid inference of the effect sizes, even at smaller sample sizes; (ii) at any fixed level of noise, there always exists a sample size above which model-misspecification leads to invalid inference; (iii) replication is a necessary but not sufficient condition for declaring valid inference because fitting separate data samples drawn from the same distribution with the same (or a similar) misspecified model twice, results in equivalent invalid inference twice; and (iv) multiple hypothesis correction methods account for the testing of multiple hypotheses only, not for false discovery in a single hypothesis, such as those due to misspecified parametric models. We conclude that it is crucial to avoid subjective

---

<sup>2</sup>The equation  $g(Y) = \alpha_0 + \underline{\alpha} \cdot \underline{W} + \beta V$  with link function  $g(Y) = Y$  and effect size  $\beta$  is used for a continuous trait  $Y$ . However, the same arguments and simulations apply to the logit link function  $g(Y) = \text{logit}(Y)$  for a binary trait  $Y$  in a case/control setting with  $\beta$  the log odds ratio of case to control.

modelling choices: model-misspecification is likely to give rise to invalid inference, especially when working with data sets as large as population biobanks which include complex population structures as well as gene-environment interactions [3].

*Non-principled choice of confounders and covariate relations.* Variables which confound the relationship between a variant  $V$  and a trait or disease  $Y$ , such as population stratification, deserve careful consideration. Usually, population genetics studies choose these variables subjectively. The extent by which such choices affect effect size estimation and hinder replication remains unclear. The UK Biobank (UKBB) project reports that inclusion of 16-20 principal components (PCs), labelled by self-reported ethnicity, is sufficient to capture sources of population stratification [1] (see also [8, 21]). However, others [7] report use of up to 40 PCs and demonstrate significant population stratification across the entire UKBB cohort. Consequently, for any particular analysis it is unclear what PCs to include, and whether by conditioning on PCAs (*e.g.*, by including them as variables in GWAS fits) the true genetic signal is faithfully revealed. Additionally, there is no consensus on what covariate combinations (*e.g.*, array batch or UKBB assessment centre or sex  $\times$  age<sup>2</sup>; for example in [8] and [Neale-UKBB-GWAS](#)) should be used.

**TarGene provides mathematical guarantees and realistic p-values.** We next introduce TarGene and show that it does not suffer from any of the above shortcomings. We also illustrate the mathematical guarantees of coverage, asymptotically normal distribution, and realistic p-values our method provides.

TarGene is based on Targeted Learning (TL), a model-independent framework of estimation integrating causal inference, machine learning, and mathematical statistics and produces powerful estimators that are provably unbiased and efficient [39, 40]. TL consists of three stages: (1) Defining the quantity of interest model-independently; (2) Employing a diverse library of learning algorithms to learn the relevant portion of the true probability distribution, which results in an initial estimate of the quantity of interest; and, (3) Applying Targeted Maximum Likelihood Estimation (TMLE) to update and target the initial fit towards the quantity of interest, thereby removing any remaining bias. Here, we briefly explain each of these steps (for more, see Methods) and show results of TarGene applied to UK Biobank data. In later sections we explain how TarGene accounts for (4) cohort population dependence, and (5) multiple hypothesis testing.

*Step 1: TarGene defines the quantity of interest model-independently.* In GWAS, whether the ground truth probability distribution  $P_0(y, v, w)$  is linear, logistic, or has a more complicated form, the effect size of  $V$  on  $Y$  can be evaluated under any distribution  $P$ . By taking advantage of such model-independent definitions of both effect size and epistasis, TarGene resolves any ambiguity in estimation due to model-misspecification. Thus TarGene does not have to resort to any particular parametric assumptions. More specifically, effect size is otherwise known as the Average Treatment Effect (ATE) [39, 40], and defined

as

$$\Psi_1(P) = \mathbb{E}_W[\mathbb{E}[Y|W, V = 1] - \mathbb{E}[Y|W, V = 0]]. \quad (2)$$

This is interpreted as the expected phenotypic change in  $Y$  for variant  $V = 1$  relative to  $V = 0$ , whilst correcting for confounders  $W$ ; see also Fig. 1C.

It is common for genetic studies to examine additive models only. Nevertheless, because epistasis is crucial to understanding complex disease [24], we similarly provide the model-independent definition of epistatic interaction between two variants  $V_1$  and  $V_2$  leading to variation in a trait or disease risk  $Y$ ,

$$I_{V_1, V_2}(P) = \mathbb{E}_W \left[ \mathbb{E}[Y|W, V_1 = 1, V_2 = 1] - \mathbb{E}[Y|W, V_1 = 0, V_2 = 1] \right. \\ \left. - (\mathbb{E}[Y|W, V_1 = 1, V_2 = 0] - \mathbb{E}[Y|W, V_1 = 0, V_2 = 0]) \right] \quad (3)$$

This definition was first introduced in [4]; see also Fig. 1D. It is interpreted as the change in effect size of variant  $V_1$  on trait  $Y$  as variant  $V_2$  changes from 0 to 1.

*Step 2A: TarGene data-adaptively incorporates confounders.* Any population genetics analysis requires accounting for sources of population stratification in the specific cohort of interest. In the current literature, this is typically done by incorporating a number of Principal Components (PC) in a trait-independent analysis. Although standard, it is important that this analysis be performed for each database and repeated if a subset of the population in the cohort is used in order to obtain a lower bound on the number of PCs required. As a first example, in [7] the first 20 PCs are selected to account for population stratification based on labelling PC plots by self-reported ethnicity and visually inspecting their symmetry. This number may well be different for other more diverse cohorts, such as the *All of Us* [2] cohort or the *Million Veterans Program* [17]. As a second example, the analysis of the population structure in the whole UK Biobank cohort [34] may require far more PCs to account for population heterogeneity than the subset of individuals of White European ancestry as used in the current work; see Fig. 5, panels A-D.

To adopt a more principled approach to choosing confounding variables, TarGene incorporates data-driven methodologies that capture sources of population stratification confounding the relationship between  $V$  and  $Y$  (Fig. 1B). Specifically we select, in a data-adaptive manner, the optimal number of PCs for a given trait and set of variants (Methods). In Fig. 5, panels A-D, we present a trait-independent PC analysis to construct a lower bound on the number of PCs required. Refining this to a trait-dependent analysis demonstrates that each trait has its own dependency on PCs, Fig. 5, panels E and F, and we recommend performing a sensitivity analysis by including higher-order PCs as part of the Super Learner library (see below). Notably, if however supernumerary PCs do not strongly confound the variant-trait relationship, then the SL sets their coefficient to zero, thereby retaining only relevant PCs even if these are not consecutive in number. In Fig. 5C, we observe that the variant rs1421085 is randomised in the cohort and so population stratification captured by PCs is not a confounder of the phenotype-genotype relationship. This

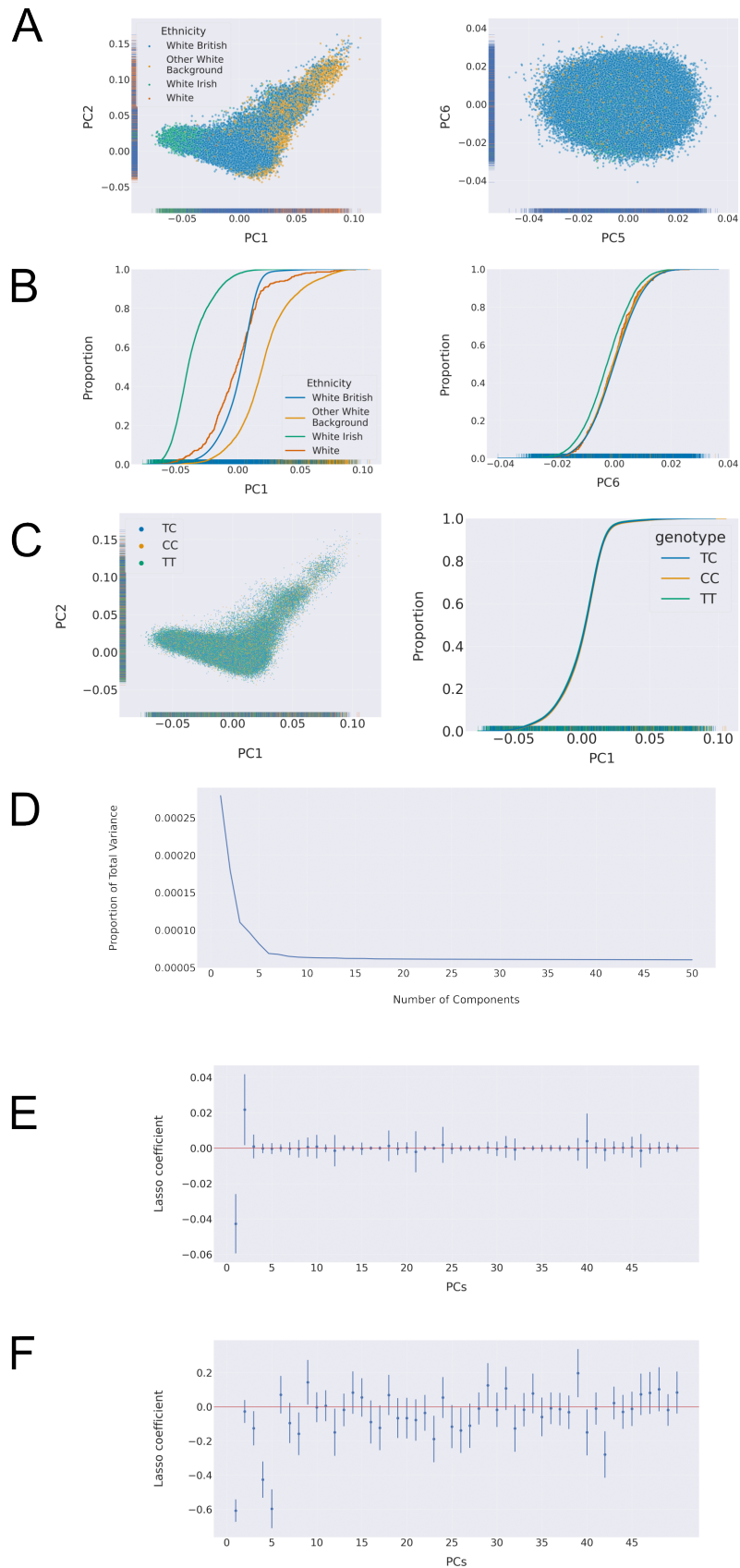


exercise is repeated for higher-order PCs and the other three variants considered in this manuscript (see Figs. 8–11). In this scenario, the only potential source of non-linearity in the functional dependence of phenotype on genotype and confounders is in the allelic copy (for a proof, see *Step 2A'* in Methods). Thus, a linear fit from  $V = 0$  to  $V = 1$  followed by a linear fit from  $V = 1$  to  $V = 2$  may yield equivalent effect sizes to TMLE up to the inference of confidence intervals. However, it is impractical to visually inspect the randomisation of each and every SNP on each and every trait in a typical statistical genetics analysis. For example, Fig. 5A and B show evidence of population stratification in the genotyped SNPs used to generate the PC plots. TarGene obviates the need for such manual inspections, irrespective of whether the variant is randomised or stratified in a cohort.

FIGURE 5. Population Structure within the UK Biobank (UKBB) cohort. **(A)** Principal component analysis labelled by ethnicity. Left: PC1 vs PC2 shows high level of population structure dependent on self-reported ethnicity. Right: PC13 vs PC14 shows a more symmetric shape indicating that there is no ethnicity structure for PCs > 13. This is more clearly visible in **(B)** via the cumulative distribution analysis of ethnicity for PC1 and PC13. Left: In PC1 self-reported ethnicity populations have different distributions indicating that ethnicity drives the first PC. Right: In PC13 this separation has disappeared. For further cumulative distributions see Supplementary Figures, Section 7. **(C)** PC Analysis of genotypes within the cohort. Left: PC1 vs PC2 labelled with the three genotypes of the *FTO* intron variant, rs1421085. There is no evidence of population structure based on this genotype. Right: A cumulative distribution shows no separation between the distributions of the genotypes indicating that the variant is randomised in the cohort. When this is the case, population stratification (and hence the PCs) is not a confounder of the phenotype-genotype relationship and the only potential source of non-linearity in the functional dependence of phenotype on genotype and confounders is in the allelic copy (*Step 2A'* in Methods). **(D)** Scree plot indicates that the proportion of variance explained by each additional PC plateaus after 8 PCs, when subset on ‘self-reported White’ UKBB population, indicating that 8 PCs is sufficient to explain the population structure of this cohort. **(E)** and **(F)** Lasso regression results on 50 PCs, in which a statistically non-zero coefficient indicates that the corresponding PC component is relevant for the trait. These results show that the choice of PCs varies according to SNP-trait pair. The error bars represent two standard deviations from the mean derived from a bootstrap ( $B = 1000$  bootstrap resamples) of the lasso coefficient. Any PCs whose error bars overlap with zero can be interpreted as adding no information to the population structure. **(E)** Lasso regression of trait “K76 other diseases of liver” (clinical\_c\_K76). This trait needs fewer than 5 PCs to explain the population structure within the cohort. **(F)** Lasso regression of trait “K20-K31 Diseases of esophagus, stomach and duodenum” (clinical\_c\_Block\_K20-K31). This trait demonstrates dependence on PCs that fluctuates more extensively than the trait in panel E. In a bespoke analysis of a trait that shows fluctuation in the Lasso regression of trait on PCs in supernumerary components, we recommend including those higher PCs in the SL in addition to the result of the trait-independent analysis as part of a sensitivity analysis.

TARGENE

13



*Step 2B: TarGene leverages a diverse combination of algorithms via Super Learning.* It is unnecessary to expend computational resources on estimating the full probability distribution  $P(y, v, w)$  in order to evaluate the target parameters of Eqs. 2 and 3. Rather, only the parts  $Q(v, w) = \mathbb{E}[Y|V = v, W = w]$  and  $g(v, w) = p(V = v|W = w)$  are required. TarGene leverages this by using a Super Learner [38], a stacking technique whereby various non-linear and/or non-parametric methods, neural networks, and tree-based algorithms in addition to parametric linear and logistic models more usually employed in GWAS, can be combined in a  $k$ -fold cross-validation scheme to find the optimal (in terms of variance explained) combination of models to fit the data. Using  $k$ -fold cross-validation, SL is mathematically guaranteed to yield the combination of models with the best predictive power as proven in [38]. This procedure is depicted in Fig. 3, second panel.

In general, the output of SL is an initial estimate  $\hat{Q}_n^0(v, w)$  of the function  $Q(v, w)$ , as well as an initial estimate of the target parameter obtained by plugging  $\hat{Q}_n^0$  into Eq. 2:

$$\Psi_1(\hat{Q}_n^0) = \frac{1}{n} \sum_{i=1}^n [\hat{Q}_n^0(1, w_i) - \hat{Q}_n^0(0, w_i)]. \quad (4)$$

The average is taken over the cohort of size  $n$ , and  $w_i$  are the covariates of participant  $i$ .

*Step 3: TarGene performs a targeted update via TMLE to remove bias.* SL is optimised to produce the best estimate,  $\hat{Q}_n^0(v, w)$ , of the function  $Q(v, w) = \mathbb{E}[Y|V = v, W = w]$ . However, it is not optimised for estimating the DNA variant's true effect size on phenotype, *i.e.*, the target parameter  $\Psi_1(P_0)$ . As a result, there may be residual bias in the initial estimate Eq. 4, *i.e.*, a discrepancy between the initial effect size estimate,  $\Psi_1(\hat{Q}_n^0)$ , and its true value,  $\Psi_1(P_0)$ . Under mild assumptions, mathematical theory [37] allows us to separate this discrepancy into three components (see Eq. 23 in Methods and Materials). The first component represents residual bias due to model-misspecification. The TMLE update is mathematically guaranteed to remove this bias and make the final estimate asymptotically normally distributed. The second component represents the variance on the estimate which is used to provide final 95% confidence intervals. The third, and final, component arises from finite sample size and is guaranteed to shrink at rate  $\sqrt{n}$  as the sample size  $n$  increases.

To perform the targeted update step of the TL framework, the TMLE step, requires an estimate of the *propensity score* or *treatment mechanism*  $g(v, w) = p(V = v|W = w)$ . In the context of GWAS, this is the probability that an individual carries variant  $v$  given that they belong to population stratum  $w$ . In particular, if  $g(v, w)$  is essentially independent of population stratum, then variant  $V$  is *not* stratified in the population. To verify this, the test for the null hypothesis  $H_0: g(v, w) = p(V = v)$  is easily incorporated into the TL framework. In practice, the propensity score is also estimated via Super Learning as discussed in the previous section.

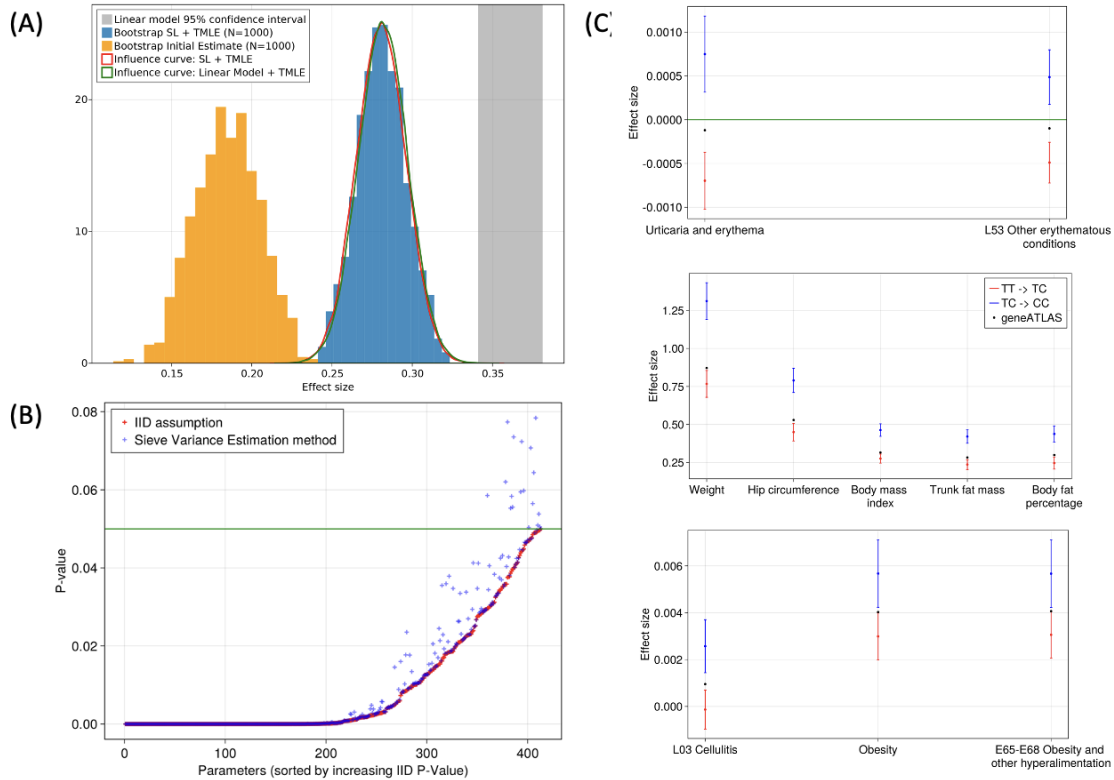


FIGURE 6. **(A) Inference results.** Comparison of various methods to estimate the effect size of rs1421085 on body mass index (BMI; X-axis; UK Biobank Data-Field 23104). The grey area indicates a 95% confidence interval on the effect size predicted by a linear-model. The orange histogram is the bootstrap distribution of the initial estimator when using a SL, without removing residual model-misspecification bias on the target quantity via TMLE. The blue histogram is also obtained from bootstrap but additionally using the TMLE. Finally, the red and green curves are the Gaussian distributions obtained via SL + TMLE and Linear Model + TMLE respectively. As can be seen, the effect sizes and corresponding p-values reported by the linear model alone, or SL alone, are overly optimistic or pessimistic, respectively. In both cases, the TMLE step brings the inference back to the same confidence region. **(B) Sieve variance correction.** P-values obtained from two variance estimation methods for rs1421085. In red, the individuals in the UK-Biobank are assumed to be i.i.d, while in blue, a sieve correction method is applied to account for the population dependence structure. Each p-value corresponds to a specific parameter of interest for which the initial i.i.d estimate was under the 0.05 threshold. **(C) Non-Linear effects.** A selection of traits for which rs1421085 TT → TC and TC → CC effect estimates are significantly different. Five continuous traits related to BMI (top), two binary traits for which effect sizes have opposite sign (middle), and three further binary traits associated with BMI (bottom). Effect sizes are reported with associated 95% confidence intervals together with estimates from GeneAtlas' LMM fits (black data points) [8]. The latter fall in-between our TT → TC and TC → CC estimates, indicative of an averaging effect.

Since the UK-Biobank contains over 450 000 samples, it is likely that the asymptotic regime is reached in this dataset. To demonstrate the behaviour of the TMLE as compared to biased estimations' methods, we performed a bootstrap analysis on the UK-Biobank (Figure 6). In this precise example, we estimated the effect of substituting TT with TC for rs1421085 on BMI. From the plot, it is notable that the effect size distributions provided by the bootstrap SL (orange) and the linear-model (grey) do not overlap. This is a real data example of the model-misspecification phenomenon illustrated in the schematic in Fig. 1 A and B. Adding an extra targeting step in addition to any of these methods, removes the residual bias and always brings estimates back to the same confidence region (see Figure legend for details).

*Step 4: TarGene accounts for population dependence structure.* Above we describe how to estimate the effect size of one DNA variant on a single phenotype measured on independent and identically distributed data. However, if the data is dependent, for example because of genetic relatedness [26, 28], care must be taken when estimating the variance on the estimates. In current practice, the genetic dependence between individuals is accounted for by incorporating the genetic relationship matrix  $\mathbf{G}$  (GRM) as random effects in a linear mixed model (LMM):

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e} \quad \text{where} \quad \mathbf{e} = \mathcal{N}(0, \sigma_g^2 \mathbf{G} + \sigma_e^2 \mathbf{I}). \quad (5)$$

Here  $\mathbf{X} \in \mathbb{R}^{N \times d}$  is the matrix of fixed effects,  $d$  the number of fixed effects, and  $\boldsymbol{\beta}$  the fixed effect sizes. There are two random components: (i)  $\sigma_g^2$  denoting the magnitude of the genetic variance, and (ii)  $\sigma_e^2$  denoting the magnitude of the residual variance. There is little justification for this model's strong restrictions. Besides model-misspecification, adopting an LMM may be invalid for three reasons: (i) The error term  $\mathbf{e}$  need not be normally distributed [16]; (ii) The variance of  $\mathbf{e}$ ,  $\text{Var}(\mathbf{e}) = \sigma_g^2 \mathbf{G} + \sigma_e^2 \mathbf{I}$ , need not decompose additively into one part due to population stratification and another due to residual environmental, technical, and other noise; and, (iii) The complexity of population stratification need not be captured by a single parameter  $\sigma_g^2$  multiplying the GRM.

In TarGene, we neither assume individuals are independent nor impose the above strong restrictive assumptions of an LMM. Instead, we adopt a network approach, drawn from [12], to incorporate the genetic dependence of individuals model-independently by taking into account ancestral diversity and familial relatedness amongst individuals as reported, *e.g.*, in Fig. 3b of [7]. The method calculates the variance of the effect size target parameter by constructing Sieve Plateau (SP) variance estimators (Methods and Materials) that incorporate genetic dependence among individuals.

In brief, the SP estimator computes a variance estimate for a range of thresholds  $\tau$ , by considering individuals to be genetically *independent* if their genetic distance exceeds  $\tau$ . The genetic distance between a pair of individuals  $(i, j)$  equals  $1 - \text{GRM}_{i,j}$ , *i.e.*, one minus

their genetic relatedness value. As the distance threshold  $\tau$  increases, fewer individuals are assumed to be genetically independent. For instance, the estimate corresponding to a distance of  $\tau = 0$  corresponds to the i.i.d. hypothesis, while a distance of  $\tau = 1$  incorporates pairs of individuals who are not genetically correlated. TarGene varies the threshold  $\tau$  from 0 to 1 and fits a curve to the corresponding variance estimates. The maximum of this curve is the most conservative estimate of the variance of the target parameter estimator and constitutes our corrected variance estimator. In Figure 6, we investigate the effect of this correction for the effect sizes obtained for rs1421085 on all parameters under investigation (see section 5). Since, the correction can only increase the variance estimate, we only correct those for which the associated p-value is under our decision threshold (0.05). The p-values resulting from both the i.i.d. (red) and the sieve variance plateau estimators (blue) are reported. For 19 parameters ( $\approx 5\%$ ), the corrected estimate changes the hypothesis test decision at the 0.05 threshold.

*Step 5: TarGene controls for multiple hypothesis testing.* When testing multiple hypotheses simultaneously in order to answer a question of interest, it is essential to state explicitly which error rate of false positives one seeks to bound, and then choose the multiple hypothesis correction procedure that affords maximal power (*i.e.*, fewest Type II errors) whilst bounding the desired Type I error rate. Since TarGene produces asymptotically normal estimators described as empirical means of efficient influence functions, the theory is sufficiently rich to be combined with any desired definition of false positives, *e.g.*, the family-wise error rate (FWER) or the false discovery rate (FDR) (see Methods). In this work, we use the marginal step-down Benjamini–Hochberg procedure of [5] to control the FDR at  $\leq 0.05$ . However, researchers can combine TarGene with any multiple testing procedures to control FWER, FDR, or other error rates, for marginal or joint multiple hypothesis testing, as described for example in [13].

**Application to the UK Biobank.** To investigate the benefits of TarGene, we performed two distinct analyses that are detailed below. The first aims at contrasting our approach with the gold standard LMM’s method on a well studied variant. The second, more exploratory, investigates whether epistatic relationships exist between multiple loci related to vitamin D receptor (VDR) biology. We present here our major findings after the targeting step, sieve variance correction and multiple hypothesis adjustment. In all analyses, we consider 776 UK Biobank traits defined by GeneATLAS [8].

*Application of TarGene to an FTO variant.* In order to demonstrate our method we performed a phenome-wide association study (PheWAS) using UK Biobank. We chose a well studied variant, rs1421085, located in the first intron of the FTO gene. The T to C nucleotide substitution has been predicted to disrupt the repression of *IRX3* and/or *IRX5*, thereby leading to a developmental shift from browning to whitening programs and loss of mitochondrial thermogenesis [11]. This variant has also been associated with several related traits such as body mass index and obesity [15]. To allow TarGene to discover



non-linear effect sizes as previously reported in [44], the analysis was not restricted to the substitution of one T allele to C. Instead, we investigated all three changes, namely  $TT \rightarrow TC$ ,  $TC \rightarrow CC$  and  $TT \rightarrow CC$ . We adjusted p-values to control for the false discovery rate using the Benjamini-Hochberg method. A summary of all estimation results is provided in Supplementary Table 1. Each change requires a dedicated targeted estimate as each corresponds to one of three different target quantities of interest. This is in contrast with a linear model assuming that these quantities are equal to  $\beta$ ,  $\beta$ , and  $2\beta$  respectively.

Of 776 traits under investigation, 20.7% are reported as significantly associated to rs1421085 by GeneATLAS. Using TarGene, we find the following for the three quantities of interest ( $FDR \leq 0.05$ , Benjamini-Hochberg):

- $TT \rightarrow TC$  *only*: 9.1% of traits are significant.
- $TT \rightarrow TC$  *or*  $TC \rightarrow CC$ : 12.5% traits are significant.
- $TT \rightarrow TC$ ,  $TC \rightarrow CC$  *or*  $TT \rightarrow CC$ : 15.4% traits are significant.

In general, the distribution of p-values obtained via TarGene is shifted towards less significant values as compared to the GeneATLAS analysis (see Supplement to Fig. 6 in Section 7). TarGene finds fewer significant results than GeneATLAS but provides mathematical guarantees of statistical coverage of ground truth on the results it finds, leading to fewer false positive whilst maximising power.

On the other hand, because we investigate each allelic change as a separate quantity, TarGene can also find associations that are undetectable by linear models such as LMMs. In Eq. 2 we define the effect size of a SNP  $V$  on a phenotype  $Y$  model-independently via the target parameter  $\Psi_1(P)$ . This is the effect on phenotype when a single allelic copy is present ( $V = 1$ ) versus when there is none ( $V = 0$ ), *e.g.*, TC vs TT. However, there is another equally valid, and potentially distinct, way of describing the effect of an additional allelic copy on  $Y$ , namely the effect on phenotype when two allelic copies are present ( $V = 2$ ) versus one ( $V = 1$ ), *e.g.*, CC vs TC. The corresponding model-independent definition is

$$\Psi_2(P) = \mathbb{E}_W [\mathbb{E}[Y|V = 2, W] - \mathbb{E}[Y|V = 1, W]]. \quad (6)$$

Conventional GWAS, such as LMMs [23], assume linearity of genetic effect on phenotype (*i.e.*, assume  $\Psi_1(P) = \Psi_2(P)$ ). More recent methods, such as KnockoffGWAS [33], do not estimate effect size altogether. However, there are no compelling biological or mathematical reasons why the two effect sizes should be equal [42].

TarGene avoids this unnecessary assumption by estimating both effect sizes separately. This does not necessarily increase the burden of multiple hypothesis testing because one can choose to query the effect size  $\Psi_1(P)$  at the TMLE step for each trait-variant pair. Alternatively, if one specifically wishes to identify non-linear effect sizes of trait-variant



pairs, or classify the type of DNA variants and/or phenotypes for which such non-linearities occur, one can directly estimate the difference  $\Psi_2(P) - \Psi_1(P)$  (Methods and Materials).

The two scenarios we are thus comparing are (i) the change TT  $\rightarrow$  TC, and (ii) the change TC  $\rightarrow$  CC. We find 54 traits for which rs1421085 displays significant non-linear effect sizes, 40 of which are highly correlated with BMI. For instance, we find that the departure from homozygous T to heterozygous TC is associated with an increase of 0.77 kg (95% CI: 0.68 – 0.85). In comparison, the departure from heterozygous TC to homozygous C is associated with a larger increase of 1.31 kg (95% CI: 1.19 – 1.43). For illustration, a subset of significant non-linear traits is presented in Figure 6; Tables 1 and 2 contain the complete list. As might be expected, estimates reported by GeneAtlas fall in-between estimates from our two scenarios, representative of an averaging effect.

Notably, for urticaria and erythema, replacing a single T with C (*i.e.*, TT  $\rightarrow$  TC) is associated with significantly ( $p = 2.9 \times 10^{-4}$ ) *lowered* risk whereas replacing the second T with C (*i.e.*, TC  $\rightarrow$  CC) is associated with significantly ( $p = 9.5 \times 10^{-3}$ ) *elevated* risk (Fig. 6, middle). Thus TarGene can capture variant-trait pairs displaying the *Heterozygote Advantage* [19]. Such patterns cannot be detected by a linear model assuming equal allelic effect sizes.

Those two results have shown that model misspecification and the curse of dimensionality (see Fig. 1) can be problematic in two ways: (i) Significant results may be reported that are false positives and may thus result in wasted follow-up studies, and (ii) significant hits may not be reported that might have otherwise informed future research (*i.e.*, they are in fact false negatives). Note that linear models are a standard part of the SL library of TarGene so if statistical inference using a linear model is correct, TarGene will choose the model in a data-driven manner. TarGene provides mathematical guarantees of coverage of ground truth thus, in combination with multiple hypothesis testing, bounding false positives whilst minimising false negatives.

*Exploration of new epistatic loci.* Epistatic interactions can be defined in a model-independent way, see Eq. 3 and adjacent discussion. Here we propose to investigate potentially interacting variants involved in VDR biology. VDR is a nuclear hormone receptor that binds to calcitriol, the active form of vitamin D, and forms a complex with the retinoid-X receptor (RXRA). This complex can then enter the nucleus and bind to specific genetic domains to regulate transcription of many genes. Because this mechanism depends on three main molecules (calcitriol, VDR and RXRA), it is a natural field of investigation for epistasis. We thus identified three genetic variants that have been associated with differential expression of each molecule in turn. From eQTLGen, A to C change of rs7971418 is associated with increased levels of VDR; G to T change of rs1045570 is associated with increased levels of RXRA; and, C to T change of rs3755967 has been associated with a

decreased level of calcitriol [20]. For each SNP, we restricted our attention to the departure from homozygosity of the major allele to heterozygosity and investigated all three pairwise interactions as well as the 3-point interaction. We provide estimation results and p-values for pairwise interactions in Supplementary Table 2, and for 3-point interactions in Supplementary Table 3. In each case, we also provide adjusted p-values to control the FDR (via Benjamini-Hochberg) across all traits being tested for each SNP combination. Although 43, 39 and 36 pairwise interactions (for rs1045570 and rs3755967, rs1045570 and rs7971418, rs3755967 and rs7971418 respectively), as well as 29 3-point interactions were found significant prior to multiple testing correction, no interaction was significant following FDR correction at the 0.05 threshold. This is not unexpected because detection of epistasis is extremely challenging [43]. However, this analysis demonstrates the opportunities that our method provides for the general investigation of higher order interactions.

### 3. DISCUSSION

We have introduced TarGene, a workflow for targeted estimation of variant effect sizes and epistatic interaction effects on polygenic traits, which dispenses with unnecessary assumptions currently widespread in the statistical genetics literature. TarGene consists of five steps: (1) Defining the quantity of interest model-independently; (2) Employing a diverse library of learning algorithms, such as data-adaptive algorithms with proven convergence properties [37], to learn the relevant portion of the true probability distribution, which results in an initial estimate of the quantity of interest; and, (3) Applying Targeted Maximum Likelihood Estimation (TMLE) to update and target the initial fit towards the quantity of interest, thereby removing any remaining bias. This is followed by (4) correction for cohort population dependence, and (5) multiple hypothesis testing. The estimators TarGene produces are context-independent, *i.e.*, are applicable beyond variant-variant interaction, to any discrete set of variables affecting an outcome of interest, *e.g.*, interactions of variant  $\times$  sex, or any other binary or discrete environmental factors [3].

TarGene offers a number of distinct advantages over current commonly employed LMM approaches, as summarised in Fig. 2. In particular, since it is firmly rooted in the mathematical estimation framework of Targeted Learning, TarGene avoids model-misspecification bias, produces asymptotically normal and efficient estimates, and is doubly-robust. Furthermore, due to the flexibility of its SL library and the TMLE step, it can be readily applied to more heterogeneous biobanks such as *All of US* [2] or the *Million Veterans Program* [17], as well as more strongly inter-related cohorts such as island communities. TarGene also allows for integration of data from multiple biobanks.

The strength of TarGene lies in bespoke analyses of effect sizes and interactions amongst targeted variants of interest, providing mathematically guaranteed coverage of the ground truth, scaling to hundreds or thousands of variants, as well as PheWAS analyses. The

run time of this workflow depends on how precise and unbiased a researcher wishes to be regarding the answer to their question of interest. For researchers interested in applying TarGene for genome-wide studies across multiple traits, we note that there is a trade-off between computational speed and guaranteeing ground truth coverage of estimates. In such a scenario we therefore recommend equipping the SL with computationally light algorithms only, such as a linear model, GLMnet or LMM, reducing the cross-validation burden significantly, before running the TMLE step. In terms of run time, with a single linear/LMM algorithm in the library, this is simply equivalent to the run-time of three linear fits for each variant-trait pair. We remark that in comparison to LMMs, TarGene does not require the memory-intensive inversion of the GRM, which for UKBB-scale population sizes may be prohibitive, depending on institutional resources.

In future work, we plan to explore non-linearities in variant allelic copies on trait, which as of yet have not been systematically explored in the literature for either homogeneous or diverse populations. We will also investigate the contribution of epistatic interactions of specific variants on various polygenic traits for a variety of biological mechanisms.

#### 4. METHODS

To enable reproducible results, all UK Biobank related analyses were made using a Nextflow pipeline that can be accessed at [TarGene pipeline](#). All runs were performed on the [Eddie](#) cluster. The configuration used for the analysis of rs1421085 and pairwise interactions is provided in Supplementary Table 2 and the configuration used for 3-points interactions in Supplementary Table 3.

**TarGene provides mathematical guarantees and realistic p-values.** Here we present the three main steps in the model-independent estimation framework of Targeted Learning (TL) [39, 40] in detail: (1) Defining the quantity of interest model-independently; (2) Super Learning; and, (3) Targeted Maximum Likelihood Estimation (TMLE). In later sections we detail how TarGene (4) incorporates population dependence in the cohort, and (5) accounts for multiple hypothesis testing.

*Step 1: TarGene defines the quantity of interest model-independently.* The effect size of a DNA variant  $V$  on trait  $Y$ , correcting for population stratification via confounders  $W$ , is defined as

$$\Psi_1(P) = \mathbb{E}_W[\mathbb{E}[Y|W, V = 1] - \mathbb{E}[Y|W, V = 0]]. \quad (7)$$

This is interpreted as the difference between the expected phenotype if the variant  $V = 1$  versus  $V = 0$ , whilst correcting for confounders. The ground truth probability distribution is called  $P_0$ . The true (but unknown) effect size is denoted by  $\psi_{1,0} = \Psi_1(P_0)$ .

The epistatic interaction between two variants in their effect on a given trait or disease  $Y$  has been defined in [4]. This model-independent definition is

$$I_{V_1, V_2}(P) = \mathbb{E}_W \left[ \mathbb{E}[Y|W, V_1 = 1, V_2 = 1] - \mathbb{E}[Y|W, V_1 = 0, V_2 = 1] \right. \\ \left. - (\mathbb{E}[Y|W, V_1 = 1, V_2 = 0] - \mathbb{E}[Y|W, V_1 = 0, V_2 = 0]) \right]. \quad (8)$$

The genomic interpretation of this definition is: “Having correctly adjusted for confounders, is the effect of variant  $V_1$  on trait modulated by the status of variant  $V_2$  and, if so, by how much and with which sign?” The ground truth interaction is denoted  $I_0 = I_{V_1, V_2}(P_0)$ . This definition of 2-point interaction, which is an extension of ATE with more than one variant, has been further generalised to higher-order interactions amongst  $n$  variants [4].

As an example, suppose the ground truth trait  $Y$  has the following expectation value:

$$\mathbb{E}[Y|W, V] = \exp(\zeta V + (\eta V + \epsilon)W), \quad (9)$$

where  $V$  represents the variant of interest, and  $W$  represents other covariates, *e.g.*, PC components.  $\zeta$  and  $\eta$  represent different parts of the parameter space. Suppose also that  $W$  is distributed according to the normal distribution  $\mathcal{N}(0, 1)$ . Since we are interested in measuring the effect size of  $V$  on trait  $Y$ , let us consider the terms appearing in the computation of the log-odds ratio or ATE effect size:

$$\mathbb{E}_w \left[ \mathbb{E}[Y|W, V = v] \right] = \int_w e^{\zeta v + (\eta v + \epsilon)w} p(w) dw = \frac{1}{\sqrt{2\pi}} \int e^{\zeta v + (\eta v + \epsilon)w - w^2/2} dw. \quad (10)$$

Completing the square and calculating the integral gives:

$$\mathbb{E}_w \left[ \mathbb{E}[Y|W, V = v] \right] = \exp \left[ \frac{1}{2}(\eta v + \epsilon)^2 + \zeta v \right]. \quad (11)$$

Therefore, the ground truth effect size is:

$$\text{ATE}_V = \mathbb{E}_w \left[ \mathbb{E}[Y|W, V = 1] \right] - \mathbb{E}_w \left[ \mathbb{E}[Y|W, V = 0] \right] = \exp \left[ \frac{1}{2}(\eta + \epsilon)^2 + \zeta \right] - \exp \left[ \frac{1}{2}\epsilon^2 \right].$$

Now, suppose that we fit the ATE using the following misspecified model for expectation value, assuming linearity:

$$\mathbb{E}[Y|W, V = v] = \beta_0 + \beta_1 V + \beta_2 W. \quad (12)$$

Then the ATE has the following expression in terms of the model parameters:

$$\text{ATE}_V = \mathbb{E}_w \left[ \mathbb{E}[Y|W, V = 1] \right] - \mathbb{E}_w \left[ \mathbb{E}[Y|W, V = 0] \right] = \beta_1. \quad (13)$$

*Simulation examples of model-misspecification.* We generated data from the ground truth, in this case Eq. 9 with  $W$  standard normally distributed as stated, at a given value of  $\epsilon = 0.3$  without loss of generality. Here  $\epsilon$  controls the levels of heteroskedasticity in the data as  $V$  varies from 0 to 1. We perform simulations in the following two scenarios:

- (1) The variant  $V$  and the source of population stratification  $W$  are independent
- (2) The variant  $V$  depends on  $W$  (*e.g.*, PC components, location, batch, ...). For simplicity, we consider the following dependence structure:

$$\begin{cases} w > 0 & \Rightarrow v = 0, \\ w < 0 & \Rightarrow v = 1. \end{cases} \quad (14)$$

When  $W \sim \mathcal{N}(0, 1)$ , this leads to  $\approx 80\%$  correlation between  $V$  and  $W$ .

We illustrate the issue of model-misspecification using simulations based on parametric models that are then fitted (i) with the correct parametric form, and (ii) with a misspecified parametric model. We perform these simulations probing different parts of the parameter space and at various sample sizes.

**Simulation 1 (top panels of Fig. 4):** We generate data from the ground truth distribution given by the exponential model explored in the previous section Eq. 9,

$$Y = \exp(\zeta V + (\eta V + \epsilon) W) \text{ , and } W \sim \mathcal{N}(0, 1) \text{ .} \quad (15)$$

under three different scenarios: (a) the variant  $V$  is generated using a binomial distribution, and completely independent of  $W$ , (b) the variant  $V$  is generated by dichotomising  $W$  using an arbitrary cut-off, (c) similar to (b) but also adding various degree of Gaussian noise when generating  $V$  from  $W$ . In scenarios (b) and (c) the variant  $V$  has various degrees of dependency on  $W$ , which is to be expected as  $W$  here represents sources of population stratification. Without loss of generality we set  $\zeta = -1$ , and vary  $\eta$  and  $\epsilon$  over a range of values to explore different part of the parameter space. More specifically, we probe the parameter space by setting  $\eta = [-2, 2]$  in steps of 0.5 and  $\epsilon = [0, 1]$  in steps of 0.125. In each case, the ground truth exponential model is fitted with the misspecified linear model:

$$Y = \alpha_0 + \alpha W + \beta V \text{ ,} \quad (16)$$

to obtain the assumed effect size  $\beta$ .

**Simulation (1A):** In the unrealistic scenario where  $V$  and  $W$  are *truly* independent, a misspecified linear model happens to coincide with the true effect size for all values  $\epsilon$  and  $\eta$  (Fig. 4, top panel, 1A). This is the case even though the goodness-of-fit measurements clearly indicate the non-linearity of the data, *e.g.*,  $R^2 \approx 0.3$  and extremely high-values of the Jarque-Bera index, indicating non-normality of the data. The previous statement holds irrespective of sample size. This behaviour has been observed in randomised control trials (RCTs) where the treatment mechanism is truly randomised with respect to known

confounders and is therefore independent of these confounders [32]. This cannot be said for GWAS where sources of population stratification, *e.g.*, genetic ancestry, are clearly visible from PCA projections in variant space.

**Simulation (1B):** The variant  $V$  is generated by dichotomising  $W$  using an arbitrary cut-off to induce dependence between  $W$  and  $V$ . Without loss of generality, the cut-off is chosen such that  $V = 0$  if  $W > -3.0$  and  $V = 1$  if  $W < -3.0$ . The conclusions below are similar irrespective of the choice of cut-off, which was tested to range from  $-3.5$  to  $3.0$  in steps of  $0.5$ , or larger values of  $\epsilon$ . The traits are again generated using Eq. 15. As an example, for sample size  $10,000$  (Fig. 4, top panel, 1B), the effect sizes estimated by the misspecified linear model are mostly incorrect. Furthermore, in about 50% of the cases, even the sign is inferred incorrectly: a positive effect size of a variant on trait is estimated to be negative, instead of positive, when fitted with a misspecified model. We also note that the dependence structure in this example is not captured by either Pearson ( $\approx 0.1$ ) or Spearman ( $\approx 0.1$ ) correlations, both indicating a weak degree of correlation.

**Simulation (1C):** This scenario is similar to (1B), following Eq. 15, but also adding various degrees of Gaussian noise when generating  $V$  from  $W$ . The noise takes on values  $[0.001, 0.01, 0.1, 0.5, 1.0, 5.0, 10, 50, 1000]$ . As observed in Fig. 4 (top panel, 1C), at fixed sample size, large levels of noise may hide the dependency between  $V$  and  $W$ , resulting in effect sizes that match the ground truth value. However, as the sample sizes increase, the level of noise has to be extreme (50 or 1000) to cover for the misspecified model.

The simulations show that (i) at any fixed level of noise, there always exists a large sample size  $N$  for which model-misspecification leads to invalid inference due to shrinking variance and, therefore, (ii) when working with large data sets, such as the UKBB, it is crucial to avoid subjective modelling choices as model-misspecification is likely to give rise to invalid inference.

**Simulation 2 (bottom panels of Fig. 4):** There is nothing special about the choice of the exponential data generating distribution in Simulation 1. To exemplify this, we perform another set of simulations with data generated from a simple polynomial model:

$$\begin{aligned} Y &= \beta V + W^2 + \mathcal{N}(1, \sigma_Y), \\ V &= Z/2 + \mathcal{N}(0, 1), \\ W &= \log(Z^2) + \mathcal{N}(0, 1), \\ Z &\sim \mathcal{N}(0.25, 1), \end{aligned} \tag{17}$$

where  $Z$  is an auxiliary variable which is used to generate dependency between  $V$  and  $W$ . The independent noise in  $Y$ , *i.e.*,  $\sigma_Y$  is varied by setting it to  $0.01, 0.1, 1, 10$  and  $100$ . Without loss of generality, the true effect size  $\beta$  is probed from  $[-2, 2]$  in steps of  $0.5$ .

Fig. 4 (bottom panel, 2A), indicates the correct estimated effect sizes when the true model is used. As expected, most results are within two standard deviations, there are only a few values ( $\approx 4\%$ ) more than two standard deviations away from the ground truth.

**Simulation (2B):** Here, the data is fitted with the following misspecified model:

$$Y = \alpha_0 + \alpha W + \beta V. \quad (18)$$

The results in Fig. 4 (bottom panel, 2B), clearly indicate that model-misspecification becomes more manifest as the sample size grows. At 10,000 samples most estimated values are incorrect. At the UKBB size, almost all estimated effect size values are incorrect.

**Simulation (2C):** Here, the data is fitted with the following misspecified model:

$$Y = \alpha_0 + \beta V. \quad (19)$$

The results are presented in Fig. 4 (bottom panel, 2C), with the same conclusions as above.

We furthermore note that replication of results, even on two different databases does not guarantee the estimated quantities and their confidence intervals span the ground truth, when a misspecified model is used. Therefore, under the assumptions that two data sets have a similar distribution, replication should be treated as a necessary but not sufficient condition for accuracy of the results: Fitting separate data samples drawn for the same distribution with the same (or similar) misspecified model twice, results in an equivalent invalid inference twice. Finally, we recall that multiple hypothesis correction methods correct for the testing of multiple hypotheses, not for false discovery in a single hypothesis, such as those predicted via misspecified parametric models.

*Step 2A: Variant preprocessing and PC analysis.* For the analysis in Fig. 5, we started with all directly genotyped variants in UKBB and applied LD block removal,  $\pm 10$  Mb around rs1421085 position, using PLINK2 [9, 30], with  $R^2 = 0.1$ . We filtered variants based on Minor Allele Frequency (MAF) threshold  $> 0.05$  and performed LD pruning using PLINK [31, 29] (1000kb window, 50 variants step size, and  $R^2 = 0.05$ ). We used the 33,483 biallelic remaining genotyped variants for 452,149 self-reported white individuals as input to FlashPCA2 to perform a partial PCA. Labelling the PCA plots by sex, age and batch, and plotting the corresponding cumulative distributions did not indicate any structure for these variables. A slight degree of separation by assessment centre was observed in PC1-2 cumulative distribution plots. Self-reported ethnicity is the dominant driver for the PCs as observed in Fig. 5.

*Step 2A': Non-confounding covariates need not be conditioned on for correct estimation.* In the estimation of the phenotype-genotype relation, covariates that do not confound this relation, *i.e.*, that do not both affect phenotype and genotype, need not be taken into account. In the causal graph of Fig. 7, the left hand side represents a confounder  $W$  of the causal effect  $V \rightarrow Y$  of genotype  $V$  on phenotype  $Y$ , whereas the right hand side  $W$  is not



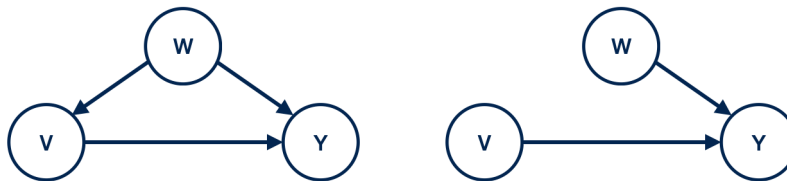


FIGURE 7. Two causal graphs from which we wish to identify the direct causal effect of  $V \rightarrow Y$ . Left:  $W$  is a confounding variable which has to be correct for, for accurate causal identification. Right:  $W$  is not a confounder.

a confounder as it merely effects phenotype, not genotype. In other words, for the graph on the right hand side, the covariate  $W$  may be ignored when estimating the causal effect of  $V$  on  $Y$  *even though it affects outcome  $Y$ ; however, the crux is that  $W$  does not affect the source  $V$* . Mathematically, this means that, for any  $v$ , we have

$$\mathbb{E}_W [\mathbb{E}_Y [Y|W, V = v]] = \mathbb{E}_Y [Y|V = v]. \quad (20)$$

This can be deduced from the fact that the directed acyclic graph on the right hand side of Fig. 7 encodes the property  $p(W|V) = p(W)$  since  $V$  and  $W$  are independent. We have:

$$\begin{aligned} \mathbb{E}_W [\mathbb{E}_Y [Y|W, V = v]] &= \int dw p(w) \int dy y p(y|w, v) \\ &= \int dw p(w) \int dy y \frac{p(y, w|v)}{p(w|v)} \\ &= \int dw p(w) \int dy y \frac{p(y, w|v)}{p(w)} \\ &= \int dy y p(y|v) = \mathbb{E}_Y [Y|V = v], \end{aligned} \quad (21)$$

where in the second equality we have used the product rule of probabilities, in the third equality we have used the independence condition of  $V$  and  $W$ , and in the fourth equality we have exchanged the integrals and used  $\int p(y, w|v)dw = p(y|v)$ .

*Step 2B: TarGene leverages a diverse combination of algorithms via Super Learning.* It is unnecessary to expend computational resources on estimating the full probability distribution  $P(y, v, w)$  in order to evaluate the target parameter Eq. 2. Indeed, only the part  $Q(v, w) = \mathbb{E}[Y|V = v, W = w]$  is required for estimating effect sizes Eq. 7 (similarly, the part  $\tilde{Q}(v_1, v_2, w) = \mathbb{E}[Y|V_1 = v_1, V_2 = v_2, W = w]$  is required for estimating the interaction Eq. 8.) Super Learner [38] applies  $k$ -fold cross-validation to a diverse library of learning algorithms to obtain an estimate of the relevant part  $Q(v, w)$  of  $P_0$ , *i.e.*, the expected phenotype given variant  $v$  and confounders  $w$ ; see Results and Fig. 3, Step 2. In  $k$ -fold cross validation, the data is split into  $k$  equally sized and disjoint folds. All algorithms are trained  $k$  times on the data, each time holding out a different ‘validation’ fold

from the training procedure. The algorithms' performance is subsequently validated on the held out fold. Finally, the best performing linear combination of algorithms is selected. The output of SL is an initial estimate  $\hat{Q}_n^0(v, w)$  of the function  $Q(v, w)$ , as well as an initial estimate of the target parameter by plugging in:

$$\Psi(\hat{Q}_n^0) = \frac{1}{n} \sum_{i=1}^n [\hat{Q}_n^0(1, w_i) - \hat{Q}_n^0(0, w_i)]. \quad (22)$$

The average is taken over the population of size  $n$ , and  $w_i$  is the covariate of participant  $i$ .

In this paper, TarGene used the following SL specifications: (i)  $k$ -fold cross-validation or stratified  $k$ -fold cross-validation based on the outcome type (continuous or binary, respectively), here  $3 \leq k \leq 20$ , selected adaptively based on the rarest class of each outcome [27], and, (ii) included the constant fit, a regularized logistic/linear regression (GLM), a gradient-boosted tree, and the Highly Adaptive Lasso (with hyper-parameters `max_degree = 1`, `smoothness_orders = 1`, `lambda = 30`) [37], as algorithms in the library. However, we note that for the optimal performance of HAL in more bespoke analyses, the parameter  $\lambda$ , tuning the total variation norm of the fit, should be left unspecified so that it is chosen by the algorithm's internal cross-validation.

*Step 3A: TarGene performs a targeted update via TMLE to remove bias.* Although SL is optimised for estimating the function  $Q(v, w)$ , it is not optimised for estimating the DNA variant's true effect size on phenotype, *i.e.*, the target parameter  $\Psi(P_0)$ . As a result, there may be residual bias in the initial estimate Eq. 22, *i.e.*, a discrepancy between the effect size estimate,  $\Psi(\hat{Q}_n^0)$ , and its true value,  $\Psi(P_0)$ . Under mild assumptions, mathematical theory (see [37]) allows us to describe and analyse this discrepancy explicitly:

$$\Psi(\hat{Q}_n^0) - \Psi(P_0) = \underbrace{-\frac{1}{n} \sum_{i=1}^n D_{\hat{Q}_n^0}^*(o_i)}_{\text{remaining bias}} + \underbrace{\frac{1}{n} \sum_{i=1}^n D_{P_0}^*(o_i)}_{\text{variance}} + \underbrace{o_P(1/\sqrt{n})}_{\text{finite sample remainder}}. \quad (23)$$

Here  $D_P^*(o_i)$  denotes the *efficient influence curve* of  $\Psi$  at  $P$  evaluated at the data point  $o_i = (y_i, v_i, w_i)$  of individual  $i$ . It quantifies the effect individual  $i$  has on the average effect size  $\Psi(P)$  across the population, *e.g.*, a large effect if  $i$  is an outlier with a rare phenotype/genotype combination (the value of  $D_P^*(o_i)$  is large), and a small effect if  $i$  is a 'typical' participant (see *Computing influence curves*). The third term  $o_P(1/\sqrt{n})$  is due to the finite sample size of the data, and shrinks at rate  $\sqrt{n}$  as sample size  $n$  increases.

Importantly, although the two averages on the right-hand side of Eq. 23 look similar, they play very different roles and have different interpretations. The first average is constructed from the initial estimate  $\hat{Q}_n^0(v, w)$  of  $Q(v, w)$ , and it quantifies the remaining bias due to the fit. We remove this bias by applying Targeted Maximum Likelihood Estimation (TMLE,

see [39]) to update the initial fit  $\hat{Q}_n^0$  repeatedly until, at the final iteration<sup>3</sup> denoted  $\hat{Q}_n^*$ , the remaining bias in Eq. 23 is approximately zero, namely  $\frac{1}{n} \sum_{i=1}^n D_{\hat{Q}_n^*}^*(o_i) = 0$ . The final TL estimate of the DNA variant's effect size on phenotype is then

$$\Psi(\hat{Q}_n^*) = \frac{1}{n} \sum_{i=1}^n [\hat{Q}_n^*(1, w_i) - \hat{Q}_n^*(0, w_i)]. \quad (24)$$

In contrast, the second average in Eq. 23, containing terms  $D_{P_0}^*$ , depends on the ground truth probability distribution,  $P_0$ , and cannot be changed by an improved analysis (apart from increasing the data size). It is responsible for the variance on the estimate  $\Psi(\hat{Q}_n^0)$ . Indeed, after removing the bias via TMLE, Eq. 23 implies

$$\sigma_n^2 \equiv \text{Var} \left[ \sqrt{n}(\Psi(\hat{Q}_n^*) - \Psi(P_0)) \right] = \frac{1}{n} \sum_{i=1}^n [D_{P_0}^*(o_i)]^2. \quad (25)$$

The right-hand side of this equation can be directly estimated from data and quantifies the variance on the final effect size estimate  $\Psi(\hat{Q}_n^*)$ . Furthermore, it allows for the direct construction of an approximate Wald-type 95% confidence interval:

$$\left[ \Psi(\hat{Q}_n^*) - 1.96 \frac{\sigma_n}{\sqrt{n}}, \Psi(\hat{Q}_n^*) + 1.96 \frac{\sigma_n}{\sqrt{n}} \right]. \quad (26)$$

This approach is valid since the TMLE step and Eq. 23 (with vanishing first term on the right-hand side) imply that  $\Psi(\hat{Q}_n^*)$  is normally distributed as  $n$  becomes large, provided the third term in Eq. 23 is indeed of order  $o_p(1/\sqrt{n})$ . The latter condition holds as long as the product of  $Q_0$  and  $g_0$  is estimated at a rate of  $n^{-1/2}$  because the effect size and interaction target quantities have the *double robust property*; see Step 3D for details. Finally, we can directly construct a p-value on the estimate of effect size from this confidence interval.

*Step 3B: Computing influence curves.* An asymptotically linear estimator behaves, for large sample sizes, like an average of independent identically distributed random variables. These random variables, called *influence curves*, are functions of the data. The Central Limit Theorem can be used to analyse the variance of asymptotically linear estimators. More precisely, an estimator  $\Psi(P_n^*)$  of a quantity of interest is asymptotically linear if

$$\Psi(P_n^*) - \Psi(P_0) = \frac{1}{n} \sum_{i=1}^n D_{P_0}(o_i) + o_p(1/\sqrt{n}). \quad (27)$$

Here  $\Psi(P_0) = \psi_0$  is the true value of the parameter,  $o_p(1/\sqrt{n})$  is a finite-sample term that shrinks to zero at rate  $\sqrt{n}$  as sample size  $n$  increases, and  $D_{P_0}(o_i)$  is the influence curve of  $\Psi$  at the probability distribution  $P_0$  evaluated at the  $i$ th data point  $o_i$ . The expectation

<sup>3</sup>For the effect size and interaction target quantities, it is a mathematical fact that convergence takes place in a single step so that no iteration is necessary.

of  $D_{P_0}$  with respect to the ground truth  $P_0$  is zero, *i.e.*,  $\mathbb{E}_{P_0}[D_{P_0}(O)] = 0$ .

The influence curve of the average treatment effect  $\Psi_1(P)$  of Eq. 7 is well known; its derivation can be found in [39, App. A.3]. It is the following function of data  $O = (Y, V, W)$ :

$$D_P^*(O) = \frac{2V - 1}{p(V|W)} [Y - Q(V, W)] + Q(1, W) - Q(0, W) - \Psi(P), \quad (28)$$

where  $Q(V, W) = \mathbb{E}_P[Y|V, W]$ . Here  $H(V, W) = (2V - 1)/p(V|W)$  is the ‘clever covariate’. Since the target parameter for interaction  $I_{V_1, V_2}$  of Eq. 8 has recently been introduced in [4], we derive its influence curve here for the first time. We record our result as a

**Proposition 4.1.** *Let  $\mathcal{M}$  be a non-parametric statistical model containing probability distributions of  $O = (Y, V_1, V_2, W)$  where  $Y$  is any outcome,  $V_1$  and  $V_2$  are binary or categorical variables, and  $W$  is any covariate. Let  $(a, b)$  be categories of  $(V_1, V_2)$ , and consider the target parameter*

$$\Psi_{a,b}(P) = \mathbb{E}_W [\mathbb{E}_P[Y|W, V_1 = a, V_2 = b]]. \quad (29)$$

The efficient influence curve of  $\Psi_{a,b}$  is given by

$$D_{a,b}^*(P)(O) = \frac{\mathbb{1}\{V_1 = a, V_2 = b\}}{p(V_1, V_2|W)} [Y - Q(V_1, V_2, W)] + Q(a, b, W) - \Psi_{a,b}(P). \quad (30)$$

Here we have defined the function  $Q(a, b, W) = \mathbb{E}_P[Y|W, V_1 = a, V_2 = b]$ .

Since the influence curve of a sum of target parameters is equal to the sum of their influence curves, we immediately deduce the influence of interaction  $I_{V_1, V_2}$ :

**Corollary 4.2.** *The influence curve of 2-point interaction  $I_{V_1, V_2}$  of Eq. 8 equals*

$$D^*(P) = [D_{1,1}^*(P) - D_{0,1}^*(P)] - [D_{1,0}^*(P) + D_{0,0}^*(P)]. \quad (31)$$

We prove Proposition 4.1 in Appendix A.

*Step 3C: TMLE updates the initial fit to obtain mathematical guarantees.* The targeted update step in TMLE removes any residual bias in Eq. 23, thus making the effect size estimator asymptotically normal and allowing for the construction of Wald-type 95% confidence intervals. This step proceeds by fluctuating the initial SL fit  $\hat{Q}_n^0(v, w)$  of the conditional phenotype  $Q(v, w) = \mathbb{E}[Y|V = v, W = w]$  in the direction of the efficient influence curve  $D_{P_0}^*$  of the quantity of interest, *e.g.*, the effect size Eq. 28 or genotype-genotype interaction Eq. 30. This fluctuation consists of a simple one-dimensional maximum likelihood estimation (MLE) of a real-valued parameter  $\epsilon$  in an auxiliary statistical model. Specifically, for binary phenotypes, the fluctuation is a logistic regression,

$$\text{logit } \hat{Q}_{n,\epsilon}^1(V, W) = \text{logit } \hat{Q}_n^0(V, W) + \epsilon H(V, W), \quad (32)$$

so that the property  $\hat{Q}_{n,\epsilon}^1(V, W) \in (0, 1)$  is preserved. For continuous phenotypes, the fluctuation is a linear regression (hence with a normally distributed noise term),

$$\hat{Q}_{n,\epsilon}^1(V, W) = \hat{Q}_n^0(V, W) + \epsilon H(V, W). \quad (33)$$

In both cases, the initial SL fit  $\hat{Q}_n^0(v, w)$  is taken as off-set and the coefficient  $\epsilon$  in front of  $H(V, W)$  is estimated; the fitted value of  $\epsilon$  in the first update is denoted by  $\epsilon_n^1$ . The clever covariate depends on the quantity of interest, satisfying  $H(V, W) = (2V - 1)/p(V|W)$  for effect sizes and  $H(V_1, V_2, W) = (2V_1 - 1)(2V_2 - 1)/p(V_1, V_2|W)$  for variant-variant interactions. Performing weighted, instead of standard, logistic (or linear) regression with weight the reciprocal of the propensity score and clever covariate  $H'(V, W) = 2V - 1$  for effect sizes and  $H'(V_1, V_2, W) = (2V_1 - 1)(2V_2 - 1)$  for interactions, results in more robust estimates when near positivity violations are present. We fit the treatment mechanism  $g(v, w) = p(V = v|W = w)$  (for effect sizes) and  $g(v_1, v_2, w) = p(V_1 = v_1, V_2 = v_2|W = w)$  (for interactions) using SL, and denote the corresponding fit by  $\hat{g}(v, w)$ . Performing MLE means solving the score equation,

$$\sum_{i=1}^n \frac{d}{d\epsilon} \log p_\epsilon(y_i|v_i, w_i) \Big|_{\epsilon=\epsilon_n^1} = 0, \quad (34)$$

where  $p_\epsilon(y|v, w)$  is the probability density of phenotype  $y$  given genotype and covariates  $v, w$ . By construction of both fluctuations, Eq. 34 evaluated at  $\epsilon = 0$  equals the empirical mean of the efficient influence curve. Thus, we iterate the TMLE step, each time taking the updated fit  $\hat{Q}_n^k(v, w)$  as off-set, until the fitted parameter  $\epsilon_n^k \approx 0$ .<sup>4</sup> Then Eq. 34 reads

$$\sum_{i=1}^n D_{\hat{Q}_n^k}(o_j) = \sum_{i=1}^n \frac{d}{d\epsilon} \log p_\epsilon(y_i|v_i, w_i) \Big|_{\epsilon=0} = 0, \quad (35)$$

and we have successfully updated  $\hat{Q}_n^0$  to  $\hat{Q}_n^* := \hat{Q}_n^k$  so as to eliminate the residual bias term in Eq. 23, thus making the estimate  $\Psi(\hat{Q}_n^*)$  asymptotically normal and unbiased.

*Step 3D: TL estimates are double-robust.* In order to obtain the final estimate,  $\Psi(\hat{Q}_n^*)$ , of a target parameter such as effect size or variant-variant interactions, the TL workflow requires the estimation of two quantities: the true conditional phenotype,  $Q_0$ , and the true treatment mechanism,  $g_0$ . This is an advantageous feature since TL estimates are double-robust in the sense that  $\Psi(\hat{Q}_n^*)$  is a consistent estimate of the ground truth  $\Psi(P_0)$  provided either (or both)  $Q_0$  and  $g_0$  are estimated consistently at a rate of convergence faster than  $n^{-1/4}$ . Put differently, if only one of  $Q_0$  and  $g_0$  is estimated *incorrectly*, *e.g.*, by a misspecified model, then the TL estimate is nevertheless consistent. If both  $Q_0$  and  $g_0$  are estimated consistently at a product of rates faster of at least  $n^{-1/2}$ , then  $\Psi(\hat{Q}_n^*)$  is

<sup>4</sup>Since the clever covariate is independent of  $Q$ , it is a mathematical fact that this algorithm converges in a single step, that is,  $\epsilon_n^1 \approx 0$  and  $\hat{Q}_n^* = \hat{Q}_n^1$ .

an efficient estimator of  $\Psi(P_0)$ , *i.e.*, the estimator with the smallest variance in its class.

The double-robustness property of a TL estimator is established separately for each target parameter. Here, we first recall the double-robustness of effect size to then establish that the interaction target parameter of Eq. 8, introduced in [4], is double-robust as well. The general approach is to recall the first-order approximation of the target parameter at a probability distribution  $P$  in terms of the influence curve as in Eq. 27,

$$\Psi(P) - \Psi(P_0) = -\mathbb{E}_{P_0}[D_P(O)] + \text{Rem}(P, P_0). \quad (36)$$

in terms of its influence curve, and analyse the second-order remainder term,

$$\text{Rem}(P, P_0) = \Psi(P) - \Psi(P_0) + \mathbb{E}_{P_0}[D_P(O)]. \quad (37)$$

For the parameter  $\Psi_1$  of Eq. 2, measuring the effect size of a DNA variant  $V$  on a phenotype  $Y$  correcting for confounders  $W$ , we have  $Q(v, w) = \mathbb{E}[Y|V = v, W = w]$  and  $g(v, w) = p(V = v|W = w)$ . It is well known that the remainder term of  $\Psi_1$  satisfies

$$\text{Rem}(P, P_0) \leq \|Q - Q_0\|_{P_0} \cdot \|(\bar{g} - \bar{g}_0)/g\|_{P_0}, \quad (38)$$

where  $\|f\|_{P_0}^2 = \mathbb{E}_{P_0}(f^2)$  for a function  $f$  of the data  $O = (Y, V, W)$ , we write  $\bar{g}(W) = g(1, W)$  and  $\bar{g}_0(W) = g_0(1, W)$ . This inequality is a special case of the Cauchy–Schwartz inequality. We show in Appendix A that the same inequality holds for the second-order exact remainder of the interaction target parameter of Eq. 8.

*In practice.* The above description of the estimation framework of Targeted Learning is predicated on the estimation of the effect size of a single SNP on a single phenotype measured on independent and identically distributed data. However, if the data are *dependent*, care must be taken in estimating the variance on the estimates; in particular, Eq. 25 needs to be generalised, see Eq. 48. Furthermore, if multiple effect sizes are estimated, multiple hypothesis correction must be incorporated in order to bound type I errors, such as FDR control, and obtain joint p-values. Both of these further steps are required when dealing with large-scale population genetics data, such as the UKBB.

**TarGene identifies non-linear effects of allelic copies on phenotype.** TarGene estimates both effect sizes  $\Psi_1(P)$  and  $\Psi_2(P)$  in Eq. 2 and Eq. 6 separately and can, thus, be leveraged to (i) determine significant non-linear effects of allelic copies on phenotypes as well as (ii) classify the type of SNPs and/or phenotypes for which such non-linearities occur. In practice, TarGene does this by combining the asymptotic description of both TL estimates as averages of independent random variables with the Central Limit Theorem (CLT). More precisely, let  $\Psi_i(P_{i,n}^*)$  be the  $i$ th effect size with efficient influence curve  $D_{i,P_0}$  for  $i = 1, 2$ . By TL theory,  $\Psi_i(P_{i,n}^*)$  is the empirical average over influence curves and thus,

by the CLT, asymptotically normally distributed:

$$\sqrt{n}[\Psi_i(P_{i,n}^*) - \Psi_i(P_0)] = \frac{1}{\sqrt{n}} \sum_{j=1}^n D_{i,P_0}(o_j) + o_P(1) \sim \mathcal{N}(0, \mathbb{E}[D_{i,P}^2]). \quad (39)$$

Here, we are interested in the difference target parameter  $\Psi_\Delta(P) = \Psi_2(P) - \Psi_1(P)$ . Taking the difference of the quantities in Eq. 39 for  $i = 1, 2$  yields a description of the difference target parameter minus the difference ground truth as an average over influence curves:

$$\begin{aligned} \sqrt{n}[\Psi_\Delta(P_n^*) - \Psi_\Delta(P_0)] &= \sqrt{n}[\Psi_2(P_{2,n}^*) - \Psi_1(P_{1,n}^*) - [\Psi_2(P_0) - \Psi_1(P_0)]] \\ &= \sqrt{n}[\Psi_2(P_{2,n}^*) - \Psi_2(P_0)] - \sqrt{n}[\Psi_1(P_{1,n}^*) - \Psi_1(P_0)] \\ &= \frac{1}{\sqrt{n}} \sum_{j=1}^n [D_{2,P_0}(o_j) - D_{1,P_0}(o_j)] + o_P(1). \end{aligned} \quad (40)$$

The asymptotic normal distribution of  $\Psi_\Delta(P_n^*)$  and its variance again follow from the CLT, the application of which is a special case of the functional delta method:

$$\sqrt{n}[\Psi_\Delta(P_n^*) - \Psi_\Delta(P_0)] \sim \mathcal{N}\left(0, \mathbb{E}[D_{1,P_0}^2] - 2\mathbb{E}[D_{1,P_0}D_{2,P_0}] + \mathbb{E}[D_{2,P_0}^2]\right) \quad (41)$$

In practice, the variance on the right-hand side is estimated by replacing the influence curves by their sample averages, followed by an update to accurately account for the population dependence structure in the cohort (see next Section and Eq. 49).

**Step 4: TarGene accounts for population dependence structure.** Biobank cohorts consist of participants that are, to some extent, related due to ancestry or kinship. TarGene accounts for this population dependence structure by appropriately adjusting variance estimates of effect sizes and interactions via Sieve Plateau (SP) variance estimators [12] which, in turn, are based on the genetic distance between participants as encoded in the genetic relationship matrix.

*Genetic Relationship Matrix.* Many statistical analyses rely on the assumption that the data are *independently* sampled<sup>5</sup> from the population. However, this assumption no longer holds for the participants in the UK Biobank since many of whom are, to some extent, genetically related. Such genetic similarity can occur on a sub-population level due to ancestry (*e.g.*, being white Irish), or on an individual level due to kinship (*e.g.*, parents, children, cousins). Moreover, genetically similar individuals may share diet and environment inducing further dependence [3].

The genetic similarity of two individuals  $i$  and  $j$  is quantified by the sample correlation coefficient  $G_{ij}$  between their (centred and scaled) SNPs. Together, these coefficients form

---

<sup>5</sup>It is also often required that the data be *identically distributed* but this assumption is typically not essential and can be circumvented.



the *Genetic Relationship Matrix* (GRM), denoted  $G$ , of size  $N \times N$  where  $N$  is the number of individuals in the population. More precisely, given a set of  $R$  SNPs, we have

$$G_{ij} = \frac{1}{R-1} \sum_{k=1}^R \frac{(s_{ik} - 2p_k)(s_{jk} - 2p_k)}{2p_k(1-p_k)}. \quad (42)$$

Here  $s_{ik} \in \{0, 1, 2\}$  denotes the number of copies of the reference allele for individual  $i$  at SNP  $k$ , and  $p_k \in (0, 1)$  denotes the frequency of the reference allele at SNP  $k$  over the population of  $N$  individuals. In particular, the population average of  $s_{ik}$  equals twice the reference allele frequency at SNP  $k$ , *i.e.*,  $2p_k$  (one for each strand copy), so

$$\frac{1}{N} \sum_{i=1}^N s_{ik} = 2p_k. \quad (43)$$

Thus  $\tilde{s}_{ik} = s_{ik} - 2p_k$  is the zero-centred count of the number of copies of the reference allele of individual  $i$  at SNP  $k$ . Considered as a random variable,  $\tilde{s}_{ik}$  takes on three values. Assuming reference alleles are sampled binomially with mean frequency  $p_k$ , the standard deviation of  $\tilde{s}_{ik}$  equals  $\sqrt{2p_k(1-p_k)}$ . This explains the additional factor in Eq. 42 that scales the variables  $\tilde{s}_{ik}$  and  $\tilde{s}_{jk}$  so as to have unit variance. Finally, note that the GRM depends on the set of  $R$  selected SNPs. These SNPs should be chosen amongst genotyped (not imputed) SNPs that, in addition, are not in linkage disequilibrium with one another.

*Sieve Plateau Variance Estimators.* In TarGene, we neither assume individuals are independent nor do we impose the strong restrictive assumptions of an LMM. Instead, we incorporate the genetic dependence of individuals model-independently in our Targeted Learning framework described in the section *TarGene provides mathematical guarantees and realistic p-values* in Methods. This approach, drawn from [12], is based on mathematical theory and addresses all the above issues. It generalises Eq. 25 for the variance of the effect size target parameter,  $\Phi(\hat{Q}_n^*)$ , by constructing *Sieve Plateau (SP) variance estimators* that incorporate genetic dependence of individuals. These estimators result in valid confidence intervals and, ultimately, realistic and valid p-values having correctly accounted for population stratification.

We now illustrate how data dependence impacts the variance estimate of Eq. 25. Since individuals  $i$  and  $j$  are in general dependent, their data  $O_i = (Y_i, V_i, W_i)$  and  $O_j$  are also in general dependent. As a result, the same holds for their corresponding influence curves  $D_{P_0}^*(O_i)$  and  $D_{P_0}^*(O_j)$ . The problem arises since for two random variables  $X_1$  and  $X_2$  the variance of their sum is *not* in general equal to the sum of their variances. The difference is exactly twice the covariance<sup>6</sup> of  $X_1$  and  $X_2$ , namely

$$\text{Var}(X_1 + X_2) = \text{Var}(X_1) + \text{Var}(X_2) + 2 \text{Cov}(X_1, X_2). \quad (44)$$

<sup>6</sup>The covariance of two random variables vanishes when they are independent.

The impact of this difference may be large depending on the size of the covariance  $\text{Cov}(X_1, X_2)$ . Thus, rather than Eq. 25, the true variance on  $\Psi(\hat{Q}_n^*)$  is given by

$$\hat{\sigma}_n^2 = \text{Var} \left[ \sqrt{n}(\Psi(\hat{Q}_n^*) - \Psi(P_0)) \right] = \frac{1}{n} \text{Var} \left[ \sum_{i=1}^n D_{P_0}^*(O_i) \right]. \quad (45)$$

Now, the distinction between Eq. 25 and Eq. 45 is only relevant for significantly genetically similar individuals. SP variance estimators define a cut-off  $\tau$  for the allowed genetic distance between individuals, and set the covariance to zero between individuals that are sufficiently genetically dissimilar. We obtain a variance estimate,  $\hat{\sigma}_n^2(\tau)$ , for each value of  $\tau$ . The true variance of the estimate is obtained where the function  $\tau \mapsto \hat{\sigma}_n^2(\tau)$  plateaus.

To construct SP variance estimators, we proceed as follows:

- (1) Using the GRM, we define a *genetic distance* between individuals  $i$  and  $j$  as

$$d(i, j) = 1 - G_{ij}, \quad (46)$$

where  $G_{ij}$  is the sample correlation coefficient of Eq. 42 quantifying the genetic dependence between individuals  $i$  and  $j$ . Since correlation is bounded,  $|G_{ij}| \leq 1$ , the genetic distance is non-negative and never larger than two, *i.e.*,  $0 \leq d(i, j) \leq 2$ . Biologically, if two individuals  $i$  and  $j$  have identical SNPs they are fully correlated,  $G_{ij} = 1$ , and thus have zero genetic distance,  $d(i, j) = 0$ , as expected.

- (2) Given a value for the cut-off  $\tau \in [0, 1]$ , we define a SP variance estimators as

$$\hat{\sigma}_n^2(\tau) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \mathbb{1}\{d(i, j) \leq \tau\} \cdot D_{\hat{Q}_n^*}^*(o_i) D_{\hat{Q}_n^*}^*(o_j). \quad (47)$$

Here, the term  $\mathbb{1}\{d(i, j) \leq \tau\}$  equals 1 if the genetic distance between individuals  $i$  and  $j$  is at most  $\tau$ , *i.e.*,  $d(i, j) \leq \tau$ , and it equals 0 otherwise.

The biological interpretation of these estimators is as follows. The correlation between the influence curves of individuals  $i$  and  $j$ , estimated by the term  $D_{\hat{Q}_n^*}^*(o_i) D_{\hat{Q}_n^*}^*(o_j)$ , is taken into account only if the genetic distance between individuals  $i$  and  $j$  is at most  $\tau$ . Thus, the SP variance estimator  $\hat{\sigma}_n^2(0)$ , *i.e.*, when  $\tau = 0$ , assumes all individuals are *independent*. By increasing  $\tau$ , we first take the covariance between strongly genetically dependent individuals into account for low  $\tau$ , and then incorporate the covariance of more weakly dependent individuals as  $\tau$  increases up to  $\tau = 1$ .

- (3) We construct the variance estimator  $\hat{\sigma}_n^2(\tau)$  for a number of values of the cut-off  $\tau$ , *e.g.*,  $\tau = 0, \frac{1}{k}, \frac{2}{k}, \dots, \frac{k-1}{k}, 1$ . Then we fit the function  $\tau \mapsto \hat{\sigma}_n^2(\tau)$  and select the value of  $\tau$  where the function plateaus, call it  $\tau_0$ .
- (4) The correct variance estimate is then  $\hat{\sigma}_n^2(\tau_0)$ .

Under mild assumptions [12, Theorem 1], the distribution of the effect size estimate  $\Psi(\hat{Q}_n^*)$  of SNP  $V$  on phenotype  $Y$  is asymptotically normal, and the SP variance estimator allows

for the construction of an approximate 95% Wald-type confidence interval for it in the usual way, namely

$$\left[ \Psi(\hat{Q}_n^*) - 1.96\sqrt{\frac{\hat{\sigma}_n^2(\tau_0)}{n}}, \Psi(\hat{Q}_n^*) + 1.96\sqrt{\frac{\hat{\sigma}_n^2(\tau_0)}{n}} \right]. \quad (48)$$

From here, we obtain realistic p-values correctly accounting for population dependence. Similarly, we need to take into account population dependence in order to obtain a realistic estimate of the variance on  $\Psi_\Delta(P_n^*) = \Psi_2(P_{2,n}^*) - \Psi_1(P_{1,n}^*)$  in Eq. 41. If this difference of effect sizes of a DNA variant  $V$  on a trait  $Y$  is significant, the effect of an additional allelic copy is non-linear. We construct an SP estimator for the variance on this difference by following steps (1)–(4) above, with the exception of appropriately generalising Eq. 47 in step (2) as follows. Given a value for the cut-off  $\tau \in [0, 1]$ , the estimator is

$$\begin{aligned} \hat{\delta}_n^2(\tau) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \mathbb{1}\{d(i, j) \leq \tau\} & \left[ \hat{D}_{1,n}^*(o_i) \hat{D}_{1,n}^*(o_j) - \hat{D}_{1,n}^*(o_i) \hat{D}_{2,n}^*(o_j) \right. \\ & \left. - \hat{D}_{2,n}^*(o_i) \hat{D}_{1,n}^*(o_j) + \hat{D}_{2,n}^*(o_i) \hat{D}_{2,n}^*(o_j) \right]. \end{aligned} \quad (49)$$

Here we have used the short-hand  $\hat{D}_{i,n}^* = D_{i, \hat{Q}_n^*}^*$  for the influence curve of effect size  $i$  evaluated at the final TMLE estimate  $\hat{Q}_n^*$  of the conditional expectation  $\mathbb{E}[Y|V, W]$ .

**Step 5: TarGene controls for multiple hypothesis testing.** In answering a typical question in population genetics, multiple hypotheses are tested simultaneously. Depending on the nature of the question, a specific multiple hypothesis correction is necessary to bound the error rate of interest. Error rates commonly employed are (i) the probability of at least one false discovery, or *family-wise error rate*,  $\text{FWER} = P(V_n > 0)$ , and (ii) the expected ratio of false discoveries to true discoveries, or *false discovery rate*,  $\text{FDR} = \mathbb{E}[V_n/R_n]$ . Here  $n$  denotes the sample size,  $R_n$  denotes the (known) number of discoveries, and  $V_n$  denotes the (unknown) false discoveries or Type I errors; see Fig. 3E. Once an error rate and bound has been chosen, *e.g.*, one seeks to bound the FDR at  $\leq 0.05$ , then a multiple hypothesis correction procedure is to be chosen. The better procedure is the one that minimises false negatives (Type II errors, denoted by  $T_n$  in Fig. 3E), *i.e.*, maximises power, at the given Type I error control.

The literature on multiple testing procedures is extensive, see for example [13]. Procedures differ mainly in that (i) they depend on the marginal distribution of the test statistics only (*marginal* procedure) or on their joint distribution (*joint* procedure), (ii) the rejection criteria of the next test is independent of the outcome of the previous tests (*single-step* procedure, such as Bonferroni correction) or the rejection criteria depend on the outcome of previous tests (*step-down* procedure, *e.g.*, Benjamini–Hochberg). Further considerations can be taken into account, all of which are aimed at maximising power at a given error

rate whilst respecting the dependence structure of the tests' null hypotheses. For example, a computationally more intensive joint procedure is unwarranted when test statistics are largely independent.

Since TarGene produces asymptotically normal estimators which are empirical means of their efficient influence curve, a vector of these estimates similarly equals an empirical mean of the vector of efficient influence curves. By the multi-variate Central Limit Theorem, their joint distribution is then again multi-variate normal. Thus, when simultaneously testing for significance of multiple (i) effect sizes, (ii) non-linearity of effect size, and/or (iii) epistatic or gene-environment interactions, the asymptotic joint distribution of the corresponding null distribution is known. As a consequence, researchers can take advantage of both marginal and joint procedures to maximise power whilst bounding their desired Type I error rate.

In this work, we use the marginal step-down Benjamini–Hochberg procedure of [5] to control the FDR at  $\leq 0.05$ .

## 5. ACKNOWLEDGEMENTS

This research has been conducted using the UK Biobank Resource under Application Number 53116. MvdL is supported by NIH grant R01AI074345. CPP is funded by the MRC (MC\_UU\_00007/15). AK was supported by the XDF Programme from the University of Edinburgh and Medical Research Council (MC\_UU\_00009/2).

## 6. COMPETING INTERESTS

No competing interests declared.

7. SUPPLEMENTARY FIGURES

**Supplement to Figure 5: Population stratification PCA analysis.** The figures 8 - 11 are empirical cumulative frequency (ECDF) diagrams.

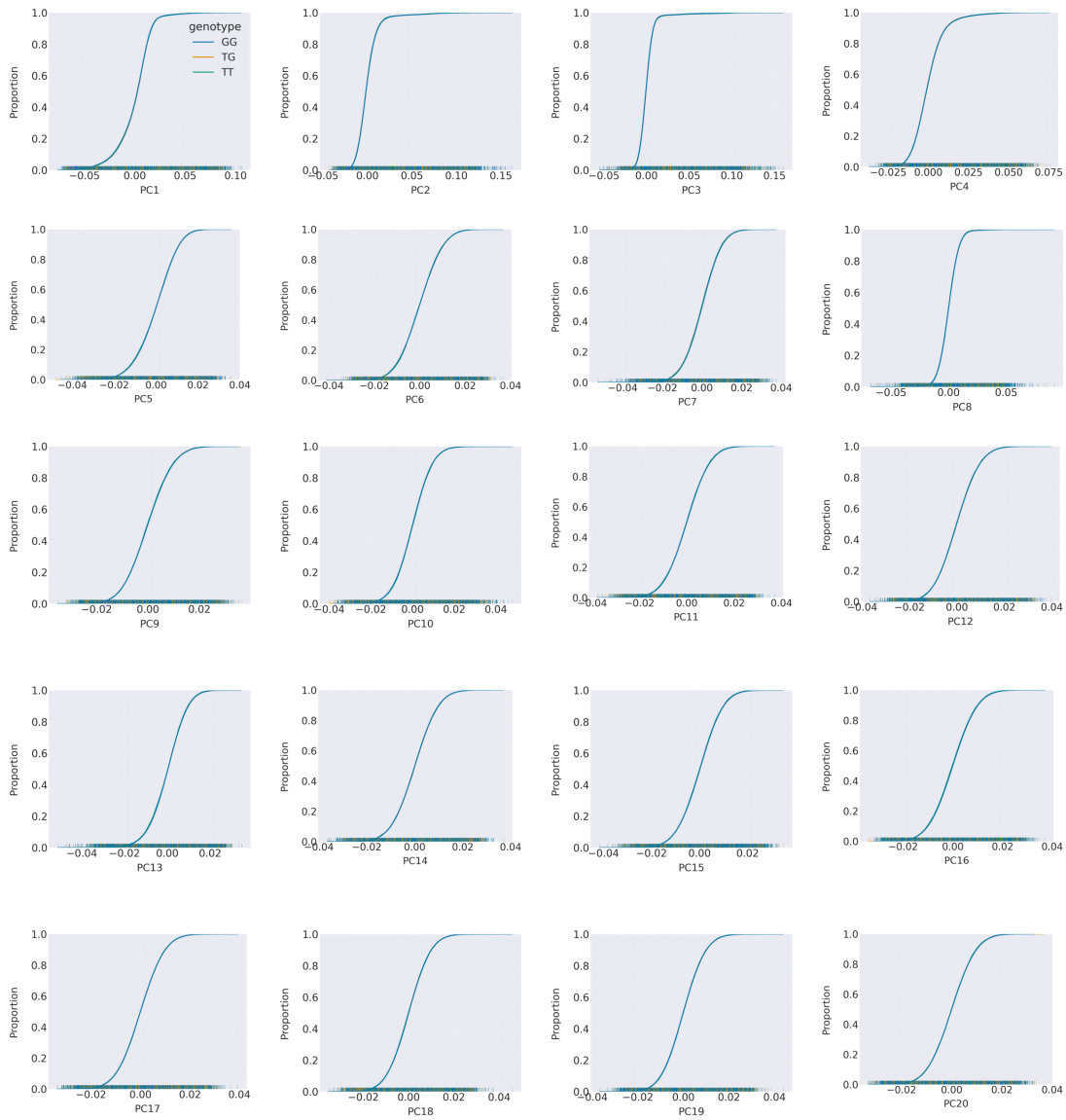


FIGURE 8. Cumulative frequency diagrams of variant rs1045570 for PCs 1-20.

TARGENE

38

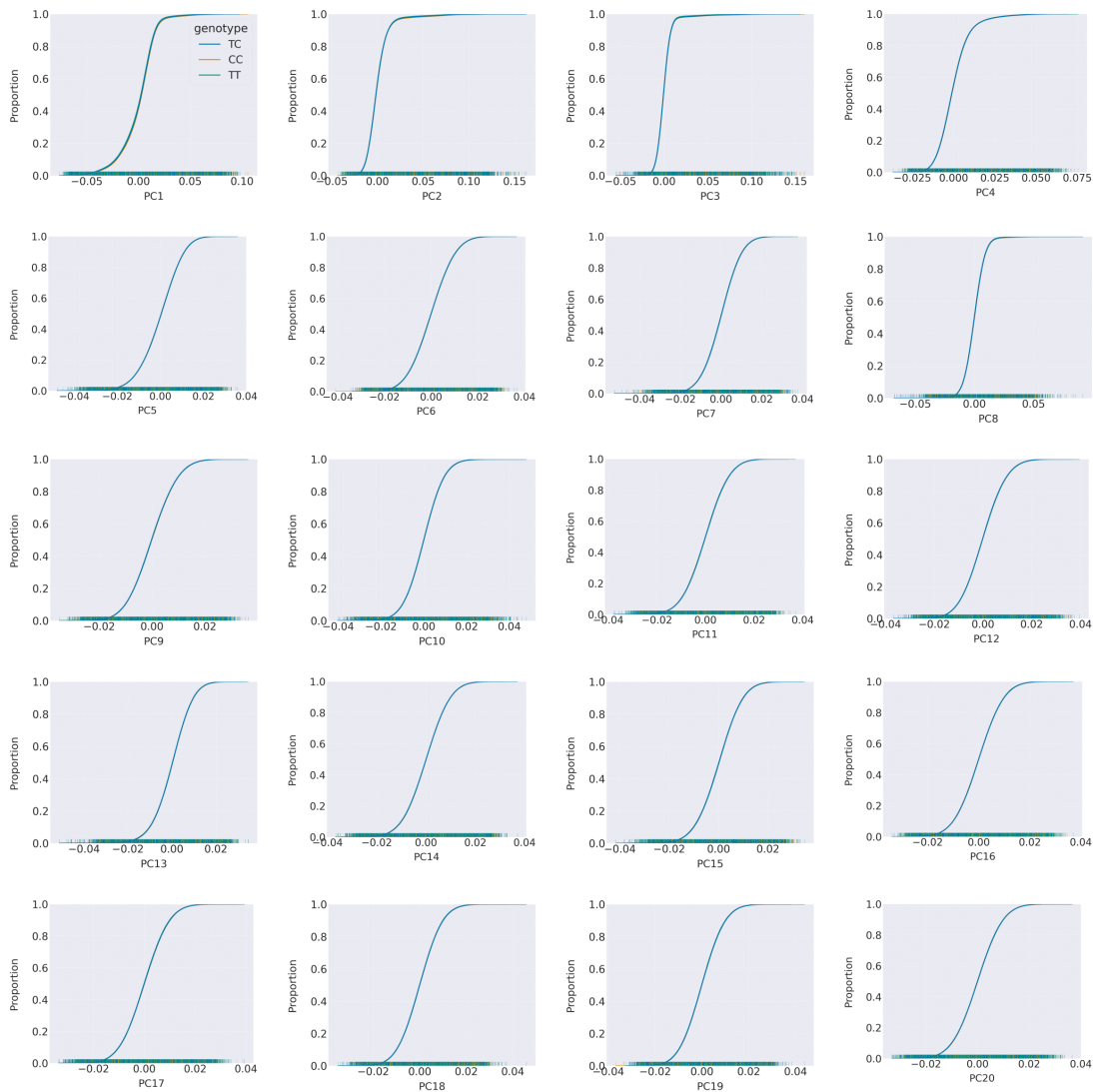


FIGURE 9. Cumulative frequency diagrams of variant rs1421085 for PCs 1-20.

The aim of these figures is to establish which of the covariates listed by UK BioBank are confounders and which PC number these should take into account. The figures can also be used to establish evidence of population structure based on variants of interest. Large differences between the separate distributions lines indicate a differences between different populations and therefore evidence that population structure is driven by this covariate. The list of variants investigated in this paper is: rs1045570, rs1421085, rs3755967, and

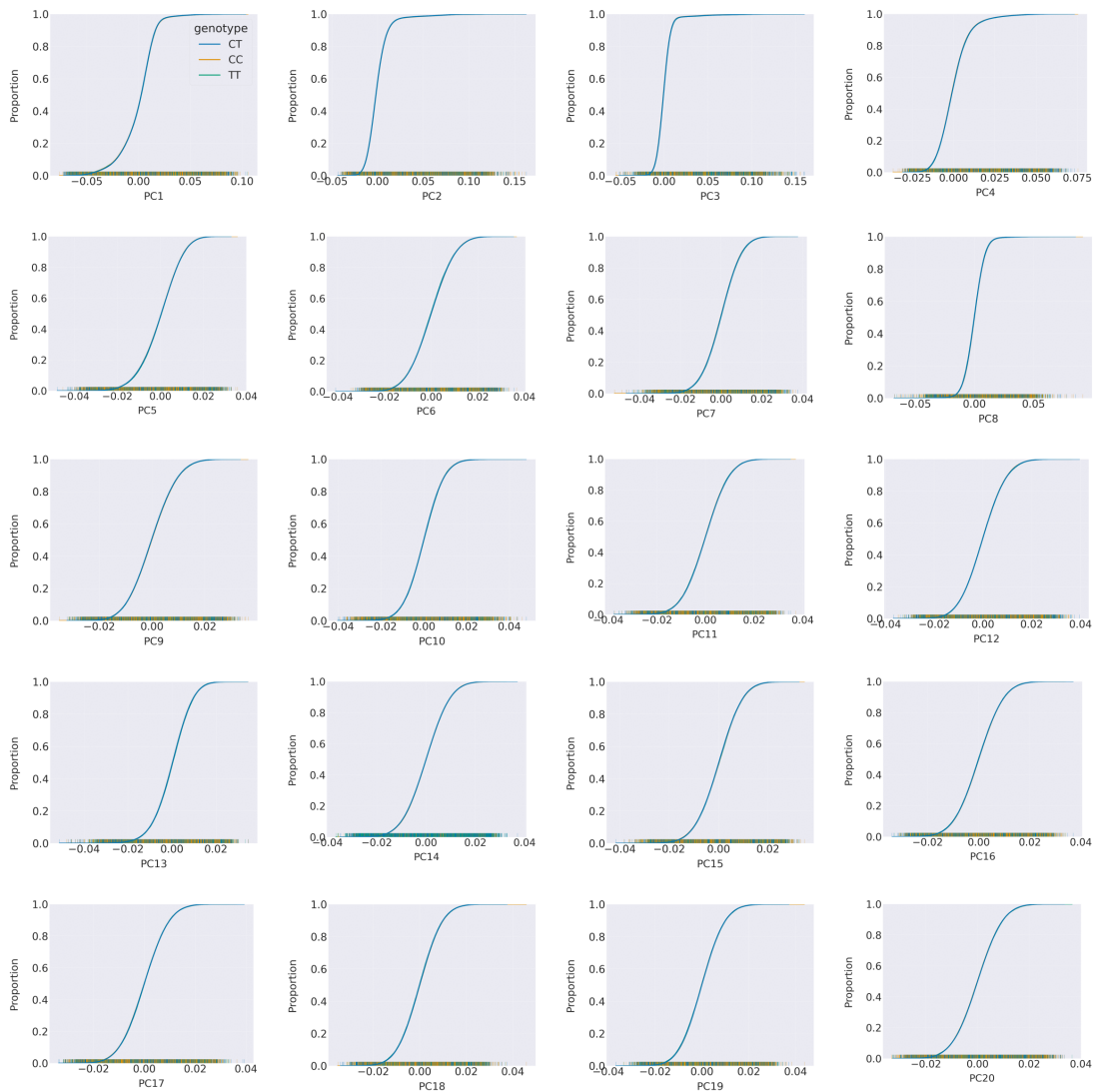


FIGURE 10. Cumulative frequency diagrams of variant rs3755967 for PCs 1-20

rs7971418. Figures 8 - 11 clearly show there is no evidence that population structure is driven by differences in alleles at these variants and therefore, the variants should not have a confounding effect in the analysis. Figure 12 shows the ECDF for self-reported ethnicity. This figure is evidence of population structure in the UK BioBank cohort among participants who self-reported as White ethnicities (including White British, White Irish, White and Any Other White Background). Here we can see that PCs 1-6 have evidence of



TARGENE

40

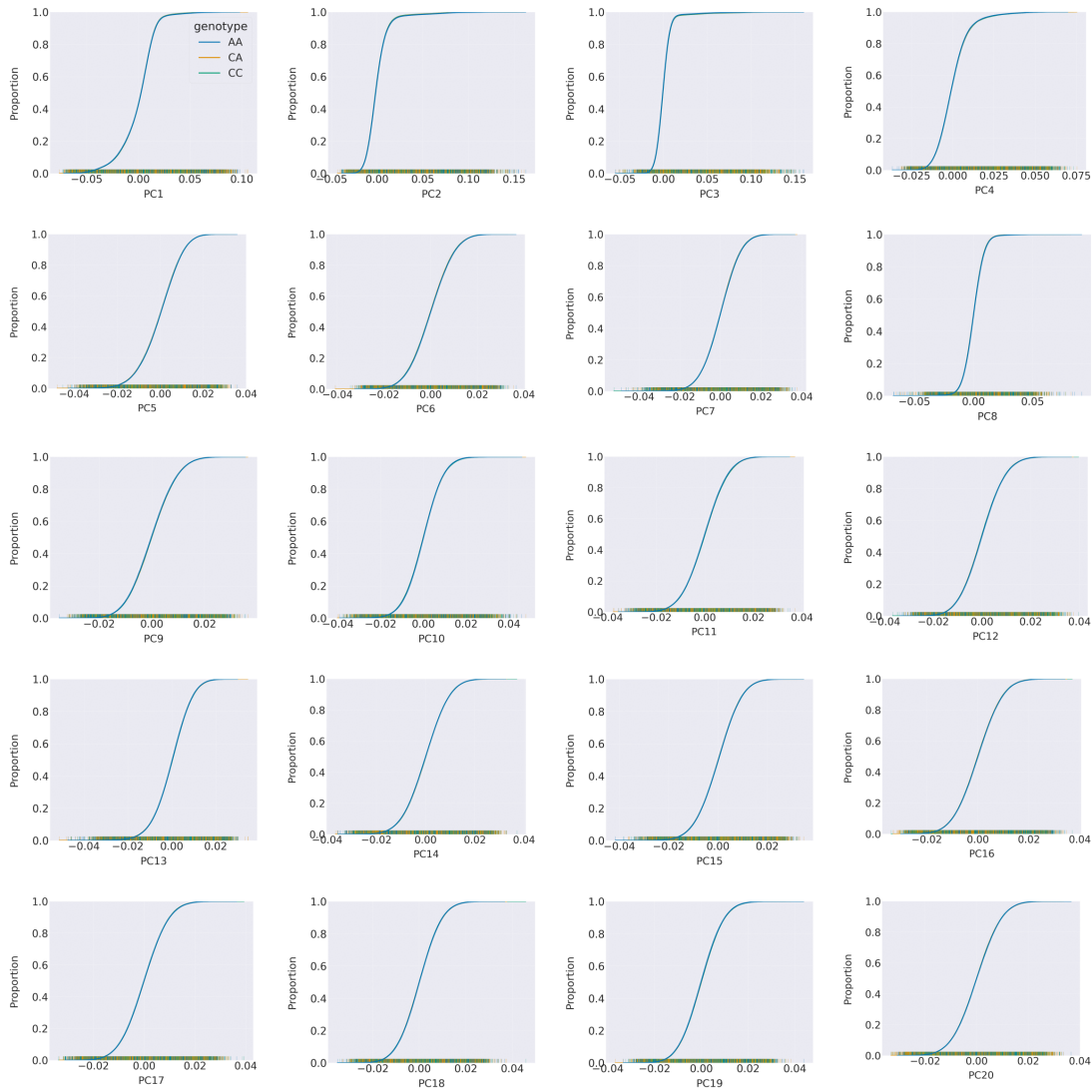


FIGURE 11. Cumulative frequency diagrams of variant rs7971418 for PCs 1-20

population structure driven by ethnicity and that this should be included as a confounder in the Super Learner.

**Supplement to Figure 6: Targeted update of our estimate.** We provide additional information regarding the effect of the targeting step on effect sizes and associated p-values. Surprisingly, in most cases, the initial estimate is shifted upward by the TMLE update (see

TARGENE

41

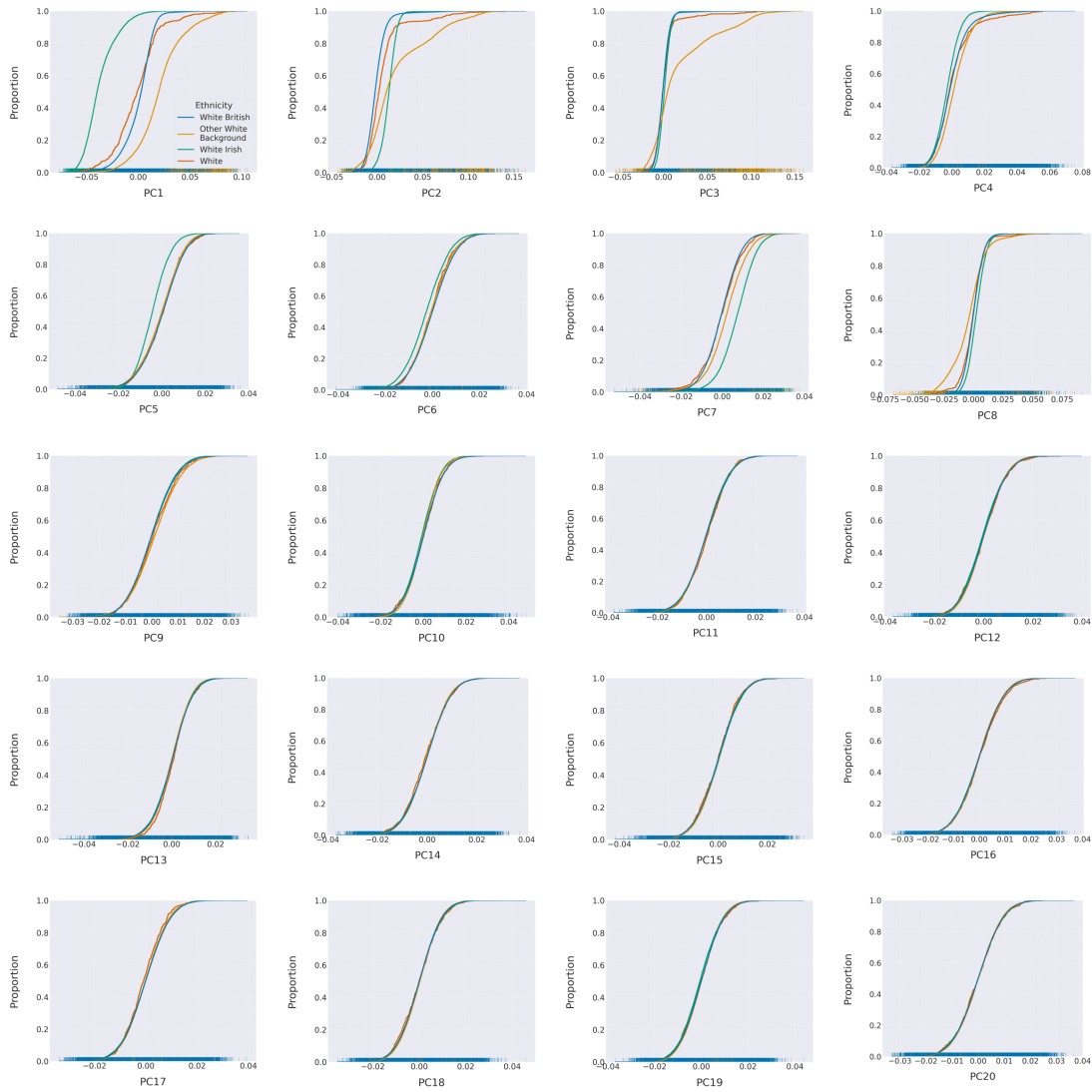


FIGURE 12. Cumulative frequency diagrams of self-reported ethnicity for PCs 1-20

figure 13). However, because the variance associated with those estimates is large, they will for the most part not be reported as significant. This behavior is confirmed by the associated p-value comparison, where the distribution of TMLE p-values is slightly shifted towards upper values as compared to GeneATLAS p-values.

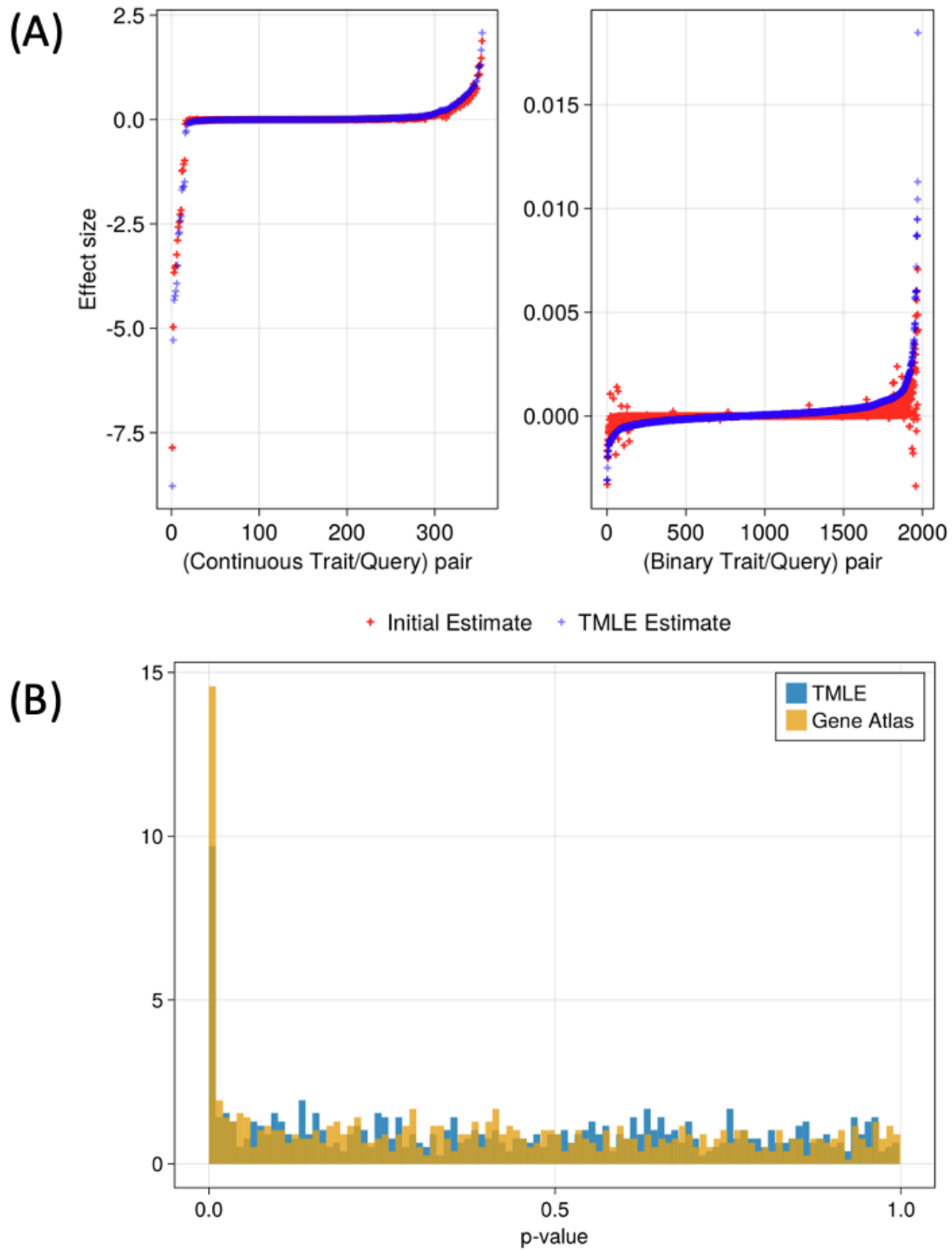


FIGURE 13. (A) Illustration of the difference between the initial estimate, reported by Super Learning, and the TMLE after the targeting step. (B) Comparative distribution of the p-values reported by TMLE vs GeneATLAS.

We present in figure 14 additional information regarding population dependence structure.

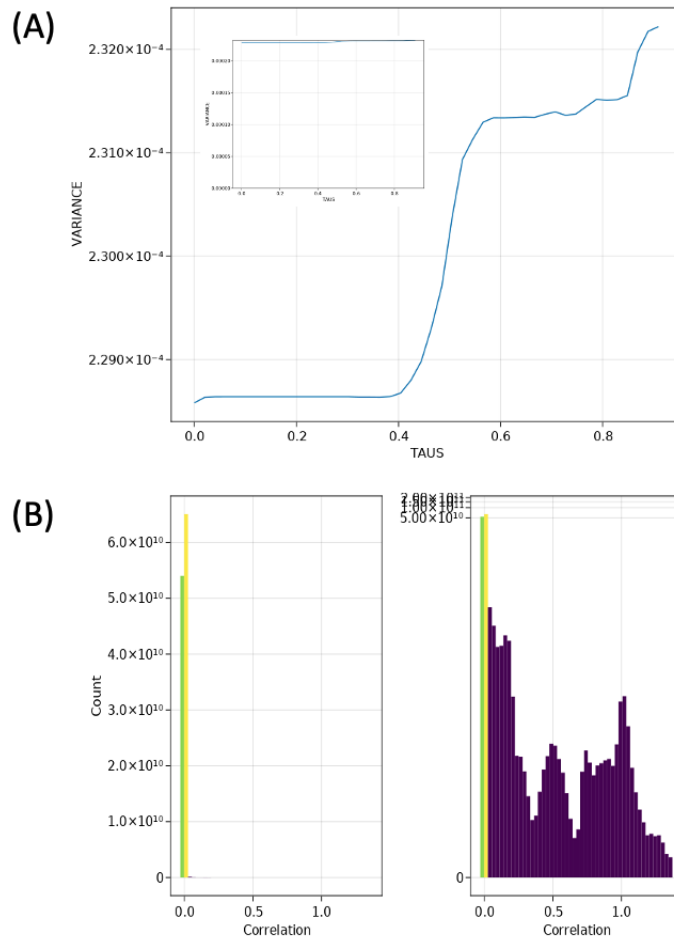


FIGURE 14. (A) A sample curve obtained via the SP estimation method (on two scales). The curve is increasing as we incorporate more dependent individuals in the estimation. (B) Histogram of the genetic relationship matrix using both a regular and log10 scale for the full UK Biobank population (488 376 individuals). As can be seen from the histogram the distribution is highly concentrated around 0. This is indicating that individuals in the UK Biobank, are to a large extent, genetically independent from one another. This potentially explains why the sieve variance correction has so little effect on the variance correction.

8. SUPPLEMENTARY TABLES

**Non-linear allelic effects.** Tables 1 and 2 present the full list of all 54 traits for which we found a non-linear allelic effect for rs1421085. That is, the effect of departing from the TT to TC genotype is significantly different from the effect of departing from the TC to CC genotype.

Description	P-value	Adjusted P-value
E66 Obesity	0.00833	0.00958
G40-G47 Episodic and paroxysmal disorders	0.0404	0.0412
I72 Other aneurysm	0.00277	0.00348
Non-oily fish intake	0.0103	0.0115
Comparative body size at age 10	2.5e-14	1.35e-12
Whole body fat mass	1.09e-9	6.55e-9
Trunk predicted mass	3.28e-5	5.54e-5
Trunk fat-free mass	4.35e-5	6.91e-5
Trunk fat mass	7.07e-9	2.39e-8
Trunk fat percentage	1.8e-6	3.6e-6
Body fat percentage	8.45e-7	1.82e-6
Arm predicted mass (left)	1.7e-7	4.18e-7
Arm fat-free mass (left)	2.42e-7	5.67e-7
Arm fat mass (left)	1.15e-10	2.15e-9
Arm fat percentage (left)	2.75e-8	7.97e-8
Arm predicted mass (right)	2.23e-6	4.29e-6
Arm fat-free mass (right)	1.25e-6	2.6e-6
Arm fat mass (right)	1.19e-10	2.15e-9
Arm fat percentage (right)	2.64e-8	7.97e-8
Leg predicted mass (left)	5.49e-9	1.98e-8
Leg fat-free mass (left)	5.36e-9	1.98e-8
Leg fat mass (left)	4.05e-9	1.82e-8
Leg fat percentage (left)	3.58e-5	5.87e-5
Impedance of arm (right)	0.00157	0.00212
Leg predicted mass (right)	6.1e-10	4.12e-9
Leg fat-free mass (right)	5.87e-10	4.12e-9
Impedance of whole body	1.19e-5	2.08e-5

TABLE 1. Table of detected non-linear allelic effects.

<b>Description</b>	<b>P-value</b>	<b>Adjusted P-value</b>
Whole body fat-free mass	3.22e-7	7.25e-7
Leg fat mass (right)	5.49e-9	1.98e-8
Impedance of arm (left)	0.00182	0.00234
Leg fat percentage (right)	0.00035	0.00054
Impedance of leg (left)	4.83e-8	1.3e-7
Impedance of leg (right)	2.59e-10	2.79e-9
Basal metabolic rate	2.8e-8	7.97e-8
Body mass index (BMI)	1.86e-10	2.51e-9
Whole body water mass	7.42e-8	1.91e-7
Weight	4.42e-10	3.97e-9
Waist circumference / Hip circumference	0.000367	0.00055
Hip circumference	3.18e-9	1.56e-8
Waist circumference	1.23e-9	6.65e-9
L03 Cellulitis	0.000915	0.00127
eye/eyelid problem	0.00181	0.00234
peripheral vascular disease	0.000447	0.000652
Fresh fruit intake	0.000504	0.000717
Number of treatments/medications taken	0.0394	0.0409
L53 Other erythematous conditions	9.81e-6	1.77e-5
O04 Medical abortion	0.00426	0.00523
L85 Other epidermal thickening	0.00562	0.0066
E65-E68 Obesity and other hyperalimentation	0.0105	0.0115
H36 Retinal disorders in diseases classified elsewhere	0.00473	0.00567
L50-L54 Urticaria and erythema	3.1e-6	5.77e-6
N00-N08 Glomerular diseases	0.0197	0.0208
C56 Malignant neoplasm of ovary	0.0181	0.0196

TABLE 2. Table of detected non-linear allelic effects (Continuation).

APPENDIX A. PROOF OF PROPOSITION

In this section, we prove Proposition 4.1. Below, we recall the content of this statement for the reader's convenience. As a direct corollary, this results in the computation of the influence curve of the interaction target parameter of Eq. 8, required for TMLE.

**Proposition A.1.** *Let  $\mathcal{M}$  be a non-parametric statistical model containing probability distributions of  $O = (Y, V_1, V_2, W)$  where  $Y$  is any outcome,  $V_1$  and  $V_2$  are binary or categorical variables, and  $W$  is any covariate. Let  $(a, b)$  be categories of  $(V_1, V_2)$ , and consider the target parameter*

$$\Psi_{a,b}(P) = \mathbb{E}_W[\mathbb{E}_P[Y|W, V_1 = a, V_2 = b]]. \quad (50)$$

The efficient influence curve of  $\Psi_{a,b}$  is given by

$$D_{a,b}^*(P)(O) = \frac{\mathbb{1}\{V_1 = a, V_2 = b\}}{p(V_1, V_2|W)} [Y - Q(V_1, V_2, W)] + Q(a, b, W) - \Psi_{a,b}(P). \quad (51)$$

Here we have defined the function  $Q(a, b, W) = \mathbb{E}_P[Y|W, V_1 = a, V_2 = b]$ .

Since the influence curve of a sum of target parameters is equal to the sum of their influence curves, we immediately deduce the influence of interaction  $I_{V_1, V_2}$ :

**Corollary A.2.** *The influence curve of 2-point interaction  $I_{V_1, V_2}$  of Eq. 8 equals*

$$D^*(P) = [D_{1,1}^*(P) - D_{0,1}^*(P)] - [D_{1,0}^*(P) + D_{0,0}^*(P)]. \quad (52)$$

A proof of a similar result, from which this Proposition can be derived via the delta method, can be found in, e.g., [39, App. A3]. For the sake of completeness, we include a proof here.

*Proof.* The probability density function of  $P$  factors as

$$p(o) = p_Y(y|v_1, v_2, w)p_V(v_1, v_2|w)p_W(w). \quad (53)$$

Let  $P_\epsilon$  be any path in  $\mathcal{M}$  through  $P$  at  $\epsilon = 0$  with score  $S(o) = (d/d\epsilon)|_{\epsilon=0} \log p_\epsilon(o)$ , where  $p_\epsilon$  is the probability density function of  $P_\epsilon$  with respect to  $\lambda$ . By [36, Thm. 3.2], a gradient  $D_{a,b}(P)$  of  $\Psi_{a,b}(P)$  can be obtained by expressing the path-wise derivative of  $\Psi_{a,b}(P_\epsilon)$  as

$$\left. \frac{d}{d\epsilon} \right|_{\epsilon=0} \Psi_{a,b}(P_\epsilon) = \langle S, D_{a,b}(P) \rangle, \quad (54)$$

where the inner product  $\langle -, - \rangle$  on the right-hand side is taken in the Hilbert space  $L_0^2(P)$  of mean-zero functions that are square-integrable with respect to  $P$ .



We now compute the path-wise derivative of the target parameter. First, we find

$$\left. \frac{d}{d\epsilon} \right|_{\epsilon=0} \Psi_{a,b}(P_\epsilon) = \left. \frac{d}{d\epsilon} \right|_{\epsilon=0} \mathbb{E}_{W,\epsilon} [\mathbb{E}_{P_\epsilon}[Y|V_1 = a, V_2 = b, W]] \quad (55)$$

$$= \int y \left. \frac{d}{d\epsilon} \right|_{\epsilon=0} \left( p_{Y,\epsilon}(y|v_1 = a, v_2 = b, w) d\lambda(y) p_{W,\epsilon}(w) \right) d\lambda(w), \quad (56)$$

where derivative and integral can be exchanged by Lebesgue's dominated convergence theorem. Next, since  $V_1, V_2$  can be considered as binary random variables, we have

$$p_{Y,\epsilon}(y|a, b, w) = \int p_{Y,\epsilon}(y|v_1, v_2, w) \frac{\mathbb{1}\{v_1 = a, v_2 = b\}}{p(v_1, v_2|W)} p(v_1, v_2|w) d\lambda(v_1, v_2). \quad (57)$$

Here  $\mathbb{1}\{v_1 = a, v_2 = b\}$  equals 1 when both  $v_1 = a$  and  $v_2 = b$  and vanishes otherwise. By another application of the dominated convergence theorem, we obtain

$$= \int y \left. \frac{d}{d\epsilon} \right|_{\epsilon=0} p_{Y,\epsilon}(y|v_1, v_2, w) \cdot d\lambda(y) \frac{\mathbb{1}\{v_1 = a, v_2 = b\} p(v_1, v_2|w)}{p(v_1, v_2|w)} d\lambda(v_1, v_2) p_W(w) d\lambda(w) \quad (58)$$

$$+ \int y p_Y(y|v_1 = 1, v_2 = 0, w) d\lambda(y) \left. \frac{d}{d\epsilon} \right|_{\epsilon=0} p_{W,\epsilon}(w) d\lambda(w) \quad (59)$$

We consider these two integrals separately, utilising the Hoeffding decomposition [41, § 11.4]. For the second integral, by the ordering  $O = (Y, V_1, V_2, W)$ , this results in

$$\left. \frac{d}{d\epsilon} \right|_{\epsilon=0} p_{W,\epsilon}(w) = \left( \mathbb{E}[S(O)|w] - \mathbb{E}[S(O)] \right) p_W(w). \quad (60)$$

We now compute the second integral  $I_2$  to be

$$I_2 = \int y p_Y(y|v_1 = a, v_2 = b, w) d\lambda(y) \left( \mathbb{E}[S(O)|w] - \mathbb{E}[S(O)] \right) p_W(w) d\lambda(w) \quad (61)$$

$$= \int Q(a, b, w) \left( \mathbb{E}[S(O)|w] - \mathbb{E}[S(O)] \right) p_W(w) d\lambda(w) \quad (62)$$

$$= \int Q(a, b, w) S(o) p_{Y,V}(y, v_1, v_2|w) d\lambda(y, v_1, v_2) p_W(w) d\lambda(w) \quad (63)$$

$$- \int S(o) p(o) d\lambda(o) \int Q(a, b, w) p_W(w) d\lambda(w) \quad (64)$$

$$= \int S(o) [Q(a, b, w) - \Psi_{a,b}(P)] p(o) d\lambda(o) \quad (65)$$

$$= \langle S, Q(a, b, W) - \Psi_{a,b}(P) \rangle. \quad (66)$$

For the first integral, again applying the Hoeffding decomposition [41, § 11.4] yields

$$\begin{aligned} \frac{d}{d\epsilon} \Big|_{\epsilon=0} p_{Y,\epsilon}(y|v_1, v_2, w) &= \left[ \mathbb{E}[S(O)|O = o] - \mathbb{E}[S(O)|v_1, v_2, w] \right] p_Y(y|v_1, v_2, w) \\ &= \left[ S(o) - \mathbb{E}[S(O)|v_1, v_2, w] \right] p_Y(y|v_1, v_2, w). \end{aligned}$$

Next, we compute the first integral  $I_1$ , which immediately splits into two pieces:

$$\begin{aligned} I_1 &= \int S(o) y \frac{\mathbb{1}\{v_1 = a, v_2 = b\}}{p(v_1, v_2|w)} p_Y(y|v_1, v_2, w) p(v_1, v_2|w) p_W(w) d\lambda(y) d\lambda(v_1, v_2) d\lambda(w) \\ &\quad - \int y \mathbb{E}[S(O)|v_1, v_2, w] p_Y(y|v_1, v_2, w) d\lambda(y) \cdot \frac{\mathbb{1}\{v_1 = a, v_2 = b\} p(v_1, v_2|w)}{p(v_1, v_2|w)} d\lambda(v_1, v_2) p_W(w) d\lambda(w) \\ &= \int S(o) y \frac{\mathbb{1}\{v_1 = a, v_2 = b\}}{p(v_1, v_2|w)} p(o) d\lambda(o) \\ &\quad - \int Q(v_1, v_2, w) \mathbb{E}[S(O)|v_1, v_2, w] \cdot \frac{\mathbb{1}\{v_1 = a, v_2 = b\} p(v_1, v_2|w)}{p(v_1, v_2|w)} d\lambda(v_1, v_2) p_W(w) d\lambda(w). \end{aligned}$$

The first piece is of the required form. In the second piece  $I_1^{(2)}$  of integral  $I_1$  we have integrated out the  $y$ -dependence, which results in the factor  $Q(v_1, v_2, w) \equiv \mathbb{E}_P[Y|v_1, v_2, w]$ . Next, we write  $\mathbb{E}[S(O)|v_1, v_2, w]$  as an integral to have  $p(o)$  under the integral:

$$\begin{aligned} I_1^{(2)} &= \int Q(v_1, v_2, w) S(o) p_Y(y|v_1, v_2, w) d\lambda(y) \cdot \frac{\mathbb{1}\{v_1 = a, v_2 = b\} p(v_1, v_2|w)}{p(v_1, v_2|w)} d\lambda(v_1, v_2) p_W(w) d\lambda(w) \\ &= \int S(o) \cdot Q(v_1, v_2, w) \frac{\mathbb{1}\{v_1 = a, v_2 = b\}}{p(v_1, v_2|w)} \cdot p(o) d\lambda(o). \end{aligned}$$

We infer that

$$I_1 = \int S(o) \left[ \frac{\mathbb{1}\{v_1 = a, v_2 = b\}}{p(v_1, v_2|w)} (y - Q(v_1, v_2, w)) \right] p(o) d\lambda(o), \quad (67)$$

and conclude that the influence curve of  $\Psi_{a,b}(P)$  is given by

$$D_{a,b}^*(P) = \frac{\mathbb{1}\{V_1 = a, V_2 = b\}}{p(V_1, V_2|W)} \left[ Y - Q(V_1, V_2, W) \right] + Q(a, b, W) - \Psi_{a,b}(P) \quad (68)$$

as claimed.  $\square$

Given the influence curve of  $\mathbb{E}_W[\mathbb{E}[Y|V, W]]$ , we can show the double-robustness of any target parameter that is a linear combination of such terms. This holds because both the exact remainder,

$$\text{Rem}(P, P_0) = \Psi(P) - \Psi(P_0) + \mathbb{E}_{P_0}[D_P(O)], \quad (69)$$

and the influence curve of a linear combination of target parameters equals the linear combination of exact remainders and influence curves respectively.

**Proposition A.3.** *Let  $\mathcal{M}$  be a non-parametric statistical model containing probability distributions of  $O = (Y, V, W)$  where  $Y$  is any outcome (e.g., trait),  $V$  a binary or categorical variable (e.g., a DNA variant), and  $W$  is any covariate. Consider the target parameter*

$$\Psi_v(P) = \mathbb{E}_W[\mathbb{E}_P[Y|W, V = v]]. \quad (70)$$

The exact remainder of  $\Psi_v$  is given by

$$\text{Rem}(P, P_0) = \mathbb{E}_{P_0} \left[ \frac{g_0(v, W) - g(v, W)}{g(v, W)} (Q_0(v, W) - Q(v, W)) \right] \quad (71)$$

We explain in what sense the TMLE of ATE or interaction is double robust. Given a function  $f$  of the data  $O = (Y, T, W)$ , we write  $\|f\|_{P_0}^2 = \mathbb{E}_{P_0}(f^2)$  for its  $L^2$ -norm. This norm is induced by the inner product  $\langle f, g \rangle_{P_0} = \mathbb{E}_{P_0}(fg)$  where  $f, g$  are two square-integrable functions of the data  $O = (Y, T, W)$ . The Cauchy–Schwarz inequality states

$$|\langle f, g \rangle_{P_0}| \leq \|f\|_{P_0} \cdot \|g\|_{P_0}.$$

Applying this inequality to Eq. 71 yields the inequality

$$\text{Rem}(P, P_0) \leq \|(g_0 - g)/g\|_{P_0} \cdot \|Q_0 - Q\|_{P_0}. \quad (72)$$

The consistency of TMLE is double-robust in the following sense. Given estimators  $(g_n, Q_n)$  of  $(g_0, Q_0)$ , e.g., constructed using a Super Learner, and assume  $g_n$  is a consistent estimator of  $g_0$ ; the same argument holds if we estimate  $Q_0$  consistently. Consistency implies

$$\|(g_0 - g_n)/g_n\|_{P_0} \longrightarrow 0, \quad \text{as } n \rightarrow \infty,$$

and so  $\text{Rem}(\hat{P}_n, P_0)$  vanishes as  $n \rightarrow \infty$ . Then by the definition of  $\text{Rem}(P, P_0)$ ,

$$\Psi(P_0) = \Psi(\hat{P}_n) + P_0 D(\hat{P}_n). \quad (73)$$

The TMLE step now achieves the following: It updates  $\hat{P}_n$  to a final estimate  $\hat{P}_n^*$ , so that

- (1)  $\hat{P}_n^*$  asymptotically has the same  $g$ -feature and  $Q$ -feature as  $P_0$ , so that it still solves the second-order remainder,  $\text{Rem}(\hat{P}_n^*, P_0) = 0$ , as  $n \rightarrow \infty$ ; and
- (2)  $\hat{P}_n^*$  solves the  $P_0$ -specific influence curve, i.e.,  $P_0 D(\hat{P}_n^*) = 0$ , so that  $\Psi(\hat{P}_n^*) = \Psi(P_0)$ .

We include the proof of the proposition for the convenience of the reader.

*Proof.* First, we recall that the influence curve of  $\Psi_v(P)$ ,

$$D_v(P) = \frac{\mathbb{1}\{V = v\}}{g(v, W)} [Y - Q(V, W)] + Q(v, W) - \Psi_v(P), \quad (74)$$

satisfies the property  $PD_v(P) \equiv \mathbb{E}_P[D_v(P)] = 0$ , *i.e.*, it is mean-centered. Second, using this property, we simplify the second-order remainder  $\text{Rem}(P, P_0)$  as follows:

$$\begin{aligned} \text{Rem}(P, P_0) &= \Psi_v(P) - \Psi_v(P_0) + (P_0 - P)D_v(P) \\ &= \Psi_v(P) - \Psi_v(P_0) + \mathbb{E}_{P_0} \left[ \frac{\mathbb{1}\{V = v\}}{g(v, W)} (Y - Q(V, W)) + Q(v, W) \right] - \Psi_v(P) \\ &= \mathbb{E}_{P_0} \left[ \frac{\mathbb{1}\{V = v\}}{g(v, W)} (Y - Q(V, W)) + Q(v, W) - Q_0(v, W) \right]. \end{aligned}$$

Here, the first equality follows by definition of  $\text{Rem}(P, P_0)$ , the second follows by  $PD_v(P) = 0$  and  $\mathbb{E}_{P_0}[\Psi_v(P)] = \Psi_v(P)$ , and the third follows by definition of the parameter  $\Psi_v(P_0)$ . Third, we examine the second half of the second-order remainder:

$$\mathbb{E}_{P_0}[Q(v, W) - Q_0(v, W)] = \mathbb{E}_{P_0} \left[ g(v, W) \cdot \frac{Q(v, W) - Q_0(v, W)}{g(v, W)} \right]$$

Fourth, we split the remaining term up into two parts,

$$\begin{aligned} \mathbb{E}_{P_0} \left[ \frac{\mathbb{1}\{V = v\}}{g(v, W)} (Y - Q(V, W)) \right] &= \mathbb{E}_{P_0} \left[ \frac{\mathbb{1}\{V = v\}}{g(v, W)} (Y - Q(v, W)) \right] \\ &= \mathbb{E}_{P_0} \left[ \frac{\mathbb{1}\{V = v\}}{g(v, W)} (Y - Q_0(v, W)) \right] + \mathbb{E}_{P_0} \left[ \frac{\mathbb{1}\{V = v\}}{g(v, W)} (Q_0(v, W) - Q(v, W)) \right] \end{aligned}$$

and treat these separately, using the tower rule for both. The first term vanishes:

$$\begin{aligned} \mathbb{E}_{P_0} \left[ \frac{\mathbb{1}\{V = v\}}{g(v, W)} (Y - Q_0(v, W)) \right] &= \mathbb{E}_{P_0} \left[ \frac{\mathbb{1}\{V = v\}}{g(v, W)} \mathbb{E}_{P_0}(Y - Q_0(v, W) | V = v, W) \right] \\ &= \mathbb{E}_{P_0} \left[ \frac{\mathbb{1}\{V = v\}}{g(v, W)} (Q_0(v, W) - Q_0(v, W)) \right] = 0, \end{aligned}$$

since  $\mathbb{E}_{P_0}[Y | V = v, W] = Q_0(v, W)$ . The second term is equal to the following expression:

$$\begin{aligned} \mathbb{E}_{P_0} \left[ \frac{\mathbb{1}\{V = v\}}{g(v, W)} \{Q_0(v, W) - Q(v, W)\} \right] &= \mathbb{E}_{P_0} \left[ \mathbb{E}_{P_0} \left\{ \frac{\mathbb{1}\{V = v\}}{g(v, W)} [Q_0(v, W) - Q(v, W)] \middle| W \right\} \right] \\ &= \mathbb{E}_{P_0} \left[ \mathbb{E}_{P_0}[\mathbb{1}\{V = v\} | W] \frac{Q_0(v, W) - Q(v, W)}{g(v, W)} \right] \\ &= \mathbb{E}_{P_0} \left[ g_0(v, W) \frac{Q_0(v, W) - Q(v, W)}{g(v, W)} \right]. \end{aligned}$$

Finally, putting all terms together yields the expression for the second-order remainder:

$$\text{Rem}(P, P_0) = \mathbb{E}_{P_0} \left[ \frac{g_0(v, W) - g(v, W)}{g(v, W)} (Q_0(v, W) - Q(v, W)) \right] \quad (75)$$

This completes the proof.  $\square$

## APPENDIX B. RESIDUALS OF ONE FIT AS DEPENDENT VARIABLE IN ANOTHER FIT

In order to provide a fast method for estimating effect sizes of large numbers of variants on traits, the 2-step residual fitting, originally proposed in GRAMMAR [14], has been used in the literature to perform GWAS fits. In this procedure, the trait  $Y$  is first fitted linearly, either using linear regression or an LMM model, as a function of the covariates and random genetic effects. The residual of the fit, *i.e.*,  $Y - Y_{\text{predict}}$ , is then used as a dependent variable in a second linear regression model, to estimate the effect size of the variants of interest. However, this procedure results in biased effect estimates and ‘conservative’ tests [35]. Here we demonstrate that the term ‘conservative’ is misleading for two reasons:

- (1) In the case where a linear model is assumed to be the true model, the estimates of the effect sizes turn out to be less significant than the ground truth value. However, this is only true when *all* relevant covariates have been regressed out in the first step, and *only* the single variant of interest is regressed in the residual fitting step (second step). In contrast, if the residuals from the step 1 are used to obtain the effect size of a variant in the presence of another variable in the fit, the quantities of interest can be under- or over-estimated, and even change sign. This happens, for example, when environmental factors and other variants of interest (*e.g.*, for epistasis quantification) are not taken into account in step 1. The extent of under- or overestimation depends on the degree of dependence (*e.g.*, correlation) between the first set and the second set of variables. The issue of 2-step residual fitting has been demonstrated using simulations in other literature, *e.g.*, economics [10].
- (2) When the ground truth is non-linear, it is unknown in general what the effect of residual fitting is on the validity of the estimates. These estimates may be under- or over-inflated, or change size.

GRAMMAR-Gamma [35] introduced a correction factor for the test-statistic and effect size estimates. However, this correction factor only applies when the (unknowable) ground truth is in fact linear and the mistake described in (1) above are not made. If the (unknowable) ground truth follows a different model, an (unknowable) different analysis is required with (unknowable) new correction factor being applied to fit of the data. For more sophisticated and powerful models, it may be far more complicated to obtain the required correction. Therefore, this is a time-consuming procedure that attempts to treat the symptom instead of the cause of the issue, namely, the 2-step residual fitting procedure. This process should thus be entirely avoided.

REFERENCES

- [1] Genotyping and quality control of UK Biobank, a large-scale, extensively phenotyped prospective resource. *Interim Data Release 2015*, v1.2, 2015.
- [2] The “All of Us” Research Program. *New England Journal of Medicine*, 381(7):668–676, 2019. PMID: 31412182.
- [3] Abdel Abdellaoui, Conor V. Dolan, Karin J. H. Verweij, and Michel G. Nivard. Gene-environment correlations across geographic regions affect genome-wide association studies. *Nature Genetics*, 2022.
- [4] Sjoerd Viktor Beentjes and Ava Khamseh. Higher-order interactions in statistical physics and machine learning: A model-independent solution to the inverse problem at equilibrium. *Phys. Rev. E*, 102:053314, Nov 2020.
- [5] Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society. Series B (Methodological)*, 57(1):289–300, 1995.
- [6] David Benkeser and Mark J. van der Laan. The highly adaptive lasso estimator. In *2016 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, pages 689–696, 2016.
- [7] Clare Bycroft, Colin Freeman, Desislava Petkova, Gavin Band, Lloyd T. Elliott, Kevin Sharp, Allan Motyer, Damjan Vukcevic, Olivier Delaneau, Jared O’Connell, Adrian Cortes, Samantha Welsh, Alan Young, Mark Effingham, Gil McVean, Stephen Leslie, Naomi Allen, Peter Donnelly, and Jonathan Marchini. The UK Biobank resource with deep phenotyping and genomic data. *Nature*, 562(7726):203–209, 2018.
- [8] Oriol Canela-Xandri, Konrad Rawlik, and Albert Tenesa. An atlas of genetic associations in UK Biobank. *Nature Genetics*, 50(11):1593–1599, 2018.
- [9] Christopher C. Chang, Carson C. Chow, Laurent C.A.M. Tellier, Shashaank Vattikuti, Shaun M. Purcell, and James J. Lee. Second-generation PLINK: rising to the challenge of larger and richer datasets. *GigaScience*, 4(1), 02 2015. s13742-015-0047-8.
- [10] Wei Chen, Paul Hribar, and Samuel Melessa. Incorrect inferences when using residuals as dependent variables. *Journal of Accounting Research*, 56(3):751–796, 2018.
- [11] Melina Claussnitzer, Simon N. Dankel, Kyoung-Han Kim, Gerald Quon, Wouter Meuleman, Christine Haugen, Viktoria Glunk, Isabel S. Sousa, Jacqueline L. Beaudry, Vijitha Puviindran, Nezar A. Abdennur, Jannel Liu, Per-Arne Svensson, Yi-Hsiang Hsu, Daniel J. Drucker, Gunnar Mellgren, Chi-Chung Hui, Hans Hauner, and Manolis Kellis. FTO Obesity Variant Circuitry and Adipocyte Browning in Humans. *The New England journal of medicine*, 373(10):895–907, September 2015.
- [12] Molly M. Davies and Mark J. van der Laan. Sieve Plateau Variance Estimators: A New Approach to Confidence Interval Estimation for Dependent Data. *U.C. Berkeley Division of Biostatistics Working Paper Series*, Working Paper 322, 2014.
- [13] Sandrine Dudoit and Mark J. van der Laan. *Multiple Testing Procedures with Application to Genomics*. Springer Series in Statistics. Springer, New York, 2008.

- [14] Yurii S. enko, Dirk-Jan de Koning, and Chris Haley. Genome-wide Rapid Association Using Mixed Model and Regression: A Fast and Simple Method For Genome-wide Pedigree-Based Quantitative Trait Loci Association Analysis. *Genetics*, 177(1):577–585, 09 2007.
- [15] Timothy M. Frayling, Nicholas J. Timpson, Michael N. Weedon, Eleftheria Zeggini, Rachel M. Freathy, Cecilia M. Lindgren, John R. B. Perry, Katherine S. Elliott, Hana Lango, Nigel W. Rayner, Beverley Shields, Lorna W. Harries, Jeffrey C. Barrett, Sian Ellard, Christopher J. Groves, Bridget Knight, Ann-Marie Patch, Andrew R. Ness, Shah Ebrahim, Debbie A. Lawlor, Susan M. Ring, Yoav Ben-Shlomo, Marjo-Riitta Jarvelin, Ulla Sovio, Amanda J. Bennett, David Melzer, Luigi Ferrucci, Ruth J. F. Loos, s Barroso, Inel, Nicholas J. Wareham, Fredrik Karpe, Katharine R. Owen, Lon R. Cardon, Mark Walker, Graham A. Hitman, Colin N. A. Palmer, Alex S. F. Doney, Andrew D. Morris, George Davey Smith, Andrew T. Hattersley, and Mark I. McCarthy. A Common Variant in the FTO Gene Is Associated with Body Mass Index and Predisposes to Childhood and Adult Obesity. *Science*, 316(5826):889–894, May 2007.
- [16] Nicolo Fusi, Christoph Lippert, Neil D. Lawrence, and Oliver Stegle. Warped linear mixed models for the genetic analysis of transformed phenotypes. *Nature Communications*, 5(1):4890, Sep 2014.
- [17] John Michael Gaziano, John Concato, Mary Brophy, Louis Fiore, Saiju Pyarajan, James Breeling, Stacey Whitbourne, Jennifer Deen, Colleen Shannon, Donald Humphries, Peter Guarino, Mihaela Aslan, Daniel Anderson, Rene LaFleur, Timothy Hammond, Kendra Schaa, Jennifer Moser, Grant Huang, Sumitra Muralidhar, Ronald Przygodzki, and Timothy J. O’Leary. Million Veteran Program: A mega-biobank to study genetic influences on health and disease. *Journal of Clinical Epidemiology*, 70:214–223, 2016.
- [18] Marta Guindo-Martínez, Ramon Amela, Silvia Bonàs-Guarch, Montserrat Puiggròs, Cecilia Salvoró, Irene Miguel-Escalada, Caitlin E. Carey, Joanne B. Cole, Sina Rüeger, Elizabeth Atkinson, Aaron Leong, Friman Sanchez, Cristian Ramon-Cortes, Jorge Ejarque, Duncan S. Palmer, Mitja Kurki, Krishna Aragam, Jose C. Florez, Rosa M. Badia, Josep M. Mercader, David Torrents, and FinnGen Consortium. The impact of non-additive genetic associations on age-related complex diseases. *Nature Communications*, 12(1):2436, 2021.
- [19] Philip W. Hedrick. What is the evidence for heterozygote advantage selection? *Trends in Ecology & Evolution*, 27(12):698–704, 2022/09/11 2012.
- [20] Xia Jiang, Paul F. O’Reilly, Hugues Aschard, Yi-Hsiang Hsu, J. Brent Richards, José Dupuis, Erik Ingelsson, David Karasik, Stefan Pilz, Diane Berry, Bryan Kestenbaum, Jusheng Zheng, Jianan Luan, Eleni Sofianopoulou, Elizabeth A. Streeten, Demetrius



- Albanes, Pamela L. Lutsey, Lu Yao, Weihong Tang, Michael J. Econs, Henri Wallaschofski, Henry Völzke, Ang Zhou, Chris Power, Mark I. McCarthy, Erin D. Michos, Eric Boerwinkle, Stephanie J. Weinstein, Neal D. Freedman, Wen-Yi Huang, Natasja M. Van Schoor, Nathalie van der Velde, Lisette C. P. G. M. de Groot, Anke Eneman, L. Adrienne Cupples, Sarah L. Booth, Ramachandran S. Vasan, Ching-Ti Liu, Yanhua Zhou, Samuli Ripatti, Claes Ohlsson, Liesbeth Vandenput, Mattias Lorentzon, Johan G. Eriksson, M. Kyla Shea, Denise K. Houston, Stephen B. Kritchevsky, Yongmei Liu, Kurt K. Lohman, Luigi Ferrucci, Munro Peacock, Christian Gieger, Marian Beekman, Eline Slagboom, Joris Deelen, Diana van Heemst, Marcus E. Kleber, Winfried März, Ian H. de Boer, Alexis C. Wood, Jerome I. Rotter, Stephen S. Rich, Cassianne Robinson-Cohen, Martin den Heijer, Marjo-Riitta Jarvelin, Alana Cavadino, Peter K. Joshi, James F. Wilson, Caroline Hayward, Lars Lind, Karl Michaëlsson, Stella Trompet, M. Carola Zillikens, Andre G. Uitterlinden, Fernando Rivadeneira, Linda Broer, Lina Zgaga, Harry Campbell, Evropi Theodoratou, Susan M. Farrington, Maria Timofeeva, Malcolm G. Dunlop, Ana M. Valdes, Emmi Tikkanen, Terho Lehtimäki, Leo-Pekka Lyytikäinen, Mika Kähönen, Olli T. Raitakari, Vera Mikkilä, M. Arfan Ikram, Naveed Sattar, J. Wouter Jukema, Nicholas J. Wareham, Claudia Langenberg, Nita G. Forouhi, Thomas E. Gundersen, Kay-Tee Khaw, Adam S. Butterworth, John Danesh, Timothy Spector, Thomas J. Wang, Elina Hyppönen, Peter Kraft, and Douglas P. Kiel. Genome-wide association study in 79,366 European-ancestry individuals informs the genetic architecture of 25-hydroxyvitamin D levels. *Nature Communications*, 9(1):260, January 2018.
- [21] Neale Lab. UK Biobank GWAS by Ben Neale et al. v2, aug 2018.
- [22] Tuuli Lappalainen and Daniel G. MacArthur. From variant to function in human disease genetics. *Science*, 373(6562):1464–1468, 2021.
- [23] Christoph Lippert, Jennifer Listgarten, Ying Liu, Carl M. Kadie, Robert I. Davidson, and David Heckerman. Fast linear mixed models for genome-wide association studies. *Nature Methods*, 8(10):833–835, 2011.
- [24] Trudy FC Mackay and Jason H. Moore. Why epistasis is important for tackling complex human disease genetics. *Genome Medicine*, 6(6):42, 2014.
- [25] Zachary R. McCaw, Thomas Colthurst, Taedong Yun, Nicholas A. Furlotte, Andrew Carroll, Babak Alipanahi, Cory Y. McLean, and Farhad Hormozdiani. DeepNull models non-linear covariate effects to improve phenotypic prediction and association power. *Nature Communications*, 13(1):241, 2022.
- [26] Ardalan Naseri, Kecong Tang, Xin Geng, Junjie Shi, Jing Zhang, Pramesh Shakya, Xiaoming Liu, Shaojie Zhang, and Degui Zhi. Personalized genealogical history of uk individuals inferred from biobank-scale ibd segments. *BMC Biology*, 19(1):32, 2021.
- [27] Rachael V. Phillips, Mark J. van der Laan, Hana Lee, and Susan Gruber. Practical considerations for specifying a super learner, April 2022.

- [28] Joseph E. Powell, Peter M. Visscher, and Michael E. Goddard. Reconciling the analysis of ibd and ibs in complex trait studies. *Nature Reviews Genetics*, 11(11):800–805, 2010.
- [29] Shaun Purcell and Christopher Chang. Plink 1.9. <https://www.cog-genomics.org/plink/1.9>, 2021.
- [30] Shaun Purcell and Christopher Chang. Plink 2.0. [www.cog-genomics.org/plink/2.0/](http://www.cog-genomics.org/plink/2.0/), 2021.
- [31] Shaun Purcell, Benjamin Neale, Kathe Todd-Brown, Lori Thomas, Manuel A.R. Ferreira, David Bender, Julian Maller, Pamela Sklar, Paul I.W. de Bakker, Mark J. Daly, and et al. Plink: A tool set for whole-genome association and population-based linkage analyses. *The American Journal of Human Genetics*, 81(3):559–575, 2007.
- [32] M. Rosenblum and M. J. van der Laan. Using regression models to analyze randomized trials: Asymptotically valid hypothesis tests despite incorrectly specified models. *Biometrics*, 65(3):937–945, 2009.
- [33] Matteo Sesia, Stephen Bates, Emmanuel Candès, Jonathan Marchini, and Chiara Sabatti. False discovery rate control in genome-wide association studies with population structure. *Proceedings of the National Academy of Sciences*, 118(40), 2021.
- [34] Cathie Sudlow, John Gallacher, Naomi Allen, Valerie Beral, Paul Burton, John Danesh, Paul Downey, Paul Elliott, Jane Green, Martin Landray, Bette Liu, Paul Matthews, Giok Ong, Jill Pell, Alan Silman, Alan Young, Tim Sprosen, Tim Peakman, and Rory Collins. Uk biobank: An open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLOS Medicine*, 12(3):1–10, 03 2015.
- [35] Gulnara R. Svishcheva, Tatiana I. Axenovich, Nadezhda M. Belonogova, Cornelia M. van Duijn, and Yurii S. Aulchenko. Rapid variance components-based method for whole-genome association analysis. *Nature Genetics*, 44(10):1166–1170, 2012.
- [36] Anastasios A. Tsiatis. *Semiparametric theory and missing data*. Springer Series in Statistics. Springer, New York, 2006.
- [37] Mark J. van der Laan. A generally efficient targeted minimum loss based estimator based on the highly adaptive Lasso. *Int. J. Biostat.*, 13(2):20150097, 35, 2017.
- [38] Mark J. van der Laan, Eric C. Polley, and Alan E. Hubbard. Super learner. *Stat. Appl. Genet. Mol. Biol.*, 6:Art. 25, 23, 2007.
- [39] Mark J. van der Laan and Sherri Rose. *Targeted learning*. Springer Series in Statistics. Springer, New York, 2011. Causal inference for observational and experimental data.
- [40] Mark J. van der Laan and Sherri Rose. *Targeted learning in data science*. Springer Series in Statistics. Springer, Cham, 2018. Causal inference for complex longitudinal studies.
- [41] A. W. van der Vaart. *Asymptotic statistics*, volume 3 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge, 1998.
- [42] Reiner A. Veitia. Nonlinear effects in macromolecular assembly and dosage sensitivity. *Journal of Theoretical Biology*, 220(1):19–25, 2003.

- [43] Wen-Hua Wei, Gibran Hemani, and Chris S. Haley. Detecting epistasis in human complex traits. *Nature Reviews Genetics*, 15(11):722–733, 2014.
- [44] Andrew R. Wood, Jessica Tyrrell, Robin Beaumont, Samuel E. Jones, Marcus A. Tuke, Katherine S. Ruth, Hanieh Yaghootkar, Rachel M. Freathy, Anna Murray, Timothy M. Frayling, Michael N. Weedon, and The GIANT consortium. Variants in the *fto* and *cdkal1* loci have recessive effects on risk of obesity and type 2 diabetes, respectively. *Diabetologia*, 59(6):1214–1221, 2016.