

Multiethnic Polygenic Risk Prediction in Diverse Populations through Transfer Learning

Peixin Tian¹, Tsai Hor Chan¹, Yong-Fei Wang², Wanling Yang², Guosheng Yin¹, Yan Dora Zhang^{1,3*}

¹ Department of Statistics and Actuarial Science, The University of Hong Kong, Pokfulam Hong Kong SAR, China;

² Department of Paediatrics and Adolescent Medicine, The University of Hong Kong, Pokfulam Hong Kong SAR, China;

³ Centre for PanorOmic Sciences, Li Ka Shing Faculty of Medicine, The University of Hong Kong, Hong Kong SAR, China.

Correspondence*:
Yan Dora Zhang
doraz@hku.hk

2 ABSTRACT

3 Polygenic risk scores (PRS) leverage the genetic contribution of an individual's genotype to a
4 complex trait by estimating disease risk. Traditional PRS prediction methods are predominantly for
5 European population. The accuracy of PRS prediction in non-European populations is diminished
6 due to much smaller sample size of genome-wide association studies (GWAS). In this article,
7 we introduced a novel method to construct PRS for non-European populations, abbreviated as
8 TL-Multi, by conducting transfer learning framework to learn useful knowledge from European
9 population to correct the bias for non-European populations. We considered non-European GWAS
10 data as the target data and European GWAS data as the informative auxiliary data. TL-Multi
11 borrows useful information from the auxiliary data to improve the learning accuracy of the target
12 data while preserving the efficiency and accuracy. To demonstrate the practical applicability of the
13 proposed method, we applied TL-Multi to predict the risk of systemic lupus erythematosus (SLE)
14 in Asian population and the risk of asthma in Indian population by borrowing information from
15 European population. TL-Multi achieved better prediction accuracy than the competing methods
16 including Lassosum and meta-analysis in both simulations and real applications.

17 **Keywords:** genome-wide association study, polygenic risk score, transfer learning, multiethnic populations

1 INTRODUCTION

18 Genetic risk prediction is an important methodology for understanding the underlying genetic architecture
19 and the inclusion of information on complex traits, such as estimating the genetic risk of complex traits
20 or diseases (e.g., coronary artery disease)(Chatterjee et al., 2016; Ge et al., 2019). Polygenic risk scores
21 (PRS) are one of the approaches to reflect a mathematical aggregation of risk by variants such as single
22 nucleotide polymorphisms (SNPs) (Peterson et al., 2019). With the application of the best linear unbiased

23 predictor to estimate PRS, some methods use summary association statistics as training data (Consortium,
24 2009; Vilhjálmsson et al., 2015; Shi et al., 2016), and others require individual-level data, such as genotype
25 data and phenotypes (De Los Campos et al., 2010; Speed and Balding, 2014; Maier et al., 2015; Moser
26 et al., 2015; Coram et al., 2017). As an implementation, PRS have become a widely used statistical tool
27 to estimate the genetic risk of certain diseases or phenotypes (Mak et al., 2017). Specifically, PRS for a
28 particular disease demonstrate the risk index for people to suffer from the disease. A remarkable study
29 of five common diseases (coronary artery disease, atrial fibrillation, type 2 diabetes, inflammatory bowel
30 disease, and breast cancer) found that people with top 8.0, 6.1, 3.5, 3.2, and 1.5% highest PRS had a
31 three-fold higher risk to develop these diseases than people with average PRS (Khera et al., 2018).

32 However, the majority of public genome-wide association study (GWAS) data has been conducted in
33 European population (Popejoy and Fullerton, 2016). Due to the limited availability of non-European
34 ancestral data and the diversity of linkage disequilibrium (LD) architectures among distinct populations,
35 previous studies showed that the genetic architectures of specific phenotypes or diseases were highly
36 consistent between populations (single-variant level and genome-wide level) (Huang et al., 2021).
37 Hence, using PRS derived from European population can result in disease associations being under-
38 or overestimated in other populations (Kim et al., 2018). Traditional approaches are insufficient to address
39 this challenge when multiple populations are involved. Recent genetic statistical studies have indicated that
40 diverse population variants share the same underlying causal variants (Brown et al., 2016; Shi et al., 2020),
41 which raises the possibility of transferability of PRS across distinct ethnic groups. However, existing studies
42 focus mostly on the application with one homogeneous population. For example, LDpred (Vilhjálmsson
43 et al., 2015) and PRS-CS (Ge et al., 2019) improve the prediction accuracy by enhancing LD modelling.
44 As an alternative, a penalized regression framework based on summary statistics, namely Lassosum was
45 proposed by Mak et al. (2017), whereas these methods are limited to GWAS data from one homogeneous
46 population. **Current multiethnic PRS construction approaches that incorporate training data from both
47 the European and target populations can leverage trans-ethnic GWAS information and stratify squared
48 trans-ethnic genetic correlation in explanation of environmental effects on genes (Mak et al., 2017; Coram
49 et al., 2017; Shi et al., 2021).** Moreover, Márquez-Luna et al. (2017) proposed PT-Multi for multiethnic
50 PRS prediction by performing LD-informed pruning and *P*-value thresholding (PT) (Consortium, 2009) on
51 each homogeneous population and linearly combining the **optimal PRS from each specific population.**

52 However, previous studies ignored the information gap among diverse populations. Li et al. (2020)
53 proposed a high-dimensional linear regression model to transfer knowledge between informative samples
54 and target samples to improve the learning performance of target samples. By using GWAS summary
55 statistics from different ancestries and incorporating the idea of transfer learning (Li et al., 2020), we
56 propose a novel statistical method called TL-Multi to enhance the transferability of polygenic risk prediction
57 across diverse populations. TL-Multi assumes most causal variants are shared among diverse populations.
58 There is a difference between the target samples and the informative auxiliary samples in the genetic
59 architecture, which causes estimation bias. TL-Multi further corrects this bias and estimates the PRS using
60 Lassosum (Mak et al., 2017). Additionally, TL-Multi inherits the advantages of Lassosum, ensuring that it
61 has more accurate performance in all circumstances than initial PT and circumvents convex optimization
62 challenges in LDpred. Moreover, TL-Multi extends the application to estimate the genetic risk from
63 unmatched ancestral populations, **and employs all available data without pruning or discarding.** For
64 practical analysis, we investigate TL-Multi prediction performance with informative auxiliary European
65 samples from UK Biobank (<https://www.ukbiobank.ac.uk>), and European summary statistics
66 and Hong Kong target samples from previous studies to predict PRS in systemic lupus erythematosus
67 (SLE) (Wang et al., 2021; Morris et al., 2016; Julià et al., 2018). We obtain a greater than 125% relative

68 improvement in prediction accuracy compared to only using GWAS data from Hong Kong population.
69 Furthermore, TL-Multi performs more accurately in PRS prediction in most scenarios in comparison with
70 the recent multiethnic methods, meta-analysis, and PT-Multi.

71 Additionally, we refer to Huang et al. (2021) to classify the PRS methods into two categories: single-
72 discovery methods and multi-discovery methods. Single-discovery methods use GWAS data from a single
73 homogeneous population, and multi-discovery methods apply the combined GWAS data of multiple
74 populations.

2 MATERIALS AND METHODS

75 2.1 Data Overview

76 In this study, we requested the individual-level genotyped data for a previous SLE GWAS in Hong Kong
77 (Wang et al., 2021) as the testing dataset, which included 1,604 SLE cases and 3,324 controls. We used
78 GWAS summary statistics of SLE from both East Asian and European populations to train the models.
79 The data for East Asians were collected from Guangzhou (GZ) and Central China (CC), including 2,618
80 SLE cases and 5,107 controls (Wang et al., 2021). The data for Europeans were obtained from previous
81 studies (Wang et al., 2021; Morris et al., 2016; Julià et al., 2018), involving a total of 4,576 cases and 8,039
82 controls. Variants with minor allele frequency greater than 1% and imputed INFO scores greater than 0.7
83 in respectively ancestral groups were reserved for the following analyses.

84 In our analysis of asthma, we requested the genotyped data of Indian and European individuals for
85 asthma from UK Biobank. The UK Biobank data consisted of 4,160 unrelated Indian samples genotyped
86 at 1,175,469 SNPs after QC and mapping HapMap 3 SNPs, and we further sampled 48,362 unrelated
87 British samples genotyped at 1,189,752 SNPs after QC and mapping HapMap 3 SNPs. We divided the
88 Indian samples into two groups: 3,160 samples as a training data set and 1,000 samples as a testing data
89 set. As stated previously, the final data set comprises of 3,160 (408 cases and 2,752 controls) unrelated
90 Indian samples for training, 1,000 (127 cases and 873 controls) unrelated Indian samples for testing, and
91 48,362 (6,555 cases and 41,807 controls) unrelated British samples for training. Variants with minor allele
92 frequency greater than 1% and P -values of Hardy-Weinberg equilibrium Fisher's exact test $< 1 \times 10^{-5}$
93 were kept. We then computed GWAS to derive the GWAS summary statistics of each population with each
94 genotypes (after quality control) and adjusting for age and gender, and the top 10 principal components.

95 2.2 Lassosum

96 Lassosum is a statistical approach introduced by Mak et al. (2017) which enables to tune parameters
97 without validation datasets and phenotype data via pseudovalidation, and outperforms PT and LDpred
98 in prediction (Consortium, 2009; Vilhjálmsón et al., 2015). It refers to the idea of Tibshirani (1996)
99 to deal with sparse matrices and calculate PRS only by using summary statistics and an external LD
100 reference panel. In this article, the ancestry-matched LD block is generally estimated by the 1000 Genome
101 project (<https://www.internationalgenome.org>). Additionally, we keep the reference panel's
102 ancestry consistent with that of our target population. Furthermore, if the SNP-wise correlation r_i is not
103 available, we can estimate r_i following Mak et al. (2017): $r_i = \frac{t_i}{\sqrt{n-1+t_i^2}}$.

104 2.3 PT-Multi

105 PT-Multi assumes the multiethnic PRS is a linear combination of the most predictive PRS from each
 106 population. First, it applies LD-pruning and P -value thresholding (PT) (Consortium, 2009) to each single
 107 ethnic summary statistics and gets the most predictive PRS. Second, it fits marginal linear regression models
 108 to get weights for each population, respectively. We apply the R package ‘*bigspnr*’ (Privé et al., 2018) to
 109 validation data for LD informed clumping with r^2 threshold of 0.1. The P -value thresholds are among:
 110 1, 0.3, 0.1, 3×10^{-2} , 10^{-2} , 3×10^{-3} , 10^{-3} , 3×10^{-4} , and 10^{-4} . We conduct 10-fold cross-validation to
 111 determine the optimal P -value threshold for each population. We use an independent validation data set to
 112 compute the final PRS and the average value of R^2 across the 10 folds.

113 This article uses single-discovery method (Lassosum) to regress European, Asian, and multi-discovery
 114 methods (meta-analysis, TL-Multi, PT-Multi) to determine the most predictive PRS with the highest R^2 .
 115 For ease of notations, let PRS_a , PRS_e , PRS_{ma} , PRS_{tl} , and PRS_{pt} represent PRS for Asian, European,
 116 meta-analysis, TL-Multi and PT-Multi, respectively.

117 2.4 Meta-analysis of two diverse ancestries

We generate the estimates of effect sizes of joint GWAS data by

$$\hat{\beta}_{ma} = \frac{\frac{\beta_a^2}{se_a} + \frac{\beta_e^2}{se_e}}{se_a^{-2} + se_e^{-2}},$$

where β_a and β_e are the effect sizes obtained from Asian and European GWAS data, respectively, and se_a and se_e are the standard errors obtained directly from ancestry-matched GWAS data. Furthermore, the estimate of the standard error in meta-analysis is defined as:

$$\hat{se}_{ma} = \sqrt{\frac{1}{se_a^{-2} + se_e^{-2}}},$$

and the estimate of z-statistic is obtained from:

$$\hat{z}_{ma} = \frac{\hat{\beta}_{ma}}{\hat{se}_{ma}}.$$

The P -value is converted from \hat{z}_{ma} following:

$$P\text{-value} = 2\Phi(-|\hat{z}_{ma}|),$$

118 where $\Phi(\cdot)$ is the cumulative distribution function of the standard normal distribution $N(0, 1)$. In this
 119 meta-analysis, the ancestry of the reference panel is consistent with the ancestry of the target population.
 120 Furthermore, due to the majority of the total sample being of European ancestry, the LD block is estimated
 121 from European population in the 1000 Genome Project.

122 2.5 Multiethnic polygenic risk scores prediction using TL-Multi

In this article, we employ European population data as our informative auxiliary data, owing to its large sample size and relative accessibility. Additionally, we treat East Asians as the target population due to the scarcity of public data (Brown et al., 2016; Shi et al., 2020). Recall the fundamental framework we

using for genetic architecture and phenotype, it is a linear combination with effect sizes β , and an n -by- p genotype matrix X , where p is the number of columns containing marker genotype codes corresponding to the number of reference alleles on the sample-specific SNP (e.g., 0, 1, 2), and n is the sample size, following as:

$$y = X\beta + \epsilon,$$

where y is a vector of clinical outcomes. Tibshirani (1996) proposed Lasso which is commonly used to estimate coefficients $\hat{\beta}$ (weights of X), when p (the columns of X or the number of elements of y) is **relative large** to result in many $\hat{\beta}$ being 0. Specifically, the optimization problem of target population is equivalent to minimizing the objective function:

$$L(\beta_a) = (y_a - X_a\beta_a)^T(y_a - X_a\beta_a) + 2\lambda\|\beta_a\|_1,$$

123 where y_a is the vector of Asian phenotypes, X_a is the genotype matrix of Asian population, $L(\cdot)$ is an
 124 optimizing function, $\|\beta_a\|_1$ is the L_1 norm of β_a , **and** λ is a data-dependent parameter determining the
 125 proportion of β_a to be estimated to 0. It can be widely extended in the scenarios in which only the summary
 126 statistics are available (Mak et al., 2017).

127 Motivated by Lassosum, we further propose a novel method, namely TL-Multi to extend its application
 128 to multiethnic polygenic prediction. We observed additional samples from auxiliary studies (e.g., European
 129 population). The estimate of the marginal effect sizes of European population, $\hat{\beta}_e$, can be generated using
 130 the auxiliary model:

$$L(\beta_e) = (y_e - X_e\beta_e)^T(y_e - X_e\beta_e) + 2\lambda\|\beta_e\|_1, \quad (1)$$

where y_e is the vector of European phenotypes, **and** X_e is the genotype matrix of European population. For illustration, we denote the auxiliary studies, in which informative auxiliary samples can be transferred, and the target model and auxiliary model are similar at certain levels (e.g., similar genetic architectures). Furthermore, we assume the difference between auxiliary samples and target samples is denoted as (Li et al., 2020):

$$\hat{\delta} = \hat{\beta}_a - \hat{\beta}_e,$$

131 where $\hat{\beta}_a$ (the weights of target population e.g., Asian population X_a) is the target regression estimator, and
 132 $\hat{\beta}_e$ (the weights of auxiliary population e.g., European population X_e) is the estimator for auxiliary study.
 133 Furthermore, the informative auxiliary set, A_q , has a requirement to ensure that the information auxiliary
 134 set includes sufficiently different information under a constrained level. Specifically, the information
 135 difference should satisfy the sufficient sparsity:

$$A_q = \{\|\hat{\delta}\|_q \leq h\}, \quad (2)$$

136 where $q \in [0, 1]$, $\|\hat{\delta}\|_q$ is the L_q norm of the information difference $\hat{\delta}$ of the informative auxiliary samples.
 137 The assumption requires the auxiliary informative population A_q to include samples in their contrast
 138 vectors with a maximum L_q -sparsity of at most h . **Moreover, we assume A_q is informative to improve
 139 the prediction performance of target population while h is relatively small compared to $\hat{\beta}_a$. Specifically,
 140 when $q = 0$, the set A_q implies that there are at most h casual variants. For $q \in (0, 1]$, this scenario may
 141 be explained that all the variants are causal variants with rapid relative amplitudes decaying effect sizes.
 142 Therefore, the smaller the h , the auxiliary samples of A_q tend to be more informative, where $|A_q|$ leverages
 143 the number of informative auxiliary samples.**

144 Our goal is to correct the bias between these populations and improve prediction performance in Asian
 145 population. First, we can estimate the marginal effect sizes of European population, $\hat{\beta}_e$ by minimizing the
 146 objective function based on equation (1):

$$L(\hat{\beta}_e) = \operatorname{argmax}_{\beta} 2\mathbf{r}_e^T \beta_e - \beta_e^T \mathbf{R}_e \beta_e - \lambda \|\beta_e\|_1, \quad (3)$$

147 where $\mathbf{r}_e = \mathbf{X}_e^T \mathbf{y}_e$ is the SNP-wise correlation between the genotype matrix of European population \mathbf{X}_e
 148 and the phenotype \mathbf{y}_e , and \mathbf{R}_e is the LD matrix indicating a matrix of correlations between SNPs. The
 149 estimates of \mathbf{r}_e can be obtained from summary statistics, and the estimates of \mathbf{R}_e can be obtained from
 150 publicly available databases, such as the 1000 Genome project. As Mak et al. (2017) indicated, the PRS
 151 can be estimated by optimizing equation (3) without extra individual-level data.

152 Specifically, TL-Multi estimates the PRS of Asian population by correcting the bias between European
 153 and Asian populations. We further denote the bias as δ which is the difference between European and
 154 Asian populations in genetic architecture. The new estimate of effect sizes of Asian population can be
 155 presented as: $\beta_{tl} = \beta_e + \delta$, in which δ is estimated by:

$$L(\hat{\delta}) = \operatorname{argmax}_{\delta} 2 \left(\mathbf{r}_a^T \delta + \delta \mathbf{R}_a \beta_e \right) - \delta^T \mathbf{R}_a \delta - \lambda_{\delta} \|\delta\|_1. \quad (4)$$

156 According to pseudovalidation proposed by Mak et al. (2017), the optimal single-discovery PRS for
 157 European and Asian populations can be determined directly by the highest R^2 without the phenotypes.
 158 The optimal estimates of effect sizes of Asian and European populations that we apply to TL-Multi
 159 are the ancestry-matched optimal PRS, respectively. The Algorithm 1 describes our proposed TL-Multi
 160 algorithm, and we further develop an R package, which is publicly available at <https://github.com/mxxptian/TL-Multi.git>.
 161

Algorithm 1 Algorithm for TL-Multi

Data: $\mathbf{r}_a, \mathbf{r}_e, \mathbf{X}_a^*$ (genotype matrix of target samples), $\mathbf{R}_a, \mathbf{R}_e, \mathbf{y}_a, (\lambda^{(1)}, \dots, \lambda^{(K)})$ (the tuning parameters for $\hat{\beta}_e$ for K -fold cross validation), $(\lambda_{\delta}^{(1)}, \dots, \lambda_{\delta}^{(K)})$ (the tuning parameters for $\hat{\delta}$ for K -fold cross validation);

Result: PRS_{tl};

- 1 Obtain $\{\hat{\beta}_e^{(k)}\}_{k=1}^K$ by solving $L(\hat{\beta}_e) = \{j \in [K] : \operatorname{argmax}_{\beta} 2\mathbf{r}_e^T \beta_e - \beta_e^T \mathbf{R}_e \beta_e - \lambda^{(j)} \|\beta_e\|_1\}$ with different tuning parameters $(\lambda^{(1)}, \dots, \lambda^{(K)})$;
 - 2 Evaluate model performance by R^2 with \mathbf{y}_a ;
 - 3 Obtain the optimal $\hat{\beta}_e$ with the maximum R^2 ;
 - 4 Obtain $\hat{\beta}_{tl}^{(j)} = \hat{\delta}^{(j)} + \hat{\beta}_e$ by solving $L(\hat{\delta}) = \{j \in [K] : \operatorname{argmax}_{\delta} 2\mathbf{r}_a^T \delta + 2\delta \mathbf{X}_a^{*T} \mathbf{X}_a \hat{\beta}_e - \delta^T \mathbf{R}_a \delta - \lambda_{\delta}^{(j)} \|\delta\|_1\}$ with different tuning parameters $(\lambda_{\delta}^{(1)}, \dots, \lambda_{\delta}^{(K)})$;
 - 5 Evaluate model performance by R^2 with \mathbf{y}_a ;
 - 6 Obtain the optimal $\hat{\beta}_{tl}$ with the maximum R^2 ;
 - 7 PRS_{tl} = $\hat{\beta}_{tl} \mathbf{X}_a^*$.
-

162 **2.6 Simulation studies**

We performed a wide range of simulation studies to evaluate the performance of TL-Multi. We used real genotypes of European population from UK Biobank and Asian population from [previous SLE study \(Wang et al., 2021\)](#). Following the quality control procedure provided in [Chang et al. \(2015\)](#), we utilized the UK Biobank and Asian lupus genotype data whose P -values of Hardy–Weinberg equilibrium Fisher’s exact test $< 1 \times 10^{-5}$ with minor allele frequency (MAF) $> 1\%$ and filtered out SNPs and missing samples. Then, we simulated the effect sizes based on the genetic architecture correlation and applied the R package ‘bigsnpr’ ([Privé et al., 2018](#)) to generate quantitative phenotypes and conduct GWAS to determine the summary statistics. Based on these estimated summary statistics, we employed the following [PRS prediction methods](#). We further extracted the common variants between European samples and Asian samples. This resulted in 69,398 SNPs in total, and 4,049 subjects in Asian population. We fixed SNP-heritability h^2 at 0.5, and further simulated genetic architectures by randomly treating 1%, 1.5%, 2%, and 5% variants as causal variants. We assumed that these causal variants were shared in multiple populations with different effect sizes. Additionally, we sampled effect sizes from a multivariate normal distribution with a wide range of cross-population genetic correlation values (0.2, 0.4, 0.6, and 0.8) ([Huang et al., 2021](#); [Bulik-Sullivan et al., 2015](#)), where for each population the variance is $\sigma^2 = \frac{h^2}{m}$ and m is the number of causal variants. There were 12 combinations in total. For each scenario, we generated 20 replicates and calculated the average values to assess the prediction accuracy. We took out the original phenotypes and generated new ones based on a linear framework:

$$y = X\beta + \epsilon,$$

163 where X is the training set of the standardized genotype matrix, and ϵ represent the random error which
164 was generated from $N(0, 1 - h^2)$. [And GWAS was implemented using the R package ‘bigsnpr’ to obtain](#)
165 [the summary statistics for each simulated phenotype.](#)

166 Due to the possibility that sample size affects performance, we investigated 25:1 and 50:1 proportions of
167 European samples to Asian samples. Additionally, we observed that the number of variants has a significant
168 influence on the prediction performance, and the majority of variants are located on chromosomes 1-
169 11. [Motivated by previous works \(Vilhjálmsón et al., 2015; Márquez-Luna et al., 2017\)](#), we further
170 [extrapolated the performance at large sample size by conducting simulations with different subsets of](#)
171 [chromosomes to increase \$\frac{N}{M}\$, where \$N\$ is the total number of samples and \$M\$ is the number of SNPs: \(1\)](#)
172 [using chromosomes 1-4; \(2\) using chromosomes 1-6; \(3\) using chromosomes 1-8; \(4\) using chromosomes](#)
173 [1-11.](#)

3 RESULTS174 **3.1 Simulations**

175 We performed simulations with real genotypes and simulated continuous phenotypes. We split the data
176 from Hong Kong population into two groups: 1,000 samples as a training data set and 3,049 samples as
177 testing data, and drew 50,000 samples from European samples. The training data set was used to simulate
178 phenotypes, and the testing data were applied to performance assessments. The prediction accuracy was
179 assessed by R^2 , which was based on the simulated phenotypes generated from the test data. Specifically,
180 LD blocks for single-discovery method were ancestry-matched as the reference panels, and they were in
181 correspondence with the ancestry of the target population for multi-discovery methods.

182 In Figure 1, we displayed the average values with a 95% upper bound of each simulation setting under
183 scenario (1) over 20 replicates. We conducted single-discovery analyses for Asian and European populations
184 by Lassosum, and multi-discovery analyses by meta-analysis, TL-Multi, and PT-Multi. Lassosum adopted
185 the PRS with the maximum R^2 by 10-fold cross-validation. We observed that meta-analysis could not
186 improve the prediction accuracy when single-discovery analysis of European population did not perform
187 better than the Asian one. **Particularly, when the genetic architecture correlation was quite low ($\rho = 0.2$),**
188 **meta-analysis and European prediction performances were comparably inferior. In this case, it was**
189 **explained that the shared information between the Asian and European populations would be limited,**
190 **preventing prediction improvement from being achieved by directly integrating the European data.** It also
191 reflected the consistent relationship between meta-analysis and single-discovery analysis of the informative
192 population. Moreover, meta-analysis could hardly outperform the European one. The performance of
193 Lassosum for European population **dominated** the performance of meta-analysis since the sample size
194 of European population is significantly larger than that of Hong Kong. Additionally, we observed that
195 TL-Multi could always improve the accuracy compared to Lassosum for Hong Kong population. **If the**
196 **genetic architecture correlation was not too high (e.g., $\rho = 0.4$ or 0.6), TL-Multi attained the highest**
197 **prediction accuracy compared to competing approaches. However, when the genetic architecture was high**
198 **($\rho = 0.8$), we noticed that TL-Multi performed slightly worse than other approaches. In this example,**
199 **the results might be explained by the remarkable similarity of the genetic architecture. When the genetic**
200 **correlation reaches 0.8, the majority of information about the Asian population would be directly explained**
201 **by that of the European population. Combining these two groups in the meta-analysis might increase the**
202 **accuracy of estimated effect sizes.**

203 In most scenarios, TL-Multi outperformed PT-Multi. **Specifically, TL-Multi substantially improved**
204 **multiethnic prediction accuracy for the instances with 1%, 1.5%, 2% causal proportions.** PT-Multi
205 conducted PT, which caused information loss in the data. However, TL-Multi could take all the data
206 information into account. We found TL-Multi performed poorly at a 5% causal proportion. We noted
207 that under this situation, the result of Lassosum for Hong Kong population was significantly inferior to
208 that of European. We referred to the assumption (2) to cast doubt on the breach of our assumption. **If**
209 **the assumption does not hold, European population could not be denoted as auxiliary informative data**
210 **because the useful information was limited. Due to it, TL-Multi would fail to borrow the information to**
211 **improve the learning performance of the target population. Alternatively, consider that the effect sizes were**
212 **simulated depending on the number of causal variations, m . As the proportion of causality rose, the effect**
213 **sizes tended to approach zero. Limited by small sample size of the Asian population, the bias between the**
214 **estimated effect sizes derived from the simulated phenotypes and the actual effect sizes would be even**
215 **larger. Some causal variants with relative small signals more likely erroneously failed to be captured which**
216 **resulted in the restricted TL-Multi's performance.** However, it is noteworthy that TL-Multi still enhanced
217 Hong Kong's prediction accuracy in this scenario. We discovered that the performance of meta-analysis
218 and PT-Multi for Hong Kong were nearly identical to that of Lassosum for Europeans, when we attributed
219 to the huge disparity in multiethnic sample sizes. To summarize, European population dominated the
220 performance of meta-analysis and PT-Multi. **In particular, TL-Multi could be employed to the moderate**
221 **genetic architecture correlations (e.g., $\rho = 0.4$, and 0.6) when the informative auxiliary population (e.g.,**
222 **European population) outperformed the target population (e.g., Hong Kong population). Referring to the**
223 **assumption (2), the performance of European population was supposed to be more accurate than that of**
224 **the target population, therefore it would be appropriate to borrow information from it. Moreover, if the**
225 **proportion of casuals increased, the estimated effect sizes of the target population would be relatively**

226 biased. We found that the precision of the effect sizes of the target population would have a substantial
227 effect on TL-Multi.

228 Alternatively, we generated phenotypes using different chromosome subsets and sample sizes of European
229 population while maintaining a fixed Hong Kong sample size. Over 20 replicates, we took the performance
230 using a fixed genetic correlation of 0.4 and 1.5% causal variants as an example. In Figure 2(A), we
231 drew 25,000 European subjects and 1,000 Hong Kong subjects. We observed that TL-Multi performed
232 much better than the competing approaches. While the performance of Hong Kong was superior to that
233 of Europeans, the performance of the meta-analysis was poor compared to that of Hong Kong. As the
234 total number of SNPs increased, the prediction accuracy of Hong Kong dramatically decreased. However,
235 the prediction accuracy of Europeans decreased relatively slowly. Specifically, under scenario (4), TL-
236 Multi was inferior to the other two multi-discovery methods. This could be explained that for this case,
237 there were 1,000 subjects from Hong Kong with 49,909 SNPs which resulted in a significant bias while
238 estimating the effect sizes by applying GWAS. In this case, TL-Multi thus failed to improve the accuracy
239 of the forecast compared to the previous scenarios, as the bias in the estimates of Hong Kong's effect
240 sizes was larger. Moreover, the consistent trend in European, meta-analysis, and PT-Multi supported
241 our previous extrapolation that the performance of European population could determine the primary
242 contribution of the other two. In Figure 2(B), we simulated 50,000 European subjects. We further observed
243 that the performance of PT-Multi was inferior to TL-Multi under the scenarios (1)-(3) and both of them
244 outperformed the single-discovery method and meta-analysis. Furthermore, the performance of meta-
245 analysis was consistent with that of European. As a result, even though the prediction accuracy of TL-Multi
246 went down, it was still better than the meta-analysis's prediction accuracy under all the scenarios.

247 3.2 Analysis of SLE in Hong Kong Population

248 We further applied the above four approaches to predict SLE risk in Hong Kong population to evaluate
249 the performance in real data analysis. We used European SLE GWAS summary statistics from previous
250 studies (Wang et al., 2021; Morris et al., 2016; Julià et al., 2018) (4,576 cases, and 8,039 controls), and the
251 ancestry-matched GWAS summary statistics (Wang et al., 2021) (2,618 cases, and 5,107 controls). The
252 validation data for Hong Kong population were from Wang et al. (2021) (1,604 cases, and 3,324 controls)
253 employing 10-fold cross-validation following Mak et al. (2017).

254 We reported the area under the receiver operating characteristic curve (AUC) to assess the prediction
255 accuracy of derived PRS. The ethnicity of the LD block is consistent with that of the majority population
256 in GWAS data, and the LD block was derived from Berisa and Pickrell (2016). Furthermore, the reference
257 panel was obtained from the 1000 Genome Project, and the ethnicity of it was consistent with the target
258 population's. We set the P -value thresholds to be the same as the values in simulation studies, and $r^2 = 0.1$.
259 In real data analysis, TL-Multi outperformed the competing methods. The optimal PRS from European
260 GWAS data yielded AUC of 0.6872 and 0.6943 from East Asian GWAS data. We further obtained the
261 optimal PRS of meta-analysis, TL-Multi and PT-Multi, with AUC values of 0.7098, 0.7131, and 0.5447,
262 and the corresponding ROC curves were depicted in Figure 3. For binary classification, we used a logistic
263 regression to obtain the mixing weights in PT-Multi. Consistent with the evaluations in simulation studies,
264 we observed that TL-Multi improved 2.7% in prediction accuracy compared to Lassosum for Hong Kong
265 population, and meta-analysis improved 2.2% compared to Lassosum. However, PT-Multi performed even
266 worse than single-discovery method in real data analysis.

267 Moreover, we reported the case prevalence of the bottom 2%, 5%, and 10% and top 2%, 5%, and 10% of
268 PRS distribution, constructed by single-discovery method, meta-analysis, and TL-Multi in Table 1. This

269 summary report demonstrated the case prevalence under different PRS conditions. For instance, the bottom
270 numbers indicate the prevalence of SLE among individuals with low PRS. We observed that TL-Multi had
271 satisfactory performance and showed 10.66, 7.50, and 5.80 fold increases comparing the top 2%, 5%, and
272 10% with bottom 2%, 5%, and 10% of the PRS distribution, respectively.

273

274 3.3 Analysis of Asthma in Indian Population

275 We applied the same methods to the Indian and European samples from UK Biobank computing associated
276 summary statistics by ‘bigsnpr’ R package. We splitted 1,000 (127 cases and 873 controls) unrelated Indian
277 samples as validation data, and 3,160 (408 cases and 2,752 controls) unrelated samples as training data,
278 and further sampled 48,362 (6,555 cases and 41,807 controls) unrelated European samples. We further
279 reported the AUC of the above four methods to evaluate the optimal prediction method. The ancestry of
280 LD blocks matches to that of the data’s predominant population. We used training data as a reference panel
281 whose ancestry was always identical to that of the target population. During pruning and clumping, the
282 P -value thresholds were set to be equal to simulation with $r^2 = 0.1$.

283 The ROC curves for binary classification are depicted in Figure 4. The optimal PRS from European and
284 Indian samples revealed AUC values of 0.5657 and 0.5441, respectively. In addition, for the multiethnic PRS
285 construction methods, the optimal PRS of meta-analysis, TL-Multi, and PT-Multi resulted in AUC values
286 of 0.5705, 0.5721, and 0.6427, respectively. We found that TL-Multi was superior to the all single-discovery
287 methods and meta-analysis. For binary classification, TL-Multi improved 5.15% in prediction accuracy
288 compared to Lassosum for Indian population, and meta-analysis improved 4.85% compared Lassosum
289 for Indian population. We noted that PT-Multi performed better than ours. However, the comparison
290 of PT-Multi method with other methods might not be fair since PT-Multi required individual level data
291 whereas other four approaches solely relied on the summary statistics. Moreover, access to individual-level
292 data was typically difficult.

293 Additionally, the case prevalence of the bottom 2%, 5%, and 10% and top 2%, 5%, and 10% of PRS
294 distribution, conducted by Lassosum for Indian and European, meta-analysis and TL-Multi was reported in
295 Table 2. We observed that TL-Multi would also perform with more accuracy in terms of case prevalence
296 than the competing methods.

4 DISCUSSION

297 In this article, we have proposed a novel approach named TL-Multi to improve the accuracy of PRS
298 prediction for non-European populations. Our proposed method leverages summary statistics and makes
299 complete use of all available data without clumping. We have shown that transferring the information
300 from the informative auxiliary populations (e.g., European) to the target populations (e.g., East Asian) can
301 indeed improve learning performance and the prediction accuracy of the target populations compared to
302 the single-discovery methods. Particularly, TL-Multi shows a higher AUC compared to meta-analysis and
303 PT-Multi in analysis of SLE in Hong Kong population. In our analysis of asthma in the Indian population,
304 TL-Multi outperforms Lassosum and meta-analysis in terms of prediction performance and case prevalence
305 prediction accuracy. Moreover we note that in the field of PRS prediction, there is no a particular method
306 outperforms all the others. It depends on the specific situation to select an appropriate method. For instance,
307 PRS-CS can always outperform PT in Huang et al. (2021), but PRS-CS may be inferior to PT in Weissbrod
308 et al. (2022) in some circumstances. Therefore, we provide some potential circumstances in which TL-Multi
309 would be an appropriate choice. First, TL-Multi is implemented using summary statistics and performs well

310 under the moderate genetic architecture correlation (e.g., $\rho = 0.4$ and 0.6) and moderate causal proportions
311 (e.g., $\frac{m}{M} = 1\%$, 1.5% , and 2%). Second, based on the assumption (2), TL-Multi would be a good alternative
312 when the single-discovery method's performance of the informative auxiliary population is superior to that
313 of the target population.

314 Compared to the single-discovery methods, we showed that the performance of TL-Multi was always
315 more accurate with an acceptable running time (e.g., 2 minutes) than the performance of Lassosum for
316 Hong Kong population, especially under moderate genetic correlation (e.g., $\rho = 0.6$). When the sample
317 size of the target data set is limited, increasing the sample size of the informative data set can enhance the
318 prediction accuracy of TL-Multi. In the simulation studies, we found that the performances of meta-analysis
319 and PT-Multi were dominated by the performance of Lassosum for European population. As the genetic
320 architecture correlation was rather high ($\rho = 0.8$), TL-Multi may perform poorly, and it would be more
321 prudent to consider approaches that integrate the whole data set across distinct populations. Therefore,
322 the performances of PT-Multi and meta-analysis were unsatisfactory, while the performance of European
323 population was worse than that of Hong Kong population.

324 Another advantage of TL-Multi is its powerful transferability, which corrects the bias in estimation
325 between European and non-European populations. De Candia et al. (2013) showed that the cross-population
326 genetic correlation could leverage the causal effect sizes in different populations. In simulation studies,
327 TL-Multi performed better when the genetic correlations were 0.4 and 0.6 . It indicated that TL-Multi
328 could be widely applied to two different populations which share some common genetic architecture
329 information. Moreover, TL-Multi retained the pseudovalidation proposed in Mak et al. (2017). It extended
330 the application of TL-Multi to fit the data without a validation data set and phenotype data.

331 Despite these advantages, some limitations of TL-Multi still remain for the future work. For example, if
332 the difference between two populations is too enormous, our proposed approach's assumptions will fail to
333 hold. It is worth bearing in mind to deal with this scenario. And in this article, we did not consider the X
334 chromosome, whose information could also contribute to prediction accuracy (Tukiainen et al., 2014). In
335 recent years, some approaches have fitted multiple diseases simultaneously (Maier et al., 2015; Turley et al.,
336 2018; Chung et al., 2019; Musliner et al., 2019; Graff et al., 2021). These studies inspire us to investigate
337 other TL-Multi extensions that bridge not only the gap between populations but also the gap between
338 illnesses in the interim.

CONFLICT OF INTEREST STATEMENT

339 The authors declare no competing interests.

AUTHOR CONTRIBUTIONS

340 Y.D.Z. conceived and supervised the study. P.T. and H.C. processed the data, implemented the software and
341 conducted the analysis. Y.W. and W.Y. provided the lupus data and offered insights of the interpretation
342 of the results. P.T. and Y.D.Z. wrote the first manuscript. All authors contributed to the revision of the
343 manuscript.

DATA AVAILABILITY STATEMENT

344 The 1000 Genome project data can be found in the <https://www.internationalgenome.org/>.
345 The UK Biobank data were accessed under Application Number 58942. Inquiries about the study's SLE
346 summary association statistics can be directed to the corresponding author.

ACKNOWLEDGMENTS

347 The data analysis were conducted using data from the UK Biobank Resource accessed under Application
348 Number 58942. This work was supported in part by Hong Kong Research Grants Council (RGC) Early
349 Career Scheme 2021/22 (project number: 27305221).

REFERENCES

- 350 Berisa, T. and Pickrell, J. K. (2016). Approximately independent linkage disequilibrium blocks in human
351 populations. *Bioinformatics* 32, 283
- 352 Brown, B. C., Ye, C. J., Price, A. L., Zaitlen, N., Consortium, A. G. E. N. T. . D., et al. (2016). Transethnic
353 genetic-correlation estimates from summary statistics. *The American Journal of Human Genetics* 99,
354 76–88
- 355 Bulik-Sullivan, B. K., Loh, P.-R., Finucane, H. K., Ripke, S., Yang, J., Patterson, N., et al. (2015). Ld score
356 regression distinguishes confounding from polygenicity in genome-wide association studies. *Nature*
357 *genetics* 47, 291–295
- 358 Chang, C. C., Chow, C. C., Tellier, L. C., Vattikuti, S., Purcell, S. M., and Lee, J. J. (2015). Second-
359 generation plink: rising to the challenge of larger and richer datasets. *Gigascience* 4, s13742–015
- 360 Chatterjee, N., Shi, J., and García-Closas, M. (2016). Developing and evaluating polygenic risk prediction
361 models for stratified disease prevention. *Nature Reviews Genetics* 17, 392–406
- 362 Chung, W., Chen, J., Turman, C., Lindstrom, S., Zhu, Z., Loh, P.-R., et al. (2019). Efficient cross-trait
363 penalized regression increases prediction accuracy in large cohorts using secondary phenotypes. *Nature*
364 *communications* 10, 1–11
- 365 Consortium, I. S. (2009). Common polygenic variation contributes to risk of schizophrenia that overlaps
366 with bipolar disorder. *Nature* 460, 748
- 367 Coram, M. A., Fang, H., Candille, S. I., Assimes, T. L., and Tang, H. (2017). Leveraging multi-ethnic
368 evidence for risk assessment of quantitative traits in minority populations. *The American Journal of*
369 *Human Genetics* 101, 218–226
- 370 De Candia, T. R., Lee, S. H., Yang, J., Browning, B. L., Gejman, P. V., Levinson, D. F., et al. (2013).
371 Additive genetic variation in schizophrenia risk is shared by populations of african and european descent.
372 *The American Journal of Human Genetics* 93, 463–470
- 373 De Los Campos, G., Gianola, D., and Allison, D. B. (2010). Predicting genetic predisposition in humans:
374 the promise of whole-genome markers. *Nature Reviews Genetics* 11, 880–886
- 375 Ge, T., Chen, C.-Y., Ni, Y., Feng, Y.-C. A., and Smoller, J. W. (2019). Polygenic prediction via bayesian
376 regression and continuous shrinkage priors. *Nature Communications* 10, 2041–1723
- 377 Graff, R. E., Cavazos, T. B., Thai, K. K., Kachuri, L., Rashkin, S. R., Hoffman, J. D., et al. (2021).
378 Cross-cancer evaluation of polygenic risk scores for 16 cancer types in two large cohorts. *Nature*
379 *communications* 12, 1–9
- 380 Huang, H., Ruan, Y., Feng, Y.-C. A., Chen, C.-Y., Lam, M., Sawa, A., et al. (2021). Improving polygenic
381 prediction in ancestrally diverse populations

- 382 Julià, A., López-Longo, F. J., Venegas, J. J. P., Bonàs-Guarch, S., Olivé, À., Andreu, J. L., et al. (2018).
383 Genome-wide association study meta-analysis identifies five new loci for systemic lupus erythematosus.
384 *Arthritis research & therapy* 20, 1–10
- 385 Khera, A. V., Chaffin, M., Aragam, K. G., Haas, M. E., Roselli, C., Choi, S. H., et al. (2018). Genome-wide
386 polygenic scores for common diseases identify individuals with risk equivalent to monogenic mutations.
387 *Nature genetics* 50, 1219–1224
- 388 Kim, M. S., Patel, K. P., Teng, A. K., Berens, A. J., and Lachance, J. (2018). Genetic disease risks can be
389 misestimated across global populations. *Genome biology* 19, 1–14
- 390 [Dataset] Li, S., Cai, T. T., and Li, H. (2020). Transfer learning for high-dimensional linear regression:
391 Prediction, estimation, and minimax optimality
- 392 Maier, R., Moser, G., Chen, G.-B., Ripke, S., Absher, D., Agartz, I., et al. (2015). Joint analysis of
393 psychiatric disorders increases accuracy of risk prediction for schizophrenia, bipolar disorder, and major
394 depressive disorder. *The American Journal of Human Genetics* 96, 283–294
- 395 Mak, T. S. H., Porsch, R. M., Choi, S. W., Zhou, X., and Sham, P. C. (2017). Polygenic scores via penalized
396 regression on summary statistics. *Genetic epidemiology* 41, 469–480
- 397 Márquez-Luna, C., Loh, P.-R., Consortium, S. A. T. . D. S., Consortium, S. T. . D., and Price, A. L. (2017).
398 Multiethnic polygenic risk scores improve risk prediction in diverse populations. *Genetic epidemiology*
399 41, 811–823
- 400 Morris, D. L., Sheng, Y., Zhang, Y., Wang, Y.-F., Zhu, Z., Tomblason, P., et al. (2016). Genome-wide
401 association meta-analysis in chinese and european individuals identifies ten new loci associated with
402 systemic lupus erythematosus. *Nature genetics* 48, 940–946
- 403 Moser, G., Lee, S. H., Hayes, B. J., Goddard, M. E., Wray, N. R., and Visscher, P. M. (2015). Simultaneous
404 discovery, estimation and prediction analysis of complex traits using a bayesian mixture model. *PLoS*
405 *genetics* 11, e1004969
- 406 Musliner, K. L., Mortensen, P. B., McGrath, J. J., Suppli, N. P., Hougaard, D. M., Bybjerg-Grauholm,
407 J., et al. (2019). Association of polygenic liabilities for major depression, bipolar disorder, and
408 schizophrenia with risk for depression in the danish population. *JAMA psychiatry* 76, 516–525
- 409 Peterson, R. E., Kuchenbaecker, K., Walters, R. K., Chen, C.-Y., Popejoy, A. B., Periyasamy, S., et al.
410 (2019). Genome-wide association studies in ancestrally diverse populations: opportunities, methods,
411 pitfalls, and recommendations. *Cell* 179, 589–603
- 412 Popejoy, A. B. and Fullerton, S. M. (2016). Genomics is failing on diversity. *Nature News* 538, 161
- 413 Privé, F., Aschard, H., Ziyatdinov, A., and Blum, M. G. (2018). Efficient analysis of large-scale genome-
414 wide data with two r packages: bigstatsr and bigsnpr. *Bioinformatics* 34, 2781–2787
- 415 Shi, H., Burch, K. S., Johnson, R., Freund, M. K., Kichaev, G., Mancuso, N., et al. (2020). Localizing
416 components of shared transethnic genetic architecture of complex traits from gwas summary data. *The*
417 *American Journal of Human Genetics* 106, 805–817
- 418 Shi, H., Gazal, S., Kanai, M., Koch, E. M., Schoech, A. P., Siewert, K. M., et al. (2021). Population-
419 specific causal disease effect sizes in functionally important regions impacted by selection. *Nature*
420 *communications* 12, 1–15
- 421 Shi, J., Park, J.-H., Duan, J., Berndt, S. T., Moy, W., Yu, K., et al. (2016). Winner's curse correction
422 and variable thresholding improve performance of polygenic risk modeling based on genome-wide
423 association study summary-level data. *PLoS genetics* 12, e1006493
- 424 Speed, D. and Balding, D. J. (2014). Multiblup: improved snp-based prediction for complex traits. *Genome*
425 *research* 24, 1550–1557

- 426 Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical*
427 *Society: Series B (Methodological)* 58, 267–288
- 428 Tukiainen, T., Pirinen, M., Sarin, A.-P., Ladenvall, C., Kettunen, J., Lehtimäki, T., et al. (2014).
429 Chromosome x-wide association study identifies loci for fasting insulin and height and evidence for
430 incomplete dosage compensation. *PLoS genetics* 10, e1004127
- 431 Turley, P., Walters, R. K., Maghzian, O., Okbay, A., Lee, J. J., Fontana, M. A., et al. (2018). Multi-trait
432 analysis of genome-wide association summary statistics using mtag. *Nature genetics* 50, 229–237
- 433 Vilhjálmsdóttir, B. J., Yang, J., Finucane, H. K., Gusev, A., Lindström, S., Ripke, S., et al. (2015). Modeling
434 linkage disequilibrium increases accuracy of polygenic risk scores. *The American journal of human*
435 *genetics* 97, 576–592
- 436 Wang, Y.-F., Zhang, Y., Lin, Z., Zhang, H., Wang, T.-Y., Cao, Y., et al. (2021). Identification of 38 novel
437 loci for systemic lupus erythematosus and genetic heterogeneity between ancestral groups. *Nature*
438 *communications* 12, 1–13
- 439 Weissbrod, O., Kanai, M., Shi, H., Gazal, S., Peyrot, W. J., Khera, A. V., et al. (2022). Leveraging fine-
440 mapping and multipopulation training data to improve cross-population polygenic risk scores. *Nature*
441 *Genetics* 54, 450–458

FIGURE CAPTIONS

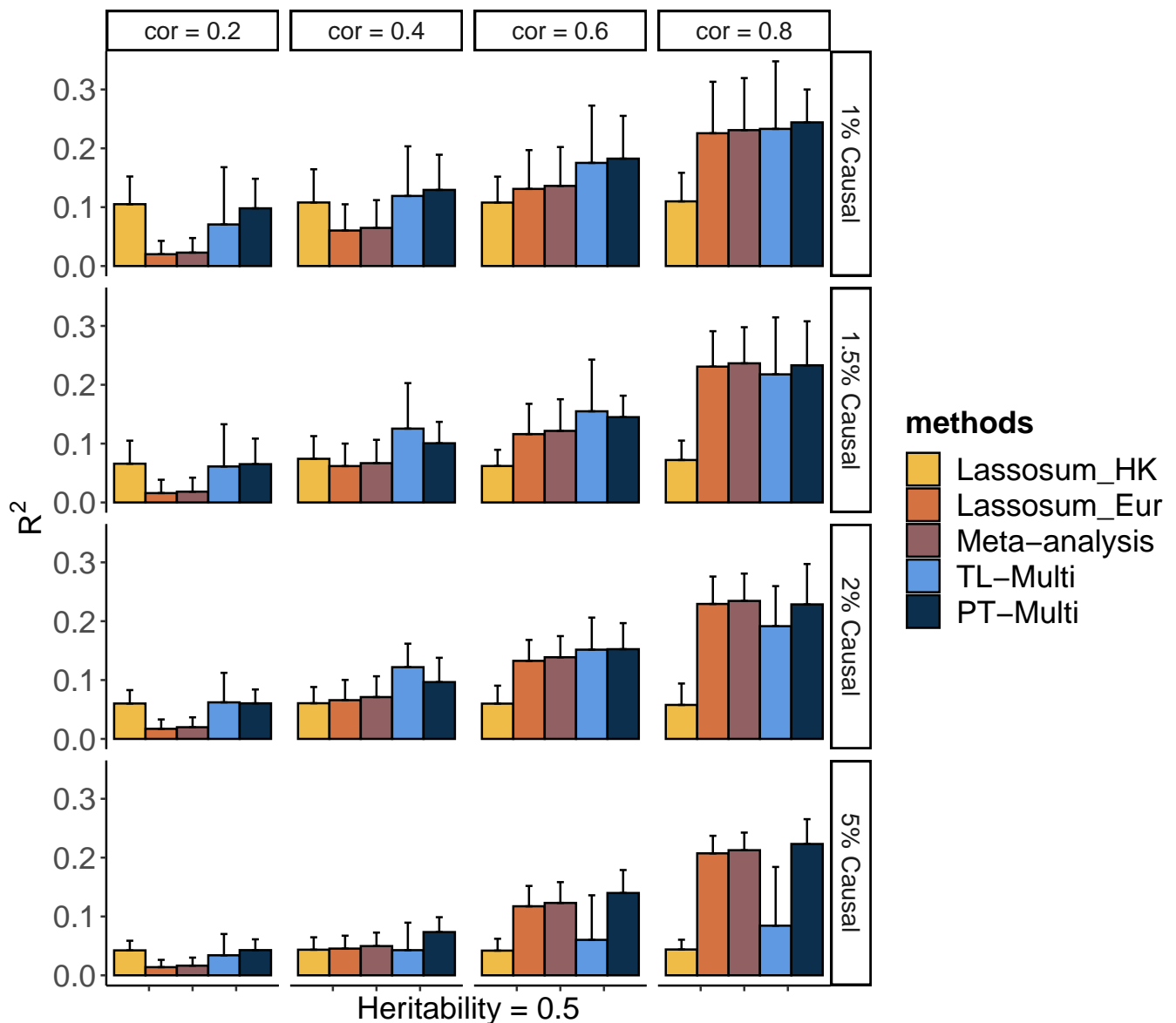


Figure 1. Prediction accuracy of Lassosum, meta-analysis, PT-Multi, and TL-Multi over 20 replications in simulations. Lassosum_HK is Lassosum for Hong Kong population, and Lassosum_Eur is Lassosum for European populations. Heritability was fixed at 0.5 and different genetic correlations (0.2, 0.4, 0.6, and 0.8) with different causal variant proportions (1%, 1.5%, 2%, and 5%) were generated. 50,000 European samples and 1,000 Hong Kong samples were simulated. The variants were generated from the common variants of the first 4 chromosomes (21,477 SNPs). The prediction accuracy was measured by R^2 between the simulated and true phenotypes. The error bar indicated the upper bound of 95% confidence interval over 20 replications.

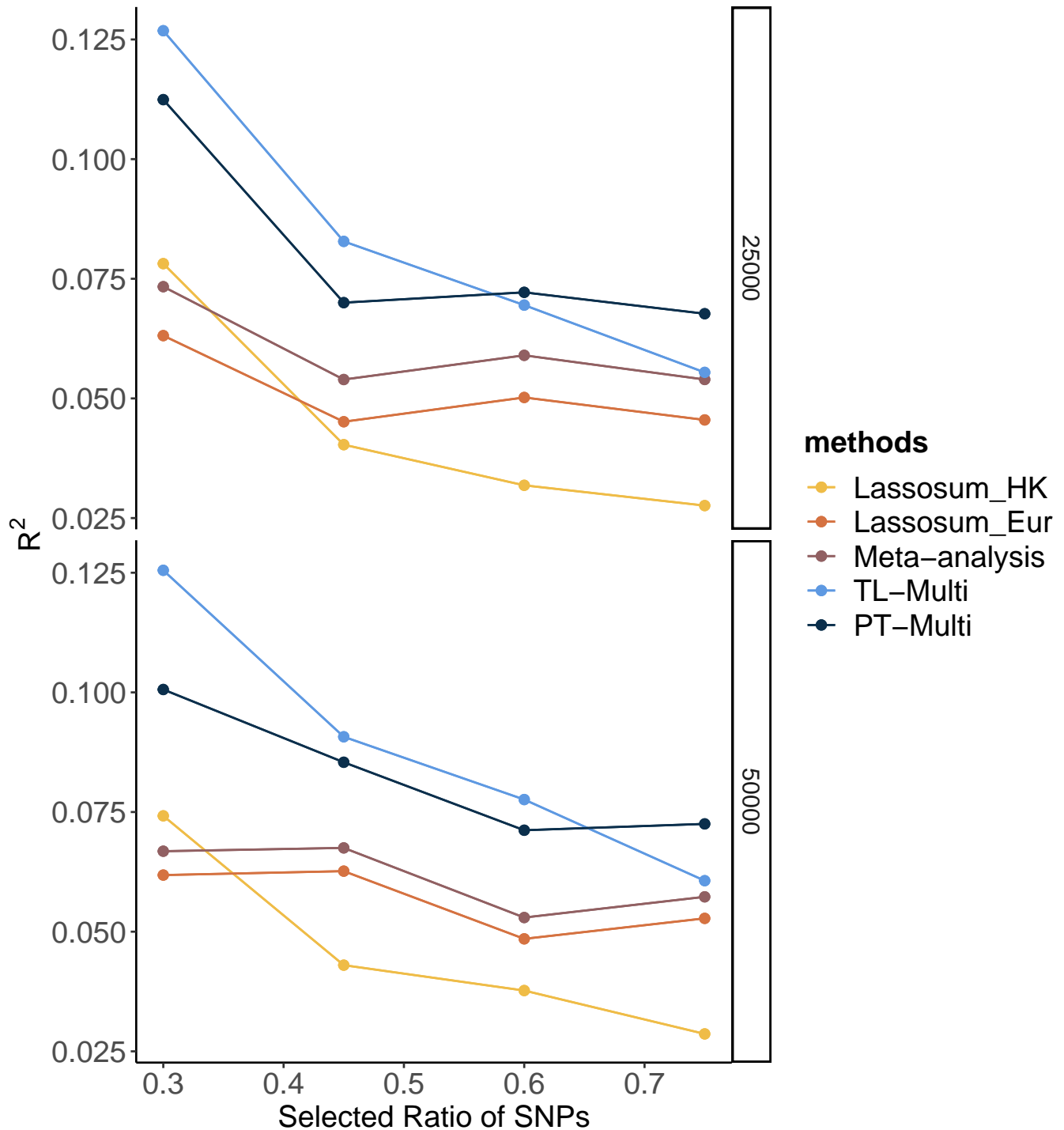


Figure 2. Prediction accuracy of Lassosum, meta-analysis, PT-Multi and TL-Multi over 20 replications in simulations. Selected ratio of SNPs is the ratio of the actual numbers of SNPs simulated to the total number of common SNPs (69,398). The actual numbers of SNPs simulated in the four scenarios are 21,477 (chromosomes 1-4), 32,151 (chromosomes 1-6), 39,682 (chromosomes 1-8), 49,909 (chromosomes 1-11) respectively. The average of R^2 are plotted. (A) The sample size of European population is 25,000, and the sample size of Hong Kong population is 1,000. (B) The sample size of European population is 50,000, and the sample size of Hong Kong population is 1,000.

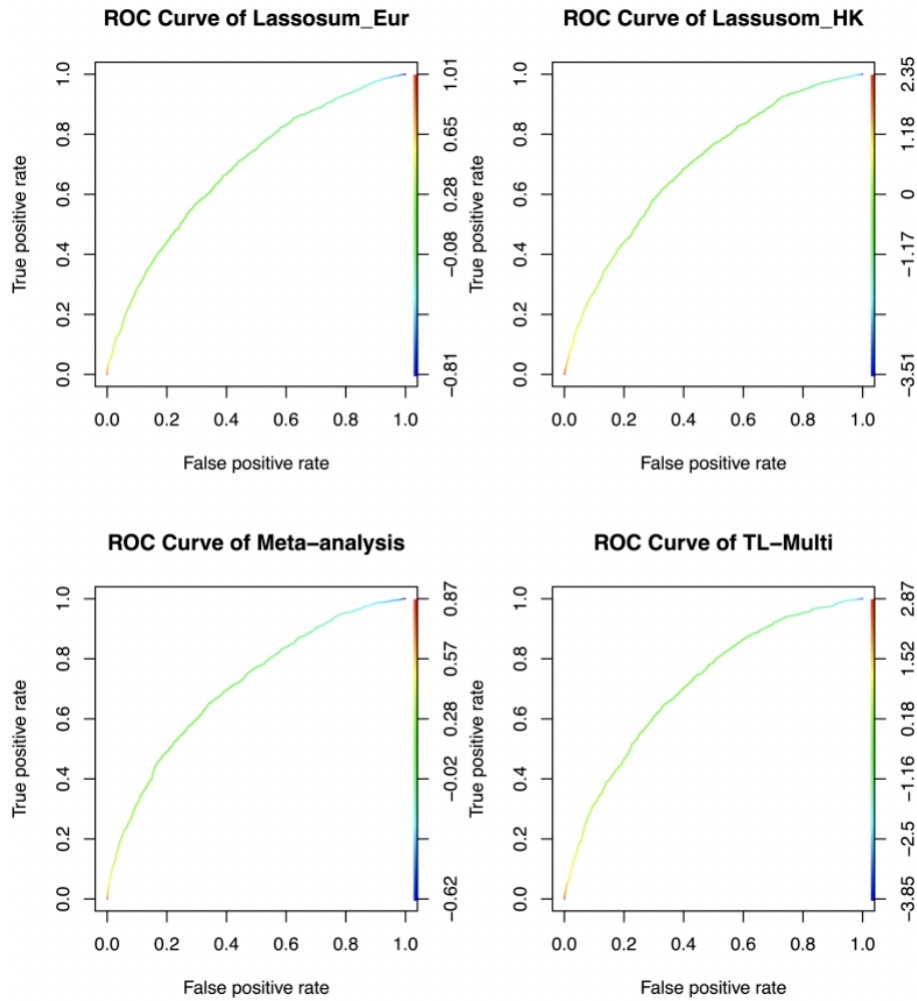


Figure 3. Receiver operating characteristic curve of Lassosum, meta-analysis, and TL-Multi in analysis of SLE study. Lassosom_HK is Lassosum for Hong Kong population, and Lassosom_Eur is Lassosum for European population. The corresponding AUC values with the optimal PRS of Lassosum for Hong Kong population and European population, meta-analysis, and TL-Multi are 0.6872, 0.6943, 0.7098 and 0.7131, respectively.

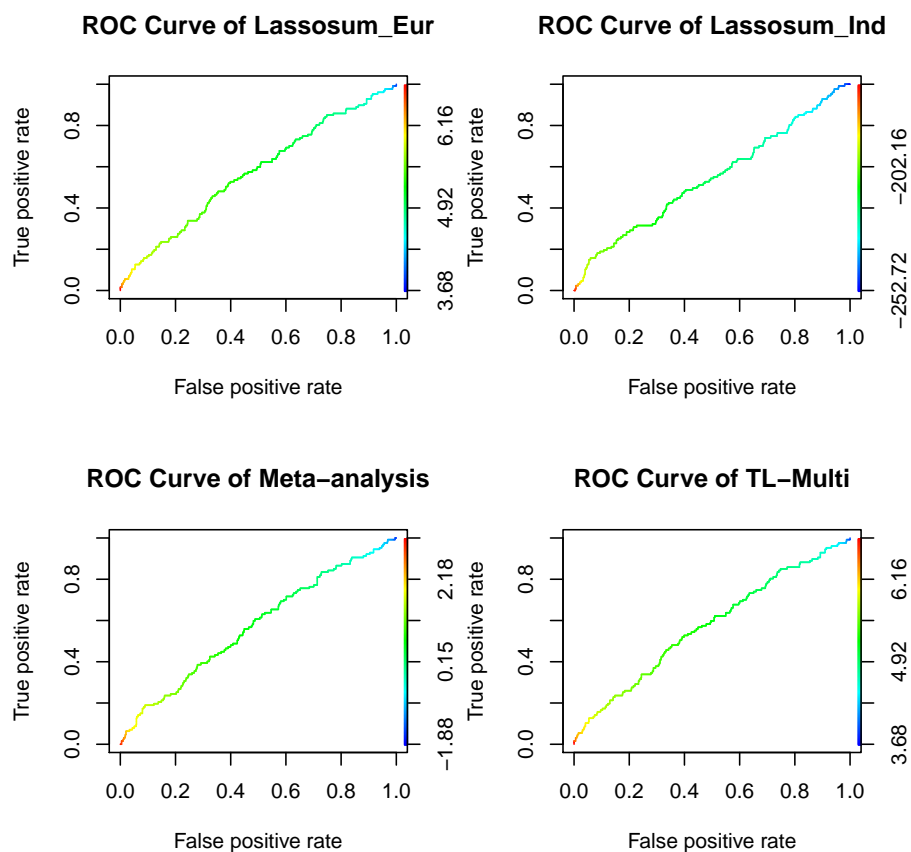


Figure 4. Receiver operating characteristic curve of Lassosum, meta-analysis, and TL-Multi in the analysis asthma study. Lassosum_Ind is Lassosum for Indian population, and Lassosum_Eur is Lassosum for European population. The corresponding AUC values with the optimal PRS of Lassosum for Indian population and European population, meta-analysis, and TL-Multi are 0.5657, 0.5441, 0.5705 and 0.5721, respectively.

Table 1. Case prevalence of 2%, 5%, and 10% for the top and bottom quantiles of the PRS distribution in analysis of SLE study with the target Indian population, generated by Lassosum, meta-analysis, and TL-Multi.

Prevalence	Bottom			Top		
	2%	5%	10%	10%	5%	2%
Lassosum_HK	0.0864	0.1133	0.1309	0.6963	0.7192	0.7407
Lassosum_Eur	0.1111	0.1281	0.1704	0.6938	0.7389	0.8148
Meta-Analysis	0.0864	0.0985	0.1309	0.7432	0.8030	0.8519
TL-Multi	0.0741	0.0985	0.1235	0.7160	0.7389	0.7901

Table 2. Case prevalence of 2%, 5%, and 10% for the top and bottom quantiles of the PRS distribution in analysis of asthma study with the target Indian population, generated by Lassosum, meta-analysis, and TL-Multi.

Prevalence	Bottom			Upper		
	2%	5%	10%	10%	5%	2%
Lassosum_Ind	0.1500	0.1800	0.1600	0.1900	0.1800	0.2000
Lassosum_Eur	0.1500	0.1200	0.1300	0.1600	0.1800	0.0000
Meta-Analysis	0.4000	0.2000	0.2100	0.1400	0.1600	0.100
TL-Multi	0.1000	0.1800	0.1700	0.2000	0.2400	0.2500