1    # Testing and controlling for horizontal pleiotropy with the probabilistic

2    # Mendelian randomization in transcriptome-wide association studies

3

4    Zhongshang Yuan[1,2], Huanhuan Zhu[2], Ping Zeng[3], Sheng Yang[2], Shiquan Sun[2], Can Yang[4], Jin

5    Liu[5], Xiang Zhou[2,6,*]

6

7    1. Department of Biostatistics, School of Public Health, Shandong University, Jinan, China.

8    2. Department of Biostatistics, University of Michigan, Ann Arbor, MI 48109, USA

9    3. Department of Epidemiology and Biostatistics, Xuzhou Medical University, Xuzhou, Jiangsu,

10    China

11    4. Department of Mathematics, Hong Kong University of Science and Technology, Hong Kong,

12    China

13    5. Centre for Quantitative Medicine, Program in Health Services and Systems Research, Duke-

14    NUS Medical School, Singapore 169857

15    6. Center for Statistical Genetics, University of Michigan, Ann Arbor, MI 48109, USA

16    *To whom correspondence should be addressed: xzhousph@umich.edu

17

18

19

## Abstract

Integrating association results from both genome-wide association studies (GWASs) and expression quantitative trait locus (eQTL) mapping studies has the potential to shed light on the molecular mechanisms underlying disease etiology. Several statistical methods have been recently developed to integrate GWASs with eQTL studies in the form of transcriptome-wide association studies (TWASs). These existing methods can all be viewed as a form of two sample Mendelian randomization (MR) analysis, which has been widely applied in various GWASs for inferring the causal relationship among complex traits. Unfortunately, most existing TWAS and MR methods make an unrealistic modeling assumption and assume that instrumental variables do not exhibit horizontal pleiotropic effects. However, horizontal pleiotropic effects have been recently discovered to be wide spread across complex traits, and, as we will show here, are also wide spread across gene expression traits. Therefore, not allowing for horizontal pleiotropic effects can be overly restrictive, and, as we will be show here, can lead to a substantial inflation of test statistics and subsequently false discoveries in TWAS applications. Here, we present a probabilistic MR method, which we refer to as PMR-Egger, for testing and controlling for horizontal pleiotropic effects in TWAS applications. PMR-Egger relies on an MR likelihood framework that unifies many existing TWAS and MR methods, accommodates multiple correlated instruments, tests the causal effect of gene on trait in the presence of horizontal pleiotropy, and, with a newly developed parameter expansion version of the expectation maximization algorithm, is scalable to hundreds of thousands of individuals. With extensive simulations, we show that PMR-Egger provides calibrated type I error control for causal effect testing in the presence of horizontal pleiotropic effects, is reasonably robust for various types of horizontal pleiotropic effect mis-specifications, is more powerful than existing MR approaches,

43    and, as a by-product, can directly test for horizontal pleiotropy. We illustrate the benefits of

44    PMR-Egger in applications to 39 diseases and complex traits obtained from three GWASs

45    including the UK Biobank. In these applications, we show how PMR-Egger can lead to new

46    biological discoveries through integrative analysis.

47

48    **Introduction**

49    Genome-wide association studies (GWASs) have identified many SNPs associated with common

50    diseases or disease related traits. Parallel expression quantitative trait loci (eQTL) mapping

51    studies have also identified many cis-acting SNPs associated with the expression level of nearby

52    genes. Integrating the existing association results from both GWASs and eQTL mapping studies

53    has the potential to shed light on the molecular mechanisms underlying disease etiology. Several

54    statistical methods have been recently proposed to integrate GWASs with eQTL mapping studies.

55    For example prediXcan[1] proposes to perform a weighted SNP set test in GWAS by inferring

56    SNP weights from eQTL studies. TWAS[2] proposes to infer the association between gene

57    expression and disease trait by leveraging the shared common set of cis-SNPs. SMR[3] or GSMR[4]

58    directly tests the causal association between gene expression and disease trait under a Mendelian

59    randomization (MR) framework through selecting a single instrument or multiple independent

60    instruments. While each of these integrative methods was originally proposed to solve a different

61    problem, all of them can be viewed as a two-sample MR method with different modeling

62    assumptions. Because of their relationship to MR, these methods effectively attempt to identify

63    genes causally associated with diseases or complex traits in the context of transcriptome-wide

64    association studies (TWAS).

65    MR analysis is a form of instrumental variable analysis that was originally developed in the

66    field of causal inference[5]. MR aims to determine the causal relationship between an exposure

67    variable (e.g. gene expression) and an outcome variable (e.g. complex trait) in observational

68    studies. MR treats SNPs as instrumental variables for the exposure variable of interest and uses

69    these SNP instruments to estimate and test the causal effect of the exposure variable on the

70    outcome variable. MR methods have been widely applied to investigate the causal relationship

71    among various complex traits[6-9], and, through a two-sample design, can be easily adapted to

72    settings where the exposure and outcome are measured on two different sets of individuals[10,11].

73    However, MR analysis for TWAS is not straightforward and requires the development of new

74    methods that can accommodate two important features of TWAS analysis.

75    First, both GWASs and eQTL mapping studies collect SNPs that are in high linkage

76    disequilibrium (LD) with each other. Traditional MR methods, such as the random effects

77    version or the fixed effect version of the inverse variance weighted regression[12], MR-Egger[13],

78    median-based regression[14], SMR[3], or GSMR[4], can only make use of a single SNP instrument or

79    multiple independent SNP instruments. Handling only independent SNPs is restrictive, as most

80    exposure variables/molecular traits are polygenic/omni-genic and are influenced by multiple

81    SNPs that are in potential LD with each other. As a result, incorporating multiple correlated

82    SNPs can often help explain a greater proportion of variance in the exposure variable than using

83    independent SNPs, and thus can help increase power and improve estimation accuracy of MR

84    analysis[5,15-17]. Due to the benefits of using multiple correlated instruments, most TWAS methods

85    (e.g. PrediXcan[1], TWAS[2], CoMM[18], DPR[19], TIGAR[20]) rely on polygenic modeling priors to

86    incorporate all cis-SNPs that are in high LD for TWAS applications. (Certainly, while the prior

87    used in PrediXcan is polygenic, the parameter estimates obtained from PrediXcan is sparse as it

88    uses posterior mode instead of posterior mean.) By incorporating all cis-SNPs, as we will show

89    below, these methods can lead to substantial power improvement over standard MR approaches

90    that use only a few independent SNPs. Unfortunately, many TWAS methods rely on a two-stage

91    MR inference procedure: they estimate SNP effect sizes in the exposure study and plug in these

92    estimates to the outcome study for causal effect inference. The two-stage inference procedure in

93    MR fails to account for the uncertainty in parameter estimates in the exposure study and can

94    often lead to biased causal effect estimates and power loss, especially in the presence of weak

95    instruments[5,16]. Indeed, similar to what have been observed in the MR filed, our previous study

96    also suggests that the likelihood based inference can substantially improve power for TWAS[18].

97    Therefore, it is important to incorporate multiple correlated instruments in a likelihood inference

98    framework for MR analysis in TWAS.

99       Second, perhaps more importantly, SNP instruments often exhibit pervasive horizontal

100    pleiotropic effects[21]. Horizontal pleiotropy occurs when a genetic variant affects the outcome

101    variable through pathways other than or in addition to the exposure variable[22]. Horizontal

102    pleiotropy is in contrast to the vertical pleiotropy, which characterizes instrument effects on the

103    outcome variable through the path of the exposure. Horizontal pleiotropy is widely distributed

104    across the genome, affects a wide spectrum of complex traits, and can be driven by LD and

105    extreme polygenicity of traits[21][23]. Despite its wide prevalence, however, only a limited number

106    of MR methods have been developed to test and control for horizontal pleiotropy; even fewer are

107    applicable for TWAS applications. For example, some existing methods (e.g. MR-PRESSO[21])

108    test for horizontal pleiotropic effects without directly controlling for them. Some methods (e.g.

109    CaMMEL[24]) control for horizontal pleiotropic effects without directly testing them[25,26]. Some

110    methods (e.g. Egger regression[13,27], GLIDE[28], GSMR[4], MR-median method[14], profile score

111    approach[29], MRMix[30] and Bayesian MR[31,32]) test and control for horizontal pleiotropic effects,

112    but can only accommodate independent instruments. As far as we are aware, there is only one

113    two-sample MR method currently developed for testing and controlling for pleiotropic effects in

114    the presence of correlated instruments: LDA MR-Egger[33]. Unfortunately, as we will show below,

115    LDA MR-Egger cannot handle realistic LD pattern among cis-SNPs for TWAS applications.

116    Here, we develop a generative two-sample MR method in a likelihood framework, which we

117    refer to as the probabilistic two-sample Mendelian randomization (PMR), to perform MR

118    analysis using multiple correlated instruments for TWAS applications. We illustrate how the

119    PMR framework can facilitate the understanding of many existing MR approaches as well as

120    many existing integrative analysis approaches. Within the PMR framework, we focus on a

121    particular horizontal pleiotropy effect modeling assumption based on the burden test assumption

122    commonly used for rare variant test. This particular horizontal pleiotropy effect, as we will show

123    later, effectively generalizes the Egger regression assumption commonly used for MR analysis to

124    correlated instruments. Our method allows us to test the causal effect in the presence of

125    horizontal pleiotropy, and, with a parameter expansion version of the expectation maximization

126    algorithm (PX-EM), is scalable to hundreds of thousands of individuals. We refer to our method

127    as PMR-Egger. With simulations, we show that PMR-Egger provides calibrated type I error for

128    causal effect testing in the presence of horizontal pleiotropic effects, is more powerful than

129    existing MR approaches, and, as a by-product, can directly test for horizontal pleiotropy. We

130    apply our method to perform TWAS for 39 diseases and complex traits obtained from three

131    GWASs with sample size ranging from 4,686 to 337,198.

132

## Methods

**PMR-Egger Overview**

We consider a probabilistic Mendelian randomization framework for performing two-sample Mendelian randomization analysis with correlated SNP instruments. Two-sample Mendelian randomization analysis aims to estimate and test for the causal effect of an exposure on an outcome in the setting where the exposure and outcome variables are measured in two separate studies with no sample overlap. In the TWAS applications we consider here, the exposure variable is gene expression level that is measured in a gene expression study, while the outcome variable is a quantitative trait or a dichotomous disease status that is measured in a GWAS. Often times, the gene expression study and GWAS are performed on two separate samples. While we mostly focus on TWAS applications in the present study, we note that the two-sample Mendelian randomization is also commonly performed in settings where both the exposure and outcome variables are complex traits that are measured in two separate GWASs. An illustrative diagram of MR analysis is displayed in Supplementary Fig. 1.

We denote $x$ as an $n_1$-vector of exposure variable (i.e. gene expression measurements) that is measured on $n_1$ individuals in the gene expression study and denote $\mathbf{Z}_x$ as an $n_1$ by $p$ matrix of genotypes for $p$ instruments (i.e. cis-SNPs) in the same study. Note that, unlike standard MR methods that select independent instruments, we follow existing TWAS approaches and use all cis-SNPs that are in LD as instruments. We denote $\mathbf{y}$ as an $n_2$-vector of outcome variable (i.e. trait) that is measured on $n_2$ individuals in the GWAS and denote $\mathbf{Z}_y$ as an $n_2$ by $p$ matrix of genotypes for the same $p$ instruments there. We consider three linear regressions to model the two studies separately

$$x = \mathbf{1}_{n_1}\mu_x + \mathbf{Z}_x\boldsymbol{\beta} + \boldsymbol{\varepsilon}_x \quad (1)$$

8

156
$$\tilde{x} = \mathbf{1}_{n_2}\mu_x + \mathbf{Z}_y\boldsymbol{\beta} + \boldsymbol{\varepsilon}_{\tilde{x}} \quad (2)$$

157
$$y = \mathbf{1}_{n_2}\mu_y + \tilde{x}\alpha + \mathbf{Z}_y\boldsymbol{\gamma} + \boldsymbol{\epsilon} \quad (3)$$

158 where the equation (1) is for the gene expression data and the equations (2)-(3) are for the

159 GWAS data. Here, $\mu_x$ and $\mu_y$ are the intercepts; $\tilde{x}$ is an unobserved $n_2$-vector of exposure

160 variable on the $n_2$ individuals in the GWAS; $\boldsymbol{\beta}$ is a $p$-vector of instrumental effect sizes on the

161 exposure variable; $\alpha$ is a scalar that represents the causal effect of the exposure variable on the

162 outcome variable; $\boldsymbol{\gamma}$ is a $p$-vector of horizontal pleiotropic effect sizes of $p$ instruments on the

163 outcome variable; $\boldsymbol{\varepsilon}_x$ is an $n_1$-vector of residual error with each element independently and

164 identically distributed from a normal distribution $N(0, \sigma_x^2)$; $\boldsymbol{\varepsilon}_{\tilde{x}}$ is an $n_2$-vector of residual error

165 with each element independently and identically distributed from the same normal distribution

166 $N(0, \sigma_x^2)$; and $\boldsymbol{\epsilon}$ is an $n_2$-vector of residual error with each element independently and identically

167 distributed from a normal distribution $N(0, \sigma_y^2)$. We note that while the above three equations are

168 specified based on two separate studies, they are joined together with the common parameter $\boldsymbol{\beta}$

169 and the unobserved gene expression measurements $\tilde{x}$. Equations (2)-(3) can also be combined

170 into

171
$$y = \mathbf{1}_{n_2}\mu_y + \mathbf{Z}_y\boldsymbol{\beta}\alpha + \mathbf{Z}_y\boldsymbol{\gamma} + \boldsymbol{\varepsilon}_y \quad (4)$$

172 where $\boldsymbol{\varepsilon}_y = \boldsymbol{\varepsilon}_{\tilde{x}}\alpha + \boldsymbol{\epsilon}$.

173 Our key parameter of interest in the above joint model is the causal effect $\alpha$. The causal

174 interpretation of $\alpha$ requires two assumptions of MR analysis to hold: (i) instruments are

175 associated with the exposure; (ii) instruments are not associated with any other confounders that

176 may be associated with both exposure and outcome. Note that our model no longer requires the

177 general exclusion restriction condition of traditional MR (i.e. instruments only influence the

178 outcome through the path of exposure), as we make explicit modeling assumptions on the

9

179    horizontal pleiotropy effects $\boldsymbol{\gamma}$. Certainly, PMR-Egger still need to satisfy the InSIDE

180    assumption that the instrument-exposure effects and instrument-outcome effects are independent

181    of each other, which is sometimes refered to as the weak exclusion restriction condition[13]. In our

182    model, we derive the causal interpretation and identification of α under the decision-theoretic

183    framework of causal inference[31,34-36] (details in Supplementary Note). Because the causal effect

184    interpretation of α depends on MR assumptions as well as other explicit modeling assumptions,

185    many of which are not easily testable in practice, MR analysis in observational studies likely

186    provides weaker causality evidence than randomized clinical trials. Therefore, while we follow

187    standard MR analysis and use the term "causal effect" through the text, we only intend to use this

188    term to emphasize the fact that α estimate from an MR analysis is more trustworthy than the

189    effect size estimate in a standard linear regression of $\boldsymbol{y}$ on $\tilde{x}$.

190       Because $p$ is often larger than $n_1$, we will need to make additional modeling assumptions on $\boldsymbol{\beta}$

191    to make the model identifiable. In addition, the two instrumental effect terms defined in equation

192    (4), the vertical pleiotropic effect $\mathbf{Z}_y\boldsymbol{\beta}\alpha$ and the horizontal pleiotropic effect $\mathbf{Z}_y\boldsymbol{\gamma}$, are also not

193    identifiable from each other, unless we make additional modeling assumptions on $\boldsymbol{\gamma}$. Here, we

194    follow standard polygenic model and assume that all elements in $\boldsymbol{\beta}$ are non-zero and that each

195    follows a normal distribution $N(0, \sigma_\beta^2)$. In addition, we follow the burden test assumption

196    commonly used for rare variant test and assume that equal horizontal pleiotropic effects across

197    SNPs $\boldsymbol{\gamma}_j = \gamma$ for $j = 1, \dots p$. With the burden test assumption on the horizontal pleiotropic

198    effects $\gamma$, our model becomes a generalization of the commonly used MR-Egger regression

199    model. In the special case where instruments are independent and treated as fixed effects and

200    where a two-stage estimation procedure is used for inference, our model reduces to MR-Egger.

201    However, our method can handle general cases where MR-Egger does not apply to. In particular,

202 unlike MR-Egger, our method can handle multiple correlated instruments and perform inference

203 in a likelihood framework.

204  In the above model, we are interested in estimating the causal effect $\alpha$ and testing the null

205 hypothesis $H_0: \alpha = 0$ in the presence of horizontal pleiotropy effects $\boldsymbol{\gamma}$. In addition, we are

206 interested in estimating the horizontal pleiotropic effect size $\gamma$ and testing the null hypothesis

207 $H_0: \gamma = 0$. We accomplish both tasks through the maximum likelihood inference framework. In

208 particular, we develop an expectation maximization (EM) algorithm for parameter inference by

209 maximizing the joint likelihood defined based on equations (1) and (4) (details in the

210 Supplementary Note). The EM algorithm allows us to obtain the maximum likelihood of the

211 joint model, together with maximum likelihood estimates for both $\alpha$ and $\gamma$. In addition, we apply

212 the EM algorithm to two reduced models, one without $\alpha$ and the other without $\gamma$, to obtain the

213 corresponding maximum likelihoods. Afterwards, we perform likelihood ratio tests for either

214 $H_0: \alpha = 0$ or $H_0: \gamma = 0$, by contrasting the maximum likelihood obtained from the joint model to

215 that obtained from each of the two reduced models, respectively. We refer to the above inference

216 procedure as probabilistic, as we place estimation and testing into a maximum likelihood

217 framework. Our inference procedure is in contrast to the commonly used two-stage estimation

218 procedure (as used in, for example, Egger regression[13,27], PrediXcan[1] and TWAS[2]), which

219 estimates $\widehat{\boldsymbol{\beta}}$ from equation (1) first and then plug in the estimates into equation (4) for inference.

220 The previous two-stage estimation procedure fails to properly account for the estimation

221 uncertainty in $\widehat{\boldsymbol{\beta}}$ and is known to lose power compared to a formal likelihood inference

222 procedure[5,16,18].

223  We refer to our model and algorithm together as the two-sample probabilistic Mendelian

224 randomization with Egger regression (PMR-Egger). As explained above, we use "probabilistic"

225     to refer to both the data generative model and the maximum likelihood inference procedure. We

226     use "Egger" to refer to the horizontal pleiotropic assumption on $\boldsymbol{\gamma}$ that effectively generalizes the

227     Egger-regression assumption to correlated instruments. We also note that the joint generative

228     Mendelian randomization model defined in equations (1) and (4) is a useful conceptual

229     framework that unifies many existing MR methods. In particular, almost all existing MR

230     methods are built upon the joint model, but with different modeling assumptions on $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$, and

231     with different inference procedures (Table 1). Compared with these existing MR approaches,

232     PMR-Egger is capable of modeling multiple correlated instruments, effectively controls for

233     horizontal pleiotropy, and places inference into a likelihood framework.

234     **Simulations**

235     We performed simulations to assess the performance of PMR-Egger and compare it with existing

236     approaches. To do so, we first obtained 556 cis-SNPs for the gene *BACE1* on chromosome 11

237     from the GEUVADIS data[37] (data processing details in the next section) and simulated gene

238     expression values. We used the gene *BACE1* because the number of cis-SNPs in this gene

239     represents the median of all genes. With the scaled genotype data $\boldsymbol{Z}_x$, we simulated SNP effect

240     sizes $\boldsymbol{\beta}$ from a normal distribution $N(0, PVE_{zx}/556)$, where the scalar $PVE_{zx}$ represents the

241     proportion of gene expression variance explained by genetic effects. We summed the genetic

242     effects across all cis-SNPs as $\boldsymbol{Z}_x\boldsymbol{\beta}$. In addition, we simulated residual errors $\boldsymbol{\varepsilon}_x$ from a normal

243     distribution $N(0, 1 - PVE_{zx})$. We then summed the genetic effects and residual errors to yield

244     the simulated gene expression level.

245       Next, we obtained genotypes for the same 556 SNPs from 2,000 randomly selected control

246     individuals in the Kaiser Permanente/UCSF Genetic Epidemiology Research Study on Adult

247     Health and Aging (GERA)[38,39] and simulated a quantitative trait. Here, we directly used $\boldsymbol{\beta}$ from

248 the gene expression data, which, when paired with the causal effect $\alpha$, yielded the vertical

249 pleiotropic effects $\alpha\boldsymbol{\beta}$. We set $\alpha = \sqrt{PVE_{zy}/PVE_{zx}}$, and we simulated residual errors $\boldsymbol{\varepsilon}_y$ from a

250 normal distribution $N(0, 1 - PVE_{zy})$. Here, the scalar parameter $PVE_{zy}$ represents the proportion

251 of phenotypic variance explained by vertical pleiotropic effects in the absence of horizontal

252 pleiotropic effects. Afterwards, we simulated horizontal pleiotropic effects $\boldsymbol{\gamma}$ for these SNPs

253 (more details below). We summed the horizontal pleiotropic effects, vertical pleiotropic effects

254 and residual errors to yield the simulated trait.

255     In the simulations, we first examined a baseline simulation setting where we set $PVE_{zx} =$

256 10%, $PVE_{zy} = 0$, with all $\gamma_j = 0$. On top of the baseline setting, we varied one parameter at a

257 time to examine the influence of various parameters. For $PVE_{zx}$, we set it to be either 1%, 5% or

258 10%, close to the median gene expression heritability estimates across genes[40,41]. For $\boldsymbol{\beta}$, we

259 examined alternative SNP effect size distributions that deviate from the polygenic assumption in

260 the baseline setting. Specifically, we randomly selected either 1 SNP, 1%, 10% or 100% of the

261 SNPs to have non-zero effect, while simulated their effects from a normal distribution to explain

262 a fixed $PVE_{zx}$. For $PVE_{zy}$, we varied its value to be either 0% (for null simulations), 0.2%, 0.4%

263 or 0.6% (for power simulations). For the horizontal pleiotropy effects $\boldsymbol{\gamma}$, we randomly assigned a

264 fixed proportion of $\gamma_j$ to be non-zero (proportion equals 10%, 30%, 50%, or 100%). Afterwards,

265 we set the absolute value of non-zero $\gamma_j$ to be the same value of $\gamma$. As a sensitivity analysis, we

266 also randomly assigned some of their signs to be positive and some of their signs to be negative,

267 with the ratio of positive effects to negative effects being either 1:9, 3:7, or 5:5. Here, we set $\gamma$ to

268 be $1 \times 10^{-4}, 5 \times 10^{-4}, 1 \times 10^{-3}$ or $2 \times 10^{-3}$, which corresponds to the 50%, 70%, 90%, 95%

269 quantiles of horizontal pleiotropic effect estimates across all genes and all traits in the WTCCC

270 data (more details below), respectively. For null simulations and type I error control examination,

13

271   we performed 10,000 simulation replicates for each simulation scenario described above. For

272   power calculation, for each scenario, we performed 1,000 alternative simulations together with

273   9,000 null simulations and calculated power based on false discovery rate (FDR).

274   While we applied PMR-Egger to analyze individual-level data from all simulations, we also

275   applied PMR-Egger to analyze summary statistics in a subset of simulations to validate the

276   implementation of the summary statistics based PMR-Egger algorithm. These results are

277   presented in the Discussion section. Here, we considered the simulation settings with a fixed

278   sample size ($n_1 = 465, n_2 = 2,000$), different causal effect sizes ($PVE_{zy} = 0$ or 0.6%) and

279   different pleiotropy effect sizes ($\gamma = 0$ or 0.0005). In the analysis, we calculated the LD matrix

280   in the eQTL data using the observed individual-level genotypes in the eQTL study. We

281   calculated the LD matrix in the GWAS data from a reference panel. The reference panel is

282   constructed in three different ways, by using individual-level genotypes from either all

283   individuals in the GWAS (n=2,000), 10% of randomly selected individuals from the GWAS

284   (n=200), or the individuals with European ancestry from the 1,000 Genomes project (n=503).

285   Besides the single gene-based simulations, we also conducted cross-gene simulations.

286   Specifically, we randomly selected 10,000 genes from GEUVADIS. We extracted cis-SNPs for

287   these 10,000 genes, obtaining a median of 576 cis-SNPs per gene (min=11; max=7,409). For

288   each gene in turn, we used its cis-SNPs to simulate its gene expression level as described above.

289   Afterwards, we applied different methods to analyze simulated data. The cross-gene based

290   simulations reflect the varying LD pattern and the varying number of cis-SNPs across genes that

291   we observe in real data, and thus are likely to be realistic than the single gene-based simulations.

292   We performed cross-gene simulations under all simulation settings described above, including

14

293   settings with varying gene expression heritability, varying genetic architectures underlying gene

294   expression, as well as varying causal and horizontal pleiotropy effects.

**Real Data Applications**

296   We applied our method to perform TWAS by integrating gene expression data with several

297   GWASs. Specifically, we obtained GEUVADIS data[37] as the gene expression data and examined

298   39 phenotypes from three GWASs. The three GWASs include the Wellcome trust case control

299   study (WTCCC)[42], the Kaiser Permanente/UCSF Genetic Epidemiology Research Study on

300   Adult Health and Aging (GERA)[38,39], and the UK Biobank[43].

301   The GEUVADIS data[37] contains gene expression measurements for 465 individuals collected

302   from five different populations that include CEPH (CEU), Finns (FIN), British (GBR), Toscani

303   (TSI) and Yoruba (YRI). In the expression data, we only focused on protein coding genes and

304   lincRNAs that are annotated in GENCODE (release 12)[44,45]. Among these genes, we removed

305   lowly expressed genes that have zero counts in at least half of the individuals to obtain a final set

306   of 15,810 genes. We performed PEER normalization to remove confounding effects and

307   unwanted variations following previous studies[19,46]. Afterwards, following[19], to remove

308   remaining population stratification, we quantile normalized the gene expression measurements

309   across individuals in each population to a standard normal distribution, and then further quantile

310   normalized the gene expression measurements to a standard normal distribution across

311   individuals from all five populations. Besides expression data, all individuals in GEUVADIS

312   also have their genotypes sequenced in the 1,000 Genomes Project. We obtained genotype data

313   from the 1,000 Genomes Project phase 3. We filtered out SNPs that have a Hardy-Weinberg

314   equilibrium (HWE) p-value $< 10^{-4}$, a genotype call rate <95%, or a minor allele frequency (MAF)

315   <0.01. We retained a total of 7,072,917 SNPs for analysis.

15

316    The WTCCC data consists of about 14,000 cases from seven common diseases and 2,938

317    shared controls[42]. The diseases include type 1 diabetes (T1D; $n$=1,963), Crohn's disease (CD;

318    $n$=1,748), rheumatoid arthritis (RA; $n$=1,861), bipolar disorder (BD; $n$=1,868), type 2 diabetes

319    (T2D; $n$=1,924), coronary artery disease (CAD; $n$=1,926), and hypertension (HT; $n$=1,952). We

320    obtained quality controlled genotypes from WTCCC and initially imputed missing genotypes

321    using BIMBAM[47] to arrive at a total of 458,868 SNPs shared across all individuals. Afterwards,

322    we further imputed SNPs using the 1,000 Genomes as the reference panel using SHAPEIT and

323    IMPUTE2[48]. We filtered out SNPs that have an HWE p-value $< 10^{-4}$, a genotype call rate <95%,

324    or an MAF<0.01 to obtain a total of 2,793,818 imputed SNPs. For each trait in turn, we first

325    regressed the phenotype on the top 10 genotype principal components (PCs) and obtained

326    phenotype residuals. We then scaled the phenotype residuals to have a mean of zero and standard

327    deviation of one and used these phenotype residuals for TWAS analysis. In addition to the main

328    analysis that uses phenotype residuals, we also performed parallel analysis with PMR-Egger

329    where we used the original phenotype as the outcome variable and the top 10 genotype PCs as

330    covariates.

331    The GERA study consists of 61,953 individuals and 675,367genotyped SNPs. We filtered out

332    SNPs that had a genotype calling rate below 0.95, MAF<0.01, or HWE p value$<10^{-4}$ to yield a

333    total of 487,609 SNPs. We phased genotypes using SHAPEIT[49] and imputed SNPs based on the

334    Haplotype Reference Consortium (HRC version r1.1) reference panel[50] on the Michigan

335    Imputation Server using Minimac3[51]. Afterwards, we further filtered out SNPs that have a HWE

336    p-value $< 10^{-4}$, a genotype call rate <95%, an MAF<0.01, or an imputation score<0.30 to arrive

337    at a total of 8,385,867 SNPs that are shared across 61,953 individuals. We examined 22 diseases

338    in GERA that include Asthma (number of cases n=10,101), Allergic Rhinitis (n=15,193),

16

339    Cardiovascular Disease (CARD, n=16,431), Cancers (n=18,714), Depressive Disorder (n=7,900),

340    Dermatophytosis (n=8,443), Type 2 Diabetes (T2D, n=7,638), Dyslipidemia (n=33,071),

341    Hypertension (HT, n=31,044), Hemorrhoids (n=9,922), Abdominal Hernia (n=6,876), Insomnia

342    (n=4,357), Iron Deficiency (n=2,706), Irritable Bowel Syndrome (n=3,367), Macular

343    Degeneration (n=4,031), Osteoarthritis (n=22,062), Osteoporosis (n=5,909), Peripheral Vascular

344    Disease (PVD, n=4,718), Peptic Ulcer (n=1,007), Psychiatric disorders (n=9408), Stress

345    Disorders (n=4,706), and Varicose Veins (n=2,714). For each trait in turn, we first regressed the

346    phenotype on the top 10 genotype principal components (PCs) and obtained phenotype residuals.

347    We then scaled the phenotype residuals to have a mean of zero and standard deviation of one and

348    used these phenotype residuals for TWAS analysis. In addition to the main analysis that uses

349    phenotype residuals, we also performed parallel analysis with PMR-Egger where we used the

350    original phenotype as the outcome and the top 10 genotype PCs as covariates.

351      The UK Biobank data consists of 487,409 individuals and 92,693,895 imputed SNPs[43]. We

352    followed        the        same        sample        QC        procedure        in        Neale        lab

353    (https://github.com/Nealelab/UK_Biobank_GWAS/tree/master/imputed-v2-gwas) to retain a

354    total of 337,198 individuals of European ancestry. We filtered out SNPs with an HWE p-value <

355    $10^{-7}$, a genotype call rate <95%, or an MAF<0.001 to obtain a total of 13,876,958 SNPs. We

356    selected 10 UK Biobank quantitative traits that have a phenotyping rate > 80%, a SNP

357    heritability > 0.2 and a low correlation among them following a previous study[52]. The 10 traits

358    include Height ($h^2 = 0.579;$), Platelet count ($h^2 = 0.404$), Bone mineral density ($h^2 = 0.401$),

359    Red blood cell count ($h^2 = 0.324$), FEV1-FVC ratio ($h^2 = 0.313$), Body mass index (BMI,

360    $h^2 = 0.308$), RBC distribution width ($h^2 = 0.288$), Eosinophils count ($h^2 = 0.277$), Forced

361    vital capacity ($h^2 = 0.277$), White blood cell count ($h^2 = 0.272$). For each trait in turn, we

17

362  regressed the resulting standardized phenotypes on sex and top 10 genotype principal

363  components (PCs) to obtain the residuals, standardized the residuals to have a mean of zero and a

364  standard deviation of one, and finally used these scaled residuals to conduct TWAS analysis. We

365  also performed parallel analysis with PMR-Egger by including the top 10 genotype PCs as

366  covariates.

367  We combined the GEUVADIS data with each of the three GWASs for TWAS analysis. To do

368  so, in the GEUVADIS data, for each gene in turn, we extracted cis-SNPs that are within either

369  100 kb upstream of the transcription start site (TSS) or 100 kb downstream of the transcription

370  end site (TES). We overlapped these SNPs in GEUVADIS with the SNPs obtained from each of

371  the three GWASs to obtain common sets of SNPs. The median number of the overlapped cis-

372  SNPs between GEUVADIS and WTCCC, GERA or UK Biobank are 200, 556 or 500,

373  respectively. Afterwards, for each pair of gene (from GEUVADIS) and trait (from GWAS) in

374  turn, we examined the causal relationship between gene expression and trait of interest while

375  testing and controlling for potential horizontal pleiotropic effects.

**Compared Methods**

377  For testing the causal effect, we compared the performance of PMR-Egger with five existing

378  methods that include: (1) SMR, which uses a single instrument and does not control for

379  horizontal pleiotropy. For SMR, we first performed a linear regression to choose the top

380  associated cis-SNP to be the instrumental variable. (2) PrediXcan, which uses multiple correlated

381  instruments but does not control for horizontal pleiotropy. For PrediXcan, we used all cis-SNPs

382  for the model and used ElasticNet implemented in the R package glmnet to obtain the coefficient

383  estimates for the cis-SNPs. (3) TWAS, which uses multiple correlated instruments but does not

384  control for horizontal pleiotropy. For TWAS, we used all cis-SNPs for the model and used

18

385    BSLMM[53] implemented in the GEMMA software[54] to obtain coefficient estimates for the cis-

386    SNPs. (4) CoMM, which uses multiple correlated instruments but does not control for horizontal

387    pleiotropy. We used all cis-SNPs for the model and used the R package CoMM for model fitting.

388    (5) LDA MR-Egger, which uses multiple correlated instruments and controls for horizontal

389    pleiotropy. We used all cis-SNPs for the model and contacted the authors of LDA MR-Egger to

390    obtain the method source code. All these methods are suitable for two-sample design and yield $p$

391    values for testing the causal effect $\alpha$. Note that PrediXcan, TWAS and CoMM are not originally

392    described as an MR method but conceptually rely on the same joint MR model based on

393    equations (1) and (4). These three methods differ in their prior assumptions on $\boldsymbol{\beta}$: PrediXcan

394    relies on ElasticNet assumption; TWAS relies on BSLMM[53] assumption; while CoMM relies on

395    the normal prior assumption. In addition, PrediXcan and TWAS rely on a two-stage regression

396    procedure while CoMM is based on maximum likelihood. We were unable to compare our

397    method with either GSRM or the standard Egger regression, as both require multiple independent

398    SNP instruments that are generally not feasible to obtain in TWAS applications.

399        Again, we used all cis-SNPs for methods that can make use of multiple correlated instruments

400    (i.e. PMR-Egger, TWAS, PrediXcan, CoMM, and LDA MR Egger). We performed a linear

401    regression to select the top associated cis-SNP as the instrumental variable for SMR, as it can

402    only use a single instrument. In all simulations and real data applications, methods that can use

403    either individual-level data or summary statistics (PMR-Egger, PrediXcan and TWAS) are

404    applied using individual-level data as input to ensure their optimal performance. Methods that

405    can only use individual-level data (CoMM) are applied using individual-level data as input.

406    Methods that can only use summary statistics (SMR and LDA MR-Egger) are applied using

19

407    summary data as input. For PMR-Egger, we used individual-level data for all main analyses and

408    used summary data for a subset of analyses that are described in the Discussion section.

409    Besides the above methods, we also compared different methods to a recently published fine-

410    mapping TWAS method, FOCUS[55]. In the FOCUS analysis, we followed[55] and obtained a set of

411    independent non-overlapping genomic regions termed as LD blocks from LDetect[56]. We

412    removed genomic regions that overlap with the MHC region due to the extensive LD structure.

413    Following[55], we also focus our analysis on a subset of regions that harbor at least one genome-

414    wide-significant SNP ($p < 5 \times 10^{-8}$; the default threshold used in FOCUS), and for each

415    TWAS/MR method (i.e. PMR-Egger, TWAS, PrediXcan, CoMM, or SMR), also harbor at least

416    one TWAS gene that is declared significant by the given method. We then applied FOCUS to

417    analyze these remaining regions and identify genes that are in the 90% credible set.

418    For testing horizontal pleiotropic effect, we compared the performance of PMR-Egger with

419    two existing methods that include (1) LDA MR-Egger; and (2) the global test in MR-PRESSO,

420    which is implemented as an R package. Both these methods examine one gene at a time and

421    output a $p$ value for testing horizontal pleiotropic effects.

422

**Results**

Our method is described in the Methods (inside the method overview subsection there), with technical details provided in the Supplementary Note. For TWAS applications, our method examines one gene at a time and estimates and tests its causal effect on a trait of interest. Our method models multiple correlated instruments, performs MR inference in a maximum likelihood inference framework, and is capable of testing and controlling for horizontal pleiotropic effects commonly encountered in TWAS. We refer to our method as the probabilistic Mendelian randomization with Egger regression (PMR-Egger), which is implemented as an R package. Our method is computationally efficient and can analyze each gene in minutes in a GWAS with a few hundred thousand individuals (Table 2).

**Simulations: Testing and estimating the causal effect**

We performed simulations to examine the effectiveness of our method and compared it with existing MR approaches. Simulation details are provided in the Methods. Briefly, we simulated gene expression values based on genotypes from 456 individuals in GEUVADIS and simulated phenotypes based on genotypes from 2,000 randomly selected individuals in GERA. In the simulations, we varied the genetic architecture underlying gene expression from sparse (one SNP or 1% of SNPs are causal) to polygenic (10% or 100% of SNPs are causal). We varied the proportion of SNPs exhibiting horizontal pleiotropic effects in a wide range (from 0%, 10%, 30%, 50% to 100%). We examined directional pleiotropy setting (the ratio of SNPs with negative vs positive horizontal pleiotropic effects is 0:10), approximately directional pleiotropy setting (1:9 or 3:7) and balanced pleiotropy settings (5:5). We varied the magnitude of horizontal pleiotropic effects $\gamma$ to be either $1\times10^{-4}$, $5\times10^{-4}$, $1\times10^{-3}$, or $2\times10^{-3}$, which corresponds to the 50%, 70%, 90%, 95% percentiles of the horizontal pleiotropic effect estimate in real data. We also

446  varied the magnitude of causal effect $\alpha$ to be either 0, 0.14, 0.2 or 0.245, which corresponds to a

447  proportion of phenotypic variance explained by vertical pleiotropic effects ($PVE_{zy}$) as 0, 0.2%,

448  0.4% and 0.6% respectively.

449  Our first set of simulations is focused on causal effect testing. Here, we compared PMR-Egger

450  with five different methods that include SMR, PrediXcan, TWAS, CoMM, and LDA MR-Egger.

451  We first examined type I error control of different methods under the null ($\alpha = 0$). In the

452  absence of horizontal pleiotropic effects, PMR-Egger, together with PrediXcan, TWAS, and

453  CoMM, all provides calibrated type I error (Fig. 1a). Consistent with previous observations[57], we

454  found that SMR produces overly-conservative/deflated p-values. The deflation of SMR p-values

455  is presumably because SMR requires the selected instrument being a true causal SNP with a

456  large effect size, which is not always guaranteed in practice. In addition, we found that LDA

457  MR-Egger produces inflated p-values, presumably because LDA MR-Egger makes a fixed effect

458  assumption on $\boldsymbol{\beta}$. The fixed effect assumption on $\boldsymbol{\beta}$ is not expected to work well in TWAS

459  settings where the number of SNPs are on the same order of the sample size in the gene

460  expression study and where the cis-SNPs are all highly correlated with each other due to LD.

461  Such fixed effect assumption on $\boldsymbol{\beta}$, when paired with the two-stage inference procedure that

462  ignores the estimation uncertainty in the first stage, makes LDA MR-Egger sensitive to the

463  collinearity induced by SNP correlations caused by LD.  Indeed, we found that the $p$-values from

464  LDA MR-Egger are well calibrated when we followed the exact same simulation setting used

465  in[33], where SNP genotypes were simulated based on an autoregressive covariance matrix with a

466  moderate correlation parameter. However, when such correlation parameter was set to be

467  realistically high (>0.9) or if we used SNPs from real data to carry out the same set of

468  simulations, then we observed p-value inflation from LDA MR-Egger (Supplementary Fig. 2).

22

469    In the presence of horizontal pleiotropic effects, PMR-Egger becomes the only method that

470    produces calibrated (or slightly conservative) p-values (Fig. 1b, c, d). In contrast, the p-values

471    from all other methods become inflated, and more so with increasingly large horizontal

472    pleiotropic effect. For example, when $\gamma$ is $5 \times 10^{-4}$, the genomic control factors from PMR-Egger,

473    SMR, PrediXcan, TWAS, CoMM, and LDA MR-Egger are 0.93, 1.30, 1.33, 1.33, 1.49 and 2.61

474    respectively. When $\gamma$ is increased to $1 \times 10^{-3}$, the genomic control factors from PMR-Egger, SMR,

475    PrediXcan, TWAS, CoMM, and LDA MR-Egger become 0.93, 2.39, 2.27, 2.46, 4.03 and 2.57

476    respectively.

477    The null p-value distributions from different methods remain largely similar regardless of the

478    genetic architecture underlying gene expression being sparse or polygenic (Supplementary Fig.

479    3). Note that, the p-values from SMR become less deflated when there is a sparse set of SNPs

480    affecting gene expression; however, such deflation is not completely abolished even when one

481    SNP has non-zero effect on gene expression, presumably because we cannot always identify the

482    true non-zero effect SNP through eQTL mapping and may supply a tagged SNP for SMR

483    analysis. In addition, the p-value distribution pattern for different methods under the null does

484    not change much with reduced the gene expression heritability value $PVE_{zx}$. When $PVE_{zx}$ is

485    either 5% or 1%, PMR-Egger still produces well-calibrated $p$ values (Supplementary Fig. 4).

486    We note that, like the standard MR-Egger regression, our PMR-Egger also makes a relatively

487    strong assumption on the horizontal pleiotropic effect and assumes that all SNPs have the same

488    horizontal pleiotropic effect. To examine the robustness of such assumption, besides the above

489    settings where either 0% or 100% SNPs have horizontal pleiotropic effects, we varied the

490    proportion of horizontal pleiotropic SNPs to be either 10%, 30%, 50%. We found that the p-

491    values from PMR-Egger remain calibrated regardless of the sparsity of the horizontal pleiotropic

23

492    SNPs (Supplementary Fig. 5). In addition, besides the above directional pleiotropy settings

493    where the ratio of SNPs with negative vs positive effects is set to be 0:10, we also examined two

494    approximately directional pleiotropy settings (1:9 or 3:7) and one balanced setting (5:5). We

495    found that the p-values from PMR-Egger remains calibrated in either the approximately

496    directional pleiotropy settings or the balanced setting when horizontal pleiotropic effect is small

497    or moderate ($\gamma = 1x10^{-4}$, $5x10^{-4}$, or $1x10^{-3}$; Supplementary Fig. 6a, b, c). However, when

498    horizontal pleiotropic effect is large ($\gamma = 2x10^{-3}$), as one would expect, the p-values from PMR-

499    Egger becomes inflated, with genomic control factor being 1.08, 1.31 and 1.37, for settings

500    where the ratio is 1:9, 3:7 and 5:5, respectively (Supplementary Fig. 6d). Finally, we repeated all

501    the above analyses with cross-gene based simulations, which provide consistent results on the

502    type I error control of different methods for testing the causal effects (Supplementary Fig. 7-12).

503        Next, we examined the power of different methods to identify the causal effect for a range of

504    possible causal effect sizes $\alpha$. Because the same p-value from different methods may correspond

505    to different type I errors, we computed power based on FDR of 0.1 instead of a nominal p-value

506    threshold to allow for fair comparison across methods. In the absence of horizontal pleiotropic

507    effects or in the presence of small horizontal pleiotropic effects, PMR-Egger, TWAS and CoMM

508    have similarly power, all outperforming the other three methods, highlighting the importance of

509    making polygenic assumptions on $\boldsymbol{\beta}$ and modeling all cis-SNPs together (Fig. 2a, b). The power

510    of PMR-Egger is slightly lower than the other two, presumably because PMR-Egger uses extra

511    parameters to model horizontal pleiotropy, which leads to a loss of degrees of freedom and

512    subsequent loss of power in the absence of horizontal pleiotropy. The power of all methods

513    increases with $\alpha$, though their relative performance rank does not change. In the presence of

514    horizontal pleiotropy, the power of all methods reduces (Fig. 2c, d). However, the power

515    reduction from PMR-Egger is substantially smaller than all other methods. For example, when

516    $PVE_{zy} = 0.006$ and $\gamma = 0.0005$, PMR-Egger reaches a power of 41%; the power of SMR,

517    PrediXcan, TWAS, CoMM, and LDA MR-Egger are 7%, 24%, 31%, 33% and 1%, respectively.

518    When $PVE_{zy} = 0.006$ but $\gamma = 0.001$, the power of PMR-Egger remains similar and is 40%; the

519    power of SMR, PrediXcan, TWAS, CoMM, and LDAMR-Egger reduces to 3%, 13%, 16%, 16%

520    and 0.9%, respectively. Besides the horizontal pleiotropic effects $\boldsymbol{\gamma}$, we examined how power is

521    influenced by the genetic architecture underlying gene expression, $\boldsymbol{\beta}$ (Supplementary Fig. 13).

522    We found that the power of different methods in the setting where 10% of SNPs have non-zero

523    effects on gene expression are similar to the baseline setting where all SNPs have non-zero

524    effects, both in the absence (Supplementary Fig. 13e vs Fig. 2a) or in the presence of horizontal

525    pleiotropic effects (Supplementary Fig. 13F vs Fig. 2d). However, the relative performance of

526    different methods changes when there is only one SNP or 1% SNPs having non-zero effect on

527    gene expression. Specifically, in the absence of horizontal pleiotropic effects, the power of both

528    PrediXcan and SMR become slightly higher than PMR-Egger, TWAS and CoMM, all of which

529    have substantially higher power than LDA MR-Egger (Supplementary Fig. 13a, c). The higher

530    power of PrediXcan and SMR in the sparse setting presumably is because the ElasticNet

531    estimation procedure employed in PrediXcan favors sparse eQTLs while SMR explicitly makes a

532    single eQTL assumption. In the presence of horizontal pleiotropic effects, however, PMR-Egger

533    remains the most powerful, even in the setting where only one SNP has non-zero effect on gene

534    expression (Supplementary Fig. 13b, d). We also found that PMR-Egger produces accurate

535    estimate of the causal effect $\alpha$, both under the null and under various alternatives, in the presence

536    or absence of horizontal pleiotropic effects (Supplementary Fig. 14). The causal effect estimates

537    remain reasonably unbiased in the two approximately directional pleiotropy settings and one

538   balanced setting (Supplementary Fig. 15a, c, e). Finally, we repeated all the above analyses in

539   cross-gene based simulations, which provide consistent results on the power of different methods

540   for detecting causal effects (Supplementary Fig. 16-17).

541   **Simulations: Testing and estimating horizontal pleiotropic effect**

542   Our second set of simulations is focused on horizontal pleiotropic effect testing. Here, we

543   compared PMR-Egger with two different methods: LDA MR-Egger and MR-PRESSO. All three

544   methods examine one gene at a time and test whether cis-SNPs within the gene exhibit non-zero

545   horizontal pleiotropic effects. Note that, unlike PMR-Egger and LDA MR-Egger, MR-PRESSO

546   requires independent instruments and uses permutation to obtain the empirical p-values. Due to

547   the heavy computational burden resulting from permutations, we restricted the number of

548   permutations in MR-PRESSO to 10,000 (the lowest possible p value from MR-PRESSO is thus

549   $10^{-4}$) and were only able to apply MR-PRESSO to a subset of simulation scenarios.

550        We first examined type I error control of different methods under the null, where there is no

551   horizontal pleiotropic effect. We found that the p-values from PMR-Egger provide calibrated

552   type I error control under a range of causal effect sizes $\alpha$ (Fig. 3). However, p-values from both

553   LDA MR-Egger and MR-PRESSO are inflated, and more so with increasingly large causal effect

554   $\alpha$. For example, when $\text{PVE}_{zy} = 0$, the genomic control factor from PMR-Egger and LDA MR-

555   Egger are 0.96 and 2.31, respectively. When $\text{PVE}_{zy}$ is increased to 0.6%, the genomic control

556   factor from PMR-Egger remains 0.96, while the genomic control factor from LDA MR-Egger

557   becomes 3.04. (We are unable to accurately compute the genomic control factor for MR-

558   PRESSO because its minimal p-value is $10^{-4}$.) The overly inflated p-values from LDA MR-Egger

559   is presumably due to its fixed effect modeling assumption on $\boldsymbol{\beta}$ and the subsequent failure to

560   control for realistic LD patterns. The inflation of MR-PRESSO p values is presumably because

26

561    MR-PRESSO can only handle independent instruments and thus does not fare well in TWAS

562    settings. Inflation of p-value on testing horizontal pleiotropy would incorrectly identify genes

563    with no pleiotropic effects, thus likely reducing the power to detect true causal effect $\alpha$.

564    Importantly, the p-values from PMR-Egger remain calibrated regardless of the genetic

565    architecture underlying gene expression (Supplementary Fig. 18). Finally, we repeated all the

566    above analyses in cross-gene based simulations, which provide consistent results on the type I

567    error control of different methods for testing pleiotropic effects (Supplementary Fig. 19-20).

568    Next, we examined the power of different methods in detecting non-zero horizontal

569    pleiotropic effect. Again, we computed power based on an FDR of 0.1 instead of the nominal p-

570    value to allow for fair comparison across methods. We dropped MR-PRESSO for comparison

571    here due to its heavy computational burden. We found that PMR-Egger outperforms LDA MR-

572    Egger in a range of possible horizontal pleiotropic effect sizes, and that the power of both

573    methods increases with increasing horizontal pleiotropy (Fig. 2e, f). For example, when $\text{PVE}_{zy} =$

574    0.6% and $\gamma = 0.0005$, PMR-Egger achieves a power of 1.6% while LDA MR-Egger achieves a

575    power of 1% (note that the power is relatively small due to the small sample size used in the

576    simulations). When $\text{PVE}_{zy} = 0.6\%$ but $\gamma = 0.001$, the power of PMR-Egger increases to 58.9%

577    while the power of LDA MR-Egger increases to 32%. In addition, the power to detect horizontal

578    pleiotropic effects is not influenced by the sparsity level of the genetic architecture underlying

579    gene expression (Supplementary Fig. 21). The power to detect horizontal pleiotropic effects does,

580    however, depend on the sparsity level of $\boldsymbol{\gamma}$ (Supplementary Fig. 22a). Specifically, power of both

581    PMR-Egger and LDA MR-Egger reduces with increasing sparsity of $\boldsymbol{\gamma}$, though the power of

582    PMR-Egger remains higher than LDA MR-Egger across a range of sparsity values. Similarly, the

583    power to detect pleiotropic effects also suffers in the absence of directional pleiotropic effect

27

584    (Supplementary Fig. 22b). In addition, PMR-Egger can estimate the horizontal pleiotropic effect

585    size accurately in the presence of directional pleiotropic effect (Supplementary Fig. 23).

586    However, in the absence of directional pleiotropic effect, as one would expect, the estimates of

587    pleiotropic effects become under-ward biased, more so in the balanced setting than in the

588    approximately directional pleiotropy settings (Supplementary Fig. 15b, d, f). Finally, we repeated

589    all the above analyses in cross-gene based simulations, which provide consistent results on the

590    power of different methods for detecting pleiotropic effects (Supplementary Fig. 24-25).

591    **Real data applications**

592    We performed TWAS to detect genes causally associated with any of the 39 phenotypes

593    collected from three GWASs (details in Methods). The examined gene expression data is

594    obtained from the GEUVADIS study and contains 15,810 genes. The examined phenotypes

595    include 7 common diseases from WTCCC, 22 diseases from GERA, and 10 quantitative traits

596    from UK Biobank. The GWAS sample size ranges from 4,686 (for Crohn's disease in WTCCC)

597    to 337,198 (for UK Biobank). We applied PMR-Egger together with five other approaches (SMR,

598    PrediXcan, TWAS, CoMM, and LDA MR-Egger) to examine pairs of gene and phenotype one at

599    a time. In the analysis, we regressed phenotypes on the top 10 genotyping PCs to obtain the

600    phenotype residuals, which we used further to conduct TWAS analysis for all compared methods.

601    The p-values for testing the causal effect of each gene on the phenotype are shown for WTCCC

602    traits (Fig. 4a, b and Supplementary Fig. 26), GERA traits (Fig. 5a, b and Supplementary Fig.

603    27), and UK Biobank traits (Fig. 6a, b and Supplementary Fig. 28); with genomic control factors

604    listed in Supplementary Table 1 and visualized in Fig. 4c, Fig. 5c and Fig. 6c. Besides these main

605    analyses, we also performed parallel analysis for PMR-Egger where we used the original

606    phenotype as the outcome and included the top 10 genotype PCs as covariates (Supplementary

607    Figures 29-31). The results from these parallel analyses are largely consistent with the main

608    results. Therefore, we will mainly report the main results in the following text. For illustration

609    purpose, we display qq-plots for two selected traits in each data, one with a relatively low

610    number of gene associations and the other with a relatively high number of gene associations, in

611    Fig. 4a, b, Fig. 5a, b and Fig. 6a, b, respectively. Among the selected six traits, the one with zero

612    number of associated genes (BD in WTCCC; Fig. 4a) and the one with one associated gene

613    (Irritable Bowel Syndrome in GERA; Fig. 5a), represent approximately null traits with no

614    apparently associated genes. For the six selected traits, consistent with simulations, we found that

615    the p-values from PMR-Egger are well calibrated, more so than the other methods. In contrast,

616    the p-values from CoMM, TWAS, PrediXcan and LDA MR-Egger are inflated and deviated

617    upward from the diagonal line, while the p-values from SMR are overly conservative and lie

618    below the diagonal line. The results observed in these exemplary traits generalize to all other

619    examined traits. For example, the genomic control factor from PMR-Egger is the lowest among

620    all methods in 25 out of the 39 traits, and ranges from 0.93 to 1.04 in WTCCC (Fig. 4c), from

621    0.92 to 1.13 in GERA (Fig. 5c), and from 1.12 to 1.34 in UK Biobank (Fig. 6c). (Note that the

622    higher genomic control factor in the large UK Biobank as compared to WTCCC and GERA is

623    expected under polygenic architecture[58] and reflects at least in part the higher power in the UK

624    Biobank as compared to GERA and WTCCC.) In contrast, the genomic control factors from

625    CoMM, TWAS, PrediXcan are often higher than that from PMR-Egger for most traits examined.

626    For example, the genomic control factor from CoMM is often the highest among all other

627    methods (except for LDA MR-Egger) in 22 out of the 39 traits, and ranges from 1.13 to 1.23 in

628    WTCCC, 0.94 to 1.62 in GERA, and 1.45 to 1.90 in UK Biobank. The genomic control factor

629    from TWAS is the highest among all other methods (except for LDA MR-Egger) in 14 out of the

630   39 traits, ranges from 1.20 to 1.31 in WTCCC, 0.98 to 1.15 in GERA, and 1.30 to 2.17 in UK

631   Biobank. The genomic control factor from PrediXcan is the highest among all other methods

632   (except for LDA MR-Egger) in 5 out of the 39 traits, and ranges from 1.21 to 1.31 in WTCCC,

633   1.00 to 1.16 in GERA, and 1.09 to 1.46 in UK Biobank. In addition, consistent with simulations,

634   we observed a substantial inflation of LDA MR-Egger p-values: its genomic control factor

635   ranges from 17.60 to 18.56 in WTCCC, 32.13 to 34.74 in GERA, and 10.48 to 16.65 in UK

636   Biobank. Also consistent with simulations, the p-value from SMR often lies underneath the

637   expected null, even though its genomic control factors are often well behaved (Fig. 4a, b, Fig. 5a,

638   b, Fig. 6a, b and Supplementary Figs. 26-28).

639      We examined the number of associated genes detected by different methods based on a

640   Bonferroni corrected genome-wide threshold (Fig. 4d, Fig. 5d and Fig. 6d; Supplementary Table

641   2). We note that the number of detected genes based on this p-value threshold may artificially

642   favors those methods that have inflated type I error control. For this analysis, we excluded LDA

643   MR-Egger for comparison, as its p-values are overly inflated. Comparing across the remaining

644   methods, we found that SMR can barely detect any genes significantly associated with traits

645   across all three data sets, much less so than that detected by the other four methods. The much

646   lower number of genes detected by SMR than the other four methods are consistent with the

647   relatively low power of SMR observed in simulations. For the other four methods, we found that

648   the number of gene-trait pairs detected by CoMM and PMR-Egger is higher than that detected by

649   TWAS and PrediXcan in all three GWASs (Fig. 4d, Fig. 5d and Fig. 6d; Supplementary Table 2).

650   The higher number of discoveries by both CoMM and PMR-Egger in the three GWASs is

651   consistent with our simulations as well as previous observations that likelihood-based inference

652   often achieves higher power than two-stage inference for MR analysis. However, we do notice

30

653      that PMR-Egger detects slightly lower number of gene-trait pairs than CoMM based on the same

654      genome-wide p-value threshold, consistent with the inflated genomic inflation factors observed

655      for CoMM. Indeed, we found that the estimated $|\frac{\alpha}{\gamma}|$ for the common set of genes detected by

656      both CoMM and PMR-Egger is higher than the set of genes only detected by CoMM across traits

657      (Supplementary Fig. 32a, b). Therefore, the genes detected by CoMM but not PMR-Egger tend

658      to have large $|\gamma|$ and small $|\alpha|$, likely reflecting false associations due to horizontal pleiotropic

659      confounding.

660      Overall, by controlling for horizontal pleiotropic effects, PMR-Egger detected many likely

661      causal genes that the other methods failed to detect. For example, the *LNK*/*SH2B3* gene

662      (111,743,752-111,989,427 on chr 12) is only identified by PMR-Egger to be associated with

663      platelet count in the UK Biobank (PMR-Egger $p = 1.17 \times 10^{-221}$; CoMM $p$=0.98; TWAS $p =$

664      $8.6 \times 10^{-5}$; PrediXcan $p$=0.68; SMR $p$=0.024). The association between *LNK* and plate count is

665      consistent with results from recent large-scale GWASs[59-61]. *LNK*/*SH2B3* encodes the lymphocyte

666      adaptor protein (LNK) that is primarily expressed in hematopoietic and endothelial cells[62]. In

667      hematopoietic cells, LNK functions as a negative regulator of cell proliferation as well as the

668      thrombopoietin-mediated cytokine signaling pathway, which is a key signaling pathway that

669      promotes megakaryocytes to form platelets[62,63]. Indeed, platelets are overproduced and

670      accumulated in *LNK* knockdown cells as well as *Lnk* knockout mouse[64-66], supporting a causal

671      role of *LNK* in platelets production. As the second example, the *NOD2* gene (50,627,514-

672      50,866,988 on chr 16) is identified by PMR-Egger to be associated with Crohn's disease (CD;

673      $p = 6.1 \times 10^{-19}$), and, with a slightly less significance, also by CoMM ($p = 7.8 \times 10^{-15}$). The

674      association between *NOD2* and CD was not identified by the other methods (TWAS $p$=0.005;

675      PrediXcan $p$=0.92; SMR $p$=0.15). *NOD2* encodes a cytosolic pattern recognition receptor that

676    acts both as a cytoplasmic sensor of microbial products and as an important mediator of innate

677    immunity and inflammatory response[67] The *NOD2* gene is a well-known susceptible gene for

678    CD and is perhaps one of the first genes ever implied for CD[68]. Multiple SNPs in *NOD2* have

679    been found to be associated with CD in both early linkage studies[69-71] and many recent

680    GWASs[72,73]. *NOD2* variants associated with CD often reside in the ligand recognition domain of

681    NOD2 and can lead to aberrant bacterial handling and antigen presentation[74]. Indeed, *NOD2*-

682    deficient mice displays dysregulated bacterial community in the ileum and *NOD2*-deficient ileal

683    epithelia exhibit impaired ability of inducing immune responses for bacteria elimination[75]. It is

684    thus hypothesized that mis-regulation of *NOD2* can causally lead to altered interactions between

685    ileal microbiota and mucosal immunity, resulting in increased disease susceptibility to CD[75]. As

686    a third example, the *TFRC* gene (195,654,054-195,909,060 on chr 3) is identified by PMR-Egger

687    to be associated with red blood cell distribution width (RDW) in the UK Biobank ($p = 3.3 \times$

688    $10^{-17}$). Such association is not identified by the other methods (CoMM *p*=0.95; TWAS *p*=0.76;

689    PrediXcan *p*=0.97; SMR *p*=0.38). *TFRC* encodes the classical transferrin receptor that is

690    involved in cellular iron uptake[76,77]. Multiple SNPs in *TFRC* have been established to be

691    associated with various erythrocyte phenotypes in GWASs[78,79]. These associated erythrocyte

692    phenotypes include the mean corpuscular hemoglobin (MCH) and mean corpuscular volume

693    (MCV, the average volume of red blood cells) which is directly related to RDW[77,78]. The variants

694    in *TFRC* likely lead to decreased iron availability for red cell precursors, as has been observed in

695    mice deficient in *TFRC*, thus resulting in a compensatory increase of red blood cell size as

696    measured by RDW[80]. The regional association plots for all these three genes are presented in the

697    Supplementary Fig.33-35.

698  We compared the results from different MR methods with a recently published TWAS fine-

699 mapping method, FOCUS[55]. The analysis details are provided in the Materials and Methods

700 section. Briefly, we follow [55] and focused on independent and non-overlapping genomic regions

701 that harbor at least one genome-wide-significant SNP and at least one TWAS gene that is

702 significant by the MR methods. The number of genes and regions analyzed by FOCUS for each

703 of the three data sets are shown in Supplementary Table 3, which also contains the number of

704 associated genes detected by FOCUS in the credible set. Due to the small number of associated

705 genes detected in WTCCC, we focus our main comparison in GERA and UK Biobank. In these

706 real data applications, we found that the results from PMR-Egger is largely consistent with that

707 of FOCUS, more so than the other methods (Supplementary Fig. 41). Specifically, the average

708 PMR-Egger -log10(p-value) for genes in the FOCUS 90% credible set is 22.43 in GERA and

709 10.67 in UK Biobank. The average -log10(p-value) of PMR-Egger is higher than CoMM (13.83

710 and 10.43), TWAS (5.71 and 7.55), PrediXcan (4.66 and 7.06) and SMR (NA for GERA, as no

711 gene in the credible set is detected by SMR; 1.78 for UKbiobank). In addition, the difference of

712 the average PMR-Egger -log10(p-value) between genes in the FOCUS credible set and genes

713 outside is large (16.61 in GERA and 7.43 in UK Biobank). The -log10(p-value) difference is

714 again larger than CoMM (8.41 and 6.02), TWAS (4.52 and 5.35), PrediXcan (3.50 and 4.74) and

715 SMR (NA and 0.28). Similarly, the proportion of significant genes detected by PMR-Egger in

716 the FOCUS credible set is 78% in GERA and 60% in UK Biobank. The proportion of significant

717 genes by PMR-Egger is higher than CoMM (75% and 53%), TWAS (50% and 47%), PrediXcan

718 (50% and 48%) and SMR (NA and 8%). In addition, the difference in the proportion of

719 significant genes detected by PMR-Egger between genes in the FOCUS credible set and genes

720 outside is high (53% in GERA and 41% in UK Biobank). This proportion difference by PMR-

721      Egger is again higher than CoMM (51% and 39%), TWAS (46% and 36%), PrediXcan (50% and

722      35%) and SMR (NA and 1%). The consistency between PMR-Egger and FOCUS validates the

723      high power of PMR-Egger.

724      Next, we shift our focus to testing horizontal pleiotropic effects. The p-values for testing the

725      causal effect of gene on phenotype are shown for WTCCC traits (Fig. 4e, f and Supplementary

726      Fig. 26), GERA traits (Fig. 5e, f and Supplementary Fig. 36), and UK Biobank traits (Fig. 6e, f

727      and Supplementary Fig. 37); with genomic control factors visualized in Fig. 4g, Fig. 5g and Fig.

728      6g. We also display qq-plots for the previously selected exemplary traits in Fig. 4e, f, Fig. 5e, f,

729      and Fig. 6e, F. Overall, consistent with simulations, the p-values from PMR-Egger are well

730      behaved while the p-value from LDA MR-Egger display substantial inflation. For example, the

731      genomic control factor from PMR-Egger ranges from 0.93 to 1.01 in WTCCC (Fig. 4g), from

732      0.92 to 1.09 in GERA (Fig. 5g), and from 1.13 to 1.71 in UK Biobank (Fig. 6g). In contrast, the

733      genomic control factor from LDA MR-Egger ranges from 34.00 to 36.00 in WTCCC, from 69.82

734      to 72.19 in GERA and from 17.75 to 29.85 in UK Biobank (Supplementary Table 1). With the

735      same Bonferroni adjusted genome-wide p-value threshold, PMR-Egger detected 33 gene-trait

736      pairs in WTCCC in which the cis-SNPs exhibit significant horizontal pleiotropy, 37 gene-trait

737      pairs in GERA, and 626 gene-trait pairs in the UK Biobank.

738      Horizontal pleiotropic effect tests can help us explain some of the discrepancy in terms of the

739      causal associations detected by PMR-Egger and the other methods. For example, for the trait of

740      red blood cell count in UK Biobank, the *MAPT* gene on chromosome 17 shows a significant

741      pleiotropy effect ($p = 2.35 \times 10^{-9}$) but displays no significant causal effect ($p=0.98$) by PMR-

742      Egger. In contrast, *MAPT* is detected to be significantly associated with red blood cell count by

743      PrediXcan ($p = 8.11 \times 10^{-10}$), and, to a much lesser extent, by TWAS ($p = 1.72 \times 10^{-3}$).

744   However, no previous evidence suggests that *MAPT* is associated with red blood cell count.

745   Indeed, we found that the genomic location of *MAPT* (43,871,748-44,205,700) is close to and

746   partially overlapped with *KANSL1* (44,007,282-44,402,733), which has been previously

747   identified to be associated with red blood cell traits[81,82]. The association between *KANSL1* and

748   red blood cell count is also detected by PMR-Egger ($p = 1.02 \times 10^{-7}$), by CoMM ($p = 2.72 \times$

749   $10^{-8}$), and, to a much lesser extent, by TWAS ($p = 1.66 \times 10^{-3}$) in the present study. By

750   controlling for the expression level of the *KANSL1* gene in the PrediXcan framework, the

751   association between the predicted *MAPT* expression level and red blood cell count is no longer

752   significant ($p = 0.10$). Therefore, the causal association between *MAPT* and red blood cell count

753   detected by PrediXcan likely reflects either the true horizontal pleiotropic effect of *MAPT* cis-

754   SNPs on red blood cell count through *KANSL1* or their tagging effects of the neighboring eQTLs

755   of *KANSL1*. As another example, for height in the UK Biobank, the pseudogene *RP11-9E13.2*

756   (70,137,755-70,340,521) on chromosome 10 has a significant pleiotropy effect ($p = 1.08 \times$

757   $10^{-13}$) but displays no significant causal effect ($p$=0.93) by PMR-Egger. In contrast, *RP11-*

758   *9E13.2* is detected to be significantly associated with height by PrediXcan ($p = 4.34 \times 10^{-10}$),

759   and, to a lesser extent, by TWAS ($p = 9.05 \times 10^{-6}$). The pseudogene *RP11-9E13.2* is in the

760   neighborhood of *MYPN* (69,765,912-70,071,774), which has been previously identified to be

761   associated with height[83]. The association between *MYPN* and height is also detected by PMR-

762   Egger ($p = 1.82 \times 10^{-7}$), CoMM ($p = 2.13 \times 10^{-14}$), and to a lesser extent, PrediXcan ($p =$

763   $3.94 \times 10^{-4}$) and TWAS ($p = 1.55 \times 10^{-3}$), in the present study. By controlling for the

764   predicted expression level of *MYPN* gene in the PrediXcan framework, the association between

765   the predicted *RP11-9E13.2* expression level and height is no longer significant at the genome-

766   wide threshold ($p = 3.37 \times 10^{-4}$). Therefore, the causal association between the pseudogene

767   *RP11-9E13.2* and height as detected by PrediXcan and TWAS likely reflects either the

768   horizontal pleiotropic effect of *RP11-9E13.2* cis-SNPs on height through *MYPN* or their tagging

769   effects of the neighboring eQTLs of *MYPN*. The results suggest the practical importance of

770   testing and controlling for pleiotropic effects in TWAS applications. Certainly, we acknowledge

771   that, both these examples are focused on the special case where the false gene association with

772   the trait disappears when conditional on a neighboring gene. We did not provide examples where

773   the apparently false gene association with the trait may be explained by horizontal pleiotropic

774   effects acted upon/through a gene far away, as it is often challenging to convincingly identify

775   trans eQTL effects. In the special case we focused on, while it is possible that SNPs display true

776   horizontal pleiotropic effects through the neighboring gene, it is equally likely that SNPs used in

777   the model are simply tagging nearby eQTLs of the neighboring causal gene[55,84] and thus display

778   apparent "horizontal pleiotropic effects" through the neighboring gene,  as also mentioned above.

779   Subsequently, the horizontal pleiotropic effect term in PMR-Egger may represent the apparent

780   "horizontal pleiotropic effects" through SNP tagging to the nearby eQTLs of the causal gene,

781   rather than the truly horizontal pleiotropic effect acted through other molecular pathways.

782   Regardless of the interpretation of the horizontal pleiotropic effect term, we found it reassuring

783   that by modeling the horizontal pleiotropic effect term in PMR-Egger can reduce false

784   discoveries in the case of SNP tagging.

785      We note that an important feature of PMR-Egger is its ability to test both causal effect and

786   horizontal pleiotropy effect simultaneously. We contrast the p-values obtained from these two

787   different tests across genes for those traits in which at least one gene is detected as significant

788   from either of the two tests (Supplementary Figs. 38-40). We found that different traits exhibit

789   different gene association patterns. For example, some traits may only contain genes with a

36

790    significant causal effect but without a significant horizontal pleiotropic effect (e.g. CD and CAD

791    in WTCCC; Allergic Rhinitis, Irritable Bowel Syndrome and Psychiatric disorders in GERA).

792    Some traits may only contain genes with a significant horizontal pleiotropic effect but without a

793    significant causal effect (e.g. Dermatophytosis in GERA). Some traits may contain genes with a

794    significant causal effect as well as genes with a significant horizontal pleiotropic effect, but with

795    the two sets of genes being non-overlapped (e.g. Asthma, Dyslipidemia, HT, Abdominal Hernia

796    and Macular Degeneration in GERA; Fored Vitral Capacity in UK Biobank). While the majority

797    of traits contain genes with both a significant causal effect and a significant horizontal

798    pleiotropic effect. The top gene which is most significant for both causal effect test and

799    pleiotropy test is highlighted in the plots. Being capable of testing both causal effect and

800    horizontal pleiotropy effect facilitates our understanding of the gene association pattern with

801    various different complex traits.

802

803

**Discussion**

We have presented a data generative model and a likelihood framework for MR analysis that unifies many existing transcriptome wide association analysis methods and many existing MR methods. Under the framework, we have presented PMR-Egger, a new method that conducts MR analysis using multiple correlated instruments while properly controlling for horizontal pleiotropic effects. By properly controlling for horizontal pleiotropic effects and making inference under a likelihood framework, PMR-Egger yields calibrated p-values across a wide range of scenarios and improves power of MR analysis over existing approaches. We have illustrated the benefits of PMR-Egger through extensive simulations and multiple real data applications of TWAS.

One important modeling assumption we made in PMR-Egger is that the horizontal pleiotropic effects of all SNPs equal to each other. The equal effect size assumption directly follows the commonly used Egger regression modeling assumption for MR analysis and is analogous to the burden effect size assumption commonly used for rare variant tests. Consistent with existing literature on applications of the Egger regression and burden test, we also found that equal effect size assumption employed in PMR-Egger works reasonably robust for causal effect estimation and testing with respect to a range of model mis-specifications and appears to be effective in several real data applications examined here. However, we do acknowledge that our equal effect size assumption in PMR-Egger can be overly restrictive in many settings. For example, as described in the Results, in the absence of direction pleiotropy, the pleiotropic effect estimate becomes down-ward biased and the pleiotropic effect test loses power. We have attempted to alleviate this restrictive modeling assumption by imposing an alternative modeling assumption on the horizontal effect sizes based on variance component assumption. In particular, we have

827    attempted to assume that the horizontal pleiotropic effect of each SNP follows a normal

828    distribution with mean zero and a certain variance component parameter, i.e. analogous to the

829    SKAT test assumption[85]. Such variance component assumption is a more flexible modeling

830    assumption than the equal effect size assumption, potentially alleviating much of the concern

831    with respect to the sensitivity and robustness of equal effect size assumption. Unfortunately,

832    under the variance component assumption, inference for the resulting PMR model becomes

833    overly complicated. In particular, due to the estimation uncertainty in the hyper-parameter

834    estimates, the p-values from the PMR variance component model becomes severely deflated

835    even under simple null simulations (Supplementary Fig. S42). Such deflation of p-values has

836    been previously observed in variance component tests for microbiome applications[86]. Only few

837    methods exist to address such p-value in-calibration issue resulting from hyper-parameter

838    estimation uncertainty[87], and it is not straightforward to adapt any of these methods to our PMR

839    variance component model. Besides the equal effect size modeling restriction, we also note that

840    neither PMR-Egger nor the PMR variance component model is capable of accounting for

841    correlation between horizontal pleiotropic effects $\gamma$ and the SNP effects on gene expression $\beta$.

842    Therefore, while we view PMR-Egger as in important first step towards effective control of

843    horizontal pleiotropic effects in TWAS applications, we emphasize that imposing more realistic

844    modeling assumptions on the horizontal pleiotropic effects in the PMR framework will likely

845    yield more fruitful results in the future.

846        We have primarily focused on modeling continuous traits with PMR-Egger. For case control

847    studies, we have followed previous approaches and directly treated binary phenotypes as

848    continuous outcomes[19,53,88,89], which appears to work well in both WTCCC and GERA data

849    applications we examined. Treating binary phenotypes as continuous outcomes can be justified

39

850     by recognizing the linear model as a first order Taylor approximation to a generalized linear

851     model[53]. However, it would be desirable to extend PMR-Egger to accommodate case control

852     data or other discrete data types in a principled way, by, for example, extending PMR-Egger into

853     the generalized linear model framework. In particular, we could use a probit or a logistic link to

854     extend PMR-Egger to directly model case control data. Extending PMR-Egger to model discrete

855     data types using the generalized linear model framework would likely lead to wider applications

856     of PMR-Egger and is thus an important avenue for future research.

857         We have primarily focused on modeling individual-level data with PMR-Egger. However, like

858     many other linear model-based methods in statistical genetics, PMR-Egger can also be easily

859     extended to make use of summary statistics. The summary statistics version of PMR-Egger is

860     described in detail in the Supplementary Text. Briefly, the summary statistics version of PMR-

861     Egger requires marginal SNP effect size estimates and their standard errors, both on the gene

862     expression and on the trait of interest. In addition, it requires a SNP by SNP correlation matrix

863     that can be constructed based on a reference panel. We validated the implementation of the

864     summary statistics-based approach of PMR-Egger in simulations (details in Materials and

865     Methods). In the comparison, we constructed the SNP by SNP correlation matrix from three

866     different reference panels, by using either all individuals from the GWAS data, 10% randomly

867     selected individuals from the GWAS data, or individuals of European ancestry from the 1,000

868     Genomes project. We applied the summary statistics-based approach of PMR-Egger to each

869     reference panel and compared results with the individual level data-based approach of PMR-

870     Egger that was applied to the complete data. The p values from both approaches for testing

871     causal effects as well as for testing pleiotropy effects are largely consistent with each other,

872     demonstrating    the    effectiveness    of    summary    statistics-based    approach    of    PMR-Egger

40

873     (Supplementary Fig. 43). The summary statistics-based approach of PMR-Egger is implemented

874     in the same software package. Being able to make use of summary statistics extends the

875     applicability of PMR-Egger to data sets where individual-level genotype or phenotype are not

876     available.

877     Finally, in addition to what we have already mentioned in the Materials and Methods, we

878     emphasize here again, that, while we have followed the previous MR literature and use "causal

879     effect" through the text, the effect is causal only when certain MR modeling assumptions hold.

880     These MR assumptions are often not straightforward to prove. For example, without measuring

881     all potential confounders, it is not straightforward to argue that the SNP instruments are not

882     associated with any other confounders that may be associated with both exposure and outcome.

883     Therefore, we caution against the over-interpretation of causal inference in observation studies

884     such as TWAS applications. However, we do believe MR is an important step that allows us to

885     move beyond standard linear regressions and is an important analysis that can provide potentially

886     more trustworthy evidence with regard to causality compared to simpler approaches.

887

888

889

890

891

892

893

894

895

896   **Code availability.** Our method is implemented in the R package PMR, freely available at

897   http://www.xzlab.org/software.html and https://cran.r-roject.org/web/packages/PPMR/index.html.

898   The code to reproduce all the analyses are available on GitHub

899   (https://github.com/yuanzhongshang/PMRreproduce).

900   **Data availability.** No data were generated in the present study. The GEUVADIS gene

901   expression data is publicly available at http://www.geuvadis.org. The WTCCC genotype and

902   phenotype data is publicly available at https://www.wtccc.org.uk. The GERA genotype and

903   phenotype data is available in dbGaP (https://www.ncbi.nlm.nih.gov/gap) with accession number

904   phs000788. The UK Biobank data is from UK Biobank resource under Application Number

905   30686.

906

## References

1    Gamazon, E. R. *et al.* A gene-based association method for mapping traits using reference transcriptome data. *Nature genetics* **47**, 1091-1098 (2015).

2    Gusev, A. *et al.* Integrative approaches for large-scale transcriptome-wide association studies. *Nature genetics* **48**, 245-252 (2016).

3    Zhu, Z. *et al.* Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. *Nature genetics* **48**, 481-487 (2016).

4    Zhu, Z. *et al.* Causal associations between risk factors and common diseases inferred from GWAS summary data. *Nature communications* **9**, 224 (2018).

5    Burgess, S., Small, D. S. & Thompson, S. G. A review of instrumental variable estimators for Mendelian randomization. *Statistical methods in medical research* **26**, 2333-2355 (2017).

6    Ference, B. A. *et al.* Variation in PCSK9 and HMGCR and risk of cardiovascular disease and diabetes. *New England Journal of Medicine* **375**, 2144-2153 (2016).

7    Helgadottir, A. *et al.* Variants with large effects on blood lipids and the role of cholesterol and triglycerides in coronary disease. *Nature genetics* **48**, 634-639 (2016).

8    Pingault, J.-B. *et al.* Using genetic data to strengthen causal inference in observational research. *Nature Reviews Genetics* **19**, 566-580 (2018).

9    Zheng, J. *et al.* Recent developments in Mendelian randomization studies. *Current epidemiology reports* **4**, 330-345 (2017).

10   Haycock, P. C. *et al.* Best (but oft-forgotten) practices: the design, analysis, and interpretation of Mendelian randomization studies. *The American journal of clinical nutrition* **103**, 965-978 (2016).

11   Lawlor, D. A. Commentary: Two-sample Mendelian randomization: opportunities and challenges. *International journal of epidemiology* **45**, 908-915 (2016).

12   Bowden, J. *et al.* A framework for the investigation of pleiotropy in two-sample summary data Mendelian randomization. *Statistics in medicine* **36**, 1783-1802 (2017).

13   Bowden, J., Smith, G. D. & Burgess, S. Mendelian randomization with invalid instruments: effect estimation and bias detection through Egger regression. *International journal of epidemiology* **44**, 512-525 (2015).

14   Bowden, J., Davey Smith, G., Haycock, P. C. & Burgess, S. Consistent estimation in Mendelian randomization with some invalid instruments using a weighted median estimator. *Genetic epidemiology* **40**, 304-314 (2016).

15   Burgess, S., Butterworth, A. & Thompson, S. G. Mendelian Randomization Analysis With Multiple Genetic Variants Using Summarized Data. *Genetic epidemiology* **37**, 658-665 (2013).

16   Burgess, S., Dudbridge, F. & Thompson, S. G. Combining information on multiple instrumental variables in Mendelian randomization: comparison of allele score and summarized data methods. *Statistics in medicine* **35**, 1880-1906 (2016).

17   Burgess, S. & Thompson, S. G. Bias in causal estimates from Mendelian randomization studies with weak instruments. *Statistics in medicine* **30**, 1312-1323 (2011).

18   Yang, C. *et al.* CoMM: a collaborative mixed model to dissecting genetic contributions to complex traits by leveraging regulatory information. *Bioinformatics*, doi:10.1093/bioinformatics/bty865 (2018).

19   Zeng, P. & Zhou, X. Non-parametric genetic prediction of complex traits with latent Dirichlet process regression models. *Nature communications* **8**, 456 (2017).

20   Nagpal, S. *et al.* TIGAR: An Improved Bayesian Tool for Transcriptomic Data Imputation Enhances Gene Mapping of Complex Traits. *The American Journal of Human Genetics* (2019).

952  21  Verbanck, M., Chen, C.-Y., Neale, B. & Do, R. Detection of widespread horizontal pleiotropy in
953      causal relationships inferred from Mendelian randomization between complex traits and
954      diseases. *Nature genetics* **50**, 693-698 (2018).
955  22  Hemani, G., Bowden, J. & Davey Smith, G. Evaluating the potential role of pleiotropy in
956      Mendelian randomization studies. *Human molecular genetics* **27**, R195-R208 (2018).
957  23  Jordan, D. M., Verbanck, M. & Do, R. The landscape of pervasive horizontal pleiotropy in human
958      genetic variation is driven by extreme polygenicity of human traits and diseases. *bioRxiv*, 311332
959      (2018).
960  24  Park, Y. *et al.* A Bayesian approach to mediation analysis predicts 206 causal target genes in
961      Alzheimer's disease. *bioRxiv*, 219428 (2017).
962  25  Kang, H., Zhang, A., Cai, T. T. & Small, D. S. Instrumental variables estimation with some invalid
963      instruments and its application to Mendelian randomization. *Journal of the American statistical
964      Association* **111**, 132-144 (2016).
965  26  Guo, Z., Kang, H., Tony Cai, T. & Small, D. S. Confidence intervals for causal effects with invalid
966      instruments by using two‐stage hard thresholding with voting. *Journal of the Royal Statistical
967      Society: Series B (Statistical Methodology)* **80**, 793-815 (2018).
968  27  Burgess, S. & Thompson, S. G. Interpreting findings from Mendelian randomization using the
969      MR-Egger method. *European journal of epidemiology* **32**, 391-392 (2017).
970  28  Dai, J. Y. *et al.* Diagnostics of Pleiotropy in Mendelian Randomization Studies: Global and
971      Individual Tests for Direct Effects. *American journal of epidemiology* **187**, 2672-2680 (2018).
972  29  Zhao, Q., Wang, J., Bowden, J. & Small, D. S. Statistical inference in two-sample summary-data
973      Mendelian randomization using robust adjusted profile score. *arXiv:1801.09652* (2018).
974  30  Qi, G. & Chatterjee, N. Mendelian Randomization Analysis Using Mixture Models (MRMix) for
975      Genetic Effect-Size-Distribution Leads to Robust Estimation of Causal Effects. *bioRxiv*, 367821
976      (2018).
977  31  Berzuini, C., Guo, H., Burgess, S. & Bernardinelli, L. A Bayesian approach to Mendelian
978      randomization with multiple pleiotropic variants. *Biostatistics*, 1-16 (2018).
979  32  Li, S. Mendelian randomization when many instruments are invalid: hierarchical empirical Bayes
980      estimation. *arXiv:1706.01389* (2017).
981  33  Barfield, R. *et al.* Transcriptome-wide association studies accounting for colocalization using
982      Egger regression. *Genetic epidemiology* **42**, 418-433 (2018).
983  34  Dawid, A. P. Causal inference without counterfactuals. *Journal of the American statistical
984      Association* **95**, 407-424 (2000).
985  35  Dawid, A. P. Statistical causality from a decision-theoretic perspective. *Annual Review of
986      Statistics and Its Application* **2**, 273-303 (2015).
987  36  Berzuini, C., Dawid, P. & Bernardinell, L. *Causality: Statistical perspectives and applications*.
988      (John Wiley & Sons, 2012).
989  37  Lappalainen, T. *et al.* Transcriptome and genome sequencing uncovers functional variation in
990      humans. *Nature* **501**, 506-511 (2013).
991  38  Banda, Y. *et al.* Characterizing race/ethnicity and genetic ancestry for 100,000 subjects in the
992      Genetic Epidemiology Research on Adult Health and Aging (GERA) cohort. *Genetics* **200**, 1285-
993      1295 (2015).
994  39  Kvale, M. N. *et al.* Genotyping informatics and quality control for 100,000 subjects in the Genetic
995      Epidemiology Research on Adult Health and Aging (GERA) cohort. *Genetics* **200**, 1051-1060
996      (2015).
997  40  Price, A. L. *et al.* Effects of cis and trans genetic ancestry on gene expression in African
998      Americans. *Plos Genetics* **4**, e1000294 (2008).

999    41    Price, A. L. *et al.* Single-tissue and cross-tissue heritability of gene expression via identity-by-descent in related or unrelated individuals. *Plos Genetics* **7**, e1001317 (2011).

1001   42    Consortium, W. T. C. C. Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature* **447**, 661-678 (2007).

1003   43    Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203-209 (2018).

1005   44    Wen, X., Luca, F. & Pique-Regi, R. Cross-population joint analysis of eQTLs: fine mapping and functional annotation. *Plos Genetics* **11**, e1005176 (2015).

1007   45    Harrow, J. *et al.* GENCODE: the reference human genome annotation for The ENCODE Project. *Genome research* **22**, 1760-1774 (2012).

1009   46    Stegle, O., Parts, L., Piipari, M., Winn, J. & Durbin, R. Using probabilistic estimation of expression residuals (PEER) to obtain increased power and interpretability of gene expression analyses. *Nature protocols* **7**, 500-507 (2012).

1012   47    Guan, Y. & Stephens, M. Practical issues in imputation-based association mapping. *Plos Genetics* **4**, e1000279 (2008).

1014   48    Howie, B. N., Donnelly, P. & Marchini, J. A flexible and accurate genotype imputation method for the next generation of genome-wide association studies. *Plos Genetics* **5**, e1000529 (2009).

1016   49    Delaneau, O., Zagury, J.-F. & Marchini, J. Improved whole-chromosome phasing for disease and population genetic studies. *Nature methods* **10**, 5-6 (2012).

1018   50    McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nature genetics* **48**, 1279-1283 (2016).

1020   51    Das, S. *et al.* Next-generation genotype imputation service and methods. *Nature genetics* **48**, 1284-1287 (2016).

1022   52    Loh, P.-R., Kichaev, G., Gazal, S., Schoech, A. P. & Price, A. L. Mixed-model association for biobank-scale datasets. *Nature genetics* **50**, 906-908 (2018).

1024   53    Zhou, X., Carbonetto, P. & Stephens, M. Polygenic Modeling with Bayesian Sparse Linear Mixed Models. *Plos Genetics* **9**, : e1003264. (2013).

1026   54    Zhou, X. & Stephens, M. Genome-wide efficient mixed-model analysis for association studies. *Nature genetics* **44**, 821-824 (2012).

1028   55    Mancuso, N. *et al.* Probabilistic fine-mapping of transcriptome-wide association studies. *Nature genetics* **51**, 675 (2019).

1030   56    Berisa, T. & Pickrell, J. K. Approximately independent linkage disequilibrium blocks in human populations. *Bioinformatics* **32**, 283 (2016).

1032   57    Barbeira, A. N. *et al.* Exploring the phenotypic consequences of tissue specific gene expression variation inferred from GWAS summary statistics. *Nature communications* **9**, 1825 (2018).

1034   58    Yengo, L. *et al.* Meta-analysis of genome-wide association studies for height and body mass index in approximately 700000 individuals of European ancestry. *Hum Mol Genet* **27**, 3641-3649, doi:10.1093/hmg/ddy271 (2018).

1037   59    Kamatani, Y. *et al.* Genome-wide association study of hematological and biochemical traits in a Japanese population. *Nature genetics* **42**, 210 (2010).

1039   60    Soranzo, N. *et al.* A genome-wide meta-analysis identifies 22 loci associated with eight hematological parameters in the HaemGen consortium. *Nature genetics* **41**, 1182 (2009).

1041   61    Auer, P. L. *et al.* Rare and low-frequency coding variants in CXCR2 and other genes are associated with hematological traits. *Nature genetics* **46**, 629 (2014).

1043   62    Bersenev, A., Wu, C., Balcerek, J. & Tong, W. Lnk controls mouse hematopoietic stem cell self-renewal and quiescence through direct interactions with JAK2. *The Journal of clinical investigation* **118**, 2832-2844 (2008).

63    Tong, W. & Lodish, H. F. Lnk inhibits Tpo–mpl signaling and Tpo-mediated megakaryocytopoiesis. *Journal of Experimental Medicine* **200**, 569-580 (2004).

64    Takizawa, H. *et al.* Lnk regulates integrin αIIbβ3 outside-in signaling in mouse platelets, leading to stabilization of thrombus development in vivo. *The Journal of clinical investigation* **120**, 179-190 (2010).

65    Viny, A. D. & Levine, R. L. Genetics of myeloproliferative neoplasms. *Cancer journal (Sudbury, Mass.)* **20**, 61 (2014).

66    Bersenev, A. *et al.* Lnk constrains myeloproliferative diseases in mice. *The Journal of clinical investigation* **120**, 2058-2069 (2010).

67    Yamamoto, S. & Ma, X. Role of Nod2 in the development of Crohn's disease. *Microbes and infection* **11**, 912-918 (2009).

68    McGovern, D., Van Heel, D., Ahmad, T. & Jewell, D. NOD2 (CARD15), the first susceptibility gene for Crohn's disease. *Gut* **49**, 752-754 (2001).

69    Hugot, J.-P. *et al.* Mapping of a susceptibility locus for Crohn's disease on chromosome 16. *Nature* **379**, 821 (1996).

70    Hugot, J.-P. *et al.* Association of NOD2 leucine-rich repeat variants with susceptibility to Crohn's disease. *Nature* **411**, 599 (2001).

71    Ogura, Y. *et al.* A frameshift mutation in NOD2 associated with susceptibility to Crohn's disease. *Nature* **411**, 603 (2001).

72    Franke, A. *et al.* Genome-wide meta-analysis increases to 71 the number of confirmed Crohn's disease susceptibility loci. *Nature genetics* **42**, 1118 (2010).

73    Franke, A. *et al.* Replication of signals from recent studies of Crohn's disease identifies previously unknown disease loci for ulcerative colitis. *Nature genetics* **40**, 713 (2008).

74    Kennedy, N. A. *et al.* The impact of NOD2 variants on fecal microbiota in Crohn's disease and controls without gastrointestinal disease. *Inflammatory bowel diseases* **24**, 583-592 (2018).

75    Sidiq, T., Yoshihama, S., Downs, I. & Kobayashi, K. S. Nod2: a critical regulator of ileal microbiota and Crohn's disease. *Front Immunol* **7**, 367 (2016).

76    Keel, S. B. *et al.* Evidence that the expression of transferrin receptor 1 on erythroid marrow cells mediates hepcidin suppression in the liver. *Experimental hematology* **43**, 469-478. e466 (2015).

77    Andrews, N. C. Genes determining blood cell traits. *Nature genetics* **41**, 1161 (2009).

78    Ganesh, S. K. *et al.* Multiple loci influence erythrocyte phenotypes in the CHARGE Consortium. *Nature genetics* **41**, 1191 (2009).

79    Lo, K. S. *et al.* Genetic association analysis highlights new loci that modulate hematological trait variation in Caucasians and African Americans. *Human genetics* **129**, 307-317 (2011).

80    Levy, J. E., Jin, O., Fujiwara, Y., Kuo, F. & Andrews, N. Transferrin receptor is necessary for development of erythrocytes and the nervous system. *Nature genetics* **21**, 396 (1999).

81    Kanai, M. *et al.* Genetic analysis of quantitative traits in the Japanese population links cell types to complex human diseases. *Nature genetics* **50**, 390-400 (2018).

82    Astle, W. J. *et al.* The allelic landscape of human blood cell trait variation and links to common complex disease. *Cell* **167**, 1415-1429 (2016).

83    Wood, A. R. *et al.* Defining the role of common variation in the genomic and biological architecture of adult human height. *Nature genetics* **46**, 1173-1186 (2014).

84    Wainberg, M. *et al.* Opportunities and challenges for transcriptome-wide association studies. *Nature genetics* **51**, 592 (2019).

85    Wu, M. C. *et al.* Rare-variant association testing for sequencing data with the sequence kernel association test. *The American Journal of Human Genetics* **89**, 82-93 (2011).

86    Zhao, N. *et al.* Testing in microbiome-profiling studies with MiRKAT, the microbiome regression-based kernel association test. *The American Journal of Human Genetics* **96**, 797-807 (2015).

1094    87    Chen, J., Chen, W., Zhao, N., Wu, M. C. & Schaid, D. J. Small sample kernel association tests for
1095          human genetic and microbiome association studies. *Genetic epidemiology* **40**, 5-19 (2016).
1096    88    Yang, J., Fritsche, L. G., Zhou, X., Abecasis, G. & Consortium, I. A.-R. M. D. G. A scalable Bayesian
1097          method for integrating functional information in genome-wide association studies. *The
1098          American Journal of Human Genetics* **101**, 404-416 (2017).
1099    89    Crawford, L., Zeng, P., Mukherjee, S. & Zhou, X. Detecting epistasis with the marginal epistasis
1100          test in genetic mapping studies of quantitative traits. *Plos Genetics* **13**, e1006869 (2017).

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

47

## Acknowledgements

## Author contributions

XZ conceived the idea and provided funding support. XZ and ZY developed the methods. ZY developed the software tool with assistance from JL and CY. ZY performed simulations and real data analysis with assistance from HZ, PZ, SY, and SS. XZ and ZY wrote the manuscript with input from all other authors. All authors reviewed and approved the final manuscript.

**Competing interests:** The authors declare no competing interests.

**Table 1** Summary of some existing MR methods

| | Design | Instrumental variable | β effect assumption | γ effect assumption | Estimation procedure |
|---|---|---|---|---|---|
| PrediXcan[1] | Two-sample | Correlated | Elastic net | N/A | Two-stage |
| TWAS[2] | Two-sample | Correlated | BSLMM | N/A | Two-stage |
| SMR[3] | Two-sample | Univariate | Fixed effect | N/A | Two-stage |
| GSMR[4] | Two-sample | independent | Fixed effect | N/A | Two-stage |
| MR-Egger[13] | Two-sample | Independent | Fixed effect | Equal effect size | Two-stage |
| CoMM[18] | Two-sample | Correlated | Normal | N/A | MLE |
| CaMMEL[24] | Two-sample | Correlated | Fixed effect | Normal | Variational Bayes |
| Kang et al.[25] | One-sample | Correlated | Fixed effect | Lasso | Two-stage |
| MRMix[30] | Two-sample | Independent | Normal Mixture | Normal Mixture | Estimating equation |
| Berzuini et al.[31] | One-sample | Correlated | Fixed effect | Horseshoe | MCMC |
| LDA MR-Egger[32] | Two-sample | Correlated | Fixed effect | Equal effect size | Two-stage |
| DPR[19] | Two-sample | Correlated | Latent Dirichlet process | N/A | Two-stage |
| TIGAR[20] | Two-sample | Correlated | Latent Dirichlet process | N/A | Two-stage |
| **PMR-Egger** | **Two-sample** | **Correlated** | **Normal** | **Equal effect size** | **MLE** |

Methods are categorized based on the experimental design (two-sample vs one-sample vs both), the characterizes of selected instrumental variables (univariate vs multiple independent vs multiple correlated), $\beta$ effect size assumption, $\gamma$ effect size assumption, estimation/inference procedure (ratio-based vs two-stage estimation vs maximum likelihood vs Bayesian), and input data type (individual-level vs summary; which is now removed per reviewer's request). The categorization of inference procedure generally follows ref [5]. In the inference procedure, the two-stage estimation procedure comprises two regression stages: the first-stage regression of the exposure on the instrumental variables, and the second-stage regression of the outcome on the fitted values of the exposure from the first stage. Some inference procedures, such as the inverse variance weighted (IVW) procedure (e.g. MR-Egger[13]) or the ratio method (e.g. for SMR[3]) are categorized as two-stage procedure here, as both are asymptotically equivalent to a two-stage estimation procedure in the case of independent instruments. We only list MR methods that directly take input instruments into the model; many MR methods that performs various selection procedures on the instruments (e.g. Guo et al[26]) are not included in the table. Some recently developed methods that only test for horizontal pleiotropy, such as GLIDE[28] and MR-PRESSO[21] are not included in the table.

**Table 2.** Mean computational time (in second) of various MR methods

| Trait | #SNP in the exemplary gene | CoMM | PMR-Egger | TWAS | LDA MR-Egger | SMR | PrediXcan | MR-PRESSO |
|---|---|---|---|---|---|---|---|---|
| T1D from WTCCC (n=4901) | 300 | 0.51(0.19) | 0.80(0.57) | 1.97(0.86) | 0.08(0.02) | 0.0003(0.0005) | 26.74(2.81) | 408.27(74.76) |
| | 500 | 1.21(0.41) | 1.42(0.77) | 3.48(1.16) | 0.14(0.03) | 0.0004(0.0005) | 11.77(0.64) | 829.04(135.79) |
| | 983 | 5.85(1.50) | 9.79(1.56) | 4.69(1.73) | 0.60(0.09) | 0.0004(0.0005) | 9.96(0.78) | 2023.77(260.43) |
| | 2106 | 111.00(12.87) | 97.33(7.63) | 5.87(2.26) | 4.18(0.59) | 0.0005(0.0005) | 22.90(2.63) | 4913.22(554.47) |
| Asthma from GERA (n=61,953) | 300 | 1.47(0.29) | 2.06(0.22) | 2.61(1.48) | 0.05(0.02) | 0.0002(0.0004) | 33.39(3.09) | 464.64(62.18) |
| | 500 | 1.21(0.33) | 4.21(0.81) | 2.54(0.87) | 0.09(0.03) | 0.0002(0.0004) | 11.71(0.70) | 919.66(102.83) |
| | 1000 | 24.37(5.13) | 21.68(1.66) | 3.07(2.55) | 0.46(0.13) | 0.0002(0.0004) | 14.29(1.30) | 2275.42(263.95) |
| | 2008 | 59.01(4.98) | 52.52(4.47) | 4.51(1.48) | 2.33(0.71) | 0.0004(0.0005) | 20.18(3.28) | 5213.73(601.46) |
| Platelet Count from UK Biobank (n=337,198) | 300 | 2.56(0.53) | 5.57(4.54) | 5.04(4.19) | 0.09(0.02) | 0.0008(0.0004) | 10.93(1.96) | 471.55(50.44) |
| | 500 | 6.82(2.75) | 7.61(2.30) | 5.44(4.30) | 0.15(0.02) | 0.0007(0.0005) | 12.17(1.04) | 876.06(92.90) |
| | 1052 | 24.92(6.28) | 23.59(3.21) | 5.91(4.79) | 0.81(0.09) | 0.0008(0.0004) | 16.05(2.38) | 2133.03(77.56) |
| | 2605 | 186.14(28.45) | 178.68(16.75) | 5.37(0.73) | 8.11(1.20) | 0.0008(0.0004) | 9.89(1.74) | 6949.72(245.75) |

Computation is carried out on a single thread of a Xeon Gold 6138 CPU. The computation time is averaged across 20 replicates, with values inside parentheses denoting the standard deviation. #SNP denotes the number of cis-SNPs for four exemplary genes in each study. The computational time for MR-PRESSO is based on 10,000 permutations.
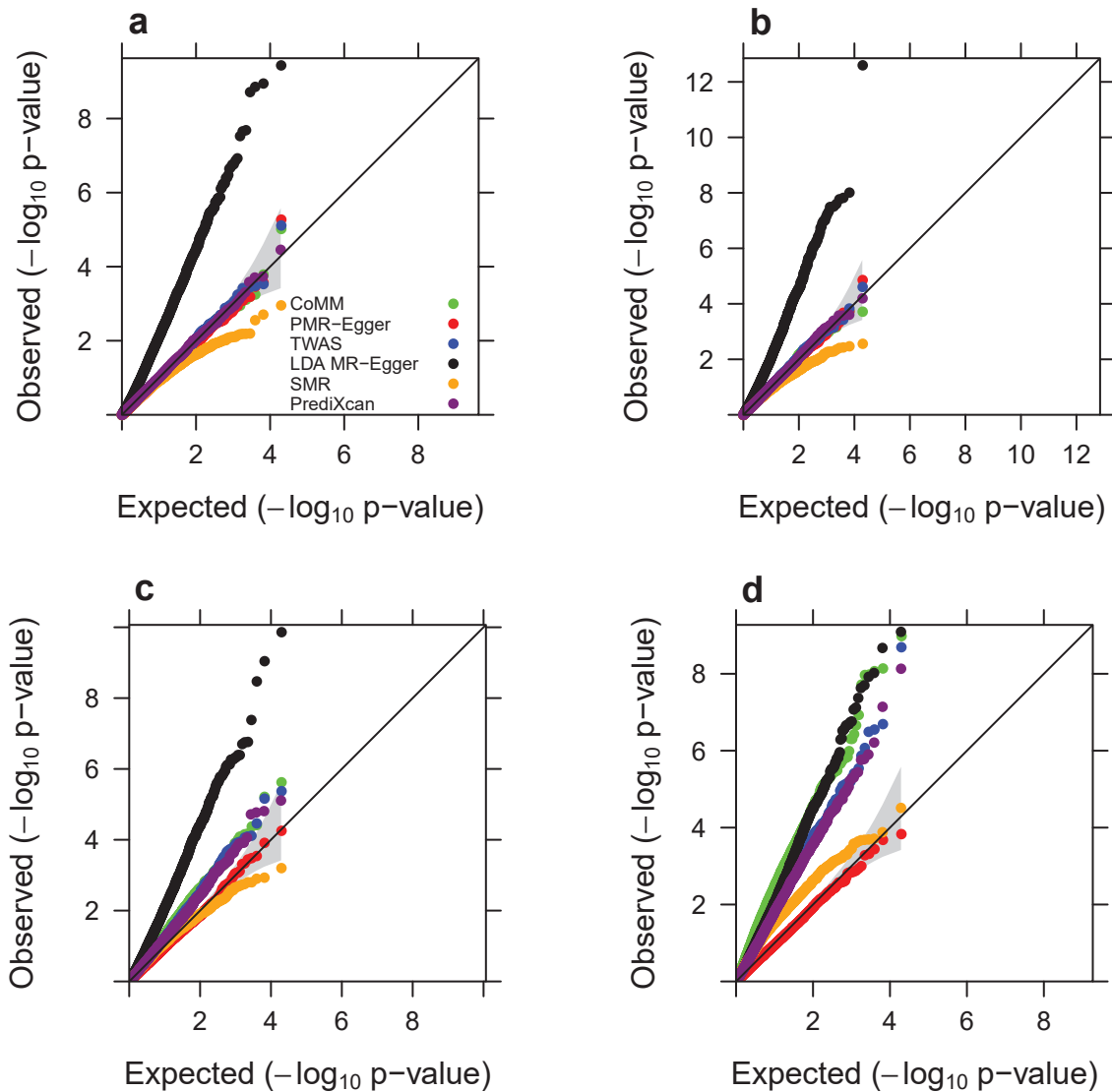
**Fig. 1** Quantile-quantile plot of -log10 p-values from different methods for testing the causal effect either in the absence or in the presence of horizontal pleiotropic effect under null simulations. Compared methods include CoMM (green), PMR-Egger (red), TWAS (blue), LDA MR-Egger (black), SMR (orange), and PrediXcan (purple). Null simulations are performed under different horizontal pleiotropic effect sizes: (a) $\gamma=0$; (b) $\gamma=0.0001$; (c) $\gamma=0.0005$; (d) $\gamma=0.001$. Only p-values from PMR-Egger adhere to the expected diagonal line across a range of horizontal pleiotropic effect sizes.
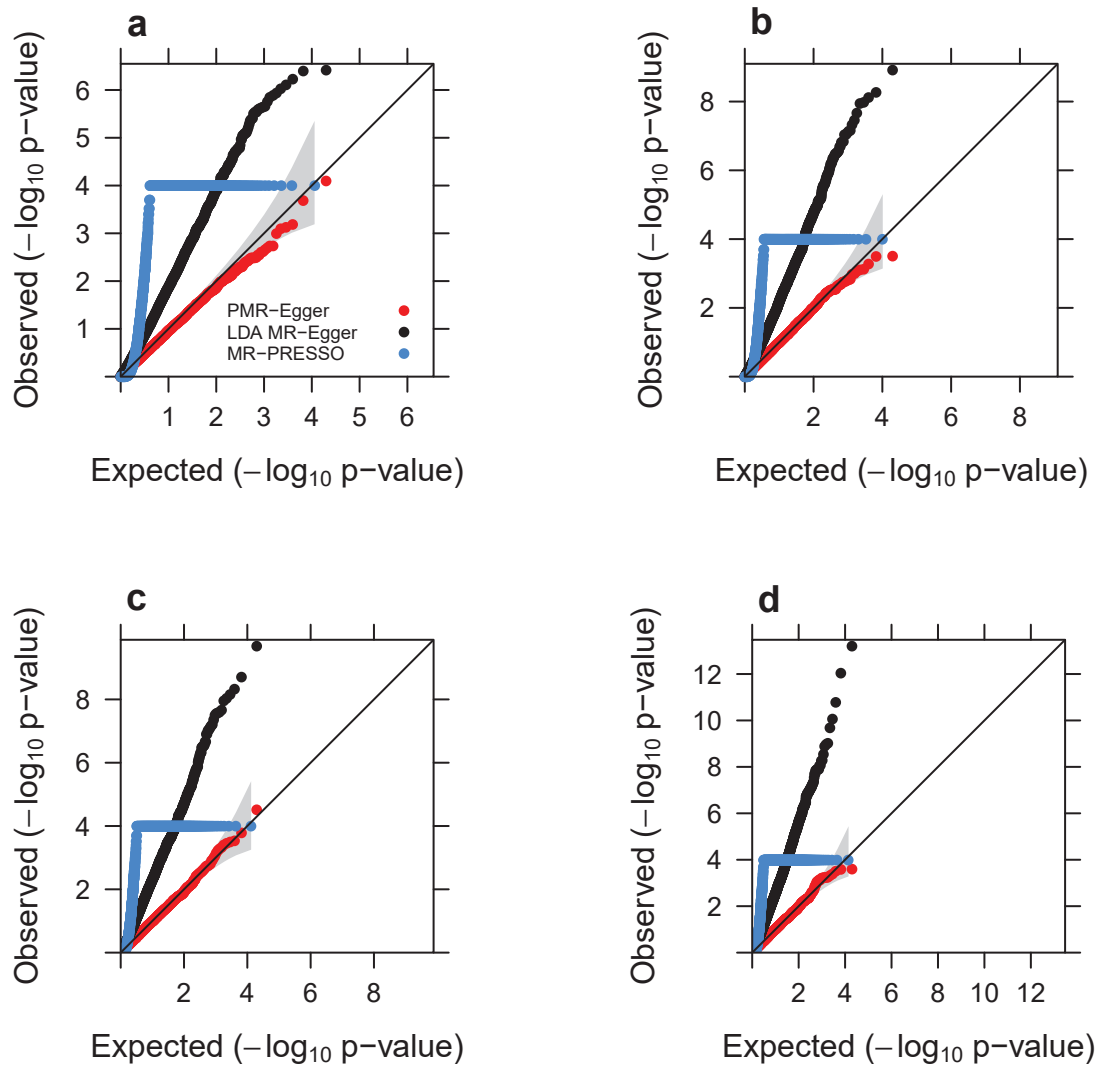
**Fig. 2** Power of different methods under various simulation scenarios. Power (y-axis) at a false discovery rate of 0.1 to detect the causal effect (a-d) or the horizontal pleiotropic effect (e-f) is plotted against different causal effect size characterized by PVE_zy (x-axis). Compared methods include CoMM (green), PMR-Egger (red), TWAS (blue), LDA MR-Egger (black), SMR (orange), and PrediXcan (purple). Simulations are performed under different horizontal pleiotropic effect sizes: (a) $\gamma=0$; (b) $\gamma=0.0001$; (c, e) $\gamma=0.0005$; (d, f) $\gamma=0.001$.

**Fig. 3** Quantile-quantile plot of -log10 p-values from different methods for testing the horizontal pleiotropic effect either in the absence or in the presence of causal effect under null simulations. Compared methods include PMR-Egger (red), LDA MR-Egger (black), and MR-PRESSO (dodger blue). Null simulations are performed under different causal effect sizes characterized by $PVE_{zy}$: **(a)** $PVE_{zy}=0$; **(b)** $PVE_{zy}=0.2\%$; **(c)** $PVE_{zy}=0.4\%$; and **(d)** $PVE_{zy}=0.6\%$. Only p-values from PMR-Egger adhere to the expected diagonal line across a range of horizontal pleiotropic effect sizes. Due to heavy computational burden, we are only able to run 10,000 permutations for MR-PRESSO. Therefore, the minimal p-value from MR-PRESSO is $10^{-4}$.
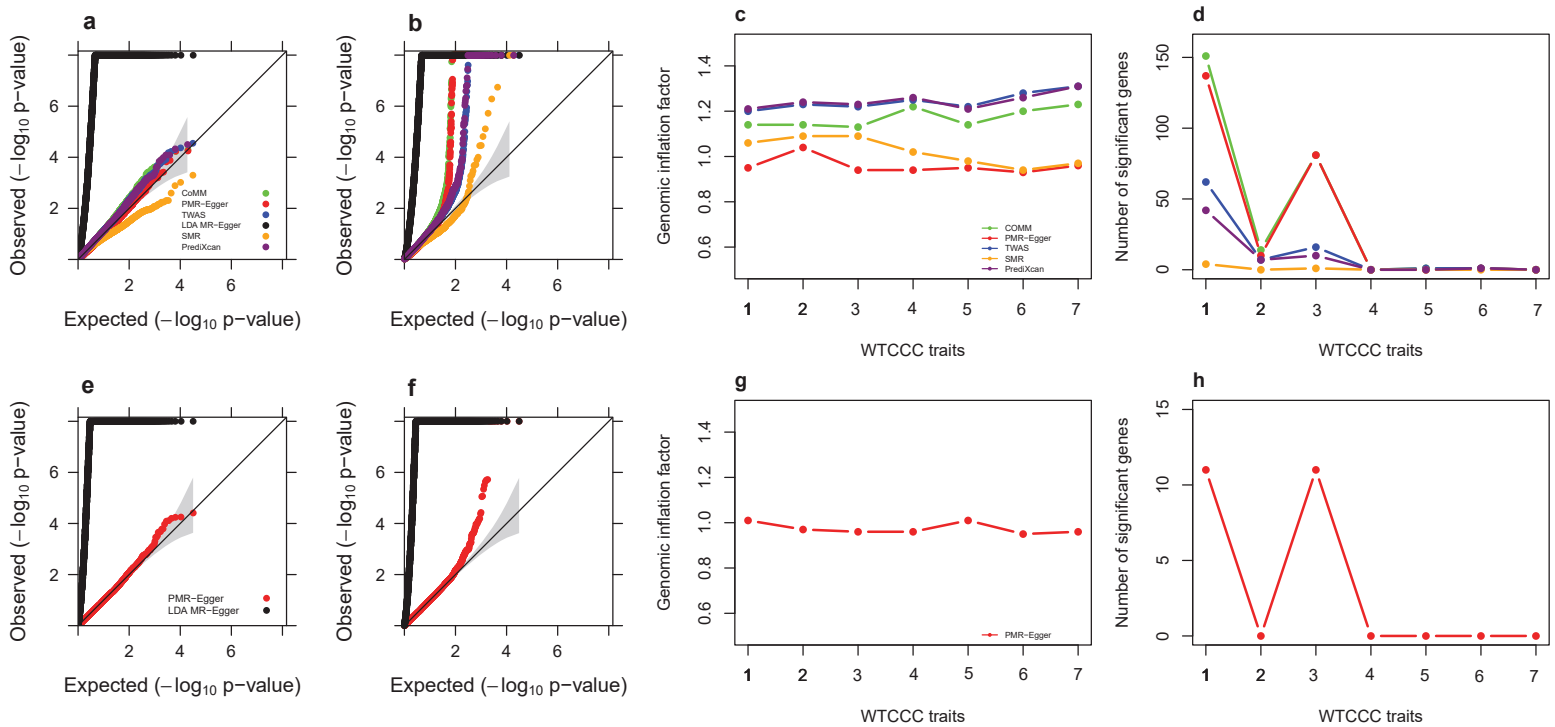
**Fig. 4** TWAS analysis results by different methods for traits in the WTCCC data. Compared methods include CoMM (green), PMR-Egger (red), TWAS (blue), LDA MR-Egger (black), SMR (orange), and PrediXcan (purple). **(a)** Quantile-quantile plot of -log10 p-values from different methods for testing the causal effect for an exemplary trait BD. **(b)** Quantile-quantile plot of -log10 p-values from different methods for testing the causal effect for another exemplary trait T1D. **(c)** Genomic inflation factor for testing the causal effect for each of the 7 traits by different methods. **(d)** Number of causal genes identified for each of the 7 traits by different methods. **(e)** Quantile-quantile plot of -log10 p-values from different methods for testing the horizontal pleiotropic effect for an exemplary trait BD. **(f)** Quantile-quantile plot of -log10 p-values from different methods for testing the horizontal pleiotropic effect for another exemplary trait T1D. **(g)** Genomic inflation factor for testing the horizontal pleiotropic effect for each of the 7 traits by different methods. **(h)** Number of genes identified to have significant horizontal pleiotropic effect for each of the 7 traits by different methods. For c, d, g, h, the number on the x-axis represents seven traits in order: T1D, CD, RA, BD, T2D, CAD, HT.
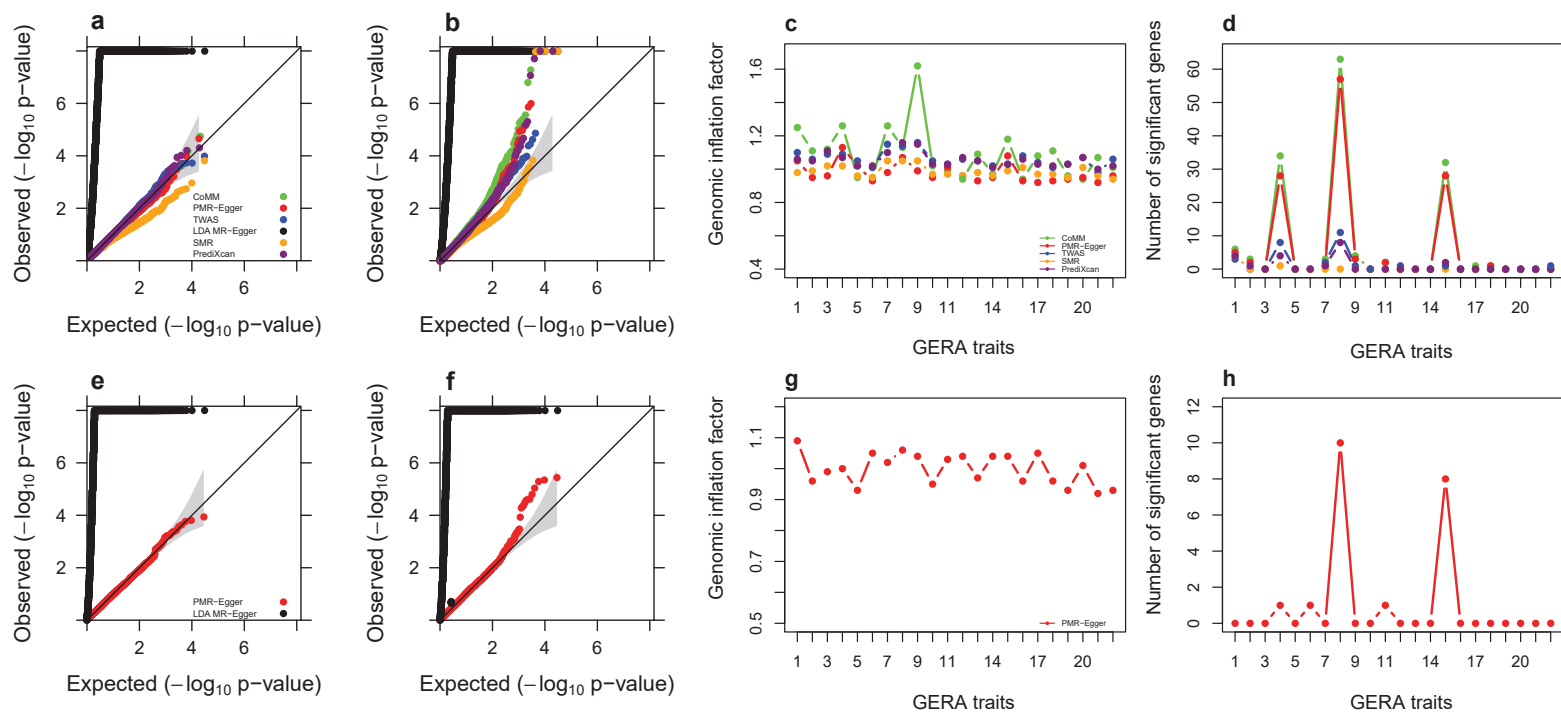
**Fig. 5** TWAS analysis results by different methods for traits in the GERA data. Compared methods include CoMM (green), PMR-Egger (red), TWAS (blue), LDA MR-Egger (black), SMR (orange), and PrediXcan (purple). **(a)** Quantile-quantile plot of -log10 p-values from different methods for testing the causal effect for an exemplary trait Irritable Bowel Syndrome. **(b)** Quantile-quantile plot of -log10 p-values from different methods for testing the causal effect for another exemplary trait Asthma. **(c)** Genomic inflation factor for testing the causal effect for each of the 22 traits by different methods. **(d)** Number of causal genes identified for each of the 22 traits by different methods. **(e)** Quantile-quantile plot of -log10 p-values from different methods for testing the horizontal pleiotropic effect for an exemplary trait Irritable Bowel Syndrome. **(f)** Quantile-quantile plot of -log10 p-values from different methods for testing the horizontal pleiotropic effect for another exemplary trait Asthma. **(g)** Genomic inflation factor for testing the horizontal pleiotropic effect for each of the 22 traits by different methods. **(h)** Number of genes identified to have significant horizontal pleiotropic effect for each of the 22 traits by different methods. For c, d, g, h, the number on the x-axis represents 22 traits in order: Asthma, Allergic Rhinitis, CARD, Cancers, Depressive Disorder, Dermatophytosis, T2D, Dyslipidemia, HT, Hemorrhoids, Abdominal Hernia, Insomnia, Iron Deficiency, Irritable Bowel Syndrome, Macular Degeneration, Osteoarthritis, Osteoporosis, PVD, Peptic Ulcer, Psychiatric disorders, Stress Disorders, Varicose Veins.
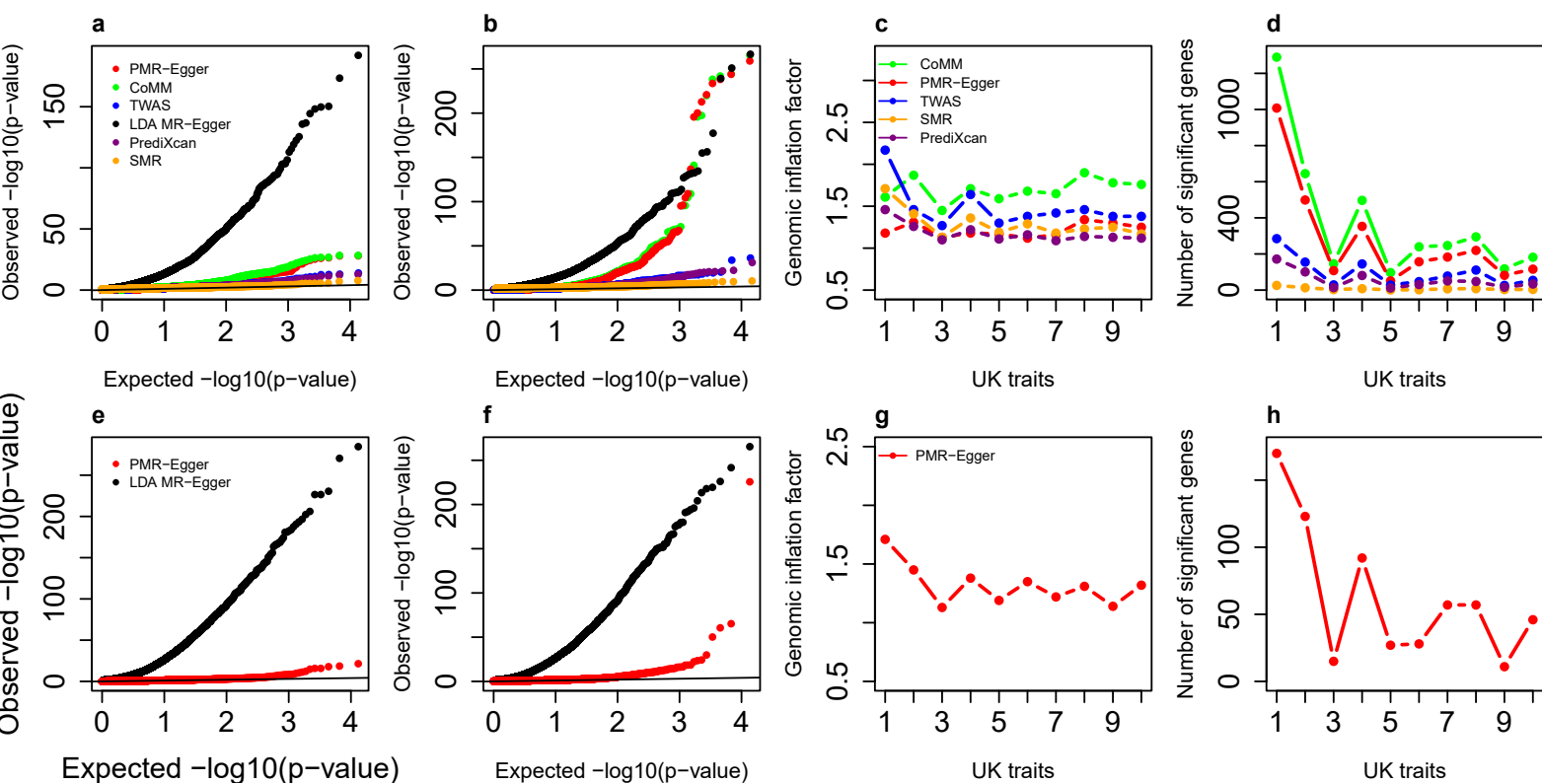
**Fig. 6** TWAS analysis results by different methods for traits in the UK Biobank data. Compared methods include CoMM (green), PMR-Egger (red), TWAS (blue), LDA MR-Egger (black), SMR (orange), and PrediXcan (purple). **(a)** Quantile-quantile plot of -log10 p-values from different methods for testing the causal effect for an exemplary trait BMI. **(b)** Quantile-quantile plot of -log10 p-values from different methods for testing the causal effect for another exemplary trait Platelet Count. **(c)** Genomic inflation factor for testing the causal effect for each of the 10 traits by different methods. **(d)** Number of causal genes identified for each of the 10 traits by different methods. **(e)** Quantile-quantile plot of -log10 p-values from different methods for testing the horizontal pleiotropic effect for an exemplary trait BMI. **(f)** Quantile-quantile plot of -log10 p-values from different methods for testing the horizontal pleiotropic effect for another exemplary trait Platelet Count. **(g)** Genomic inflation factor for testing the horizontal pleiotropic effect for each of the 10 traits by different methods. **(h)** Number of genes identified to have significant horizontal pleiotropic effect for each of the 10 traits by different methods. For c, d, g, h, the number on the x-axis represents 10 traits in order: Height, Platelet count, Bone mineral density, Red blood cell count, FEV1-FVC ratio, BMI, RDW, Eosinophils count, Forced vital capacity, White blood cell count.