

Nonlinear Ridge Regression Improves Robustness of Cell-Type-Specific Differential Expression Studies

Fumihiko Takeuchi (✉ fumihiko@takeuchi.name)

Kokuritsu Kenkyu Kaihatsu Hojin Kokuritsu Kokusai Iryo Kenkyu Center <https://orcid.org/0000-0003-3185-5661>

Norihiro Kato

Kokuritsu Kenkyu Kaihatsu Hojin Kokuritsu Kokusai Iryo Kenkyu Center

Methodology article

Keywords: Epigenome-wide association study, Differential gene expression analysis, Cell type, Ridge regression, Nonlinear regression, mQTL, eQTL

Posted Date: July 1st, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-39226/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.
[Read Full License](#)

Version of Record: A version of this preprint was published on March 22nd, 2021. See the published version at <https://doi.org/10.1186/s12859-021-03982-3>.

Nonlinear ridge regression improves robustness of cell-type-specific differential expression studies

Fumihiko Takeuchi and Norihiro Kato

Department of Gene Diagnostics and Therapeutics, Research Institute,
National Center for Global Health and Medicine (NCGM), Tokyo, Japan

Correspondence

Fumihiko Takeuchi

Department of Gene Diagnostics and Therapeutics, Research Institute,
National Center for Global Health and Medicine (NCGM)

1-21-1 Toyama, Shinjuku-ku, Tokyo, 162-8655, Japan

Email: fumihiko@takeuchi.name

1 **Abstract**

2 **Background:** Epigenome-wide association studies (EWAS) and differential
3 gene expression analyses are generally performed on tissue samples, which
4 consist of multiple cell types. Cell-type-specific effects of a trait, such as
5 disease, on the omics expression are of interest but difficult or costly to
6 measure experimentally. By measuring omics data for the bulk tissue, cell
7 type composition of a sample can be inferred statistically. Subsequently, cell-
8 type-specific effects are estimated by linear regression that includes terms
9 representing the interaction between the cell type proportions and the trait.
10 This approach involves two issues, scaling and multicollinearity.

11 **Results:** First, although cell composition is analyzed in linear scale,
12 differential methylation/expression is analyzed suitably in the logit/log scale.
13 To simultaneously analyze two scales, we developed nonlinear regression.
14 Second, we show that the interaction terms are highly collinear, which is
15 obstructive to ordinary regression. To cope with the multicollinearity, we
16 applied ridge regularization. In simulated and real data, the improvement was
17 modest by nonlinear regression and substantial by ridge regularization.

18 **Conclusion:** Nonlinear ridge regression performed cell-type-specific
19 association test on bulk omics data more robustly than previous methods.
20 The omicwas package for R implements nonlinear ridge regression for cell-
21 type-specific EWAS, differential gene expression and QTL analyses. The
22 software is freely available from <https://github.com/fumi-github/omicwas>
23

24 **Keywords**

25 Epigenome-wide association study, Differential gene expression analysis, Cell
26 type, Ridge regression, Nonlinear regression, mQTL, eQTL
27

28 **Background**

29 Epigenome-wide association studies (EWAS) and differential gene expression
30 analyses elucidate the association of disease traits (or conditions) with the
31 level of omics expression, namely DNA methylation and gene expression.
32 Thus far, tissue samples, which consist of heterogeneous cell types, have
33 mainly been examined, because cell sorting is not feasible in most tissues
34 and single-cell assay is still expensive. Nevertheless, the cell type
35 composition of a sample can be quantified statistically by comparing omics
36 measurement of the target sample with reference data obtained from sorted
37 or single cells [1,2]. By utilizing the composition, the disease association
38 specific to a cell type was statistically inferred for gene expression [3-10] and
39 DNA methylation [11-14].

40 For the imputation of cell type composition, omics markers are usually
41 analyzed in the original linear scale, which measures the proportion of mRNA
42 molecules from a specific gene or the proportion of methylated cytosine
43 molecules among all cytosines at a specific CpG site [15]. The proportion can
44 differ between cell types, and the weighted average of cell-type-specific
45 proportions becomes the proportion in a bulk tissue sample. Using the fact
46 that the weight equals the cell type composition, the cell type composition of
47 a sample is imputed. In contrast, gene expression analyses are performed in
48 the log-transformed scale because the signal and noise are normally
49 distributed after log-transformation [16]. In DNA methylation analysis, the
50 logit-transformed scale, which is called the M-value, is statistically valid [17].
51 Consequently, the optimal scales for analyzing differential gene expression
52 or methylation can differ from the optimal scale for analyzing cell type
53 composition.

54 Aiming to perform cell-type-specific EWAS or differential gene expression
55 analyses by using unsorted tissue samples, we study two issues that have
56 been overlooked. Whereas previous studies were performed in linear scale,
57 we develop a nonlinear regression, which simultaneously analyzes cell type
58 composition in linear scale and differential expression/methylation in log/logit
59 scale. The second issue is multicollinearity. Cell-type-specific effects of a trait,
60 such as disease, on omics expression are usually estimated by linear

61 regression that includes terms representing the interaction between the cell
62 type proportions and the trait. We show that the interaction terms can
63 mutually be highly correlated, which obstructs ordinary regression. To cope
64 with the multicollinearity, we implement ridge regularization. Our methods
65 and previous ones are compared in simulated and real data.
66

67 **Results**

68 **Multicollinearity of interaction terms**

69 Typically, cell-type-specific effects of a trait on omics marker expression is
70 analyzed by the linear regression in equation (2). The goal is to estimate $\beta_{h,k}$
71 the effect of trait k on the expression level in cell type h . This is estimated
72 based on the relation between the bulk expression level Y_i of a sample and
73 the regressor $W_{h,i}X_{i,k}$, which is an interaction term defined as the product of
74 the cell type proportion $W_{h,i}$ and the trait value $X_{i,k}$ of the sample. The
75 variable W_h for cell type composition cannot be mean-centered, and
76 interaction terms involving uncentered variables cause multicollinearity [18].
77 We first survey the extent of multicollinearity in real data for cell-type-specific
78 association.

79 In peripheral blood leukocyte data from a rheumatoid arthritis study
80 (GSE42861), the proportion of cell types ranged from 0.59 for neutrophils to
81 0.01 for eosinophils (Table 1A). The proportion of neutrophils was negatively
82 correlated with the proportion of other cell types (apart from monocytes) with
83 correlation coefficient of -0.68 to -0.46 , whereas the correlation was weaker
84 for other pairs (Table 1B). Rheumatoid arthritis status was modestly
85 correlated with proportions of cell types. The product of the disease status
86 X_k , centered to have zero mean, and the proportion of a cell type becomes
87 an interaction term. The correlation coefficients between the interaction
88 terms are mostly >0.8 , apart from eosinophils (Table 1C). The ratio of mean
89 to SD of the proportion is high for all cell types apart from eosinophils (Table
90 1A). The interaction terms for high-ratio cell types are strongly correlated

91 with X_k , which in turn causes strong correlation between the relevant
92 interaction terms.

93 The situation was the same for the interaction with age in GTEx data. The
94 granulocytes (which include neutrophils and eosinophils) were the most
95 abundant (Table 2A). The proportion of granulocytes was negatively
96 correlated with other cell types (apart from monocytes) with correlation
97 coefficient of -0.89 to -0.41 , and the correlation between other pairs was
98 generally weaker (Table 2B). Age was modestly correlated with proportions
99 of cell types. In this dataset, the ratio of mean to SD of the proportion was
100 high in all cell types (Table 2A), which caused strong mutual correlation
101 between interaction terms (Table 2C).

102 In the above empirical data, multicollinearity between interaction terms
103 seemed to arise not due to the correlation between cell type proportions or
104 X_k , but due to the high ratio of mean to SD in the cell type proportions.
105 Subsequently, this property was derived mathematically. As we derived in
106 equation (17), the correlation between interaction terms $W_h X_k$ and $W_{h'} X_k$
107 approaches to one, when the ratios $E[W_h]/SD[W_h]$ and $E[W_{h'}]/SD[W_{h'}]$ are
108 high, irrespective of $Cor[W_h, W_{h'}]$ (Figure 1). The ratio was 1.6 to 5.3 (apart
109 from eosinophils) in the rheumatoid arthritis dataset and ≥ 4.3 in the GTEx
110 dataset. We looked up datasets of several ethnicities and found the ratio to
111 be ≥ 1.5 in majority of cell types (Additional file 1: Table S1). Thus,
112 multicollinearity can be a common problem for cell-type-specific association
113 analyses.

114 **Evaluation in simulated data**

115 By using simulated data, we evaluated previous methods and new
116 approaches of the omicwas package. In order to simultaneously analyze two
117 scales, the linear scale for heterogeneous cell mixing and the log/logit scale
118 for trait effects, we applied nonlinear regression in omicwas (equations (4)
119 and (5)). To cope with the multicollinearity of interaction terms, we applied
120 ridge regularization (equations (9) and (10)).

121 Previous regression type methods are based either on the full model of
122 linear regression (equation (2)) or the marginal model (equation (3)). The

123 full model fits and tests cell-type-specific effects for all cell types
124 simultaneously, and its derivatives include TOAST, csSAM.lm,
125 CellDMC.unfiltered and CellDMC.filtered. The marginal model fits and tests
126 cell-type-specific effect for one cell type at a time, and its derivatives include
127 csSAM.monovariate and TCA.

128 The simulation data was generated from real datasets of DNA methylation
129 and gene expression. The original cell type composition was retained for all
130 samples, and the case-control status was randomly assigned. In each sample,
131 expression level in each cell type was randomly determined according to a
132 scenario, and then averaged according to the sample's cell type composition.
133 Under each statistical algorithm, the disease association in the target cell type
134 was assessed by a Z-score, comparing cases vs controls.

135 In scenario A for DNA methylation, expression of all cell types had
136 identical distribution, irrespective of the case/control status (Figure 2A). The
137 type I error rate was controlled (≤ 0.05) in all algorithms. In scenario B, cases
138 had higher expression level in one randomly selected cell type, and that cell
139 type was tested (Figure 2B). Here, the most appropriate algorithm is the
140 marginal test applied to the perturbed cell type, which indeed attained the
141 highest power. For the most abundant neutrophils, the Z-score was in the
142 high range of 9.9 to 14.9 for the marginal test. With regards to the power,
143 the ridge regression methods (omicwas.identity.ridge and
144 omicwas.logit.ridge) came next. The algorithms based on full model, without
145 ridge regularization, (Full, TOAST, CellDMC.unfiltered, omicwas.identity and
146 omicwas.logit) gained modest power. TCA, which is similar to the marginal
147 test, detected neutrophil-specific association with high Z-score, but the power
148 over all cell types was modest. In scenario C, the expression level of cases
149 was lower in one cell type, which was not the tested cell type (Figure 2C).
150 Since the expression of the tested cell type is identical between cases and
151 controls, a correct algorithm should detect no signal. The type I error rate
152 was inflated, being highest for the marginal test, followed by the ridge
153 regression methods and TCA. Extremely strong spurious signals of Z-score $<$
154 -6 were detected in marginal and TCA. Scenario D combined scenarios B and
155 C, where the tested cell type had higher expression in cases, and one non-

156 tested cell type had lower expression in cases (Figure 2D). The distribution
157 of neutrophil Z-score was similar to scenario B, and the spurious signals with
158 low Z-scores were similar to scenario C. Over all scenarios, the similarity in
159 performance of omicwas.identity vs omicwas.logit, as well as
160 omicwas.identity.ridge vs omicwas.logit.ridge, indicates that the scaling was
161 not influential in DNA methylation data.

162 The results for simulated gene expression data were similar. In scenario
163 A with no true signal, type I error rate was controlled (≤ 0.05) in all algorithms
164 (Figure 3A). In scenario B, where true signal exists only for the tested cell
165 type, the power was the highest in marginal and relatively high in
166 csSAM.monovariate (Figure 3B). The power was in decreasing order,
167 omicwas.log.ridge > omicwas.identity.ridge > omicwas.log >
168 omicwas.identity; proper scaling modestly improved performance. In
169 scenario C, where cases have lower expression in one non-target cell type,
170 the type I error was inflated in the negative direction, with the largest inflation
171 in marginal, and moderate inflation in ridge regression methods and
172 csSAM.monovariate (Figure 3C). Extremely strong false signals of Z-score <
173 -6 occurred in marginal and csSAM.monovariate. In scenario D, where the
174 tested cell type has higher expression in cases, while one non-tested cell type
175 has lower expression, we could observe the overlay of power gain of scenario
176 B and type I error inflation of scenario C (Figure 3D).

177 Although we roughly grouped previous algorithms into derivatives of full
178 or derivatives of marginal, some implement treatments beyond simple linear
179 models. For example, the TCA algorithm tends to detect neutrophil signals
180 similarly as the marginal test (Fig. 2B), yet had smaller type I error rate (Fig.
181 2C).

182 **Cell-type-specific association with rheumatoid arthritis and age**

183 The cell-type-specific association of DNA methylation with rheumatoid
184 arthritis was predicted using bulk peripheral blood leukocyte data and was
185 evaluated in sorted monocytes (Figure 4A) and B cells (Figure 4B). Whereas
186 the full model (and its derivatives) performed the best and the marginal
187 model (and its derivatives) performed the worst in monocytes, the

188 performance ranking was opposite in B cells. A robust algorithm would
189 consistently achieve high performance relative to the best algorithm in each
190 instance. Nonlinear ridge regression (omicwas.logit.ridge) was the most
191 robust, performing 65% to 93% relative to the best method.

192 The cell-type-specific association of gene expression with age was
193 predicted using whole blood data and was evaluated in sorted CD4⁺ T cells
194 (Figure 4C) and monocytes (Figure 4D). All algorithms performed poorly in
195 CD4⁺ T cells, and the marginal model performed the best in monocytes.
196 Overall, nonlinear ridge regression (omicwas.log.ridge) was next to the
197 marginal model, performing 21% to 47% to the marginal.

198 For dataset GSE42861 and for GTEx whole blood, the omicwas.logit.ridge
199 and omicwas.log.ridge models of the omicwas package was computed in 8.1
200 and 0.7 hours respectively, using 8 cores of a 2.5 GHz Xeon CPU Linux server.
201

202 Discussion

203 Aiming to elucidate cell-type-specific trait association in DNA methylation and
204 gene expression, this article explored two aspects, multicollinearity and scale.
205 We observed multicollinearity in real data and derived mathematically how it
206 emerges. To cope with the multicollinearity, we proposed ridge regression. To
207 properly handle multiple scales simultaneously, we developed nonlinear
208 regression. By testing in simulated and real data, we found proper scaling to
209 modestly improve performance. In contrast, ridge regression achieved
210 performance that was more robust than previous methods.

211 The statistical methods discussed in this article are applicable, in principle,
212 to any tissue. For validation of the methods, we need datasets for bulk tissue
213 as well as sorted cells, ideally of >100 samples. Currently, the publicly
214 available data is limited to peripheral blood. By no means, we claim the
215 rheumatoid arthritis EWAS datasets [19-21] or the datasets for age
216 association of gene expression [22,23] to be representative. Nevertheless,
217 we think verification in real data is important, which has not been performed
218 previously in large sample size.

219 By the performance in simulated and real data, we can roughly divide

220 algorithms into three groups: full (and its derivatives), marginal (and its
221 derivatives) and ridge models. In marginal models, we test one cell type at a
222 time. If we knew in advance that one particular cell type is associated with
223 the trait, which would be a rare situation, testing that cell type in the marginal
224 model is the most simple and correct approach. Indeed, under such a
225 simulated scenario, the marginal test attained highest power (Figs. 2B, 3B).
226 However, when the test target cell type is not associated, but instead another
227 cell type is associated, the marginal tests can pick up false signals due to the
228 collinearity between regressor variables (Figs. 2C, 3C). The high power and
229 high error rate of the marginal tests can lead to unstable performance; in real
230 data, the marginal tests were the most powerful for detecting B cell specific
231 association with rheumatoid arthritis (Fig. 4B) but were the least powerful for
232 monocytes (Fig. 4A). The full model tests all cell types together, and its
233 performance was the opposite of the marginal. By fitting all cell types
234 simultaneously, the full model adjusts for the effects of other cell types. The
235 full models did not detect false association coming indirectly from non-target
236 cell types (Figs. 2C, 3C), yet their power was relatively low (Figs. 2B, 3B).
237 The ridge tests (omicwas.identity.ridge, omicwas.logit.ridge and
238 omicwas.log.ridge) were in the middle between full and marginal tests with
239 regards to the power (Figs. 2B, 3B, 4). The false positives of ridge tests were
240 modest compared to the marginal tests (Figs. 2C, 3C).

241 We mathematically modeled and implemented the logit scale for DNA
242 methylation and log scale for gene expression. It turns out that the
243 improvement by formulating the nonlinear scale was negligible for DNA
244 methylation (Fig. 2B) and modest for gene expression (Fig. 3B;
245 omicwas.identity vs omicwas.log, and omicwas.identity.ridge vs
246 omicwas.log.ridge). This implies that previous works, which were almost
247 exclusively in linear scale, were not losing much power due to scaling.

248

249 **Conclusions**

250 For cell-type-specific differential expression analysis by using unsorted tissue
251 samples, we recommend trying ridge regression as a first choice because it

252 balances power and type I error. Although marginal tests can be powerful
253 when the tested cell type actually is the only one associated with the trait,
254 caution is needed due to its high type I error rate. For a signal detected by
255 the marginal test, reanalysis in full model could be valuable. Ridge regression
256 is preferable compared to the full model without ridge regularization because
257 ridge estimator of the effect size has smaller MSE (equation (13)). Nonlinear
258 regression, which models scales properly, is recommended more than the
259 linear regression, yet the difference can be modest. We do not claim the ridge
260 model to substitute previous models. Indeed, we think none of the current
261 algorithms is superior to others in all aspects, indicating possibility for future
262 improvement.

263

264 **Methods**

265 **Linear regression**

266 We begin by describing the linear regressions used in previous studies. Let
267 the indexes be h for a cell type, i for a sample, j for an omics marker (CpG
268 site or gene), k for a trait that has cell-type-specific effects on marker
269 expression, and l for a trait that has a uniform effect across cell types. The
270 input data is given in four matrices. The matrix $W_{h,i}$ represents cell type
271 composition. The matrices $X_{i,k}$ and $C_{i,l}$ represent the values of the traits that
272 have cell-type-specific and uniform effects, respectively. We assume the two
273 matrices are centered: $\sum_i X_{i,k} = \sum_i C_{i,l} = 0$. The matrix $Y_{i,j}$ represents the
274 omics marker expression level in tissue samples.

275 The parameters we estimate are the cell-type-specific trait effect $\beta_{h,j,k}$,
276 tissue-uniform trait effect $\gamma_{j,l}$, and basal marker level $\alpha_{h,j}$ in each cell type.
277 For the remaining of the first five sections (up to “Multicollinearity of
278 interaction terms”), we focus on one marker j , and omit the index for
279 readability. For cell type h , the marker level of sample i is

$$280 \quad \alpha_h + \sum_k \beta_{h,k} X_{i,k}. \quad (1)$$

281 This is a representative value rather than a mean because we do not model
 282 a probability distribution for cell-type-specific expression. By averaging the
 283 value over cell types with weight $W_{h,i}$, and combining with the tissue-uniform
 284 trait effects, we obtain the mean marker level in bulk tissue of sample i ,

$$285 \quad \mu_i = \sum_h \alpha_h W_{h,i} + \sum_{h,k} \beta_{h,k} W_{h,i} X_{i,k} + \sum_l \gamma_l C_{i,l}.$$

286 With regards to the statistical model, we assume the error of the marker
 287 level to be normally distributed with variance σ^2 , independently among
 288 samples, as

$$289 \quad Y_i = \mu_i + \varepsilon_i,$$

$$290 \quad \varepsilon_i \sim N(0, \sigma^2).$$

291 The statistical significance of all parameters is tested under the *full* model of
 292 linear regression,

$$293 \quad Y_i = \sum_h \alpha_h W_{h,i} + \sum_{h,k} \beta_{h,k} W_{h,i} X_{i,k} + \sum_l \gamma_l C_{i,l} + \varepsilon_i, \quad (2)$$

294 or its derivatives [5,10,13]. Alternatively, the cell-type-specific effects of
 295 traits can be fitted and tested for one cell type h at a time by the *marginal*
 296 model,

$$297 \quad Y_i = \sum_{h'} \alpha_{h'} W_{h',i} + \sum_k \beta_{h,k} W_{h,i} X_{i,k} + \sum_l \gamma_l C_{i,l} + \varepsilon_i, \quad (3)$$

298 or its derivatives [7-9,11,14].

299 **Nonlinear regression**

300 Aiming to simultaneously analyze cell type composition in linear scale and
 301 differential expression/methylation in log/logit scale, we develop a nonlinear
 302 regression model. The differential analyses are performed after applying
 303 normalizing transformation. The normalizing function is the natural logarithm
 304 $f = \log$ for gene expression, and $f = \text{logit}$ for methylation (see Background).
 305 Conventional linear regression can be formulated by defining f as the identity
 306 function. We denote the inverse function of f by g ; $g = \exp$ for gene
 307 expression, and $g = \text{logistic}$ for methylation. Thus, f converts from the linear
 308 scale to the normalized scale, and g does the opposite.

309 The marker level in a specific cell type (formula (1)) is modeled in the
 310 normalized scale. The level is linearized by applying function g , then averaged
 311 over cell types with weight $W_{h,i}$, and normalized by applying function f .
 312 Combined with the tissue-uniform trait effects, the mean normalized marker
 313 level in bulk tissue of sample i becomes

$$314 \quad \mu_i = f\left(\sum_h W_{h,i} g\left(\alpha_h + \sum_k \beta_{h,k} X_{i,k}\right)\right) + \sum_l \gamma_l C_{i,l}. \quad (4)$$

315 We assume the normalized marker level to have an error that is normally
 316 distributed with variance σ^2 , independently among samples, as

$$317 \quad f(Y_i) = \mu_i + \varepsilon_i, \quad (5)$$

$$318 \quad \varepsilon_i \sim N(0, \sigma^2).$$

319 We obtain the ordinary least squares (OLS) estimator of the parameters by
 320 minimizing the residual sum of squares,

$$321 \quad \text{RSS} = \sum_i (f(Y_i) - \mu_i)^2, \quad (6)$$

322 and then estimate the error variance as

$$323 \quad \widehat{\sigma^2} = \frac{1}{n-p} \text{RSS}, \quad (7)$$

324 where n is the number of samples and p is the number of parameters [[24],
 325 section 6.3.1].

326 **Ridge regression**

327 The parameters $\beta_{h,k}$ for cell-type-specific effect cannot be estimated
 328 accurately by ordinary linear regression because the regressors $W_{h,i} X_{i,k}$ in
 329 equation (2) are highly correlated between cell types (see below).
 330 Multicollinearity also occurs to the nonlinear case in formula (4) because of
 331 local linearity. To cope with the multicollinearity, we apply ridge regression
 332 with a regularization parameter $\lambda \geq 0$, and obtain the ridge estimator of the
 333 parameters that minimizes

$$334 \quad \text{RSS} + \lambda \sum_{h,k} \beta_{h,k}^2, \quad (8)$$

335 where the second term penalizes $\beta_{h,k}$ for taking large absolute values. The
 336 ridge estimator $\hat{\boldsymbol{\theta}}(\lambda)$ is asymptotically normally distributed (see Additional file
 337 2: [Supplementary note](#)) with

$$338 \quad \text{Mean}[\hat{\boldsymbol{\theta}}(\lambda)] = Q(\lambda)^{-1} Q(0) \boldsymbol{\theta}, \quad (9)$$

$$339 \quad \text{Var}[\hat{\boldsymbol{\theta}}(\lambda)] = \sigma^2 Q(\lambda)^{-1} \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^T \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right) Q(\lambda)^{-1}, \quad (10)$$

$$340 \quad Q(\lambda) = \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^T \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right) + \lambda \begin{pmatrix} 0 & 0 & 0 \\ 0 & I & 0 \\ 0 & 0 & 0 \end{pmatrix} - (f(Y) - \boldsymbol{\mu}(\boldsymbol{\theta})) \cdot \left(\frac{\partial^2 \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \right),$$

341 where $\boldsymbol{\mu}$ is the vector form of μ_i , $\boldsymbol{\theta}$ is the vector form of the parameters α_h ,
 342 $\beta_{h,k}$ and γ_l combined, $(\partial \boldsymbol{\mu} / \partial \boldsymbol{\theta})$ is the Jacobian matrix, $(\partial^2 \boldsymbol{\mu} / \partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T)$ is the
 343 array of Hessian matrices for μ_i taken over samples, and T indicates matrix
 344 transposition. The product of $f(Y) - \boldsymbol{\mu}(\boldsymbol{\theta})$ and the Hessian is taken by
 345 multiplying for each sample and then summing up over samples. The matrix
 346 after λ has one only in the diagonal corresponding to $\beta_{h,k}$. The assigned value
 347 $\boldsymbol{\theta}$ is the true parameter value. By taking the expectation of Q , we obtain a
 348 rougher approximation [25] as

$$349 \quad \text{Mean}[\hat{\boldsymbol{\theta}}(\lambda)] = Q^*(\lambda)^{-1} Q^*(0) \boldsymbol{\theta}, \quad (11)$$

$$350 \quad \text{Var}[\hat{\boldsymbol{\theta}}(\lambda)] = \sigma^2 Q^*(\lambda)^{-1} \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^T \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right) Q^*(\lambda)^{-1}, \quad (12)$$

$$351 \quad Q^*(\lambda) = E[Q(\lambda)] = \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^T \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right) + \lambda \begin{pmatrix} 0 & 0 & 0 \\ 0 & I & 0 \\ 0 & 0 & 0 \end{pmatrix}.$$

352 The matrices Q and Q^* are the observed and expected Fisher matrices
 353 multiplied by σ^2 and adapted to ridge regression, respectively.

354 Since our objective is to predict the cell-type-specific trait effects, we
 355 choose the regularization parameter λ that can minimize the mean squared
 356 error (MSE) of $\beta_{h,k}$. Our methodology is based on [26]. To simplify the
 357 explanation, we assume the Jacobian matrices $(\partial \boldsymbol{\mu}(\boldsymbol{\theta}) / \partial \boldsymbol{\alpha})$, $(\partial \boldsymbol{\mu}(\boldsymbol{\theta}) / \partial \boldsymbol{\beta})$ and
 358 $(\partial \boldsymbol{\mu}(\boldsymbol{\theta}) / \partial \boldsymbol{\gamma})$ to be mutually orthogonal, where $\boldsymbol{\alpha}$, $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$ are the vector
 359 forms of α_h , $\beta_{h,k}$ and γ_l , respectively. Then, from formulae (11) and (12), the
 360 ridge estimator $\hat{\boldsymbol{\beta}}(\lambda)$ is asymptotically normally distributed with

$$361 \quad \text{Mean}[\hat{\boldsymbol{\beta}}(\lambda)] = \left[\left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\beta}} \right)^T \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\beta}} \right) + \lambda I \right]^{-1} \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\beta}} \right)^T \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\beta}} \right) \boldsymbol{\beta},$$

362
$$\text{Var}[\widehat{\boldsymbol{\beta}}(\lambda)] = \sigma^2 \left[\left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\beta}} \right)^T \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\beta}} \right) + \lambda I \right]^{-1} \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\beta}} \right)^T \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\beta}} \right)$$

363
$$\left[\left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\beta}} \right)^T \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\beta}} \right) + \lambda I \right]^{-1},$$

364 where the assigned values $\boldsymbol{\theta}$ and $\boldsymbol{\beta}$ are the true parameter values. We apply
365 singular value decomposition

366
$$\left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\beta}} \right) = U D V^T,$$

367 where U and V are orthogonal matrices, the columns of V are $\mathbf{v}_1, \dots, \mathbf{v}_M$, and
368 the diagonals of diagonal matrix D are sorted $d_1 \geq \dots \geq d_M \geq 0$. The bias,
369 variance and MSE of the ridge estimator are decomposed as

370
$$\text{Bias}[\widehat{\boldsymbol{\beta}}(\lambda)] = \text{E}[\widehat{\boldsymbol{\beta}}(\lambda) - \boldsymbol{\beta}]$$

371
$$= -\lambda \left[\left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\beta}} \right)^T \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\beta}} \right) + \lambda I \right]^{-1} \boldsymbol{\beta}$$

372
$$= \left\{ \sum_{m=1}^M \mathbf{v}_m \frac{-\lambda}{d_m^2 + \lambda} \mathbf{v}_m^T \right\} \boldsymbol{\beta},$$

373
$$\text{Var}[\widehat{\boldsymbol{\beta}}(\lambda)] = \sigma^2 \sum_{m=1}^M \mathbf{v}_m \frac{d_m^2}{(d_m^2 + \lambda)^2} \mathbf{v}_m^T,$$

374
$$\text{MSE}[\widehat{\boldsymbol{\beta}}(\lambda)] = \text{E}[\|\widehat{\boldsymbol{\beta}}(\lambda) - \boldsymbol{\beta}\|^2]$$

375
$$= \|\text{Bias}[\widehat{\boldsymbol{\beta}}(\lambda)]\|^2 + \text{tr}(\text{Var}[\widehat{\boldsymbol{\beta}}(\lambda)])$$

376
$$= \sum_{m=1}^M \left(\frac{\lambda}{d_m^2 + \lambda} \right)^2 (\mathbf{v}_m^T \boldsymbol{\beta})^2 + \left(\frac{d_m^2}{d_m^2 + \lambda} \right)^2 \left(\frac{\sigma^2}{d_m^2} \right). \quad (13)$$

377 For each m in the summation of (13), the minimum of the summand is
378 attained at $\lambda_m = \sigma^2 / (\mathbf{v}_m^T \boldsymbol{\beta})^2$. To minimize MSE, we need to find some
379 "average" of the optimal λ_m over the range of m . Hoerl et al. [27] proposed
380 to take the harmonic mean $\lambda = M\sigma^2 / \|\boldsymbol{\beta}\|^2$. However, if an OLS estimator $\widehat{\boldsymbol{\beta}}(0)$
381 is plugged in, $\|\boldsymbol{\beta}\|^2$ is biased upwards, and λ is biased downwards. Indeed,
382 with regards to the estimator of $1/\sqrt{\lambda_m}$, we notice that

383
$$\frac{1}{\sigma} \mathbf{v}_m^T \widehat{\boldsymbol{\beta}}(0) \sim N\left(\frac{1}{\sigma} \mathbf{v}_m^T \boldsymbol{\beta}, \frac{1}{d_m^2}\right),$$

384 where the terms with larger m have larger variance. Thus, we take the
385 average of $(\mathbf{v}_m^T \widehat{\boldsymbol{\beta}}(0))^2 / \sigma^2$, weighted by $d_m^2 / \sum_{m=1}^M d_m^2$, and also subtract the
386 upward bias as,

387
$$\kappa = \frac{1}{\sum_{m=1}^M d_m^2} \sum_{m=1}^M \left\{ \frac{d_m^2 (\mathbf{v}_m^T \hat{\boldsymbol{\beta}}(0))^2}{\sigma^2} - 1 \right\}. \quad (14)$$

388 The weighting and subtraction were mentioned in [26], where the subtraction
 389 term was dismissed, under the assumption of large effect-size $\boldsymbol{\beta}$. Since the
 390 effect-size could be small in our application, we keep the subtraction term.
 391 The statistic κ can be nonpositive, and is unbiased in the sense that

392
$$E[\kappa] = \frac{1}{\sum_{m=1}^M d_m^2} \sum_{m=1}^M \frac{d_m^2 (\mathbf{v}_m^T \boldsymbol{\beta})^2}{\sigma^2} = \frac{1}{\sum_{m=1}^M d_m^2} \sum_{m=1}^M \frac{d_m^2}{\lambda_m}.$$

393 Our choice of regularization parameter is

394
$$\lambda = \begin{cases} 1/\kappa & \text{if } \kappa > 0, \\ d_1^2 & \text{otherwise,} \end{cases} \quad (15)$$

395 where d_1^2 is taken instead of positive infinity.

396 **Implementation of omicwas package**

397 For each omics marker, the parameters $\boldsymbol{\alpha}$, $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$ (denoted in combination
 398 by $\boldsymbol{\theta}$) are estimated and tested by nonlinear ridge regression in the following
 399 steps. As we assume the magnitude of trait effects $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$ to be much
 400 smaller than that of basal marker level $\boldsymbol{\alpha}$, we first fit $\boldsymbol{\alpha}$ alone for numerical
 401 stability.

- 402 1. Compute OLS estimator $\hat{\boldsymbol{\alpha}}(0)$ by minimizing formula (6) under $\boldsymbol{\beta} = \boldsymbol{\gamma} = \mathbf{0}$.
- 403 Apply Wald test.
- 404 2. Calculate $\hat{\sigma}^2$ by formula (7). Use it as a substitute for σ^2 . The residual
 405 degrees of freedom $n - p$ is the number of samples minus the number of
 406 parameters in $\boldsymbol{\alpha}$.
- 407 3. Compute OLS estimators $\hat{\boldsymbol{\beta}}(0)$ and $\hat{\boldsymbol{\gamma}}(0)$ by minimizing formula (6) under
 408 $\boldsymbol{\alpha} = \hat{\boldsymbol{\alpha}}(0)$. Let $\hat{\boldsymbol{\theta}}(0) = (\hat{\boldsymbol{\alpha}}(0)^T, \hat{\boldsymbol{\beta}}(0)^T, \hat{\boldsymbol{\gamma}}(0)^T)^T$.
- 409 4. Apply singular value decomposition $(\partial \boldsymbol{\mu}(\hat{\boldsymbol{\theta}}(0))/\partial \boldsymbol{\beta}) = UDV^T$.
- 410 5. Calculate κ and then the regularization parameter λ by formulae (14)
 411 and (15).
- 412 6. Compute ridge estimators $\hat{\boldsymbol{\beta}}(\lambda)$ and $\hat{\boldsymbol{\gamma}}(\lambda)$ by minimizing formula (8) under
 413 $\boldsymbol{\alpha} = \hat{\boldsymbol{\alpha}}(0)$. Let $\hat{\boldsymbol{\theta}}(\lambda) = (\hat{\boldsymbol{\alpha}}(0)^T, \hat{\boldsymbol{\beta}}(\lambda)^T, \hat{\boldsymbol{\gamma}}(\lambda)^T)^T$.
- 414 7. Approximate the variance of ridge estimator, according to formula (10),

415 by

$$416 \quad \text{Var} \left[\begin{pmatrix} \hat{\boldsymbol{\beta}}(\lambda) \\ \hat{\boldsymbol{\gamma}}(\lambda) \end{pmatrix} \right] = \hat{\sigma}^2 Q(\lambda)^{-1} \begin{pmatrix} \frac{\partial \boldsymbol{\mu}(\hat{\boldsymbol{\theta}}(\lambda))}{\partial \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{pmatrix}} \end{pmatrix}^T \begin{pmatrix} \frac{\partial \boldsymbol{\mu}(\hat{\boldsymbol{\theta}}(\lambda))}{\partial \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{pmatrix}} \end{pmatrix} Q(\lambda)^{-1},$$

$$417 \quad Q(\lambda) = \begin{pmatrix} \frac{\partial \boldsymbol{\mu}(\hat{\boldsymbol{\theta}}(\lambda))}{\partial \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{pmatrix}} \end{pmatrix}^T \begin{pmatrix} \frac{\partial \boldsymbol{\mu}(\hat{\boldsymbol{\theta}}(\lambda))}{\partial \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{pmatrix}} \end{pmatrix} + \lambda \begin{pmatrix} I & O \\ O & O \end{pmatrix} - \left(f(Y) - \boldsymbol{\mu}(\hat{\boldsymbol{\theta}}(\lambda)) \right) \cdot \begin{pmatrix} \frac{\partial^2 \boldsymbol{\mu}(\hat{\boldsymbol{\theta}}(\lambda))}{\partial \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{pmatrix} \partial \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{pmatrix}^T} \end{pmatrix}.$$

418 8. Apply the "non-exact" t -type test [28]. For the s -th coordinate,

$$420 \quad \frac{\begin{pmatrix} \hat{\boldsymbol{\beta}}(\lambda) \\ \hat{\boldsymbol{\gamma}}(\lambda) \end{pmatrix}_s}{\sqrt{\text{Var} \left[\begin{pmatrix} \hat{\boldsymbol{\beta}}(\lambda) \\ \hat{\boldsymbol{\gamma}}(\lambda) \end{pmatrix} \right]_{s,s}}} \sim t_{n-p}, \quad (16)$$

419 under the null hypothesis $\begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{pmatrix}_s = 0$.

421 The formula (16) is the same as a Wald test, but the test differs, because the
 422 ridge estimators are not maximum-likelihood estimators. The algorithm was
 423 implemented as a package for the R statistical language. We used the NL2SOL
 424 algorithm of the PORT library [29] for minimization.

425 In analyses of quantitative trait locus (QTL), such as methylation QTL
 426 (mQTL) and expression QTL (eQTL), an association analysis that takes the
 427 genotypes of a single nucleotide polymorphism (SNP) as $X_{i,k}$ is repeated for
 428 many SNPs. In order to speed up the computation, we perform rounds of
 429 linear regression. First, the parameters $\hat{\boldsymbol{\alpha}}(0)$ and $\hat{\boldsymbol{\gamma}}(0)$ are fit by ordinary
 430 linear regression under $\boldsymbol{\beta} = \mathbf{0}$, which does not depend on $X_{i,k}$. By taking the
 431 residuals, we practically dispense with $\hat{\boldsymbol{\alpha}}(0)$ and $\hat{\boldsymbol{\gamma}}(0)$ in the remaining steps.
 432 Next, for $X_{i,k}$ of each SNP, $\hat{\boldsymbol{\beta}}(0)$ is fit by ordinary linear regression under $\boldsymbol{\alpha} =$
 433 $\hat{\boldsymbol{\alpha}}(0)$, $\boldsymbol{\gamma} = \hat{\boldsymbol{\gamma}}(0)$. The regularization parameter λ is computed according to
 434 steps 4 and 5 above. Finally, $\hat{\boldsymbol{\beta}}(\lambda)$ is fitted and tested by linear ridge
 435 regression under $\boldsymbol{\alpha} = \hat{\boldsymbol{\alpha}}(0)$, $\boldsymbol{\gamma} = \hat{\boldsymbol{\gamma}}(0)$.

436 **Multicollinearity of interaction terms**

437 The regressors for cell-type-specific trait effects in the full model (equation
 438 (2)) are the interaction terms $W_{h,i}X_{i,k}$. To assess multicollinearity, we
 439 mathematically derive the correlation coefficient between two interaction
 440 terms $W_{h,i}X_{i,k}$ and $W_{h',i}X_{i,k}$. In this section, we treat $W_{h,i}$, $W_{h',i}$ and $X_{i,k}$ as
 441 sampled instances of random variables W_h , $W_{h'}$ and X_k , respectively. For
 442 simplicity, we assume W_h and $W_{h'}$ are independent of X_k . Let $E[\bullet]$, $\text{Var}[\bullet]$,
 443 $\text{Cov}[\bullet]$ and $\text{Cor}[\bullet]$ denote the expectation, variance, covariance and
 444 correlation, respectively. Since X_k is centered, $E[W_h X_k] = E[W_{h'} X_k] = 0$. The
 445 correlation coefficient between interaction terms becomes

$$\begin{aligned}
 446 \quad \text{Cor}[W_h X_k, W_{h'} X_k] &= \frac{E[W_h X_k W_{h'} X_k]}{\sqrt{E[W_h^2 X_k^2]} \sqrt{E[W_{h'}^2 X_k^2]}} \\
 447 \quad &= \frac{E[W_h W_{h'}]}{\sqrt{E[W_h^2]} \sqrt{E[W_{h'}^2]}} \\
 448 \quad &= \frac{\text{Cov}[W_h, W_{h'}] + E[W_h] E[W_{h'}]}{\sqrt{\text{Var}[W_h] + E[W_h]^2} \sqrt{\text{Var}[W_{h'}] + E[W_{h'}]^2}} \\
 449 \quad &= \frac{\text{Cor}[W_h, W_{h'}] + \frac{E[W_h]}{\sqrt{\text{Var}[W_h]}} \frac{E[W_{h'}]}{\sqrt{\text{Var}[W_{h'}]}}}{\sqrt{1 + \frac{E[W_h]^2}{\text{Var}[W_h]}} \sqrt{1 + \frac{E[W_{h'}]^2}{\text{Var}[W_{h'}]}}}. \quad (17)
 \end{aligned}$$

450 If the ratios $E[W_h]/\sqrt{\text{Var}[W_h]}$ and $E[W_{h'}]/\sqrt{\text{Var}[W_{h}]}$ are high, the correlation of
 451 interaction terms approaches to one, irrespective of $\text{Cor}[W_h, W_{h'}]$.

452 **EWAS of rheumatoid arthritis**

453 EWAS datasets for rheumatoid arthritis were downloaded from the Gene
 454 Expression Omnibus (GEO). Using the RnBeads package (version 2.2.0) [30]
 455 of R, IDAT files of HumanMethylation450 array were preprocessed by
 456 removing low quality samples and markers, by normalizing methylation level,
 457 and by removing markers on sex chromosomes and outlier samples. The
 458 association of methylation level with disease status was tested with

459 adjustment for sex, age, smoking status and experiment batch; the
460 covariates were assumed to have uniform effects across cell types. After
461 quality control, dataset GSE42861 included bulk peripheral blood leukocyte
462 data for 336 cases and 322 controls [20]. GSE131989 included sorted CD14⁺
463 monocyte data for 63 cases and 31 controls [21]. By meta-analysis of
464 GSE131989 and GSE87095 [19], we obtained sorted CD19⁺ B cell data for
465 108 cases and 95 controls. The cell type composition of bulk samples was
466 imputed using the Houseman algorithm [31] in the GLINT software (version
467 1.0.4) [32].

468 **Differential gene expression by age**

469 Whole blood RNA-seq data of GTEx v7 was downloaded from the GTEx
470 website [22]. Genes of low quality or on sex chromosomes were removed,
471 expression level was normalized, outlier samples were removed, and 389
472 samples were retained. The association of read count with age was tested
473 with adjustment for sex. From GEO dataset GSE56047 [23], we obtained
474 sorted CD14⁺ monocyte data for 1202 samples and sorted CD4⁺ T cell data
475 for 214 samples. The cell type composition of bulk samples was imputed using
476 the DeconCell package (version 0.1.0) [9] of R.

477 **Simulation of cell-type-specific disease association**

478 Bulk tissue sample data for case-control comparison were simulated based
479 on real data. We generated four scenarios. Each omics marker was simulated
480 independently. The mean expression level was defined for each cell type,
481 separately in cases and controls. The standard deviation (SD) was set to be
482 the same for each combination. We tested disease association specific to one
483 cell type, which we call the target cell type. In each scenario, the mean
484 expression level was set as follows.

- 485 A. The mean was equal for all cell types both in cases and controls (null
486 scenario).
- 487 B. The mean in cases was higher by 1 SD for the target cell type. Other
488 combinations had the same mean value.

489 C. The mean in cases was lower by 1 SD for one non-target cell type. Other
490 combinations had the same mean value.

491 D. The mean in cases was higher by 1 SD for the target cell type, and lower
492 by 1 SD for one non-target cell type. Other combinations had the same
493 middle mean value.

494 The target and non-target cell types were randomly chosen for each marker.
495 For each sample, the cell-type-specific expression level was randomly
496 sampled from a normal distribution that was specified in the scenario. The
497 cell-type-specific expression levels were converted to the linear scale, and
498 then averaged across cell types, according to the predefined cell type
499 composition. The result becomes the bulk expression level of the sample in
500 linear scale.

501 We used the above-mentioned bulk tissue data, namely DNA methylation
502 data for 658 peripheral blood leukocyte samples (GSE42861) and gene
503 expression data for 389 whole blood samples (GTEx). We applied the same
504 simulation procedure to each dataset. The cell type composition in the original
505 data was retained for all samples. Half of the samples were randomly
506 assigned as cases, and the other half were assigned as controls. Normalizing
507 transformation (i.e., logit or log) was applied to the bulk expression data, and
508 500 omics markers were randomly selected. For each marker, we measured
509 the average μ and the standard deviation σ of the expression level. For
510 control samples, the expression level in each cell type was sampled from
511 $N(\mu, \sigma^2)$. For case samples, the expression level in each cell type was sampled
512 from $N(\mu, \sigma^2)$, $N(\mu + \sigma, \sigma^2)$ or $N(\mu - \sigma, \sigma^2)$ according to the scenario.

513 **Evaluation of statistical methods**

514 Cell-type-specific effects of traits was statistically tested by using bulk tissue
515 data as input. We applied the omicwas package with the normalizing function
516 $f = \log$, logit , identity without ridge regularization (omicwas.log,
517 omicwas.logit, omicwas.identity) or under ridge regression
518 (omicwas.log.ridge, omicwas.logit.ridge, omicwas.identity.ridge). The
519 omicwas package was used also for conventional linear regression under the
520 full and marginal models.

521 Among previous methods, we evaluated those that accept cell type
522 composition as input and compute test statistics for cell-type-specific
523 association. For DNA methylation data, we applied TOAST (version 1.2.0)
524 [10], CellDMC (version 2.4.0) [13] and TCA (version 1.1.0) [14]. CellDMC
525 first tests association for all combinations, and then filters out those not
526 differentially methylated. We took all of the initial results as
527 CellDMC.unfiltered; in CellDMC.filtered, Z-score was set to zero for those
528 filtered out. For gene expression data, we applied TOAST and csSAM (version
529 1.4) [5]. For csSAM, we either fitted all cell types together or one cell type at
530 a time, and denoted the results as csSAM.lm and csSAM.monovariate,
531 respectively. The csSAM method is applicable to binomial traits but not to
532 quantitative traits.

533 For simulated data, we adopted the nominal significance level $P < 0.05$
534 (two-sided). In scenario B, the power was defined as the frequency of Z-score
535 > 1.96 .

536 For the association with rheumatoid arthritis and age, "true" association
537 was determined from the measurements in physically sorted blood cells,
538 under the nominal significance level $P < 0.05$ (two-sided). The significant
539 markers were "up-regulated" (in rheumatoid arthritis cases or elders) or
540 "down-regulated." For a set of differentially expressed markers in a cell type
541 (e.g., up-regulated in monocytes), the prediction performance of an
542 algorithm was measured by the area under the curve (AUC) of receiver
543 operating characteristic (ROC). Standard error of AUC was computed by the
544 jackknife estimator by splitting the markers into 100 groups by chromosomal
545 position. The relative performance of an algorithm was evaluated by its AUC
546 $- 0.5$ divided by that for the best algorithm in each scenario.

547

548 **Supplementary information**

549 **Additional file 1: Table S1.** Blood cell type proportion in Tsimane
550 Amerindians, Caucasians and Hispanics.

551 **Additional file 2: Supplementary note.** Asymptotic distribution of ridge
552 estimator.

553 **Abbreviations**

554 AUC : area under the curve, eQTL: expression QTL, EWAS: epigenome-wide
555 association study, GEO: Gene Expression Omnibus, mQTL: methylation QTL,
556 MSE: mean squared error, OLS: ordinary least squares, QTL: quantitative trait
557 locus, ROC: receiver operating characteristic, SD: standard deviation, SNP:
558 single nucleotide polymorphism

559 **Declarations**

560 **Ethics approval and consent to participate**

561 Not applicable.

562 **Consent for publication**

563 Not applicable.

564 **Availability of data and materials**

565 The datasets generated and analyzed during the current study are available
566 in the figshare repository, <https://dx.doi.org/10.6084/m9.figshare.10718282>

567 **Competing interests**

568 The authors declare that they have no competing interests.

569 **Funding**

570 This work was supported by JSPS KAKENHI [grant number JP16K07218] and
571 by the NCGM Intramural Research Fund [grant numbers 19A2004, 20A1013].
572 The funding body had no role in the design and collection of the study,
573 experiments, analyses and interpretations of data, and in writing the
574 manuscript.

575 **Author's contributions**

576 FT developed the methodology, wrote the software, implemented the study,
577 and wrote the manuscript. NK revised the manuscript. All authors read and
578 approved the final manuscript.

579 **Acknowledgements**

580 Not applicable.

581

582 **References**

583 1. Teschendorff AE, Zheng SC. Cell-type deconvolution in epigenome-wide
584 association studies: a review and recommendations. *Epigenomics*.
585 2017;9:757–68.

586 2. Sturm G, Finotello F, Petitprez F, Zhang JD, Baumbach J, Fridman WH, et
587 al. Comprehensive evaluation of transcriptome-based cell-type
588 quantification methods for immuno-oncology. *Bioinformatics*.
589 2019;35:i436–45.

590 3. Ghosh D. Mixture models for assessing differential expression in complex
591 tissues using microarray data. *Bioinformatics*. 2004;20:1663–9.

592 4. Stuart RO, Wachsman W, Berry CC, Wang-Rodriguez J, Wasserman L,
593 Klacansky I, et al. In silico dissection of cell-type-associated patterns of
594 gene expression in prostate cancer. *Proc Natl Acad Sci USA*. National
595 Academy of Sciences; 2004;101:615–20.

596 5. Shen-Orr SS, Tibshirani R, Khatri P, Bodian DL, Staedtler F, Perry NM, et
597 al. Cell type-specific gene expression differences in complex tissues. *Nat*
598 *Meth*. Nature Publishing Group; 2010;7:287–9.

599 6. Erkkilä T, Lehmusvaara S, Ruusuvaara P, Visakorpi T, Shmulevich I,
600 Lähdesmäki H. Probabilistic analysis of gene expression measurements from
601 heterogeneous tissues. *Bioinformatics*. 2010;26:2571–7.

602 7. Kuhn A, Thu D, Waldvogel HJ, Faull RLM, Luthi-Carter R. Population-
603 specific expression analysis (PSEA) reveals molecular changes in diseased
604 brain. *Nat Meth*. Nature Publishing Group; 2011;8:945–7.

605 8. Westra H-J, Arends D, Esko T, Peters MJ, Schurmann C, Schramm K, et
606 al. Cell Specific eQTL Analysis without Sorting Cells. Pastinen T, editor. *PLoS*

607 Genet. Public Library of Science; 2015;11:e1005223–17.

608 9. Aguirre-Gamboa R, de Klein N, di Tommaso J. Deconvolution of bulk
609 blood eQTL effects into immune cell subpopulations. bioRxiv. 2019.

610 10. Li Z, Wu Z, Jin P, Wu H. Dissecting differential signals in high-
611 throughput data from complex tissues. Hancock J, editor. Bioinformatics.
612 2019;35:3898–905.

613 11. Montañó CM, Irizarry RA, Kaufmann WE, Talbot K, Gur RE, Feinberg AP,
614 et al. Measuring cell-type specific differential methylation in human brain
615 tissue. Genome Biol. BioMed Central; 2013;14:R94–9.

616 12. White N, Benton M, Kennedy D, Fox A, Griffiths L, Lea R, et al.
617 Accounting for cell lineage and sex effects in the identification of cell-
618 specific DNA methylation using a Bayesian model selection algorithm.
619 Sawalha AH, editor. PLoS ONE. 2017;12:e0182455–18.

620 13. Zheng SC, Breeze CE, Beck S, Teschendorff AE. Identification of
621 differentially methylated cell types in epigenome-wide association studies.
622 Nat Meth. Nature Publishing Group; 2018;15:1059–66.

623 14. Rahmani E, Schweiger R, Rhead B, Criswell LA, Barcellos LF, Eskin E, et
624 al. Cell-type-specific resolution epigenetics without the need for cell sorting
625 or single-cell biology. Nature Communications. Nature Publishing Group;
626 2019;10:3417–11.

627 15. Cobos FA, Vandesompele J, Mestdagh P. Computational deconvolution
628 of transcriptomics data from mixed cell populations. Bioinformatics.
629 2018;34:1969–79.

630 16. Hoyle DC, Rattray M, Jupp R, Brass A. Making sense of microarray data
631 distributions. Bioinformatics. 2002;18:576–84.

632 17. Du P, Zhang X, Huang C-C, Jafari N, Kibbe WA, Hou L, et al.
633 Comparison of Beta-value and M-value methods for quantifying methylation

- 634 levels by microarray analysis. BMC Bioinformatics. BioMed Central;
635 2010;11:1–9.
- 636 18. Aiken LS, West SG. Multiple Regression: Testing and Interpreting
637 Interactions. Sage Publications; 1991.
- 638 19. Julià A, Absher D, López-Lasanta M, Palau N, Pluma A, Waite Jones L, et
639 al. Epigenome-wide association study of rheumatoid arthritis identifies
640 differentially methylated loci in B cells. Hum Mol Genet. 2017;26:2803–11.
- 641 20. Liu Y, Aryee MJ, Padyukov L, Fallin MD, Hesselberg E, Runarsson A, et
642 al. Epigenome-wide association data implicate DNA methylation as an
643 intermediary of genetic risk in rheumatoid arthritis. Nat Biotechnol. Nature
644 Publishing Group; 2013;31:142–7.
- 645 21. Rhead B, Holingue C, Cole M, Shao X, Quach HL, Quach D, et al.
646 Rheumatoid Arthritis Naive T Cells Share Hypermethylation Sites With
647 Synoviocytes. Arthritis & Rheumatology. John Wiley & Sons, Ltd;
648 2017;69:550–9.
- 649 22. GTEx Consortium. Genetic effects on gene expression across human
650 tissues. Nature. Nature Publishing Group; 2017;550:204–13.
- 651 23. Reynolds LM, Taylor JR, Ding J, Lohman K, Johnson C, Siscovick D, et
652 al. Age-related variations in the methylome associated with gene expression
653 in human monocytes and T cells. Nature Communications. 2014;5:5366.
- 654 24. Riazoshams H, Midi H, Ghilagaber G. Robust Nonlinear Regression: with
655 Applications using R. John Wiley & Sons; 2019.
- 656 25. Lim C. Robust ridge regression estimators for nonlinear models with
657 applications to high throughput screening assay data. Statist. Med.
658 2014;34:1185–98.
- 659 26. Lawless JF, Wang P. A simulation study of ridge and other regression
660 estimators. Communications in Statistics - Theory and Methods.

661 1976;5:307–23.

662 27. Hoerl AE, Kannard RW, Baldwin KF. Ridge regression: some simulations.
663 Communications in Statistics - Theory and Methods. 1975;4:105–23.

664 28. Halawa AM, Bassiouni El MY. Tests of regression coefficients under ridge
665 regression models. Journal of Statistical Computation and Simulation.
666 2000;65:341–56.

667 29. Dennis JE, Gay DM, Welsch RE. An adaptive nonlinear least-squares
668 algorithm. ACM Transactions on Mathematical Software. 1981;7:348–68.

669 30. Müller F, Scherer M, Assenov Y, Lutsik P, Walter J, Lengauer T, et al.
670 RnBeads 2.0: comprehensive analysis of DNA methylation data. Genome
671 Biol. BioMed Central; 2019;20:55–12.

672 31. Houseman EA, Accomando WP, Koestler DC, Christensen BC, Marsit CJ,
673 Nelson HH, et al. DNA methylation arrays as surrogate measures of cell
674 mixture distribution. BMC Bioinformatics. 2012;13:86.

675 32. Rahmani E, Yedidim R, Shenhav L, Schweiger R, Weissbrod O, Zaitlen
676 N, et al. GLINT: a user-friendly toolset for the analysis of high-throughput
677 DNA-methylation array data. Hancock JM, editor. Bioinformatics.
678 2017;33:1870–2.

679

680 TABLES

681

Table 1A Blood cell type proportion in rheumatoid arthritis dataset

Cell type	Neu	CD4+T	CD8+T	NK	Mono	Bcells	Eos
Mean	0.59	0.10	0.08	0.08	0.07	0.07	0.01
SD	0.11	0.06	0.05	0.04	0.02	0.03	0.02
Mean/SD	5.3	1.8	1.6	1.9	3.2	2.6	0.4

Table 1B Correlation between blood cell type proportion and rheumatoid arthritis (X_k)

r	Neu	CD4+T	CD8+T	NK	Mono	Bcells	Eos	X_k =Disease
Neu	1	-0.68	-0.60	-0.46	-0.06	-0.49	-0.48	0.44
CD4+T	-0.68	1	0.14	0.05	-0.17	0.38	0.26	-0.33
CD8+T	-0.60	0.14	1	0.08	-0.05	0.19	0.13	-0.27
NK	-0.46	0.05	0.08	1	-0.04	0.01	0.11	-0.27
Mono	-0.06	-0.17	-0.05	-0.04	1	-0.17	0.05	0.10
Bcells	-0.49	0.38	0.19	0.01	-0.17	1	0.11	-0.22
Eos	-0.48	0.26	0.13	0.11	0.05	0.11	1	-0.10

Table 1C Correlation between interaction terms

r	Neu* X_k	CD4+T* X_k	CD8+T* X_k	NK* X_k	Mono* X_k	Bcells* X_k	Eos* X_k
Neu*X_k	1	0.83	0.80	0.85	0.93	0.90	0.27
CD4+T*X_k	0.83	1	0.78	0.78	0.83	0.88	0.42
CD8+T*X_k	0.80	0.78	1	0.77	0.82	0.83	0.35
NK*X_k	0.85	0.78	0.77	1	0.85	0.83	0.35
Mono*X_k	0.93	0.83	0.82	0.85	1	0.88	0.35
Bcells*X_k	0.90	0.88	0.83	0.83	0.88	1	0.36
Eos*X_k	0.27	0.42	0.35	0.35	0.35	0.36	1

Neu, neutrophils; Mono, monocytes; Eos, eosinophils.

682

Table 2A Blood cell type proportion in GTEx dataset

Cell type	Gran	CD4+T	CD8+T	Mono	NK	Bcells
Mean	0.53	0.22	0.10	0.07	0.05	0.03
SD	0.037	0.020	0.013	0.004	0.012	0.003
Mean/SD	14.4	10.7	7.7	16.6	4.3	8.3

Table 2B Correlation between blood cell type proportion and age (X_k)

r	Gran	CD4+T	CD8+T	Mono	NK	Bcells	$X_k=Age$
Gran	1	-0.89	-0.83	0.56	-0.76	-0.41	-0.23
CD4+T	-0.89	1	0.59	-0.64	0.50	0.51	0.14
CD8+T	-0.83	0.59	1	-0.40	0.59	0.15	0.15
Mono	0.56	-0.64	-0.40	1	-0.44	-0.42	0.02
NK	-0.76	0.50	0.59	-0.44	1	0.13	0.31
Bcells	-0.41	0.51	0.15	-0.42	0.13	1	-0.03

Table 2C Correlation between interaction terms

r	Gran* X_k	CD4+T* X_k	CD8+T* X_k	Mono* X_k	NK* X_k	Bcells* X_k
Gran*X_k	1	0.99	0.98	1.00	0.96	0.99
CD4+T*X_k	0.99	1	1.00	0.99	0.98	1.00
CD8+T*X_k	0.98	1.00	1	0.99	0.98	0.99
Mono*X_k	1.00	0.99	0.99	1	0.96	0.99
NK*X_k	0.96	0.98	0.98	0.96	1	0.97
Bcells*X_k	0.99	1.00	0.99	0.99	0.97	1

Gra, granulocytes; Mono, monocytes.

685 **FIGURE LEGENDS**

686 **Figure 1**

687 Contour plot of the correlation coefficient between interaction terms $W_h X_k$
688 and $W_{h'} X_k$. W_h and $W_{h'}$ represent proportions of cell types h and h' , and X_k
689 represents the value of trait k . For this plot, we assume the ratios
690 $E[W_h]/SD[W_h]$ and $E[W_{h'}]/SD[W_{h'}]$ to be equal. As the ratio increases 1.5, 2 to
691 3, the correlation coefficient raises >0.5 , >0.7 to >0.8 , over most range of
692 $Cor[W_h, W_{h'}]$. SD stands for standard deviation.

693 **Figure 2**

694 Detection of cell-type-specific association in simulated data for DNA
695 methylation. (A), (B), (C) and (D) correspond to the respective scenarios.
696 Results from different algorithms are aligned horizontally. Vertical axis
697 indicates the Z-score for the disease effect (cases vs controls) specific to the
698 target cell type. Points are colored according to the target cell type. The
699 middle bar of the box plot indicates the median, and the lower and upper
700 hinges correspond to the first and third quartiles. The whiskers extend to the
701 value no further than $1.5 * \text{inter-quartile range}$ from the hinges. Neu,
702 neutrophils; Mono, monocytes; Eos, eosinophils.

703 **Figure 3**

704 Detection of cell-type-specific association in simulated data for gene
705 expression. (A), (B), (C) and (D) correspond to the respective scenarios.
706 Results from different algorithms are aligned horizontally. Vertical axis
707 indicates the Z-score for the disease effect (cases vs controls) specific to the
708 target cell type. Points are colored according to the target cell type. The
709 middle bar of the box plot indicates the median, and the lower and upper
710 hinges correspond to the first and third quartiles. The whiskers extend to the
711 value no further than $1.5 * \text{inter-quartile range}$ from the hinges. Gra,
712 granulocytes; Mono, monocytes.

713 **Figure 4**

714 Performance of cell-type-specific association prediction. For rheumatoid
715 arthritis association of DNA methylation in monocytes (A) and B cells (B);
716 Age association of gene expression in CD4⁺ T cells (C) and monocytes (D).
717 The prediction is evaluated separately for up-regulated and down-regulated
718 markers. The AUC of ROC and its 95% confidence interval are plotted for each
719 statistical algorithm.

Figures

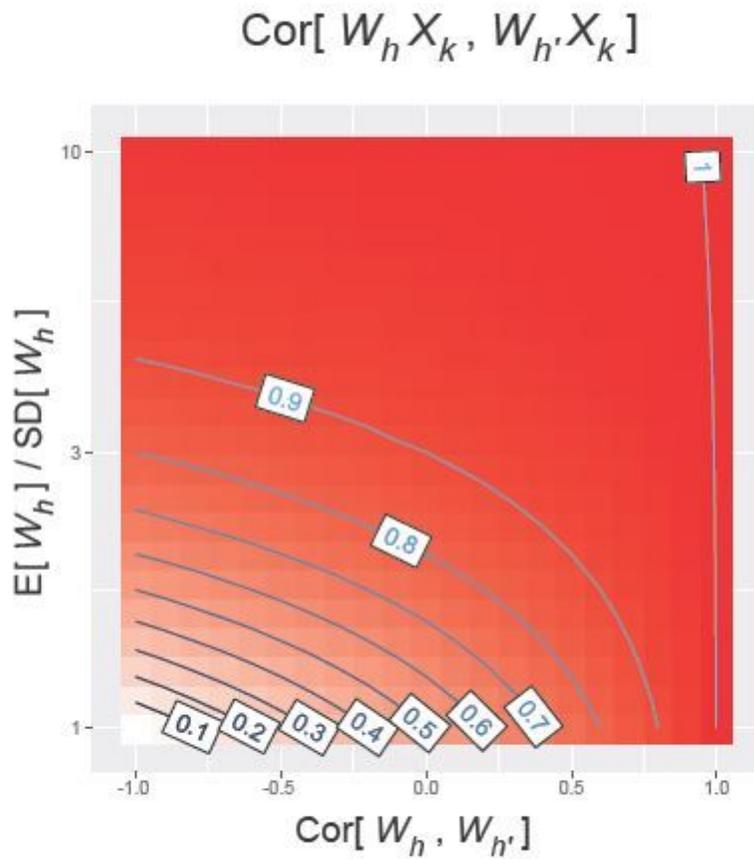


Figure 1

Contour plot of the correlation coefficient between interaction terms $W_h X_k$ and $W_{(h')} X_k$. W_h and $W_{(h')}$ represent proportions of cell types h and h' , and X_k represents the value of trait k . For this plot, we assume the ratios $E[W_h] / \text{SD}[W_h]$ and $E[W_{(h')}] / \text{SD}[W_{(h')}]$ to be equal. As the ratio increases 1.5, 2 to 3, the correlation coefficient raises >0.5 , >0.7 to >0.8 , over most range of $\text{Cor}[W_h, W_{(h')}]$. SD stands for standard deviation.

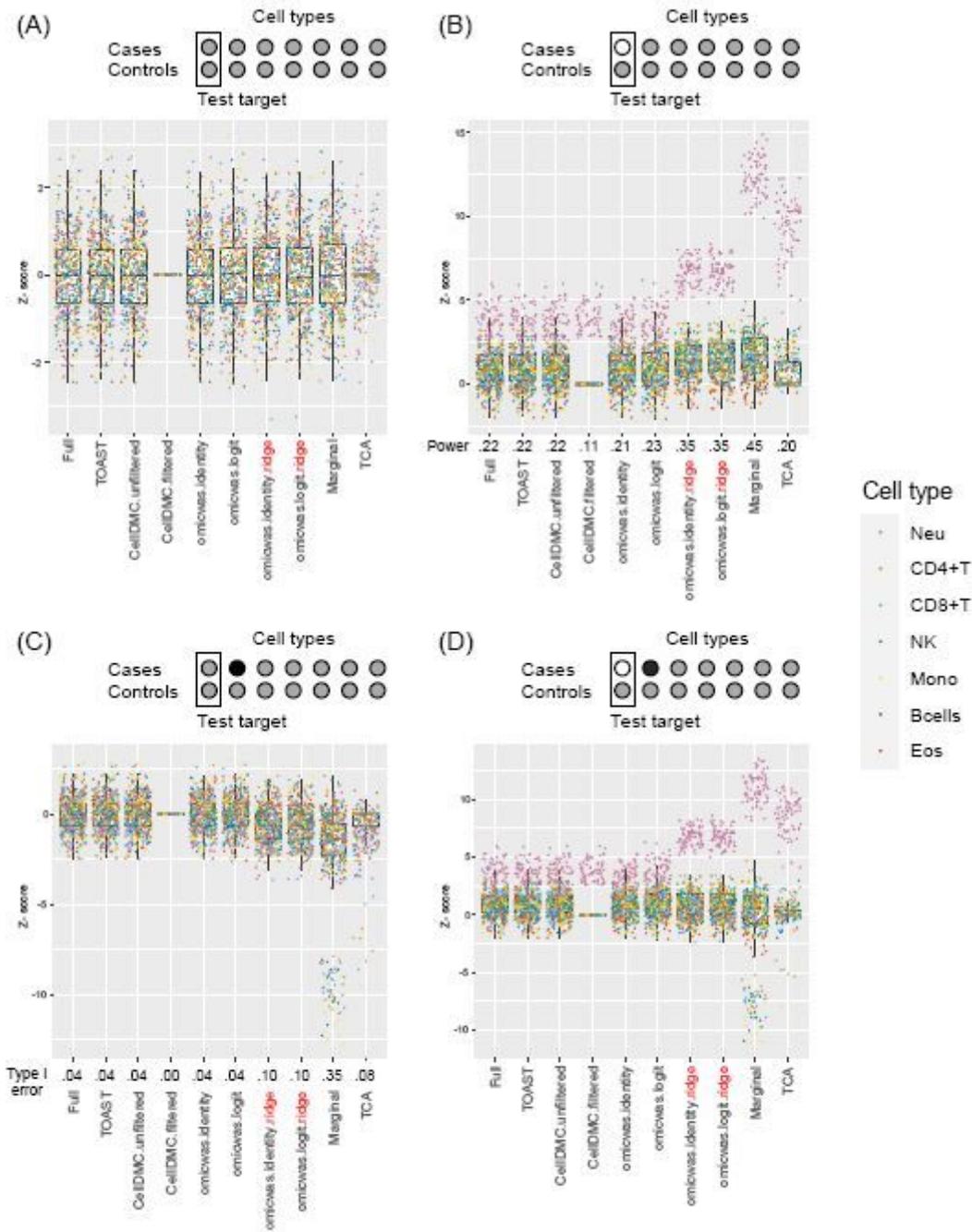


Figure 2

Detection of cell-type-specific association in simulated data for DNA methylation. (A), (B), (C) and (D) correspond to the respective scenarios. Results from different algorithms are aligned horizontally. Vertical axis indicates the Z-score for the disease effect (cases vs controls) specific to the target cell type. Points are colored according to the target cell type. The middle bar of the box plot indicates the median, and the lower and upper hinges correspond to the first and third quartiles. The whiskers extend to the value no further than $1.5 \times$ inter-quartile range from the hinges. Neu, neutrophils; Mono, monocytes; Eos, eosinophils.

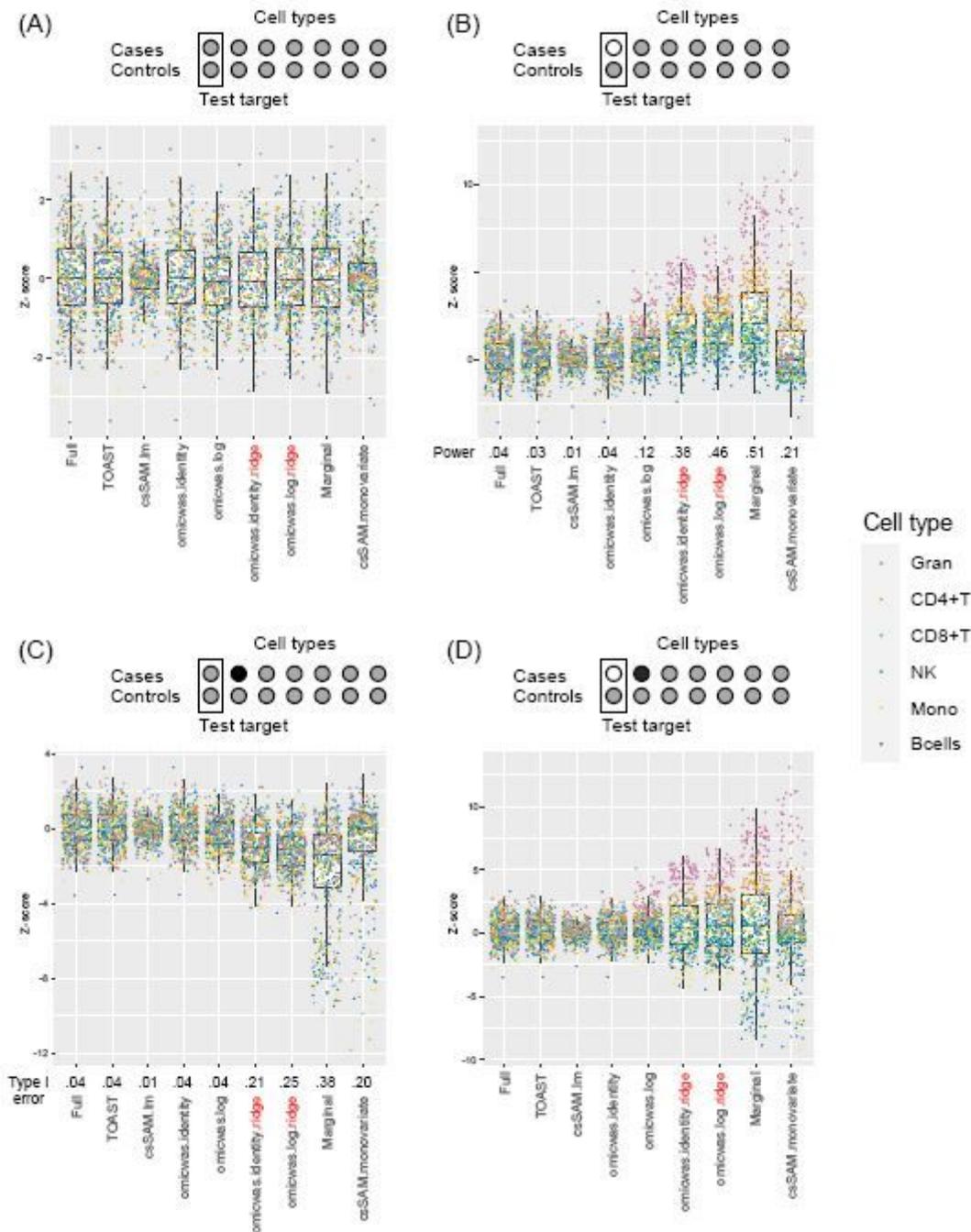


Figure 3

Detection of cell-type-specific association in simulated data for gene expression. (A), (B), (C) and (D) correspond to the respective scenarios. Results from different algorithms are aligned horizontally. Vertical axis indicates the Z-score for the disease effect (cases vs controls) specific to the target cell type. Points are colored according to the target cell type. The middle bar of the box plot indicates the median, and the lower and upper hinges correspond to the first and third quartiles. The whiskers extend to the value no further than 1.5 * inter-quartile range from the hinges. Gra, granulocytes; Mono, monocytes.

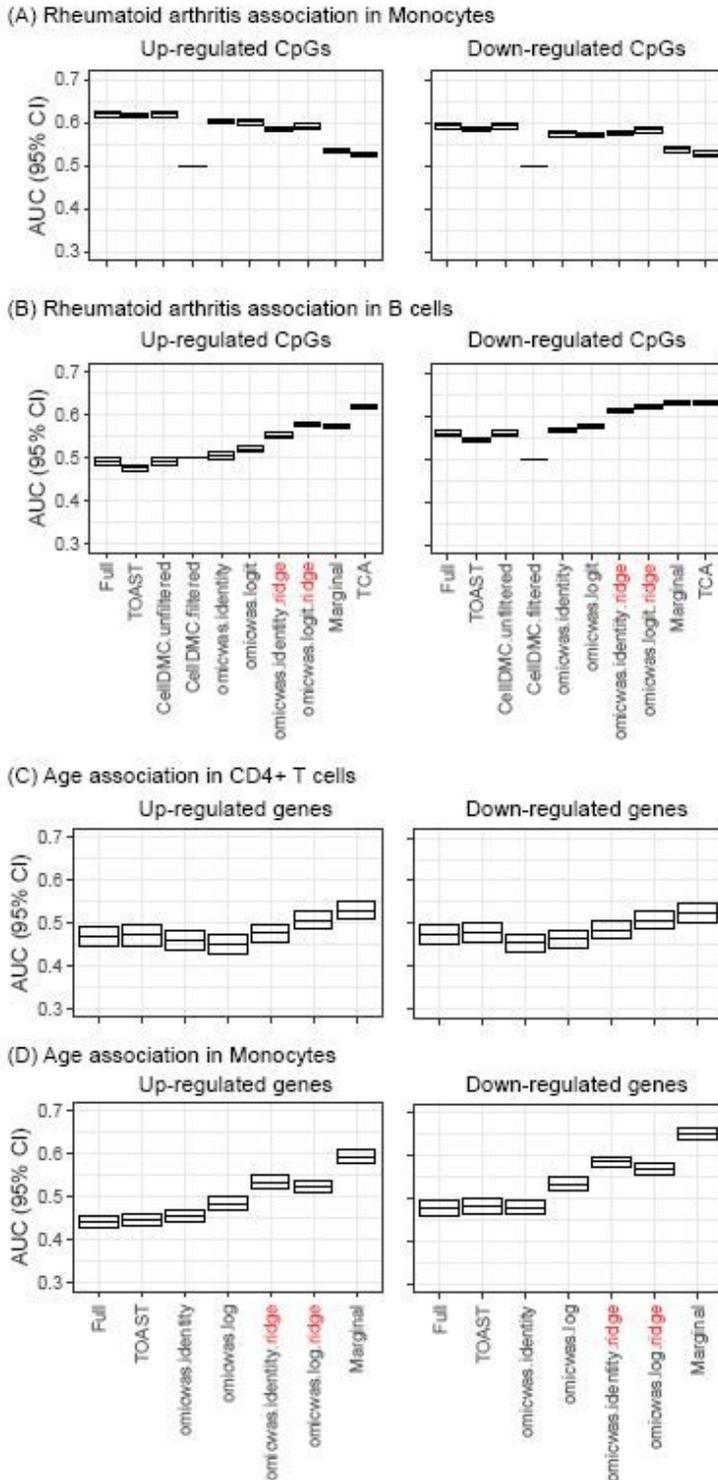


Figure 4

Performance of cell-type-specific association prediction. For rheumatoid arthritis association of DNA methylation in monocytes (A) and B cells (B); Age association of gene expression in CD4+ T cells (C) and monocytes (D). The prediction is evaluated separately for up-regulated and down-regulated markers. The AUC of ROC and its 95% confidence interval are plotted for each statistical algorithm.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SuppNote.pdf](#)
- [TableS1.xls](#)