

Optimal Genomic Control in Large-scale Genetic Associations for Binary Diseases

Yuxin Song¹, Li Jiang¹, Zhiyu Hao² and Runqing Yang^{1,2}

1. Research Centre for Aquatic Biotechnology, Chinese Academy of Fishery Sciences, Beijing 100141, People's Republic of China

2. College of Animal Scientific Technology, Northeast Agricultural University, Harbin 150030, China

Correspondence to Runqing Yang (runqingyang@cafs.ac.cn)

Complex computation and approximate solution hinder the application of generalized linear mixed models (GLMM) into genome-wide association studies. We extended GRAMMAR to handle binary diseases by considering genomic breeding values (GBVs) estimated in advance as a known predictor in genomic logit regression, and then controlled polygenic effects by regulating downward genomic heritability. Using simulations and case analyses, we showed in optimizing GRAMMAR, polygenic effects and genomic controls could be evaluated using the fewer sampling markers, which extremely simplified GLMM-based association analysis in large-scale data. In addition, joint analysis for quantitative trait nucleotide (QTN) candidates chosen by multiple testing offered significant improved statistical power to detect QTNs over existing methods.

Introduction

Although usually expressed as binary phenotypes, many disease traits in plants and animals are thought to be controlled by a number of loci each having a small effect^{1,2}. Thus, random polygenic effects excluding the tested markers should be considered in genome-wide association study (GWAS) for disease traits to correct the population stratification and cryptic relatedness, as linear mixed model (LMM) does for quantitative traits^{3,4}. Logistic regression, as a kind of generalized linear model (GLM)^{5,6}, has been used earlier for genome-wide association analysis with the disease traits⁷. Despite correction for fixed-effect covariates⁸⁻¹⁰, logistic regression still produces inflation of association test statistics. Therefore, it is necessary to introduce a generalized linear mixed model (GLMM)¹¹ that considers random polygenic effects to increase the power to map QTNs for disease traits.

However, genome-wide GLMM association consumes much more computing time than mixed model association with either restricted maximum likelihood estimation

34 (REML) ¹² or Markov Chain Monte Carlo (MCMC) iteration ¹³. If using the maximum
35 likelihood in estimation and approximations to avoid numerical integration ¹⁴, the
36 GLMM yields a problem of serious bias induced by the approximations ¹⁵, especially
37 solutions tending toward the positive/negative infinity. In the GWAS for case-control
38 studies, a non-random sampling of cases from the population results in biased
39 estimation for genomic heritability. Analyzed by the LMM for binary phenotypes,
40 genomic heritability estimate is transformed to a liability scale by adjusting both for
41 scale and for ascertainment of the case samples ¹⁶. The genomic heritability of liability
42 is estimated in a biased manner for disease traits, even though it is done by GLMM via
43 MCMC iteration. Based on the calibrated genomic heritability for case-control
44 ascertainment, a Chi-squared score statistic for GWAS of disease traits is computed
45 from posterior mean liabilities (PMLs) under the liability-threshold model ¹⁷. The
46 individual PMLs are estimated with a multivariate Gibbs sampler, which increases the
47 computational demand. For simplification to GLMM-based association analysis, the
48 GMMAT ¹⁸ and SAIGE ¹⁹ separately extend the EMMAX ²⁰ and BOLT-LMM ²¹,
49 respectively, for normally distributed traits to binary phenotypes.

50 Owing to the computational intensity and approximate solutions obtained, GLMM
51 can hardly be employed in GWAS for disease traits. Moreover, genomic heritability
52 cannot be accurately estimated for complex diseases, especially in ascertained case-
53 control studies. Motivated by the optimal genomic control for mixed model association
54 analysis for quantitative traits distributed normally, we extended GRAMMAR ²² to
55 handle binary traits by considering genomic breeding values (GBVs) estimated in
56 advance as a known predictor in genomic logit regression, and then, optimized the
57 genomic control for GRAMMAR for binary traits by regulating the downward genomic
58 heritability to estimate the residual phenotypes. The complicated GLMM does not need
59 to be directly solved by the Optim-GRAMMAR for binary traits, and it only repeatedly
60 estimates GBVs with genomic best linear unbiased prediction (GBLUP) ²³ for GLMM,
61 achieving genome-wide GLMM association analysis rapidly. Finally, we jointly
62 analyzed the candidate quantitative trait nucleotides (QTNs) chosen by multiple testing
63 to improve the statistical power to detect QTNs.

64 **Results**

65 *Statistical properties of Optim-GRAMMAR for binary traits*

66 Based on the two genomic datasets, we simulated phenotypes controlled by 40, 200,

67 and 1,000 QTNs at the low (0.2), moderate (0.5), and high (0.8) genomic heritability,
68 respectively. The statistical properties of Optim-GRAMMAR using a test at once for
69 binary traits were investigated by comparing it with GRAMMAR, GMMAT, LTMLM,
70 and SAIGE. The Q-Q and ROC profiles are displayed selectively in Figure 1 and Figure
71 2, respectively, and in Supplementary Figure 1S and Figure 2S, respectively, in detail.
72 The genomic controls are estimated in Table 1S. Making genomic control infinitely
73 close to 1.0, Optim-GRAMMAR achieved almost the same statistical power to detect
74 QTNs as the GMMAT which approximates the exact GLMM, irrespective of how many
75 QTNs and heritabilities are simulated. Among Optim-GRAMMAR and the four
76 competing methods, GRAMMAR had the lowest genomic controls and statistical
77 power, and for GRAMMAR, the population structure was more complex and the false
78 negative rate was larger. Although LTMLM achieved the highest statistical power to
79 detect QTNs for all simulated phenotypes for the maize dataset, and SAIGE
80 demonstrated a higher statistical power for the dataset controlled by 1,000 QTNs at the
81 genomic heritability of 0.2. A strong false positive error rate was observed for SAIGE.
82 In the human dataset, there were no distinct differences in the statistical properties
83 between GRAMMAR-lambda and the four competing methods, although GRAMMAR
84 provided some false negative errors.

85 After optimization for genomic control, Optim-GRAMMAR jointly analyzed
86 multiple QTN candidates chosen from a test at once at a significance level of 0.05. For
87 convenience for comparison, we analyzed the statistical powers obtained with one test
88 at a time and joint analyses together. By backward regression analysis, Optim-
89 GRAMMAR evidently exhibited improved statistical power. In contrast, LTMLM was
90 inferior to joint analysis of Optim-GRAMMAR in the terms of statistical power, even
91 with the highest false positive rates.

92 ***Calculation of GRMs and GCs with the sampling markers***

93 To investigate the effects of sampling markers on Optim-GRAMMAR, we randomly
94 took 3 K, 5 K, 10 K, 20 K, and 25 K SNPs from the entire genomic markers to calculate
95 GRM. Changes in the genomic control at the varied sampling levels of SNPs are
96 depicted in Figure 3 for Optim-GRAMMAR, GRAMMAR, GMMAT, and SAIGE.
97 Because LTMLM cannot sample SNPs, it was not included in the comparison. No
98 competing method stably controlled the positive/negative false errors using less than 25
99 K sampling SNPs, besides SAIGE for human phenotypes. Specifically, GMMAT

100 gradually controlled the positive false errors as the sampling markers increased;
101 GRAMMAR controlled the negative false rate by sampling less markers, while SAIGE
102 produced serious false negative errors in the complex maize population. In comparison,
103 Optim-GRAMMAR still retained a high statistical power to detect QTNs through almost
104 perfect genomic control, even using less than 3000 sampling markers (see
105 Supplementary Figure 3S and Figure 4).

106 *Application of Optim-GRAMMAR to WTCCC study*

107 We were authorized to re-analyze the Wellcome trust case-control consortium (WTCCC)
108 study 1²⁴. There were the 11,985 cases from six common diseases and 3,004 shared
109 controls, genotyped at a total of 490,032 SNPs. For each dataset, a standard quality
110 control (QC) procedure was performed: SNPs with MAFs < 0.01 and HWE > 0.05 were
111 excluded, and individuals with missing rates > 0.01 were also excluded. After the QC
112 process, the number of samples and SNPs used for generalized mixed model association
113 analyses were 5002 individuals (1998 cases and 3004 controls) and 409,642 SNPs for
114 bipolar disorder (BD), 4992 individuals (1988 cases and 3004 controls) and 409,516
115 SNPs for coronary artery disease (CAD), 5003 individuals (1999 cases and 3004
116 controls) and 409,924 SNPs for rheumatoid arthritis (RA), 5005 individuals (2001 cases
117 and 3004 controls) and 409,742 SNPs for hypertension (HT), 5004 individuals (2000
118 cases and 3004 controls) and 40,9674 SNPs for type I diabetes (T1D), and 5003
119 individuals (1999 cases and 3004 controls) and 409,805 SNPs for type II diabetes (T2D).
120 All data analyses were performed in a CentOS Linux sever with 2.60 GHz Intel(R)
121 Xeon(R) 40 CPUs E5-2660 v3, and 512 GB memory.

122 For the six common diseases, we implemented Optim-GRAMMAR using entire
123 genomic markers and 5,000 sampling SNPs, respectively, to estimate the GRM. The Q-
124 Q and Manhattan profiles for the six common diseases are depicted in Figure 4S and
125 Figure 5S obtained with the Optim-GRAMMAR using a test at once and the four
126 competing methods used in simulations, while in Figure 5S with the Optim-
127 GRAMMAR using joint association analyses. The association analyses illustrated that
128 (1) under perfect genomic control, Optim-GRAMMAR found the QTNs for each
129 disease on each chromosome, and the numbers of detected QTNs were not less than all
130 the competing methods; and (2) in Optim-GRAMMAR, joint association analyses
131 detected more QTNs than a test at once. As compared to Optim-GRAMMAR,
132 GRAMMAR detected less QTNs with the lowest genomic control among all the

133 methods, while GMMAT yielded more SNPs whose $-\log(p)$ exceeded the Bonferroni
134 corrected thresholds for CAD, T1D, T2D, and HT, but it obtained the highest genomic
135 control. Additionally, LTMMLM estimated the abnormal genomic heritabilities for CAD,
136 BD, T2D, and HT, producing unstable genomic controls.

137 Further, we conducted strict QC for each dataset, as done in ¹⁶ for estimating genomic
138 heritability. Despite this, the missing heritabilities could not be normally estimated for
139 BD and HT. As shown in Figure 6S and Figure 7S, all methods exhibited clear and
140 comparable association results, except for GRAMMAR. Interestingly, both LTMMLM
141 and GMMAT seriously underestimated the genomic heritability for each disease after
142 strict QC. In summary, Optim-GRAMMAR could efficiently and robustly map QTNs
143 for binary diseases and did not depend on estimation of genomic heritability and QC
144 for genomic datasets. For each dataset with standard QC, we recorded the running times
145 from input of genotypes and phenotype to output of mapping QTNs for all the methods.
146 Table 2S shows that Optim-GRAMMAR reduced the computing time by dozens of
147 times with the lowest memory footprint.

148 **Discussion**

149 Development of the GRAMMAR for GLMM association was an essential prerequisite
150 for extending the Optim-GRAMMAR to rapidly optimize mixed model association
151 analysis for binary traits. For genomic GLMM, however, no binary residuals could be
152 produced because of the scale difference between binary phenotype and predictors.
153 Thus, we considered the GBVs estimated in advance as a known predictor in genomic
154 logit regression and then executed association tests for candidate markers. This ensured
155 that GRAMMAR had the lowest computing complexity for association tests for binary
156 traits among the existing GLMM-based association methods ¹⁷⁻¹⁹. Because
157 GRAMMAR for binary diseases produces high false negative rates for quantitative
158 traits distributed normally, we optimized the genomic control for GRAMMAR by
159 regulating downward genomic heritability to underestimate the GBVs with GBLUP
160 equations for GLMM. Thus, optim-GRAMMAR solved the GBLUP equations and
161 performed association tests with simple logit regression only for several iterations, thus
162 improving the computational efficiency for genome-wide GLMM association analysis.

163 Several GLMM-based association methods such LTMMLM, GMMAT, and SAIGE
164 have simplified genome-wide mixed model association analysis for binary traits to a
165 certain extent, but they are more appropriate to handle the less complex populations

166 such as human datasets^{18,19}. Moreover, the heritability for binary diseases could not be
167 robustly and precisely estimated using genomic markers^{16 14,15}, which also limited the
168 efficient application of these association methods. In contrast, because optim-
169 GRAMMAR does not need to directly estimate genomic heritability, it can powerfully
170 and robustly map QTNs for binary traits in complexly structured populations. Within
171 the framework of Optim-GRAMMAR, further, joint analysis for the candidate
172 quantitative trait nucleotides chosen by multiple testing significantly improved
173 statistical power to detect QTNs with almost perfect genomic control.

174 The Optim-GRAMMAR extremely simplified genome-wide GLMM association
175 analysis for binary traits in large-scale population. For a genomic dataset containing m
176 SNPs genotyped on n individuals, Optim-GRAMMAR for binary traits took only the
177 computing complexity of $O(mn^2)$ to build the relationship matrix and $O(imn)$ for
178 association tests with i rounds to optimize genomic controls. If we solved genetic effects
179 of the m_1 sampled markers using ridge regression²⁵ with given heritability and then
180 estimated GBVs, then the computing complexity to build the information matrix would
181 reduce to $O(m_1^2n)$, as in FaST-LMM-Select²⁶. At the same time, if we evaluated
182 genomic control using m_1 sampled markers at each iteration, then the computing
183 complexity for association tests would reduce to $O(im_2n)$. For the simulated 8 million
184 SNPs on 400,000 individuals, Optim-GRAMMAR required only 5.1hr to analyze
185 single binary phenotype by sampling 5,000 SNPs to calculate GRM, 0.9hr of which to
186 build the information matrix and optimize genomic controls, while SAIGE did about
187 534 hr¹⁵. A user friendly GRL-Binary software was developed, which is freely available
188 at <https://github.com/RunKingProgram/Binary-Optim-GRAMMAR>.

189 **Acknowledgements**

190 The research is financially supported by the National Natural Science Foundations of
191 China (32072726) and the Special Scientific Research Funds for Central Non-profit
192 Institutes, Chinese Academy of Fishery Sciences (2017A001).

193 **Competing interests**

194 The authors declare no competing financial interests.

195 **References**

196 1. Bulmer, M.G. The Effect of Selection on Genetic Variability. *American Naturalist* **105**, 201-211

- 197 (1971).
- 198 2. Falconer, D.S. *Introduction to Quantitative Genetics, 2nd ed.*, (Longman, London, 1981).
- 199 3. Henderson, C.R. *Applications of linear models in animal breeding*, (University of Guelph, Guelph,
200 1984).
- 201 4. Yu, J.M. *et al.* A unified mixed-model method for association mapping that accounts for multiple
202 levels of relatedness. *Nature Genetics* **38**, 203-208 (2006).
- 203 5. Wedderburn, R.W.M. Quasi-likelihood functions, generalized linear models, and the gauss-newton
204 method. *Biometrika* **61**, 439-447 (1974).
- 205 6. McCullagh, P. & Nelder, J.A. *Generalized linear models, 2nd ed.*, (Chapman and Hall, New York,
206 1989).
- 207 7. Chang, C.C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets.
208 *GigaScience* **4**, 7 (2015).
- 209 8. Mefford, J. & Witte, J.S. The Covariate's Dilemma. *PLoS Genet* **8**, e1003096 (2012).
- 210 9. Zaitlen, N. *et al.* Informed conditioning on clinical covariates increases power in case-control
211 association studies. *PLoS Genet* **8**, e1003032 (2012).
- 212 10. Zaitlen, N. *et al.* Analysis of case-control association studies with known risk variants.
213 *Bioinformatics* **28**, 1729-1737 (2012).
- 214 11. Breslow, N.E. & Clayton, D.G. Approximate inference in generalized linear mixed models. *Journal*
215 *of the American statistical Association* **88**, 9-25 (1993).
- 216 12. Patterson, H.D. & Thompson, R. Recovery of inter-block information when block sizes are unequal.
217 *Biometrika* **58**, 545-554 (1971).
- 218 13. Sorensen, D. & Gianola, D. *Likelihood, Bayesian, and MCMC methods in quantitative genetics*,
219 (Springer, New York, 2002).
- 220 14. Schall, R. Estimation in generalized linear models with random effects. *Biometrika* **78**, 719-727
221 (1991).
- 222 15. Gilmour, A.R., Anderson, R.D. & Rae, A.L. The Analysis of Binomial Data by a Generalized Linear
223 Mixed Model. *Biometrika* **72**, 593-599 (1985).
- 224 16. Lee, S.H., Wray, N.R., Goddard, M.E. & Visscher, P.M. Estimating missing heritability for disease
225 from genome-wide association studies. *Am J Hum Genet* **88**, 294-305 (2011).
- 226 17. Hayeck, T.J. *et al.* Mixed Model with Correction for Case-Control Ascertainment Increases
227 Association Power. *American Journal of Human Genetics* **96**, 720-730 (2015).
- 228 18. Chen, H. *et al.* Control for Population Structure and Relatedness for Binary Traits in Genetic
229 Association Studies via Logistic Mixed Models. *Am J Hum Genet* **98**, 653-66 (2016).
- 230 19. Zhou, W. *et al.* Efficiently controlling for case-control imbalance and sample relatedness in large-
231 scale genetic association studies. *Nat Genet* **50**, 1335-1341 (2018).
- 232 20. Kang, H.M. *et al.* Variance component model to account for sample structure in genome-wide
233 association studies. *Nature Genetics* **42**, 348-354 (2010).
- 234 21. Loh, P.R. *et al.* Efficient Bayesian mixed-model analysis increases association power in large
235 cohorts. *Nature Genetics* **47**, 284-290 (2015).
- 236 22. Aulchenko, Y.S., de Koning, D.J. & Haley, C. Genomewide rapid association using mixed model
237 and regression: a fast and simple method for genomewide pedigree-based quantitative trait loci
238 association analysis. *Genetics* **177**, 577-585 (2007).
- 239 23. Vanraden, P.M. *et al.* Invited review: reliability of genomic predictions for North American Holstein
240 bulls. *Journal of Dairy Science* **92**, 16-24 (2009).

- 241 24. Consortium, W.T.C.C. Genome-wide association study of 14,000 cases of seven common diseases
242 and 3,000 shared controls. *Nature* **447**, 661-78 (2007).
- 243 25. Hoerl, A.E. & Kennard, R.W. Ridge Regression: Biased Estimation for Nonorthogonal Problems.
244 *Technometrics* **12**, 55-67 (1970).
- 245 26. Jennifer, L. *et al.* Improved linear mixed models for genome-wide association studies. *Nature*
246 *Methods* **9**, 525-526 (2012).
- 247

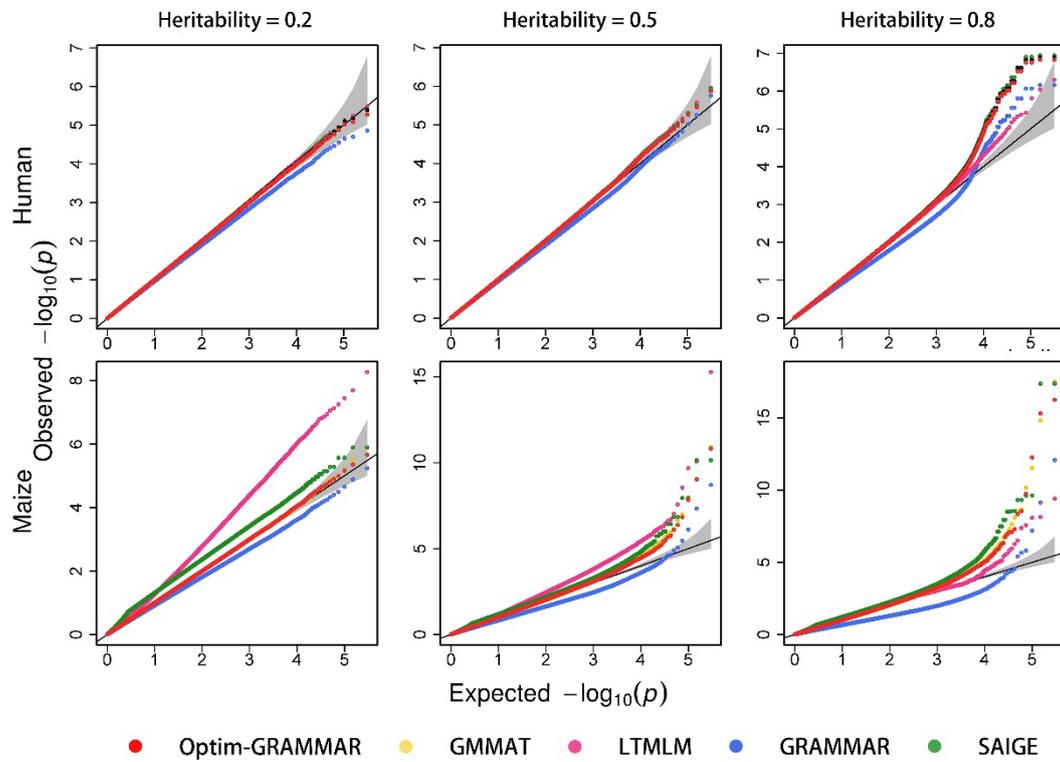


Figure 1: Comparison in the Q-Q profiles between Optim-GRAMMAR and the four competing methods. The simulated phenotypes are controlled by 200 QTNs with the low, moderate and high heritabilities in human and maize. The Q-Q profiles for all simulated phenotypes are reported in Supplementary Figure 1S.

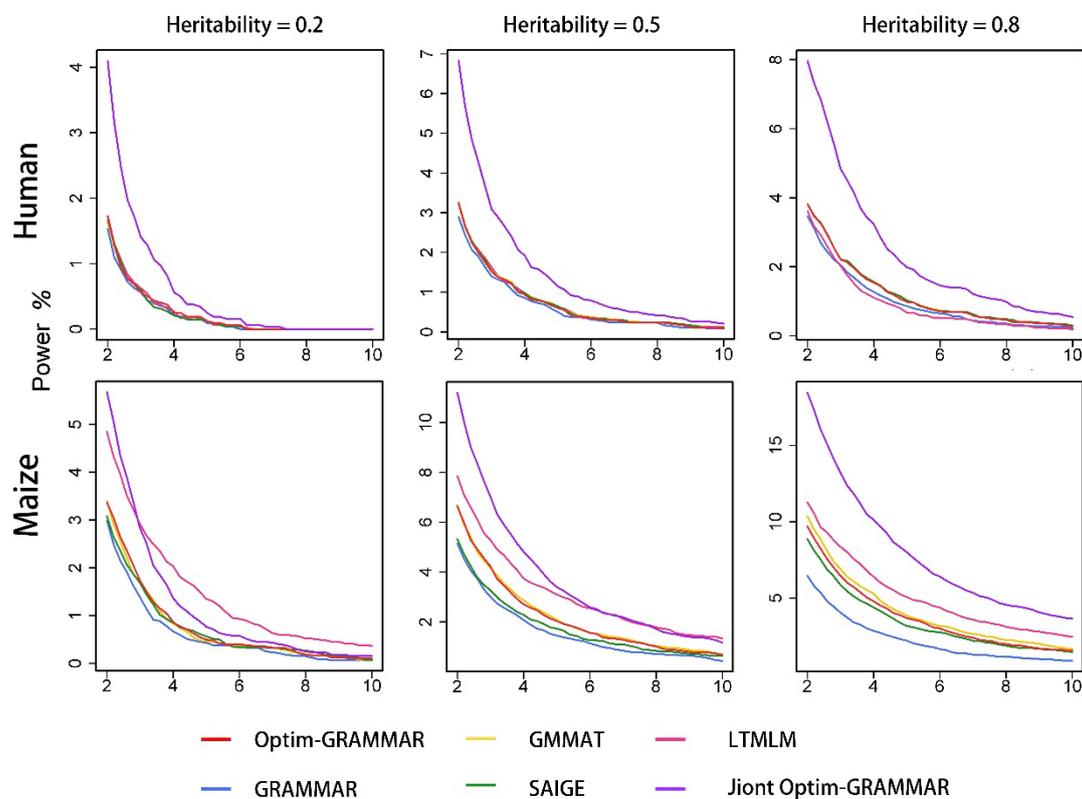


Figure 2: Comparison in the ROC profiles between Optim-GRAMMAR and the four competing methods. The ROC profiles are plotted using the statistical powers to detect QTNs relative to the given series of Type I errors. Here, the simulated phenotypes are controlled by 200 QTNs with the low, moderate and high heritabilities in human and maize. The ROC profiles for all simulated phenotypes are reported in Supplementary Figure 2S.

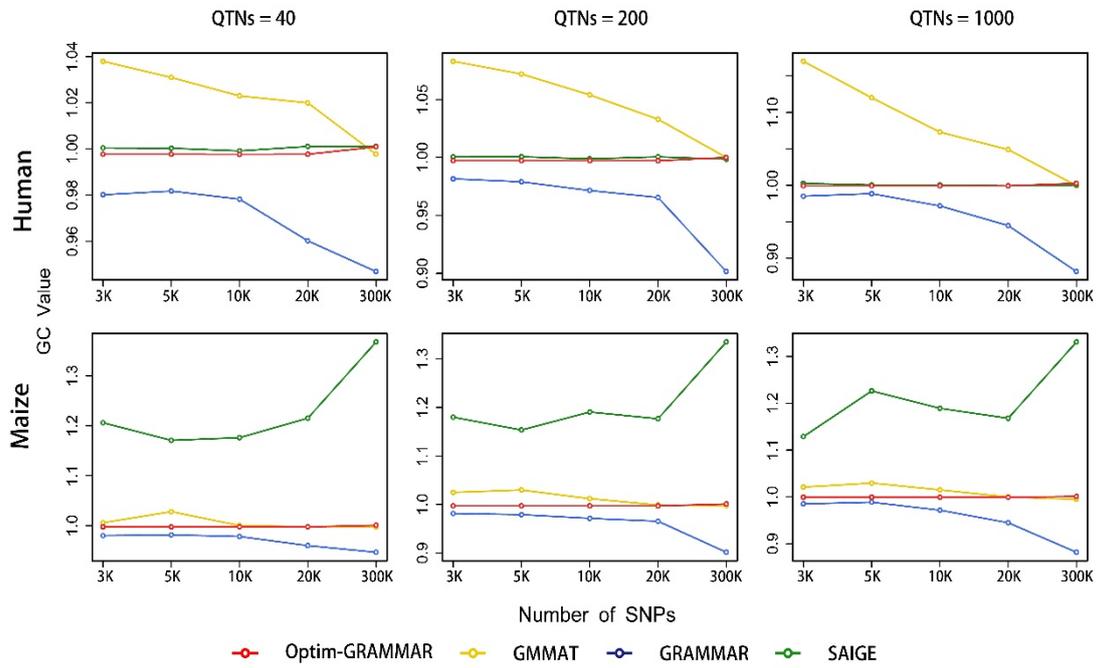


Figure 3: Changes in genomic controls with the number of sampling SNPs for Optim-GRAMMAR and the three competing methods. Genomic control is calculated by averaging genome-wide test statistics. The simulated phenotypes are controlled by 40, 200 and 1000 QTNs with the moderate heritability in human and maize.

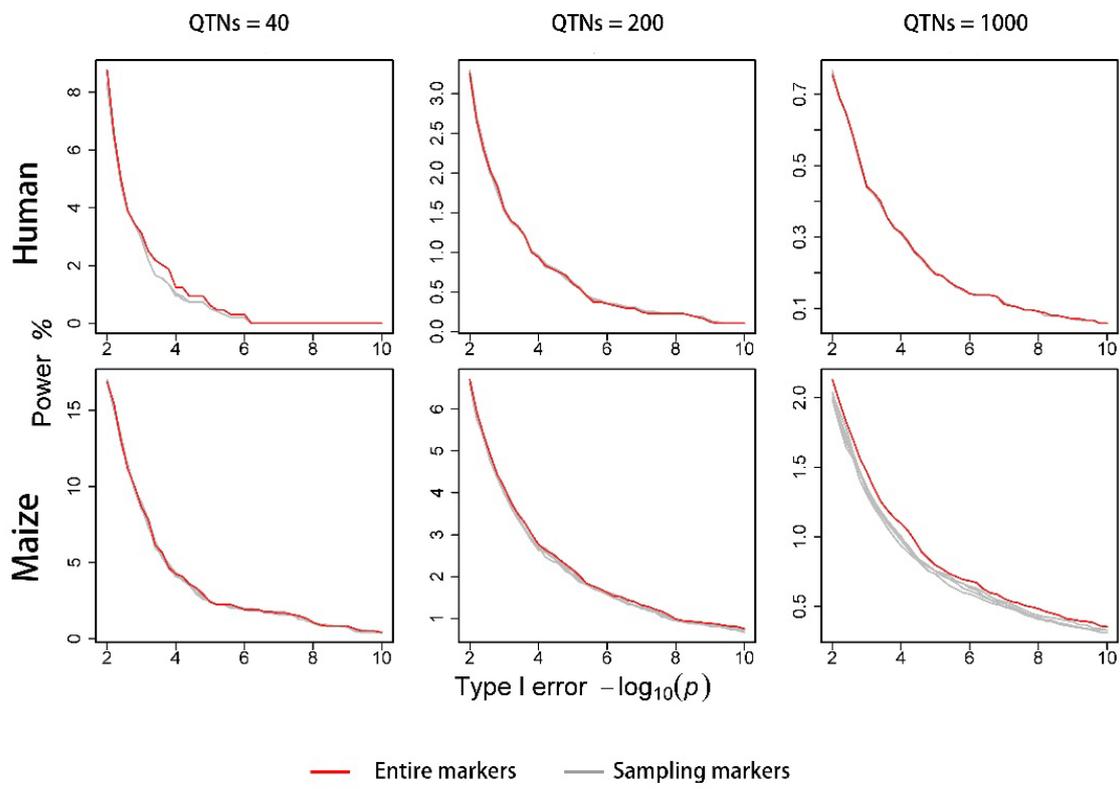


Figure 4: Changes in ROC profiles with the number of sampling SNPs for Optim-GRAMMAR. The simulated phenotypes are controlled by 40, 200, 1,000 QTNs with the moderate heritability in human and maize.

Online Methods

Genomic logit regression

Complex disease traits, as binary ones, usually follow binomial or Poisson distributions, so the generalized linear model (GLM) ^{1,2} is used to map QTLs controlling the traits. Assume that n individuals are recorded for phenotypic values and genotyped for m genetic markers. Distinguishing these markers from major and common alleles in a magnitude of effects of the markers on quantitative traits, we describe the relationship between all markers (predictors) and the mean of the exponential distribution family in the following logit regression:

$$\ln\left(\frac{\mu}{1-\mu}\right) = \mathbf{X}\mathbf{b} + \mathbf{Z}_1\mathbf{a}_1 + \mathbf{Z}_2\mathbf{a}_2 \quad (1)$$

where μ denotes the expectations of phenotypic distribution, \mathbf{b} is the systematic environment effect; the population structure (stratification) which results in phenotypic differences among subpopulations is always considered here, except for sex, age, and some initial experimental conditions. \mathbf{a}_1 is the large genetic effect of q markers on phenotype, \mathbf{a}_2 is the minor or zero effect of the $m-q$ markers on phenotype, and \mathbf{X} , \mathbf{Z}_1 , and \mathbf{Z}_2 are the corresponding design matrices of \mathbf{b} , \mathbf{a}_1 , and \mathbf{a}_2 , respectively.

GRAMMAR for binary disease traits

We define the GBVs as

$$\mathbf{g} = \mathbf{Z}_1\mathbf{a}_1 + \mathbf{Z}_2\mathbf{a}_2 \quad (2)$$

Then, model (1) becomes

$$\ln\left(\frac{\mu}{1-\mu}\right) = \mathbf{X}\mathbf{b} + \mathbf{g} \quad (3)$$

which is a the GLMM ³. Under the assumption that $(\mathbf{a}_1, \mathbf{a}_2) : N_m(\mathbf{0}, \mathbf{I}\sigma_a^2)$ with minor σ_a^2 for each marker, the GBVs are turned into random effects and $\mathbf{g} : N_n(\mathbf{0}, \mathbf{K}\sigma_g^2)$ with genomic variance of traits $\sigma_g^2 = m\sigma_a^2$ and the genomic relationship matrix (GRM) \mathbf{K} ⁴. Based on the model (3), the GBVs can be estimated with the following GBLUP equations:

$$\begin{bmatrix} \mathbf{X}^T \mathbf{W} \mathbf{X} & \mathbf{X}^T \mathbf{W} \\ \mathbf{W} \mathbf{X} & \mathbf{Z}^T \mathbf{W} \mathbf{Z} + \frac{1-h^2}{h^2} \mathbf{K}^{-1} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{b}} \\ \hat{\mathbf{g}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}^T \mathbf{W} \mathbf{y}^* \\ \mathbf{W} \mathbf{y}^* \end{bmatrix} \quad (4)$$

with

$$\mathbf{W} = \boldsymbol{\mu}(\mathbf{1} - \boldsymbol{\mu}) \quad \text{and} \quad \mathbf{y}^* = \ln\left(\frac{\boldsymbol{\mu}}{1-\boldsymbol{\mu}}\right) + \frac{\mathbf{y} - \boldsymbol{\mu}}{\boldsymbol{\mu}(1-\boldsymbol{\mu})}.$$

where \mathbf{y} is a binary phenotype, and $h^2 = \frac{\sigma_g^2}{\sigma_g^2 + 1}$ is the unknown genomic heritability of liability with the residual variance of 1 to be assumed in GLMM, which can be estimated in advance using the REML for GLMM^{3,5}.

Unlike the normally distributed quantitative traits, the residuals for binary traits cannot be directly obtained due to the difference in scale between the predictors and response variables. Within the framework of GRAMMAR, thus, we eliminated polygenic effects on the binary phenotype by regarding the estimated GBVs $\hat{\mathbf{g}}$ as a known predictor in the following GLM:

$$\ln\left(\frac{\boldsymbol{\mu}}{1-\boldsymbol{\mu}}\right) = \mathbf{z}_{\text{SNP}} a_{\text{SNP}} + \hat{\mathbf{g}} \quad (5)$$

with a regression item $\mathbf{z}_{\text{SNP}} a_{\text{SNP}}$ of the SNP tested.

With the iteratively re-weighted least square method¹, we obtained the maximum likelihood estimate for the SNP effect as:

$$\hat{a}_{\text{SNP}} = (\mathbf{z}_{\text{SNP}}^T \mathbf{W} \mathbf{z}_{\text{SNP}})^{-1} \mathbf{z}_{\text{SNP}}^T \mathbf{W} (\mathbf{y}^* - \hat{\mathbf{g}}) \quad (6)$$

The test statistic to infer the association of the SNP with binary traits is generally formulated by

$$\chi^2 = \frac{\hat{a}_{\text{SNP}}^2}{\mathbf{z}_{\text{SNP}}^T \mathbf{W} \mathbf{z}_{\text{SNP}}} \quad (7)$$

which is subject to the chi-squared distribution with the 1 degree of freedom.

Optimal Genomic control

In GRAMMAR for binary traits, replacement of polygenic effects excluding QTNs with GBVs deflates the test association statistics, which yields a high false negative rate. By regulating the downward genomic heritability, we can more accurately estimate the polygenic effects with the GBLUP Equation (4). The polygenic heritability less than genomic heritability is determined by optimizing genomic control for association tests.

Such an optimization for GRAMMAR can be summarized in the following steps:

- 1) Set the searching open interval of h^2 to (0, 1)
- 2) Estimate the GBVs $\hat{\mathbf{g}}$ using Equation (4);
- 3) Statistically infer the genetic effect for each SNP by the chi-squared statistic (7);
- 4) Calculate the genome-wide chi-squared mean or statistical probability for each SNP;
- 5) Plot the quantile-quantile (Q-Q) profile for genome-wide statistical probabilities;
- 6) Update h^2 with Brent's method ⁶;
- 7) Repeat step (2)-(6) until the genome-wide chi-squared mean reaches 1.0 plus or yields a satisfactory Q-Q plot.

Joint association analysis

After optimizing genomic control for GRAMMAR, we jointly analyzed multiple QTN candidates to improve the statistical power to detect QTNs for binary disease traits. Multiple QTN candidates were chosen within the interval of significance level 0.05 to the Bonferroni corrected criterion ⁷, so that the number of QTN candidates was limited to be no greater than the population size. Backward regression was adopted to optimize the multiple GLM with known optimized polygenic effects in a stepwise manner:

$$\ln\left(\frac{\mu}{1-\mu}\right) = \mathbf{X}\mathbf{b} + \mathbf{Z}_1\mathbf{a}_1 + \hat{\mathbf{g}} \quad (8)$$

Given the Bonferroni corrected significance level, the significant QTN effects remained in the model (8) according to the corrected statistic (7).

Simulations

Two genomic datasets of human ⁸ and maize ³ samples were used to simulate the adaptability of GRAMMAR for binary traits to population structure. The maize population has a more complex structure than the human population. Then, 300,000 SNPs for both 12000 people and 2640 maize were extracted through higher quality control. In whole simulations, control and case samples were constrained to 1:1 for the maize population, and 3000 cases were selected from the human population with low incidence rate of 5% simulated in advance. QTNs were distributed randomly over the entire SNPs, whose additive effects were sampled from a gamma distribution with

shape = 1.66 and scale = 0.4. Given the genomic heritability of liability, phenotypes of control (0) and case (1) can be generated from the genomic logit model (1).

In addition to population structure, the number of QTNs, genome heritability, and sampling number of SNPs were considered as experimental factors in the simulations. Under the optimized genomic control infinitely close to 1.0, the ROC profiles can be plotted by statistical powers to detect the QTNs relative to a given series of Type I errors. Statistical powers are defined as the percentage of identified QTNs that have the maximum test statistic among their 20 closest neighbors over the total number of simulated QTNs. Simulations were repeated 50 times, and in each simulation, the positions and effects of QTNs simulated were varied and the average results were recorded.

Method References

1. Wedderburn, R.W.M. Quasi-likelihood functions, generalized linear models, and the gauss-newton method. *Biometrika* **61**, 439-447 (1974).
2. McCullagh, P. & Nelder, J.A. *Generalized linear models, 2nd ed.*, (Chapman and Hall, New York, 1989).
3. Breslow, N.E. & Clayton, D.G. Approximate inference in generalized linear mixed models. *Journal of the American statistical Association* **88**, 9-25 (1993).
4. Vanraden, P.M. Efficient methods to compute genomic predictions. *Journal of Dairy Science* **91**, 4414-4423 (2008).
5. Chen, H. *et al.* Control for Population Structure and Relatedness for Binary Traits in Genetic Association Studies via Logistic Mixed Models. *Am J Hum Genet* **98**, 653-66 (2016).
6. Brent, R.P. *Algorithms for minimization without derivatives*, (Prentice-Hall, New Jersey, 1973).
7. Hochberg, Y. & Tamhane, A.C. *Multiple Comparison Procedures*, (John Wiley & Sons, Inc., New York, 1987).
8. Romay, M.C. *et al.* Comprehensive genotyping of the USA national maize inbred seed bank. *Genome biology* **14**, R55 (2013).