

1 Efficient implementation of penalized regression
2 for genetic risk prediction

3 Florian Privé ^{1,*}, Hugues Aschard ² and Michael G.B. Blum ^{1,*}

4
5 ¹Université Grenoble Alpes, CNRS, Laboratoire TIMC-IMAG, UMR 5525, France,

6 ²Centre de Bioinformatique, Biostatistique et Biologie Intégrative (C3BI), Institut Pasteur, Paris,
7 France.

8 *To whom correspondence should be addressed.

Abstract

Polygenic Risk Scores (PRS) consist in combining the information across many single-nucleotide polymorphisms (SNPs) in a score reflecting the genetic risk of developing a disease. PRS might have a major public health impact, possibly allowing for screening campaigns to identify high-genetic risk individuals for a given disease. The “Clumping+Thresholding” (C+T) approach, which is the most common method to derive PRS, uses only univariate genome-wide association studies (GWAS) summary statistics, which makes it fast and easy to use. However, previous work showed that jointly estimating SNP effects for computing PRS has the potential to significantly improve the predictive performance of PRS as compared to C+T.

In this paper, we present an efficient method to jointly estimate SNP effects, allowing for practical application of penalized logistic regression on modern datasets including hundreds of thousands of individuals. Moreover, our implementation of penalized logistic regression directly includes automatic choices for hyper-parameters. The choice of hyper-parameters for a predictive model is very important since it can dramatically impact its predictive performance. As an example, AUC values range from less than 60% to 90% in a model with 30 causal SNPs, depending on the p-value threshold in C+T.

We compare the performance of penalized logistic regression to the C+T method and to a derivation of random forests. Penalized logistic regression consistently achieves higher predictive performance than the two other methods while being very fast. We find that improvement in predictive performance is more pronounced when there are few effects located in nearby genomic regions with correlated SNPs; AUC values increase from 83% with the best prediction of C+T to 92.5% with penalized logistic regression. We confirm these results in a data analysis of a case-control study for celiac disease where penalized logistic regression and the standard C+T method achieve AUC of 89% and of 82.5%.

In conclusion, our study demonstrates that penalized logistic regression is applicable to large-scale individual-level data and can achieve more discriminative polygenic risk scores. Our implementation is publicly available in R package `bigstatsr`.

Contact: florian.prive@univ-grenoble-alpes.fr & michael.blum@univ-grenoble-alpes.fr

Supplementary information:

1 Introduction

Polygenic Risk Scores (PRS) consist in combining the information across many single-nucleotide polymorphisms (SNPs) in a score reflecting the genetic risk of developing a disease. PRS are useful for genetic epidemiology when testing the polygenicity of one disease and finding a common genetic contribution between two diseases (Purcell *et al.* 2009). Personalized medicine is another major application of PRS. Personalized medicine envisions to use PRS in screening campaigns in order to identify high-risk individuals for a given disease (Chatterjee *et al.* 2016). As an example of practical application, targeting screening to men at higher polygenic risk could reduce the problem of overdiagnosis and lead to a better benefit-to-harm balance in screening for prostate cancer (Pashayan *et al.* 2015). Yet, PRS would have to show a high discriminative power between cases and controls in order to be used for helping in the diagnosis of diseases. For screening high-risk individuals and for presymptomatic diagnosis of the general population, it is suggested that the AUC must be greater than 75% and 99% respectively (Janssens *et al.* 2007).

Several methods have been developed to predict disease status, or more generally any phenotype, based on SNP information. A commonly used method, called “P+T” or “C+T” (which stands for “Clumping and Thresholding”) – or even just PRS – is used to derive PRS from results of Genome-Wide Association Studies (GWAS) (Chatterjee *et al.* 2013; Dudbridge 2013; Evans *et al.* 2009; Purcell *et al.* 2009; Wray *et al.* 2007). This technique uses GWAS summary statistics only, allowing for fast implementation. However, the “C+T” approach also has several limitations. Previous studies have shown that predictive performance of the C+T method is very sensitive to the threshold of inclusion of SNPs, depending on the disease architecture (Ware *et al.* 2017). Linear Mixed-Models (LMMs) are another widely-used method in fields such as plant and animal breeding or for predicting highly heritable quantitative human phenotypes such as height (Lello *et al.* 2017; Yang *et al.* 2010). Yet, models resulting from LMM, known e.g. as “gBLUP”, are not optimal for predicting disease status based on genotypes (Abraham *et al.* 2013). Moreover, these methods and their derivatives are often computationally demanding, both in terms of memory and time required, which makes them

67 unlikely to be used for prediction on very large datasets (Golan and Rosset 2014; Maier *et al.*
68 2015; Speed and Balding 2014; Zhou *et al.* 2013). Finally, statistical learning methods have
69 also been used to derive PRS for complex human diseases by jointly estimating SNP effects.
70 Such methods include joint logistic regression, Support Vector Machine (SVM) and random
71 forests (Abraham *et al.* 2012, 2014; Botta *et al.* 2014; Okser *et al.* 2014; Wei *et al.* 2009).

72 We recently developed two R packages, *bigstatsr* and *bigsnpr*, for efficiently analyzing
73 large-scale genome-wide data (Privé *et al.* 2018). Package *bigstatsr* now includes an efficient
74 algorithm with a new implementation for computing sparse linear and logistic regressions on
75 huge datasets as large as the UK Biobank (Bycroft *et al.* 2017). In this paper, we present a com-
76 prehensive comparative study of our implementation of penalized logistic regression against
77 the C+T method and the T-Trees algorithm, a derivation of random forests that has shown high
78 predictive performance (Botta *et al.* 2014). In this comparison, we do not include any LMM
79 method for the reasons mentioned before and do not include any SVM method because it is ex-
80 pected to give similar results to logistic regression (Abraham *et al.* 2012). For the C+T model,
81 we report results for a large grid of hyper-parameters. For the penalized logistic regression, the
82 choice of hyper-parameters is included in the algorithm so that we report only one model for
83 each simulation. We also use this penalized logistic regression on a feature-augmented dataset
84 in order to capture not only linear effects, but also recessive and dominant effects.

85 To perform simulations, we use real genotype data and simulate new phenotypes (Zhou
86 *et al.* 2013). In order to make our comparison as comprehensive as possible, we compare
87 different disease architectures by varying the number, size and location of causal effects as well
88 as the heritability. We also compare different models for simulating phenotypes, one with only
89 linear effects, and one that combines linear, dominant and interaction-type effects. Overall, we
90 find that the penalized logistic regression consistently achieves higher predictive performance
91 than the C+T and T-Trees methods while being very fast. This demonstrates the feasibility and
92 relevance of this approach for PRS computation on large modern datasets.

2 Methods

2.1 Genotype data

We use real genotypes of European individuals from a case-control celiac disease study (Dubois *et al.* 2010). The composition of this dataset is presented in table S1. Details of quality control and imputation for this dataset are available in Privé *et al.* (2018). For simulations presented later, we first restrict this dataset to controls in order to remove the genetic structure induced by the celiac disease status. Then, we decided to remove population structure because it can affect the predictive performance of methods (Martin *et al.* 2017). In order to alleviate population structure, we keep people from the UK only and we further remove outliers based on a robust Mahalanobis distances computed using the first 10 principal components of the remaining individuals. This 3-step filtering process results in a sample of 7100 individuals with minimal population structure (see supplementary notebook “preprocessing”). We also use this dataset for real data application, in this case keeping all 15,155 individuals (4496 cases and 10,659 controls). Both datasets contain 281,122 SNPs.

2.2 Simulations of phenotypes

We simulate binary phenotypes using a Liability Threshold Model (LTM) with a prevalence of 30% (Falconer 1965). We vary parameters of the simulations in order to match a range of genetic architecture from low to high polygenicity. This is achieved by varying the number of causal variants and their location (30, 300, or 3000 anywhere in all 22 chromosomes or 30 in the HLA region of chromosome 6), and the heritability (50% or 80%). In a second phase, in order to further increase the proportion of causal variants, we restrict the dataset to chromosome 6 only (18,941 SNPs) instead of using all 22 autosomal chromosomes (281,122 SNPs). We also consider deviation from the standard normal additive model, drawing effects of causal SNPs from either a Normal or from a Laplace distribution, and computing liability scores either from a “simple” model with linear effects only or a “fancy” model that combines linear, dominant and interaction-type effects. For the “simple” model, we compute the liability score of the i -th

individual

$$y_i = \sum_{j \in S_{\text{causal}}} w_j \cdot \widetilde{G}_{i,j} + \epsilon_i,$$

108 where w_j are weights generated from a Gaussian or a Laplace distribution, $G_{i,j}$ is the allele
109 count of individual i for SNP j , $\widetilde{G}_{i,j}$ corresponds to its standardized version (zero mean and
110 unit variance for all SNPs), ϵ follows a Gaussian distribution $N(0, 1 - h^2)$ and S_{causal} is the
111 set of causal SNPs. For the “fancy” model, we simulate phenotypes using linear, dominant and
112 interaction-type effects (see Supplementary Materials).

113 We implement 3 different simulation scenarios, summarized in table 2. Scenario №1 uses
114 the whole dataset (all 22 autosomal chromosomes) and a training set of size 6000. It com-
115 pares all methods described in section 2.4. For each combination of the remaining parameters,
116 results are based on 100 simulations excepted in the first simulation comparing penalized lo-
117 gistic regression with T-Trees; this simulation relies on 5 simulations only because of a much
118 higher computational burden (several hours of computation for a single simulation) of T-Trees
119 as compared to other approaches. Scenario №2 consists of 100 simulations per combination of
120 parameters on a dataset composed of chromosome 6 only. Reducing the number of SNPs aims
121 at increasing the polygenicity (the proportion of causal SNPs) of the simulated models and at
122 virtually increasing the sample size (Dudbridge 2013; Márquez-Luna *et al.* 2017; Vilhjálmsson
123 *et al.* 2015). For this scenario, we use the additive (“simple”) model only, but continue to
124 compare all previous different values of the other parameters. Finally, scenario №3 reuses the
125 whole dataset but this time varying the size of the training set in order to assess how the sample
126 size affects predictive performance of methods. A total of 100 simulations per combination of
127 parameters are run using 300 causal SNPs randomly chosen anywhere on the genome.

128 2.3 Predictive performance measures

129 In this study, we use two different measures of predictive accuracy. First, we use the Area Un-
130 der the Receiver Operating Characteristic (ROC) Curve (AUC) (Fawcett 2006; Lusted 1971).
131 In the case of our study, the AUC is the probability that the PRS of a case is greater than the
132 PRS of a control. This measure indicates the extent to which we can distinguish between cases

133 and controls using PRS. As a second measure, we also report the partial AUC for specificities
134 between 90% and 100% (Dodd and Pepe 2003; McClish 1989). This measure is similar to
135 the AUC, but focuses on high specificities, which is the most useful part of the ROC curve
136 in clinical settings. When reporting AUC results of simulations, we use estimates of maxi-
137 mum achievable AUC of 84% and 94% for heritabilities of respectively 50% and 80%. These
138 estimates are based on three different yet consistent estimations (see Supplementary Materials).

139 2.4 Methods compared

140 In this study, we compare three different types of methods: the C+T method, T-Trees and
141 penalized logistic regression.

The C+T (Clumping + Thresholding) method directly derives a Polygenic Risk Score (PRS) from the results of Genome-Wide Associations Studies (GWAS). In GWAS, a coefficient of regression (i.e. the estimated effect size $\hat{\beta}_j$) is learned independently for each SNP j along with a corresponding p-value p_j . The SNPs are first clumped (C) so that there remain only loci that are weakly correlated with one another (this set of SNPs is denoted S_{clumping}). Then, thresholding (T) consists in removing SNPs with p-values larger than a threshold p_T to be determined. Finally, a PRS is defined as the sum of allele counts of the remaining SNPs weighted by the corresponding effect coefficients

$$\text{PRS}_i = \sum_{\substack{j \in S_{\text{clumping}} \\ p_j < p_T}} \hat{\beta}_j \cdot G_{i,j},$$

142 where $\hat{\beta}_j$ (p_j) are the effect sizes (p-values) learned from the GWAS. In this study, we mostly
143 report scores for a clumping threshold at $r^2 > 0.2$ within regions of 500kb, but we also investi-
144 gate thresholds of 0.05 and 0.8. We report three different scores of prediction: one including all
145 the SNPs remaining after clumping (denoted “PRS-all”), one including only SNPs remaining
146 after clumping and that have a p-value under the GWAS threshold of significance ($p < 5 \cdot 10^{-8}$,
147 “PRS-stringent”), and one that maximizes the AUC (“PRS-max”) for these two thresholds (0
148 and $5 \cdot 10^{-8}$) and a sequence of 100 values of thresholds ranging from $10^{-0.1}$ to 10^{-100} and

149 equally spaced on the log-log-scale (Table S2). As we report the optimal threshold based on
150 the test set, the AUC for “PRS-max” is an upper bound of the AUC for the C+T method.

151 T-Trees (*Trees inside Trees*) is an algorithm derived from random forests (Breiman 2001)
152 that takes into account the correlation structure among the genetic markers implied by linkage
153 disequilibrium in GWAS data (Botta *et al.* 2014). We use the same parameters as reported in
154 Table 4 of Botta *et al.* (2014), except that we use 100 trees instead of 1000 because using 1000
155 trees provides a minimal increase of AUC while requiring a disproportionately long processing
156 time (e.g. AUC of 81.5% instead of 81%, data not shown).

Finally, for the penalized logistic regression, we find regression coefficients β_0 and β that minimize the following regularized loss function

$$L(\lambda, \alpha) = \underbrace{-\sum_{i=1}^n (y_i \log(p_i) + (1 - y_i) \log(1 - p_i))}_{\text{Loss function}} + \lambda \underbrace{\left((1 - \alpha) \frac{1}{2} \|\beta\|_2^2 + \alpha \|\beta\|_1 \right)}_{\text{Penalization}},$$

157 where $p_i = 1 / (1 + \exp(-(\beta_0 + x_i^T \beta)))$, x is denoting the genotypes and covariables (e.g.
158 principal components), y is the disease status to predict, λ and α are two regularization hyper-
159 parameters that need to be chosen. Different regularizations can be used to prevent overfitting,
160 among other benefits: the L2-regularization (“ridge”, Hoerl and Kennard (1970)) shrinks coeffi-
161 cients and is ideal if there are many predictors drawn from a Gaussian distribution (corresponds
162 to $\alpha = 0$ in the previous equation); the L1-regularization (“lasso”, Tibshirani (1996)) forces
163 some of the coefficients to be equal to zero and can be used as a means of variable selection,
164 leading to sparse models (corresponds to $\alpha = 1$); the L1- and L2-regularization (“elastic-net”,
165 Zou and Hastie (2005)) is a compromise between the two previous penalties and is particularly
166 useful in the $m \gg n$ situation (m : number of SNPs), or any situation involving many corre-
167 lated predictors (corresponds to $0 < \alpha < 1$) (Friedman *et al.* 2010). In this study, we use an
168 embedded grid search over $\alpha \in \{1, 0.5, 0.05, 0.001\}$.

169 To fit this penalized logistic regression, we use a very efficient algorithm (Friedman *et al.*
170 2010; Tibshirani *et al.* 2012; Zeng *et al.* 2017) from which we derived our own implementation
171 in R package bigstatsr. This type of algorithm builds predictions for many values of λ , which is
172 called a “regularization path”. To obtain an algorithm free of the choice of this hyper-parameter

173 λ , we developed a procedure that we call Cross-Model Selection and Averaging (CMSA, figure
174 S1). Because of L1-regularization, the resulting vectors of coefficients are sparse and can be
175 used to make a PRS based on a *linear* combination of allele counts. We refer to this method as
176 “logit-simple” in the results section.

177 In order to capture recessive and dominant effects in addition to linear effects, we use
178 feature engineering: we construct a separate dataset with, for each SNP variable, two more
179 variables coding for recessive and dominant effects: one variable is coded 1 if homozygous
180 variant and 0 otherwise, and the other is coded 0 for homozygous referent and 1 otherwise.
181 This results in a dataset with 3 times as many variables as the initial one, on which we can
182 apply penalized logistic regression. We refer to this method “logit-triple” in the results.

183 **2.5 Evaluating predictive performance for Celiac data**

184 We use Monte Carlo cross-validation to compute AUC, partial AUC, the number of predictors
185 and execution time for the original Celiac dataset with real phenotypes: we randomly split 100
186 times the dataset in a training set of 12,000 individuals and a test set composed of the remaining
187 3155 individuals.

188 **2.6 Reproducibility**

189 All the code used in this paper along with results such as figures and tables, are available as
190 HTML R notebooks in the Supplementary Materials.

191 **3 Results**

192 **3.1 Joint estimation improves predictive performance**

193 We compared penalized logistic regression (“logit-simple”) with the C+T method (“PRS”) us-
194 ing whole-genome simulations of scenario N°1 (Table 2).

195 When simulating a model with 30 causal SNPs and an heritability of 80%, penalized logis-
196 tic regression provides AUC greater than 93%, nearly reaching the maximum achievable AUC

197 of 94%, whereas AUC values obtained with C+T method range between 83% and 90% (Figures
198 1 and 2). Moreover, penalized logistic regression consistently provides higher predictive per-
199 formance than the C+T method across all scenarios we considered, excepted in some cases of
200 high polygenicity or small sample size where all methods perform poorly (AUC values below
201 60% – figures 3 and S3).

202 Method “logit-simple” provides particularly higher predictive performance than “PRS-max”
203 when there are correlations between predictors, i.e. when we choose causal SNPs to be in the
204 HLA region. In this situation, the mean AUC reaches 92.5% with the “logit-simple” approach
205 and 84% with “PRS-max”, while the maximum achievable AUC is 94% (Figure 1).

206 Note that for the simulations we do not report results in terms of partial AUC because partial
207 AUC values have a Spearman correlation of 98% with the AUC results for all methods (Figure
208 S2).

209 **3.2 Importance of hyper-parameters**

210 In practice, a particular value of the threshold of inclusion of SNPs should be chosen for the
211 C+T method and this choice can dramatically impact the predictive performance of C+T. For
212 example, in a model with only 30 causal SNPs, AUC ranges from less than 60% when using all
213 SNPs passing clumping to 90% if choosing the optimal p-value threshold (Figures 2 and S4).

214 Concerning the r^2 threshold of the clumping step in C+T, we mostly used the common
215 value of 0.2. Yet, using a more stringent value of 0.05 provides higher predictive performance
216 than using 0.2 in most of the cases tested in this paper (Figures S5, 3 and S6)

217 Method “logit-simple” that automatically chooses hyper-parameter λ provides similar pre-
218 dictive performance than the best predictive performance of the implementation of R package
219 biglasso (Zeng *et al.* 2017), only slightly better for biglasso, which is likely due to over-fitting
220 when reporting the best prediction (Figure S10).

221 **3.3 Non-linear effects**

222 We tested the T-Trees method in scenario №1. As compared to “logit-simple”, T-Trees perform
223 worse in terms of predictive ability, while taking much longer to run and making more complex
224 predictive models because T-Trees use more predictors and non-linear effects (Figure S7). Even
225 when simulating a “fancy” model in which there are dominant and interaction-type effects that
226 T-Trees should be able to handle, AUC is still lower when using T-Trees than when using
227 “logit-simple” (Figure S7).

228 We also compared the two penalized logistic regressions in scenario №1, “logit-simple”
229 and “logit-triple” that uses additional features (variables) coding for recessive and dominant
230 effects. Predictive performance of “logit-triple” are nearly as good as “logit-simple” when
231 there are only linear effects (differences of AUC are always smaller than 2%) and can lead to
232 significantly greater results when there are also dominant and interactions effects (Figures S8
233 and S9). For the “fancy model”, “logit-triple” provides AUC values at least 3.5% higher than
234 “logit-simple”, excepted when there are 3000 causal SNPs. Yet, the “triple” solution takes 2-3
235 times as much time to run and requires 3 times as much disk storage as the “simple” solution.

236 **3.4 Simulations varying number of SNPs and training size**

237 First, when reproducing simulations of scenario №1 using chromosome 6 only (scenario №2),
238 the predictive performance of “logit-simple” always increase (Figure S6). There is particularly
239 a large increase when simulating 3000 causal SNPs: AUC from the “logit-simple” increases
240 from 60% to nearly 80% for Gaussian effects and a heritability of 80%. On the contrary, when
241 simulating only 30 or 300 causal SNPs on the corresponding dataset, AUC of the “PRS-max”
242 does not increase, and even decreases for an heritability of 80% (Figure S6). Secondly, when
243 varying the training size (scenario №3), we report an increase of AUC when increasing the
244 training size, with a faster increase of AUC provided by “logit-simple” as compared to “PRS-
245 max” (Figure 3).

3.5 Polygenic scores for the celiac disease

Joint logistic regressions also provide higher AUC values for the Celiac data: 88.7% with “logit-simple” and 89.1% with “logit-triple” as compared to 82.5% with the C+T method. The relative increase in partial AUC, for specificities larger than 90%, is even larger (42% and 47%) with partial AUC values of 0.0411, 0.0426 and 0.0289 obtained with “logit-simple”, “logit-triple” and the C+T method, respectively. Moreover, logistic regressions use less predictors, respectively 1570, 2260 and 8360 (Table 1, figures 4 and supplementary notebook “results-celiac”). Note that for the C+T method, we still report the best result among 102 p-value thresholds. In terms of computation time, we report only the GWAS computation for the C+T method and we show that the “logit-simple” method, while learning jointly on all SNPs at once and testing different hyper-parameter values, is almost as fast as the C+T method (190 vs 130 seconds), and the “logit-triple” takes less than twice as long as the “logit-simple” (296 vs 190 seconds).

Table 1: Results for the real Celiac dataset. The results are averaged over 100 runs where the training step is randomly composed of 12,000 individuals. In the parentheses is reported the standard deviation of 10^5 bootstrap samples of the mean of the corresponding variable. Results are reported with 3 significant digits.

Method	AUC	pAUC	# predictors	Execution time (s)
PRS-max	0.825 (0.000664)	0.0289 (0.000187)	8360 (744)	130 (0.143)
logit-simple	0.887 (0.00061)	0.0411 (0.000224)	1570 (46.4)	190 (1.21)
logit-triple	0.891 (0.000628)	0.0426 (0.000219)	2260 (56.1)	296 (2.03)

4 Discussion

4.1 Joint estimation improves predictive performance

In this comparative study, we present a computationally efficient implementation of penalized logistic regression. This model can be used to build polygenic risk scores based on very large SNP datasets such as the UK biobank (Bycroft *et al.* 2017). In agreement with previous work (Abraham *et al.* 2013), we show that jointly estimating SNP effects has the potential to sub-

265 stantially improve predictive performance as compared to the standard C+T approach in which
266 SNP effects are learned independently. Penalized logistic regression nearly always outperform
267 the C+T method, and the benefits of using it are more pronounced with an increasing sample
268 size or when causal SNPs are correlated with one another.

269 **4.2 Importance of hyper-parameters**

270 The choice of hyper-parameter values is very important since it can greatly impact method
271 performance. In the C+T method, there are two main hyper-parameters: the r^2 and the p_T
272 thresholds that control how stringent are the clumping and thresholding steps, respectively.
273 The choice of the r^2 threshold of the clumping step is important. Indeed, on the one hand,
274 choosing a low value for this threshold may discard independently predictive SNPs that are in
275 Linkage Disequilibrium; yet, on the other hand, when choosing a high value for this threshold,
276 too much redundant information would be included in the model, which would bias SNP ef-
277 fects. Based on the simulations, we find that using a stringent threshold ($r^2 = 0.05$) leads to
278 higher predictive performance, even when causal SNPs are correlated. It means that accurately
279 estimating SNP effects is more important than including all causal SNPs. Moreover, in this pa-
280 per, we reported the maximum AUC of 102 different p-value thresholds, a threshold that should
281 normally be learned on the training set only. The choice of this threshold is very important as
282 it can greatly impact the predictive performance of the C+T method, which we confirm in this
283 study (Ware *et al.* 2017).

284 On the contrary, in the penalized logistic regression presented here, we developed an auto-
285 matic procedure called Cross-Model Selection and Averaging (CMSA) that releases investiga-
286 tors from the burden of choosing hyper-parameter λ that accounts for the amount of regular-
287 ization used in the model. Not only this procedure provides near-optimal results (as compared
288 to the best prediction when using R package biglasso), but it also accelerates the training of
289 the model thanks to the development of an early stopping criterion. Usually, cross-validation is
290 used to choose hyper-parameter values and then the model is trained again with these particular
291 hyper-parameter values (Hastie *et al.* 2008; Wei *et al.* 2013). Yet, performing cross-validation
292 and retraining the model is computationally demanding; CMSA offers a less burdensome alter-

293 native. Concerning hyper-parameter α that accounts for the relative importance of the L1 and
294 L2 regularizations, we use a grid search directly embedded in the CMSA procedure.

295 **4.3 Non-linear effects**

296 In this paper, we also explored how to capture non-linear effects. For this, we introduced a
297 simple feature engineering technique that enables logistic regression to detect and learn not
298 only additive effects, but also dominant and recessive effects. This technique improves the
299 predictive performance of logistic regression when there are some non-linear effects in the
300 simulations, while providing nearly the same predictive performance when there are only linear
301 effects. Moreover, it also improves predictive performance for the celiac disease.

302 Yet, this approach is not able to detect interaction-type effects. In order to capture interaction-
303 type effects, we tested T-Trees, a method that is able to exploit SNP correlations thanks to
304 special decision trees (Botta *et al.* 2014). However, predictive performance of T-Trees were
305 consistently lower than with penalized logistic regression, even when simulating a model with
306 dominant and interaction-type effects that T-Trees should be able to handle.

307 **4.4 Limitations**

308 Our approach has one major limitation: the main advantage of the C+T method is that it is
309 applicable directly to summary statistics, allowing to leverage the largest GWAS sample size
310 to date, even when individual cohort data cannot be merged because of practical and ethical
311 reasons (e.g. consortium data including many cohorts). As of today, the proposed penalized
312 logistic regression does not allow for the analysis of summary data, but this represents an
313 important future direction of our work. The current version is of particular interest for the
314 analysis of modern SNP dataset including hundreds of thousands of individuals.

315 Finally, in this comparative study, we did not consider the problem of population structure
316 (Márquez-Luna *et al.* 2017; Martin *et al.* 2017; Vilhjálmsón *et al.* 2015) and also did not
317 consider non-genetic data such as environmental and clinical data (Dey *et al.* 2013; Van Vliet
318 *et al.* 2012). In next study, we will assess how can we use models and effects learned in one

319

population to improve learning and prediction in another population.

Table 2: Summary of all simulations. Where there is symbol ‘-’ in a box, it means that the parameters are the same as the ones in the upper box.

Numero of scenario	Dataset	Size of training set	Causal SNPs (number and location)	Distribution of effects	Heritability	Simulation model	Methods
1	All 22 chromosomes	6000	30 in HLA 30 in all 300 in all 3000 in all	Gaussian Laplace	0.5 0.8	simple fancy	PRS logit-simple logit-triple (T-Trees)
2	Chromosome 6 only	-	-	-	-	simple	PRS logit-simple
3	All 22 chromosomes	1000 2000 3000 4000 5000	300 in all	-	-	-	-

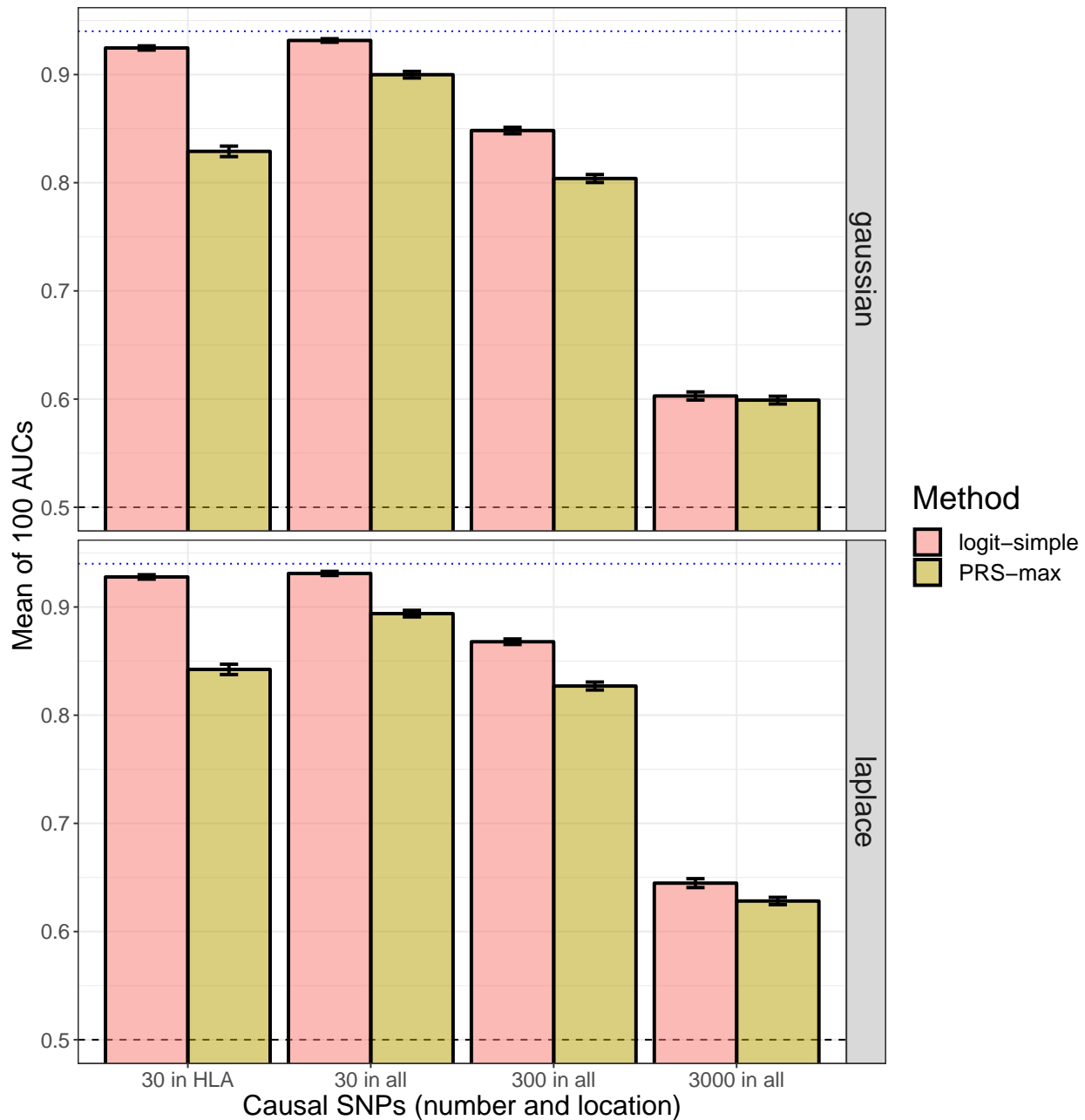


Figure 1: Main comparison of the C+T and “logit-simple” methods in scenario №1 for the “simple” model and an heritability of 80%. Mean of AUC over 100 simulations for “logit-simple” and the maximum AUC reported with the C+T method (“PRS-max”). Upper (lower) panel is presenting results for effets following a Gaussian (Laplace) distribution. Error bars are representing $\pm 2SD$ of 10^5 non-parametric bootstrap of the mean of AUC. The blue dotted line represents the maximum achievable AUC.

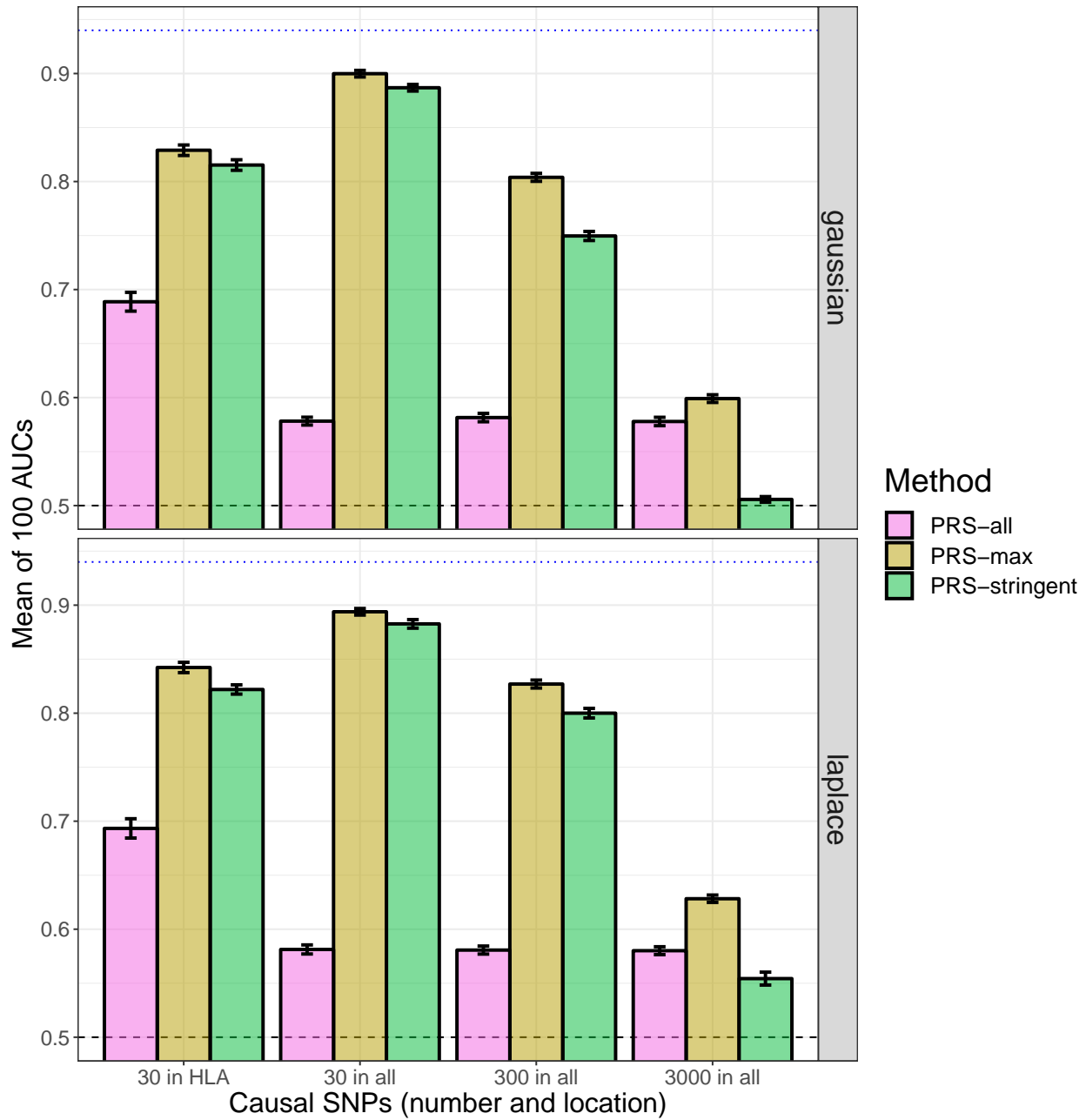


Figure 2: Comparison of three different p-value thresholds used in the C+T method in scenario N°1 for the “simple” model and an heritability of 80%. Mean of AUC over 100 simulations. Upper (lower) panel is presenting results for effets following a Gaussian (Laplace) distribution. Error bars are representing $\pm 2SD$ of 10^5 non-parametric bootstrap of the mean of AUC. The blue dotted line represents the maximum achievable AUC.

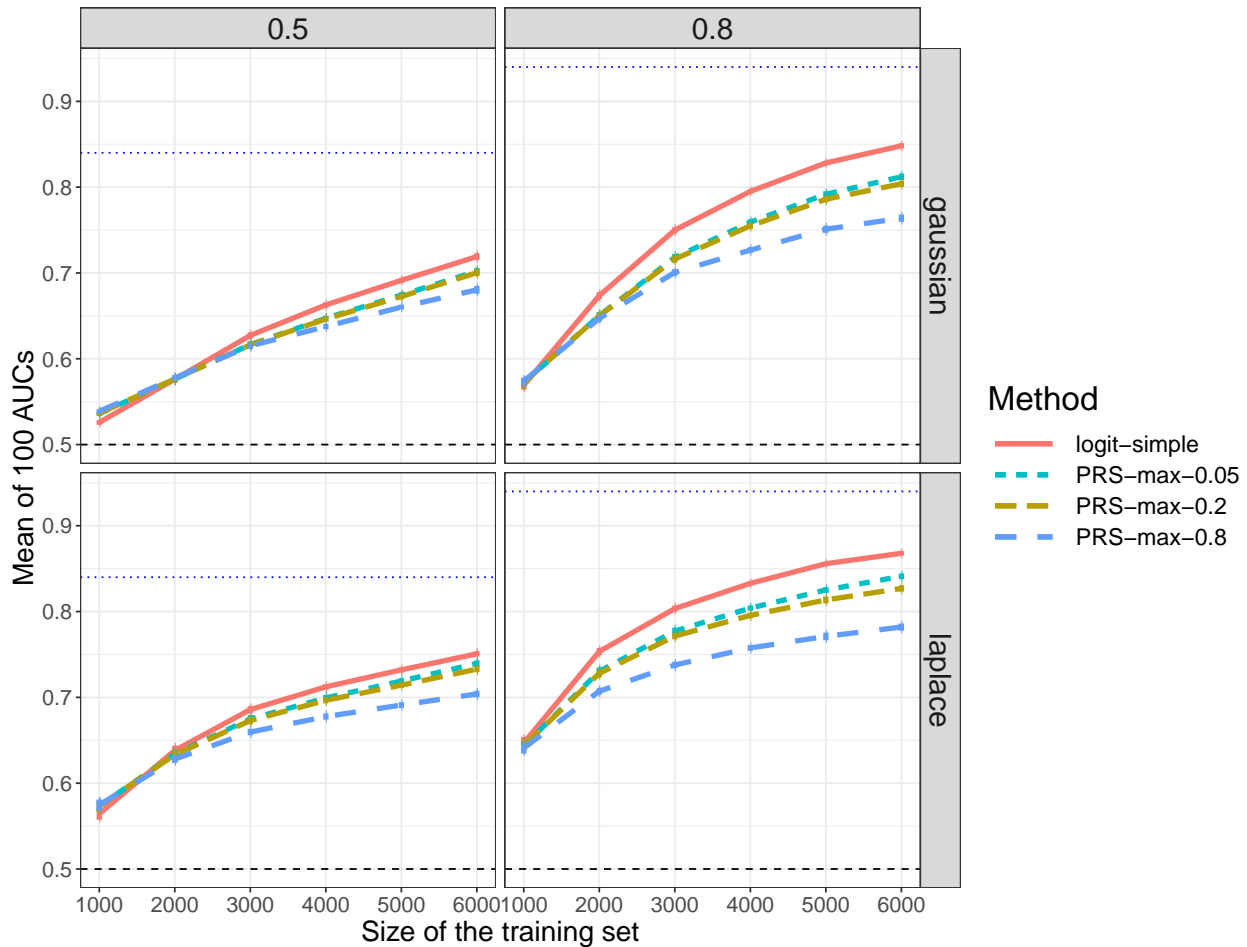


Figure 3: Comparison of models when varying sample size in scenario №3 for the “simple” model with 300 causal SNPs sampled anywhere on the genome. Mean of AUC over 100 simulations for the maximum values of PRS for three different r^2 thresholds (0.05, 0.2 and 0.8) and “logit-simple” as a function of the training size. Upper (lower) panels are presenting results for effects following a Gaussian (Laplace) distribution and left (right) panels are presenting results for an heritability of 0.5 (0.8). Error bars are representing $\pm 2SD$ of 10^5 non-parametric bootstrap of the mean of AUC. The blue dotted line represents the maximum achievable AUC.

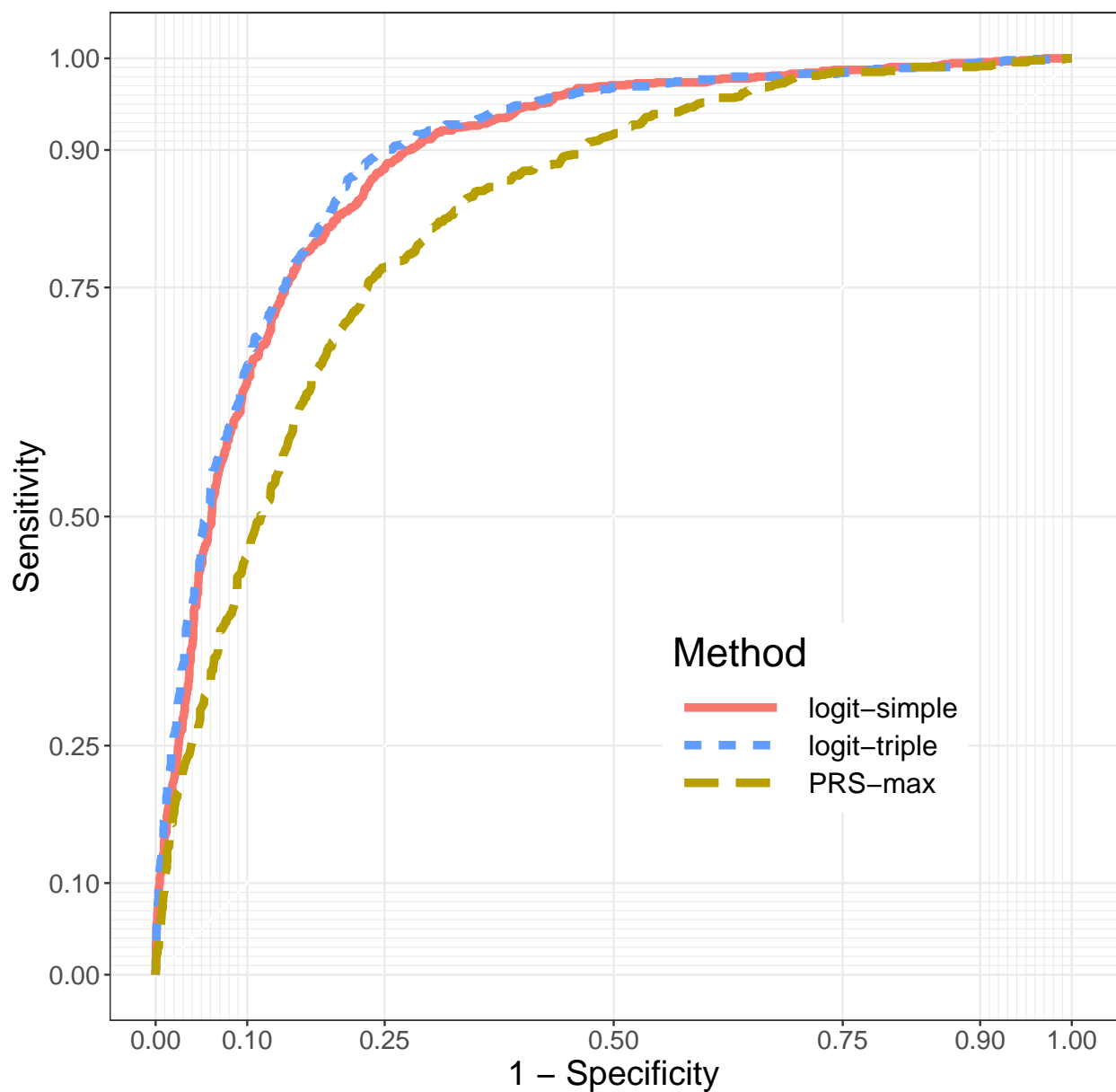


Figure 4: ROC Curves for the “C+T”, “logit-simple” and “logit-triple” methods for Celiac disease dataset. Models were trained using 12,000 individuals. These are results projecting these models on the remaining 3155 individuals. The figure is plotted using R package plotROC (Sachs *et al.* 2017).

Acknowledgements

Authors acknowledge LabEx PERSYVAL-Lab (ANR-11-LABX-0025-01). Authors also acknowledge the Grenoble Alpes Data Institute that is supported by the French National Research Agency under the “Investissements d’avenir” program (ANR-15-IDEX-02). We are also grateful to Félix Balazard for useful discussions about T-Trees, and to Yaohui Zeng for useful discussions about R package biglasso.

References

- Abraham, G., Kowalczyk, A., Zobel, J., and Inouye, M. (2012). Sparsnp: Fast and memory-efficient analysis of all snps for phenotype prediction. *BMC bioinformatics*, **13**(1), 88.
- Abraham, G., Kowalczyk, A., Zobel, J., and Inouye, M. (2013). Performance and robustness of penalized and unpenalized methods for genetic prediction of complex human disease. *Genetic Epidemiology*, **37**(2), 184–195.
- Abraham, G., Tye-Din, J. A., Bhalala, O. G., Kowalczyk, A., Zobel, J., and Inouye, M. (2014). Accurate and robust genomic prediction of celiac disease using statistical learning. *PLoS genetics*, **10**(2), e1004137.
- Botta, V., Louppe, G., Geurts, P., and Wehenkel, L. (2014). Exploiting snp correlations within random forest for genome-wide association studies. *PLoS one*, **9**(4), e93379.
- Breiman, L. (2001). Random forests. *Machine learning*, **45**(1), 5–32.
- Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L. T., Sharp, K., Motyer, A., Vukcevic, D., Delaneau, O., O’Connell, J., *et al.* (2017). Genome-wide genetic data on ~500,000 uk biobank participants. *bioRxiv*, page 166298.
- Chatterjee, N., Wheeler, B., Sampson, J., Hartge, P., Chanock, S. J., and Park, J.-H. (2013). Projecting the performance of risk prediction based on polygenic analyses of genome-wide association studies. *Nature genetics*, **45**(4), 400–405.
- Chatterjee, N., Shi, J., and García-Closas, M. (2016). Developing and evaluating polygenic risk prediction models for stratified disease prevention. *Nature Reviews Genetics*, **17**(7), 392.
- Dey, S., Gupta, R., Steinbach, M., and Kumar, V. (2013). Integration of clinical and genomic data: a methodological survey. *Briefings in Bioinformatics*.
- Dodd, L. E. and Pepe, M. S. (2003). Partial auc estimation and regression. *Biometrics*, **59**(3), 614–623.
- Dubois, P. C., Trynka, G., Franke, L., Hunt, K. A., Romanos, J., Curtotti, A., Zhernakova, A., Heap, G. A., Ádány, R., Aromaa, A., *et al.* (2010). Multiple common variants for celiac disease influencing immune gene expression. *Nature genetics*, **42**(4), 295–302.
- Dudbridge, F. (2013). Power and predictive accuracy of polygenic risk scores. *PLoS genetics*, **9**(3), e1003348.
- Evans, D. M., Visscher, P. M., and Wray, N. R. (2009). Harnessing the information contained within genome-wide association studies to improve individual prediction of complex disease risk. *Human molecular genetics*, **18**(18), 3525–3531.

- 350 Falconer, D. S. (1965). The inheritance of liability to certain diseases, estimated from the incidence among relatives. *Annals of human genetics*,
351 **29**(1), 51–76.
- 352 Fawcett, T. (2006). An introduction to roc analysis. *Pattern recognition letters*, **27**(8), 861–874.
- 353 Friedman, J., Hastie, T., and Tibshirani, R. (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of statistical*
354 *software*, **33**(1), 1.
- 355 Golan, D. and Rosset, S. (2014). Effective genetic-risk prediction using mixed models. *The American Journal of Human Genetics*, **95**(4), 383–393.
- 356 Hastie, T., Tibshirani, R., and Friedman, J. (2008). Model assessment and selection. In *The Elements of Statistical Learning*, pages 219–259.
357 Springer New York.
- 358 Hoerl, A. E. and Kennard, R. W. (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, **12**(1), 55–67.
- 359 Janssens, A. C. J., Moonesinghe, R., Yang, Q., Steyerberg, E. W., van Duijn, C. M., and Khoury, M. J. (2007). The impact of genotype frequencies
360 on the clinical validity of genomic profiling for predicting common chronic diseases. *Genetics in Medicine*, **9**(8), 528–535.
- 361 Lello, L., Avery, S. G., Tellier, L., Vazquez, A., Campos, G. d. I., and Hsu, S. D. (2017). Accurate genomic prediction of human height. *arXiv*
362 *preprint arXiv:1709.06489*.
- 363 Lusted, L. B. (1971). Signal detectability and medical decision-making. *Science*, **171**(3977), 1217–1219.
- 364 Maier, R., Moser, G., Chen, G.-B., Ripke, S., Absher, D., Agartz, I., Akil, H., Amin, F., Andreassen, O. A., Anjorin, A., *et al.* (2015). Joint analysis
365 of psychiatric disorders increases accuracy of risk prediction for schizophrenia, bipolar disorder, and major depressive disorder. *The American*
366 *Journal of Human Genetics*, **96**(2), 283–294.
- 367 Márquez-Luna, C., Loh, P.-R., and Price, A. L. (2017). Multiethnic polygenic risk scores improve risk prediction in diverse populations. *Genetic*
368 *epidemiology*, **41**(8), 811–823.
- 369 Martin, A. R., Gignoux, C. R., Walters, R. K., Wojcik, G. L., Neale, B. M., Gravel, S., Daly, M. J., Bustamante, C. D., and Kenny, E. E. (2017).
370 Human demographic history impacts genetic risk prediction across diverse populations. *The American Journal of Human Genetics*, **100**(4),
371 635–649.
- 372 McClish, D. K. (1989). Analyzing a portion of the roc curve. *Medical Decision Making*, **9**(3), 190–195.
- 373 Okser, S., Pahikkala, T., Airola, A., Salakoski, T., Ripatti, S., and Aittokallio, T. (2014). Regularized machine learning in the genetic prediction of
374 complex traits. *PLoS genetics*, **10**(11), e1004754.
- 375 Pashayan, N., Duffy, S. W., Neal, D. E., Hamdy, F. C., Donovan, J. L., Martin, R. M., Harrington, P., Benlloch, S., Al Olama, A. A., Shah, M., *et al.*
376 (2015). Implications of polygenic risk-stratified screening for prostate cancer on overdiagnosis. *Genetics in Medicine*, **17**(10), 789–795.
- 377 Privé, F., Aschard, H., Ziyatdinov, A., and Blum, M. G. B. (2018). Efficient analysis of large-scale genome-wide data with two R packages: bigstatsr
378 and bigsnpr. *Bioinformatics*, **34**(16), 2781–2787.
- 379 Purcell, S. M., Wray, N. R., Stone, J. L., Visscher, P. M., O'donovan, M. C., Sullivan, P. F., Sklar, P., Ruderfer, D. M., McQuillan, A., Morris, D. W.,
380 *et al.* (2009). Common polygenic variation contributes to risk of schizophrenia and bipolar disorder. *Nature*, **460**(7256), 748–752.
- 381 Sachs, M. C. *et al.* (2017). plotroc: A tool for plotting roc curves. *Journal of Statistical Software*, **79**(c02).

- 382 Speed, D. and Balding, D. J. (2014). Multiblup: improved snp-based prediction for complex traits. *Genome research*, **24**(9), 1550–1557.
- 383 Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages
384 267–288.
- 385 Tibshirani, R., Bien, J., Friedman, J., Hastie, T., Simon, N., Taylor, J., and Tibshirani, R. J. (2012). Strong rules for discarding predictors in lasso-type
386 problems. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **74**(2), 245–266.
- 387 Van Vliet, M. H., Horlings, H. M., Van De Vijver, M. J., Reinders, M. J., and Wessels, L. F. (2012). Integration of clinical and gene expression data
388 has a synergetic effect on predicting breast cancer outcome. *PLoS one*, **7**(7), e40358.
- 389 Vilhjálmsson, B. J., Yang, J., Finucane, H. K., Gusev, A., Lindström, S., Ripke, S., Genovese, G., Loh, P.-R., Bhatia, G., Do, R., *et al.* (2015).
390 Modeling linkage disequilibrium increases accuracy of polygenic risk scores. *The American Journal of Human Genetics*, **97**(4), 576–592.
- 391 Ware, E. B., Schmitz, L. L., Faul, J. D., Gard, A., Mitchell, C., Smith, J. A., Zhao, W., Weir, D., and Kardia, S. L. (2017). Heterogeneity in polygenic
392 scores for common human traits. *bioRxiv*, page 106062.
- 393 Wei, Z., Wang, K., Qu, H.-Q., Zhang, H., Bradfield, J., Kim, C., Frackleton, E., Hou, C., Glessner, J. T., Chiavacci, R., *et al.* (2009). From disease
394 association to risk assessment: an optimistic view from genome-wide association studies on type 1 diabetes. *PLoS genetics*, **5**(10), e1000678.
- 395 Wei, Z., Wang, W., Bradfield, J., Li, J., Cardinale, C., Frackleton, E., Kim, C., Mentch, F., Van Steen, K., Visscher, P. M., *et al.* (2013). Large sample
396 size, wide variant spectrum, and advanced machine-learning technique boost risk prediction for inflammatory bowel disease. *The American
397 Journal of Human Genetics*, **92**(6), 1008–1012.
- 398 Wray, N. R., Goddard, M. E., and Visscher, P. M. (2007). Prediction of individual genetic risk to disease from genome-wide association studies.
399 *Genome research*, **17**(10), 1520–1528.
- 400 Wray, N. R., Yang, J., Goddard, M. E., and Visscher, P. M. (2010). The genetic interpretation of area under the roc curve in genomic profiling. *PLoS
401 genetics*, **6**(2), e1000864.
- 402 Yang, J., Benyamin, B., McEvoy, B. P., Gordon, S., Henders, A. K., Nyholt, D. R., Madden, P. A., Heath, A. C., Martin, N. G., Montgomery, G. W.,
403 *et al.* (2010). Common snps explain a large proportion of the heritability for human height. *Nature genetics*, **42**(7), 565–569.
- 404 Zeng, Y., Breheny, P., and Yang, T. (2017). Efficient feature screening for lasso-type problems via hybrid safe-strong rules. *arXiv preprint
405 arXiv:1704.08742*.
- 406 Zhou, X., Carbonetto, P., and Stephens, M. (2013). Polygenic modeling with bayesian sparse linear mixed models. *PLoS genetics*, **9**(2), e1003264.
- 407 Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical
408 Methodology)*, **67**(2), 301–320.

409

Supplementary Materials

410

“Fancy” model

For the “fancy” model, we separate the causal SNPs in three equal sets $S_{\text{causal}}^{(1)}$, $S_{\text{causal}}^{(2)}$ and $S_{\text{causal}}^{(3)}$; $S_{\text{causal}}^{(3)}$ is further separated in two equal sets, $S_{\text{causal}}^{(3.1)}$ and $S_{\text{causal}}^{(3.2)}$. We then compute

$$y_i = \underbrace{\sum_{j \in S_{\text{causal}}^{(1)}} w_j \cdot \widetilde{G}_{i,j}}_{\text{linear}} + \underbrace{\sum_{j \in S_{\text{causal}}^{(2)}} w_j \cdot \widetilde{D}_{i,j}}_{\text{dominant}} + \underbrace{\sum_{\substack{k=1 \\ j_1=e_k^{(3.1)} \\ j_2=e_k^{(3.2)}}}^{|S_{\text{causal}}^{(3.1)}|} w_{j_1} \cdot \widetilde{G}_{i,j_1} \widetilde{G}_{i,j_2}}_{\text{interaction}} + \epsilon_i,$$

411

where $D_{i,j} = \mathbb{1} \{G_{i,j} \neq 0\}$ and $S_{\text{causal}}^{(q)} = \left\{ e_k^{(q)}, k \in \left\{ 1, \dots, |S_{\text{causal}}^{(q)}| \right\} \right\}$.

412

Maximum AUCs

413

We used three different ways to estimate the maximum achievable AUC for our simulations.

414

First, we used the estimation from equation (3) of Wray *et al.* (2010). For a prevalence fixed at

415

30% and an heritability of 50% (respectively 80%), the approximated theoretical values of AUC

416

are 84.1% (respectively 93.0%). Note that this approximation is reported to be less accurate

417

for high heritabilities. Secondly, if we assume that the genetic part of the liabilities follows a

418

Gaussian distribution $N(0, h^2)$ and that the environmental part follows a Gaussian distribution

419

$N(0, 1 - h^2)$, we can estimate the theoretical value of the AUC that can be achieved given

420

the heritability h^2 through Monte Carlo simulations. We report AUCs of 84.1% and 94.1% for

421

an heritability of 50% and 80%, respectively. Thirdly, we reproduce the exact same procedure

422

of simulations and, for each combination of parameters (Table 2), we estimate the AUC of

423

the “oracle”, i.e. the true simulated genetic part of the liabilities through 100 replicates. For

424

every combination of parameters, AUC of oracles are comprised between 83.2% and 84.2%

425

for an heritability of 50% and between 93.2% and 94.1% for an heritability of 80%. Given

426

all these estimates of the maximal achievable AUC and for the sake of simplicity, we report

427

maximum AUCs of 84% (94%) for heritabilities of 50% (80%) whatever are the parameters of

428

the simulations.

Population	UK	Finland	Netherlands	Italy	Total
Cases	2569	637	795	495	4496
Controls	7492	1799	828	540	10659
Total	10061	2436	1623	1035	15155

Table S1: Number of individuals by population and disease status in the celiac disease case-control study (after quality control, genotyped on 281,122 SNPs).

1.00e+00	7.22e-01	5.87e-01	4.20e-01	2.43e-01	1.00e-01	2.35e-02	2.21e-03	4.69e-05	8.81e-08	3.18e-12	1.83e-19	2.89e-31	1.70e-50	7.71e-82
5.00e-08	7.05e-01	5.65e-01	3.95e-01	2.20e-01	8.47e-02	1.79e-02	1.42e-03	2.28e-05	2.73e-08	4.69e-13	8.08e-21	1.80e-33	4.30e-54	1.06e-87
7.94e-01	6.87e-01	5.42e-01	3.69e-01	1.97e-01	7.08e-02	1.34e-02	8.83e-04	1.05e-05	7.74e-09	6.03e-14	2.86e-22	7.73e-36	5.97e-58	5.49e-94
7.81e-01	6.69e-01	5.19e-01	3.43e-01	1.75e-01	5.85e-02	9.79e-03	5.31e-04	4.61e-06	2.01e-09	6.69e-15	7.92e-24	2.24e-38	4.37e-62	1.00e-100
7.67e-01	6.50e-01	4.95e-01	3.18e-01	1.54e-01	4.76e-02	7.01e-03	3.08e-04	1.90e-06	4.72e-10	6.32e-16	1.70e-25	4.26e-41	1.61e-66	
7.53e-01	6.30e-01	4.70e-01	2.93e-01	1.35e-01	3.82e-02	4.90e-03	1.72e-04	7.31e-07	1.00e-10	5.04e-17	2.75e-27	5.16e-44	2.83e-71	
7.38e-01	6.09e-01	4.46e-01	2.68e-01	1.17e-01	3.02e-02	3.33e-03	9.18e-05	2.63e-07	1.89e-11	3.35e-18	3.31e-29	3.84e-47	2.26e-76	

Table S2: The 102 thresholds used for the C+T method for this study.

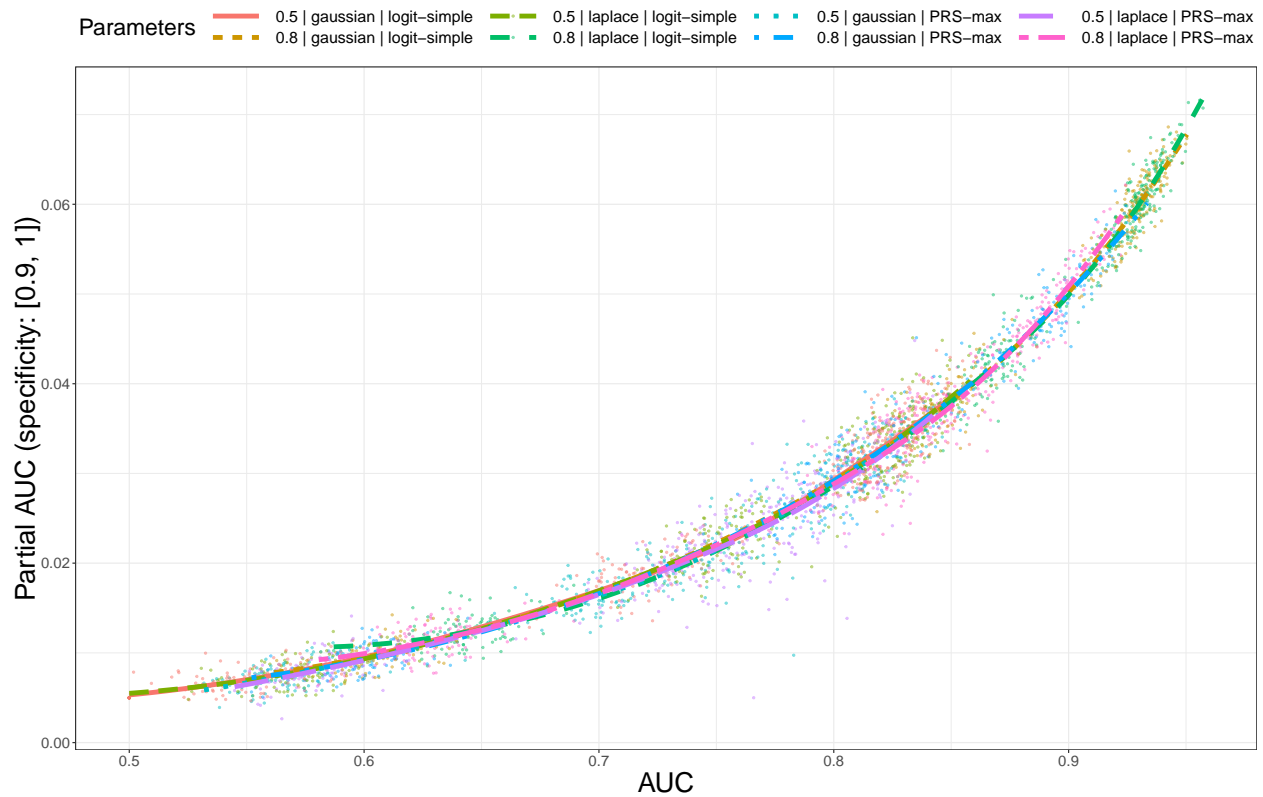


Figure S2: Correlation between AUC and partial AUC values in scenario N^o1. There is a Spearman correlation of 98% between values of AUC and partial AUC. The relation between the two values are the same whatever are the heritability, distribution of effects and method used.

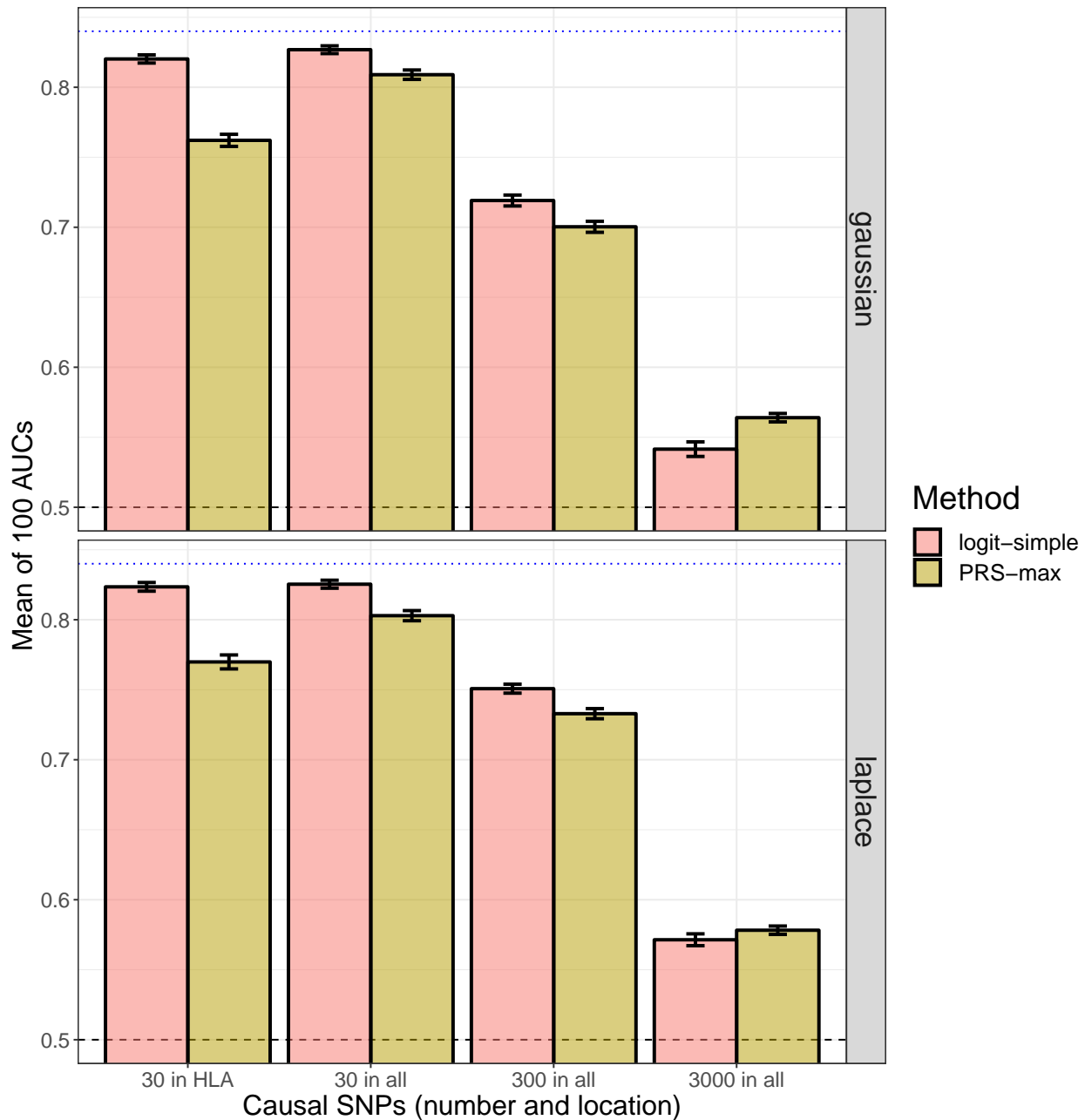


Figure S3: Comparison of the C+T and “logit-simple” methods in scenario №1 for the “simple” model and an heritability of 50%. Mean of AUC over 100 simulations for “logit-simple” and the maximum AUC reported with the C+T method (“PRS-max”). Upper (lower) panel is presenting results for effects following a Gaussian (Laplace) distribution. Error bars are representing $\pm 2SD$ of 10^5 non-parametric bootstrap of the mean of AUC. The blue dotted line represents the maximum achievable AUC.

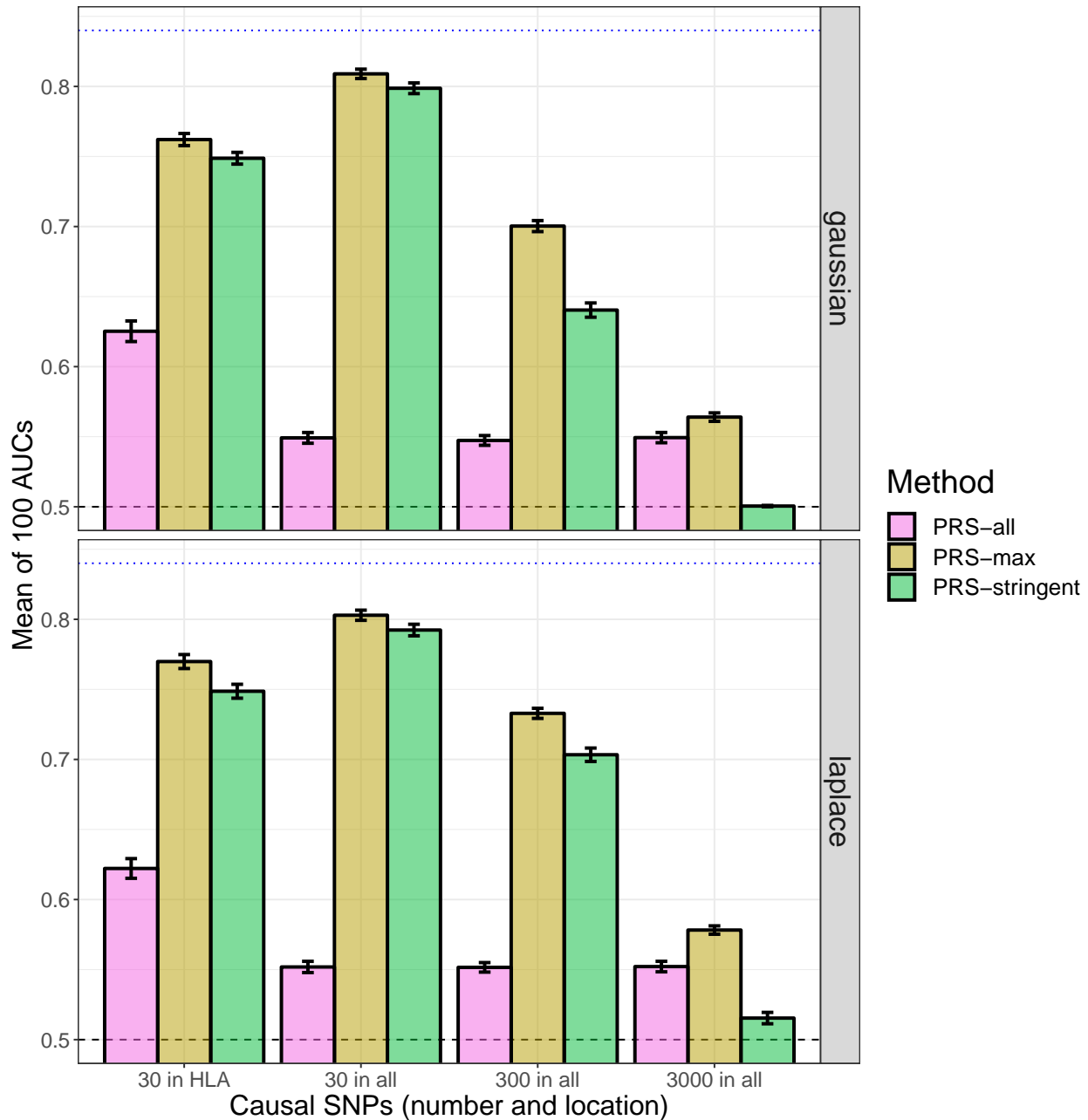


Figure S4: Comparison of three different p-value thresholds used in the C+T method in scenario №1 for the “simple” model and an heritability of 50%. Mean of AUC over 100 simulations. Upper (lower) panel is presenting results for effets following a Gaussian (Laplace) distribution. Error bars are representing $\pm 2SD$ of 10^5 non-parametric bootstrap of the mean of AUC. The blue dotted line represents the maximum achievable AUC.

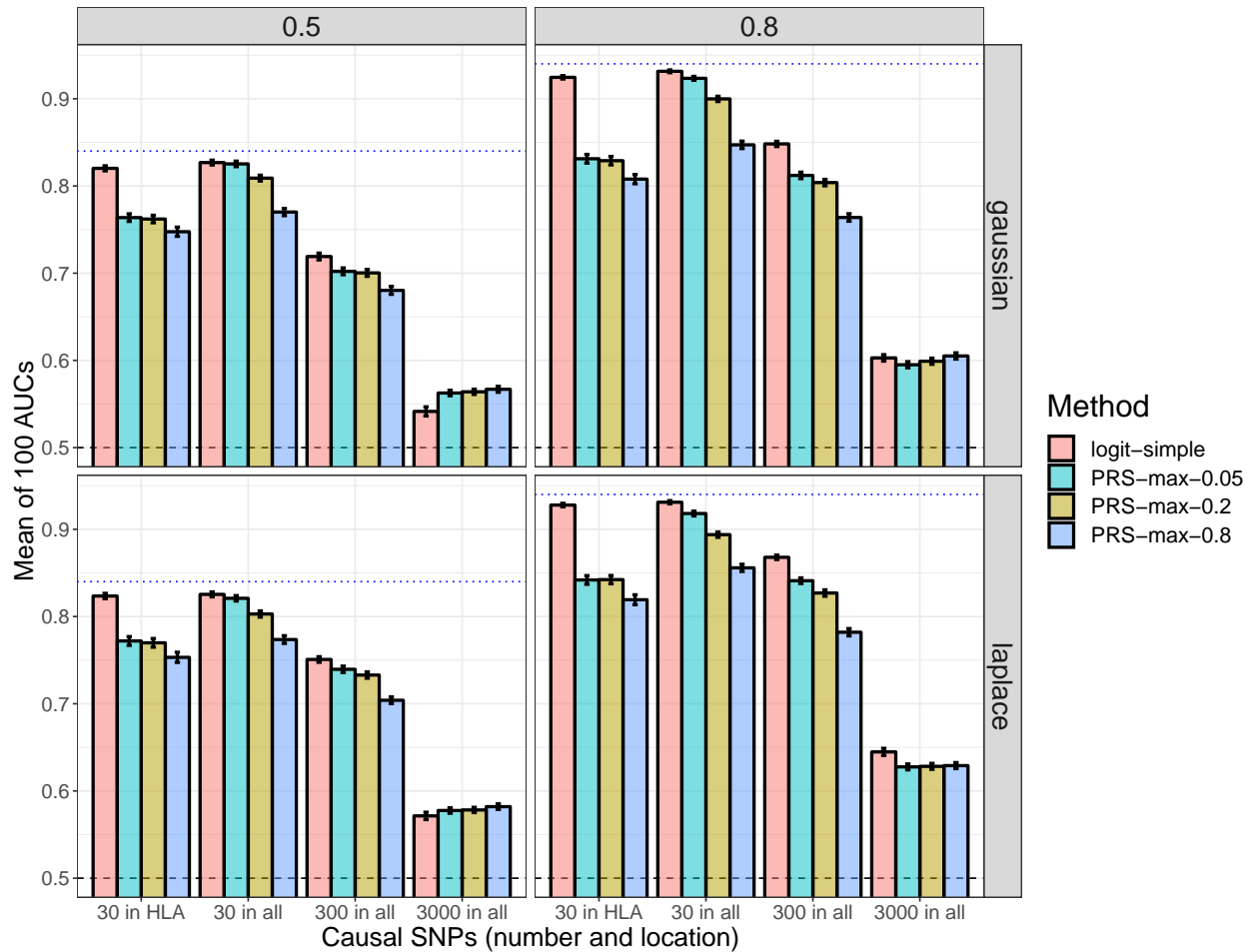


Figure S5: Comparison of models in scenario №1 for the “simple” model. Mean of AUC over 100 simulations for the maximum values of PRS for three different r^2 thresholds (0.05, 0.2 and 0.8) and “logit-simple” as a function of the number and location of causal SNPs. Upper (lower) panels are presenting results for effects following a Gaussian (Laplace) distribution and left (right) panels are presenting results for an heritability of 0.5 (0.8). Error bars are representing $\pm 2SD$ of 10^5 non-parametric bootstrap of the mean of AUC. The blue dotted line represents the maximum achievable AUC.

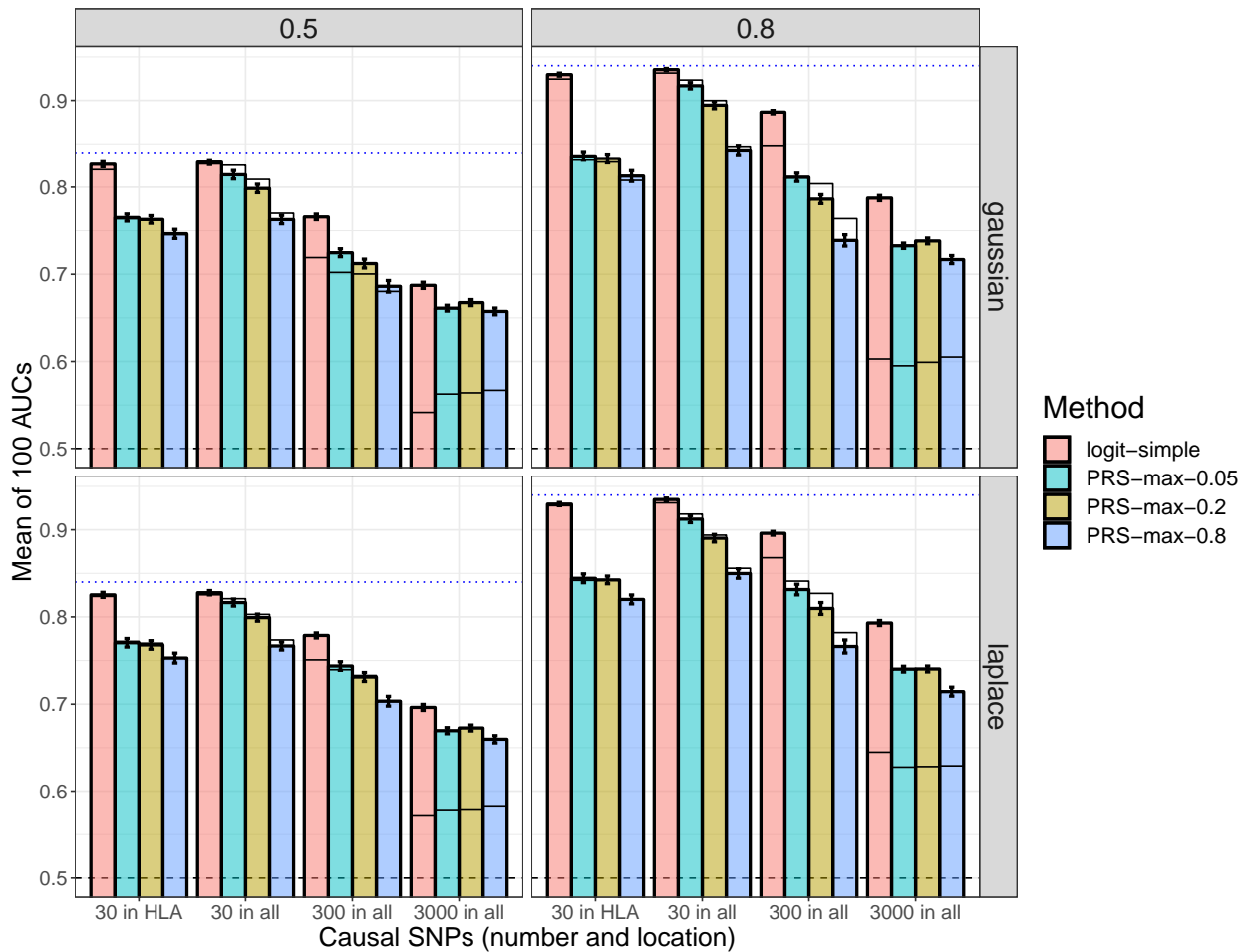


Figure S6: Comparison of models in scenario N^o2 (using chromosome 6 only) for the “simple” model. Thinner lines represents results in scenario N^o1 (Figure S5). Mean of AUC over 100 simulations for the maximum values of PRS for three different r^2 thresholds (0.05, 0.2 and 0.8) and “logit-simple” as a function of the number and location of causal SNPs. Upper (lower) panels are presenting results for effects following a Gaussian (Laplace) distribution and left (right) panels are presenting results for an heritability of 0.5 (0.8). Error bars are representing $\pm 2SD$ of 10^5 non-parametric bootstrap of the mean of AUC. The blue dotted line represents the maximum achievable AUC.

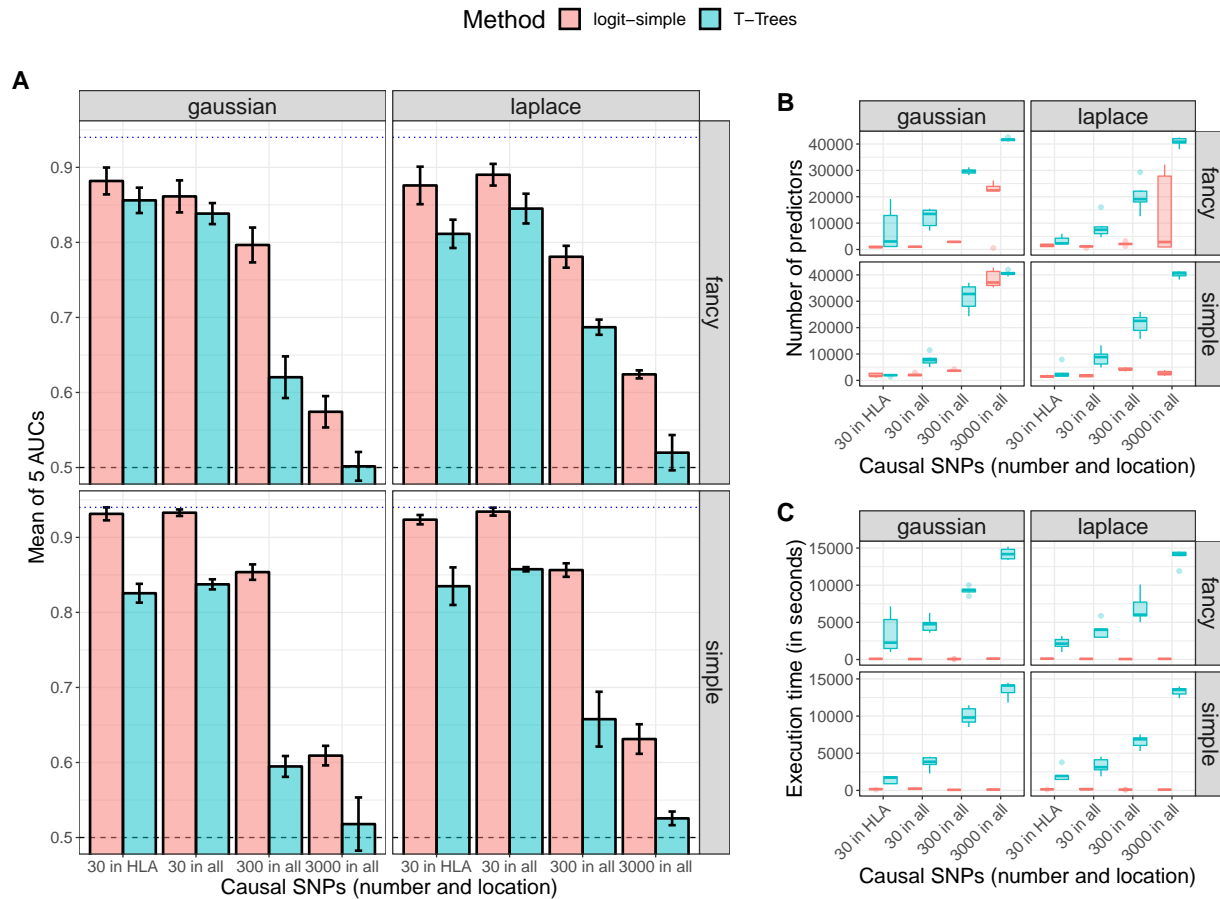


Figure S7: Comparison of T-Trees and “logit-simple” in scenario N°1 for an heritability of 80%. Vertical panels are presenting results for effects following a Gaussian or Laplace distribution. Horizontal panels are presenting results for the “simple” and “fancy” models for simulating phenotypes. **A:** Mean of AUC over 5 simulations. Error bars are representing $\pm 2SD$ of 10^5 non-parametric bootstrap of the mean of AUC. The blue dotted line represents the maximum achievable AUC. **B:** Boxplots of numbers of predictors used by the methods for 5 simulations. **C:** Boxplots of execution times for 5 simulations.

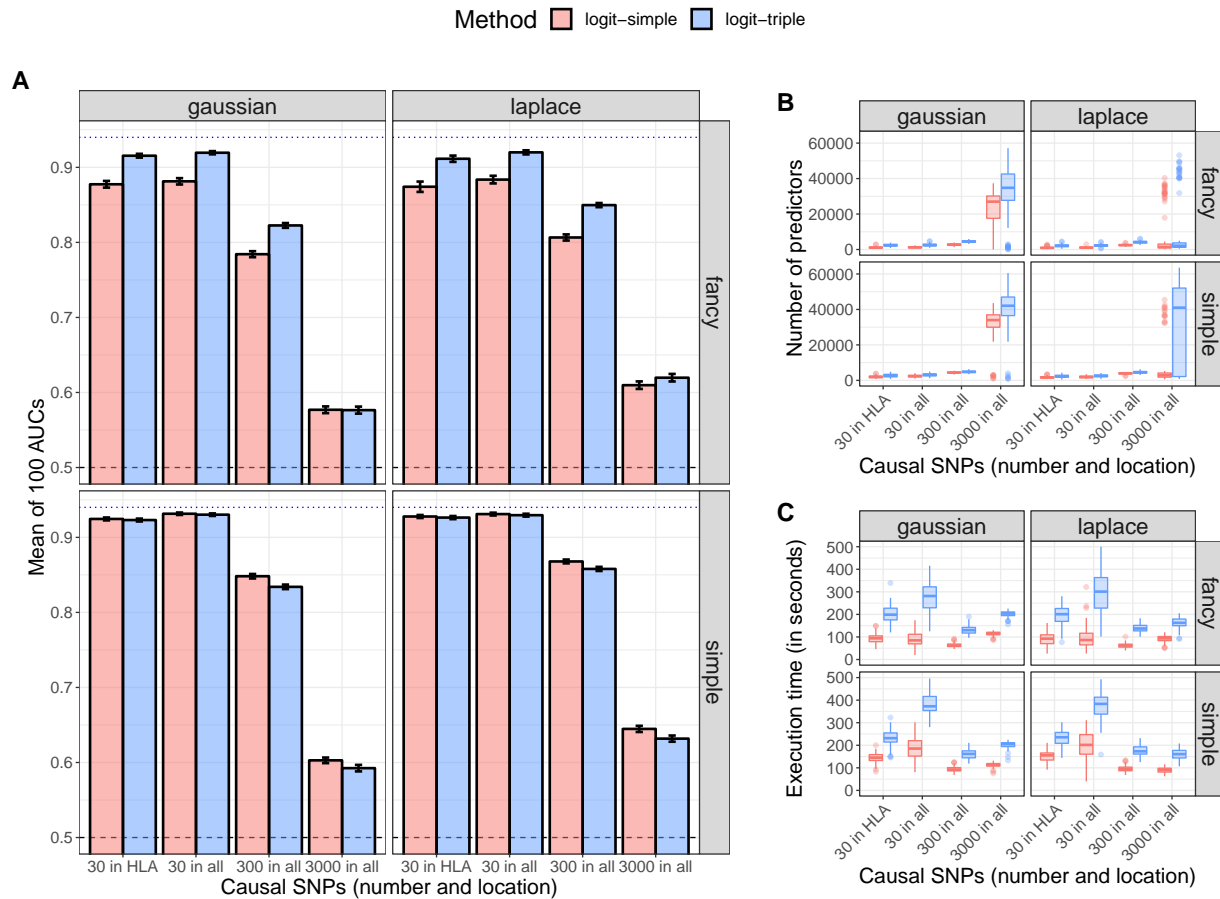


Figure S8: Comparison of “logit-triple” and “logit-simple” in scenario №1 for an heritability of 80%. Vertical panels are presenting results for effects following a Gaussian or Laplace distribution. Horizontal panels are presenting results for the “simple” and “fancy” models for simulating phenotypes. **A:** Mean of AUC over 100 simulations. Error bars are representing $\pm 2SD$ of 10^5 non-parametric bootstrap of the mean of AUC. The blue dotted line represents the maximum achievable AUC. **B:** Boxplots of numbers of predictors used by the methods for 100 simulations. **C:** Boxplots of execution times for 100 simulations.

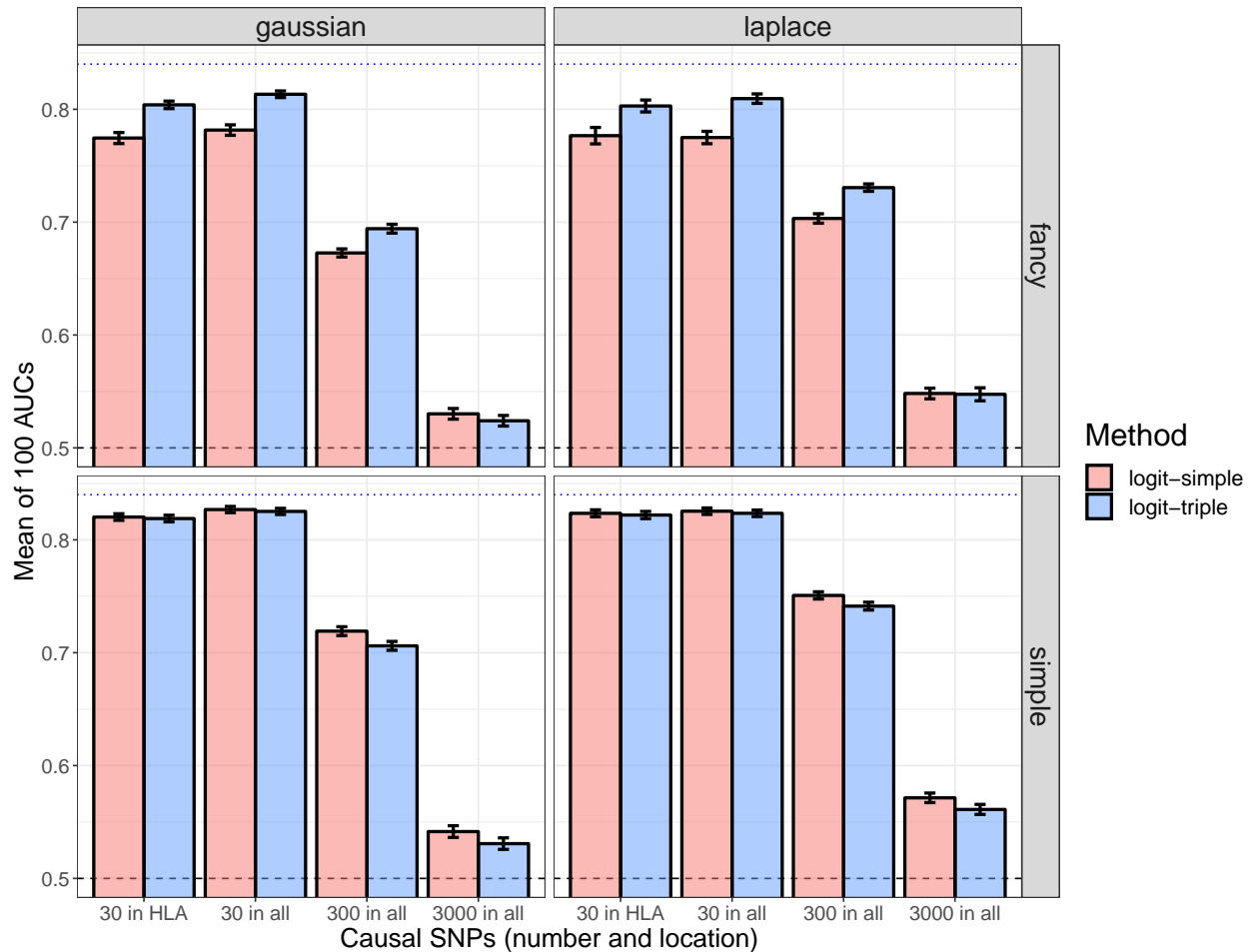


Figure S9: Comparison of “logit-triple” and “logit-simple” in scenario №1 for an heritability of 50%. Vertical panels are presenting results for effects following a Gaussian or Laplace distribution. Horizontal panels are presenting results for the “simple” and “fancy” models for simulating phenotypes. **A:** Mean of AUC over 100 simulations. Error bars are representing $\pm 2SD$ of 10^5 non-parametric bootstrap of the mean of AUC. The blue dotted line represents the maximum achievable AUC.

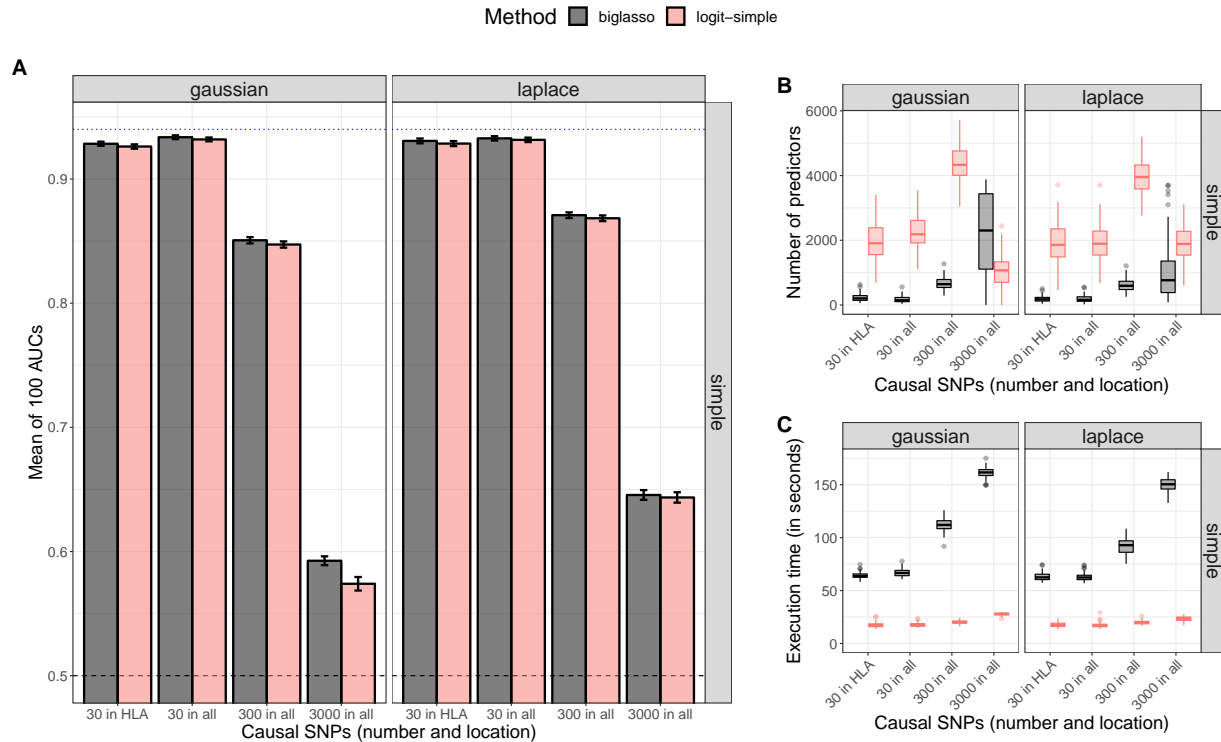


Figure S10: Comparison of “logit-triple” and the best prediction (among 100 tested λ values) for “biglasso” (another implementation of penalized logistic regression) in scenario №1. Simulations use the “simple” model, an heritability of 80% and $\alpha = 1$. Vertical panels are presenting results for effects following a Gaussian or Laplace distribution. **A**: Mean of AUC over 100 simulations. Error bars are representing $\pm 2SD$ of 10^5 non-parametric bootstrap of the mean of AUC. The blue dotted line represents the maximum achievable AUC. **B**: Boxplots of numbers of predictors used by the methods for 100 simulations. **C**: Boxplots of execution times for 100 simulations.