

# LDpred2: better, faster, stronger

Florian Privé,<sup>1,\*</sup> Julyan Arbel,<sup>2</sup> and Bjarni J. Vilhjálmsson<sup>1,3,\*</sup>

<sup>1</sup>National Centre for Register-Based Research, Aarhus University, Aarhus, 8210, Denmark.

<sup>2</sup>Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LJK, Grenoble, 38000, France.

<sup>3</sup>Bioinformatics Research Centre, Aarhus University, Aarhus, 8000, Denmark.

\*To whom correspondence should be addressed.

## Contacts:

- `florian.prive.21@gmail.com`
- `bjv@econ.au.dk`

## Abstract

Polygenic scores have become a central tool in human genetics research. LDpred is a popular and powerful method for deriving polygenic scores based on summary statistics and a matrix of correlation between genetic variants. However, LDpred has limitations that may result in limited predictive performance. Here we present LDpred2, a new version of LDpred that addresses these issues. We also provide two new options in LDpred2: a “sparse” option that can learn effects that are exactly 0, and an “auto” option that directly learns parameters from data. We benchmark LDpred2 against the previous version on simulated and real data, demonstrating substantial improvements in robustness and efficiency, as well as providing much more accurate polygenic scores. LDpred2 is implemented in R package `bigsnpr`.

# 1 Introduction

In recent years the use of polygenic scores (PGS) has become widespread. A PGS aggregates (risk) effects across many genetic variants into a single predictive score. These scores have proven useful for studying the genetic architecture and relationships between diseases and traits (Purcell *et al.* 2009; Kong *et al.* 2018). Moreover, there are high hopes for using these scores in clinical practice to improve disease risk estimates and prognostic accuracies. The heritability, i.e. the proportion of phenotypic variance that is attributable to genetics, determines an upper limit on the predictive performance of PGS and thus their value as a diagnostic tool. Nevertheless, a number of studies have explored the use of PGS in clinical settings (Pashayan *et al.* 2015; Willoughby *et al.* 2019; Abraham *et al.* 2019). PGS are also extensively used in epidemiology and economics as predictive variables of interest (Horsdal *et al.* 2019). For example, a recently derived PGS for education attainment has been one of the most predictive variables in behavioural sciences so far (Allegrini *et al.* 2019).

LDpred is a popular and powerful method for deriving polygenic scores based on summary statistics and a Linkage Disequilibrium (LD) matrix only (Vilhjálmsdóttir *et al.* 2015). It assumes there is a proportion  $p$  of variants that are causal. However, LDpred has several limitations that may result in limited predictive performance. The non-infinitesimal version of LDpred, a Gibbs sampler, is particularly sensitive to model misspecification when applied to summary statistics with large sample sizes. It is also unstable in long-range LD regions such as the human leukocyte antigen (HLA) region of chromosome 6. This issue has led to the removal of such regions from analyses (Marquez-Luna *et al.* 2018; Lloyd-Jones *et al.* 2019), which is unfortunate since this region of the genome contains many known disease-associated variants, particularly with autoimmune diseases and psychiatric disorders (Mokhtari and Lachman 2016; Matzaraki *et al.* 2017).

Here, we present LDpred2, a new version of LDpred that addresses these issues while markedly improving its computational efficiency. We provide this faster and more robust implementation of LDpred in R package bigsnpr (Privé *et al.* 2018). We also provide two new options in LDpred2. First, we provide a “sparse” option, where LDpred2 truly fits some effects to zero, therefore providing a sparse vector of effects. Second, we also provide an “auto” option, where LDpred2 automatically estimates the sparsity  $p$  and the SNP heritability  $h^2$ , and therefore does not require validation data to tune hyper-parameters. We show that LDpred2 provides higher predictive performance than LDpred1 (LDpred v1.0.0), especially when there are causal variants in long-range LD regions, or when the proportion of causal variants is small. We also show that the new sparse option performs equally well as the non-sparse version, enabling LDpred2 to provide sparse effects without losing prediction accuracy.

## 2 Results

### Overview of methods

Here we present LDpred2, a new version of LDpred (Vilhjálmsdóttir *et al.* 2015). LDpred2 has 4 options: 1) LDpred2-inf, which provides an analytical solution under the infinitesimal model of LDpred1; 2) LDpred2-grid (or simply LDpred2) that is the main LDpred model, where a grid of values for hyper-parameters  $p$  (the

proportion of causal variants) and  $h^2$  (the SNP heritability) are tuned using a validation set; 3) LDpred2-sparse, which is similar to LDpred2-grid but where effects can be exactly 0, offering a version of LDpred that can provide sparse effects; 4) LDpred2-auto, which automatically estimate  $p$  and  $h^2$  and therefore is free of hyper-parameters to tune. As a recall, LDpred v1 has two options: LDpred1-grid where only  $p$  is optimized, and LDpred1-inf.

We compare the two versions of LDpred using six simulation scenarios to understand the expected impact of using the new version of LDpred. We also compare these two versions of LDpred using eight case-control phenotypes of major interest and for which there are published external summary statistics available and a substantial number of cases in the UK Biobank data. Area Under the ROC Curve (AUC) values are reported.

## Simulations

Figure 1 presents the simulation results comparing LDpred1 (v1.0.0 as implemented by Vilhjálmsson *et al.* (2015)) with the new LDpred2 (as implemented in R package bigsnpr). Six simulation scenarios are used, each repeated 10 times. In the first four simulation scenarios, a heritability  $h^2$  of 40% and a prevalence of 15% are used. We simulate 300, 3000, 30,000 or 300,000 causal variants anywhere on the genome. In these scenarios, infinitesimal models perform similarly. When testing a grid of hyper-parameters, LDpred2 performs substantially better than LDpred1 in the cases where the number of causal variants is small, i.e. in the case of 300 or 3000 causal variants ( $p=2.5e-4$  and  $2.5e-3$ ). For example, in simulations with 300 causal variants, a mean AUC of 73.5% is obtained with LDpred1 while a value of 81.9% is obtained with LDpred2. In these scenarios, all 3 non-infinitesimal models of LDpred2 perform equally well. As expected, LDpred1 performs poorly in HLA scenarios, i.e. when causal variants are sampled in a region with variants in high LD. In contrast, all LDpred2 models perform well in these two HLA scenarios, although LDpred2-auto performs slightly worse than other LDpred2 models.

## Real data

Figure 2 presents the results of real data applications comparing LDpred1 (v1.0.0 as implemented by Vilhjálmsson *et al.* (2015)) with the new LDpred2 (as implemented in R package bigsnpr). Eight case-control phenotypes are used, summarized in table 1. For T1D, T2D, BRCA and PRCA, all main LDpred2 models perform much better than LDpred1. These improvements are also significant for MDD and CAD, albeit to a lesser extent. For example, for BRCA, AUC improves from 58.5% with LDpred1 to 64.4% with LDpred2, and from 54.9% to 74.3% for T1D. For Asthma and RA, predictive performances of LDpred1 and LDpred2 are similar.

As in simulations, the sparse version of LDpred2 performs equally well as the original version. For PRCA and T1D, LDpred2-auto perform much worse than LDpred2, otherwise it performs similarly.

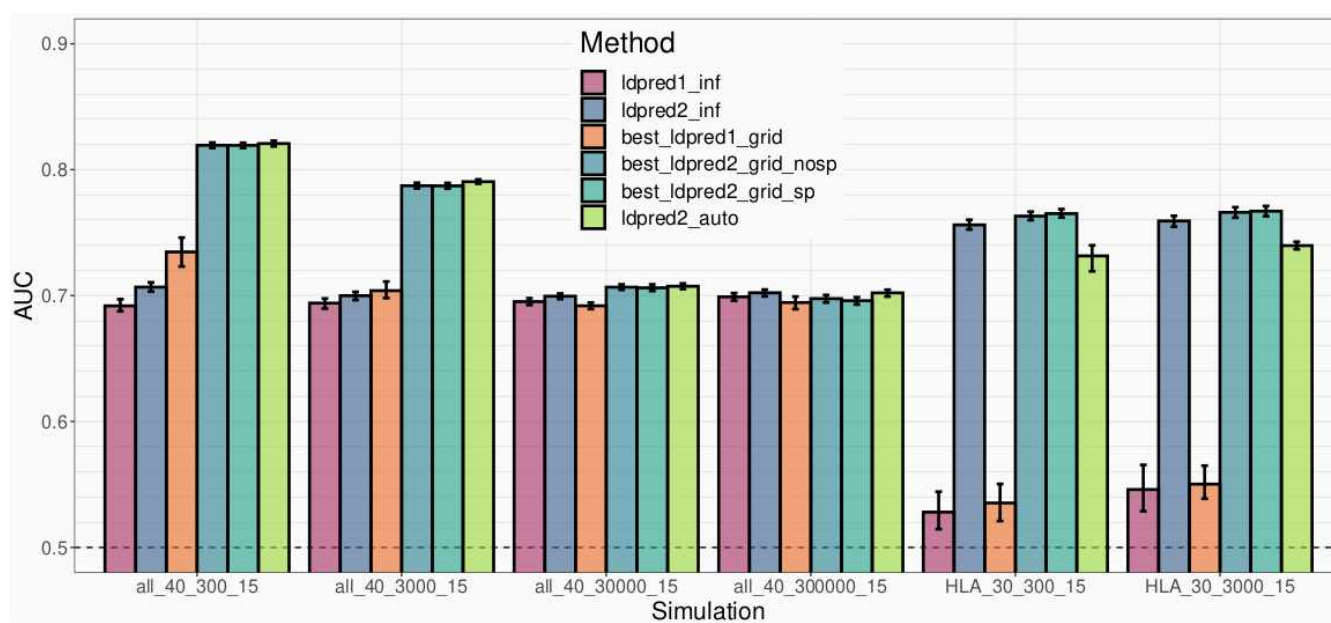


Figure 1: Results of the six simulation scenarios. Scenarios are named in 4 parts separated by underscores: 1) “all” means that causal variants are randomly sampled anywhere on the genome while “HLA” means that they are sampled in the HLA region of chromosome 6 (25.5-33.5 Mb); 2) the heritability (in %); 3) the number of causal variants; 4) the prevalence (in %). Bars present the mean and 95% CI of 10,000 non-parametric bootstrap replicates of the mean AUC of 10 simulations for each scenario.

### 3 Discussion

Here we present LDpred2, a new version of LDpred. LDpred is widely used and has the potential to provide polygenic models with good predictive performance (Khera *et al.* 2018). Yet, it has some instability issues that have been pointed out (Marquez-Luna *et al.* 2018; Lloyd-Jones *et al.* 2019) and likely contributed to discrepancies in reported prediction accuracies (Choi and O’Reilly 2019; Ge *et al.* 2019). We therefore implemented a new version of LDpred to solve these instability issues and improve its computational efficiency.

We show that LDpred2 is much more stable and provides higher prediction than LDpred1. We demonstrate that LDpred1 has defects, particularly when handling long-range LD regions such as the HLA region. Indeed, LDpred1 performs poorly in the simulations where causal variants are in the HLA region. In contrast, LDpred2 performs very well. We hypothesize that LDpred1 does not use a window size that is large enough to account for long-range LD such as in the HLA region. In LDpred2, we use a window size of 3 cM, which is much larger than the default value used in LDpred1 and which enables LDpred2 to work well even when causal variants are in long-range LD regions. We strongly discourage against removing these regions as sometimes suggested in the literature. Indeed, these regions, especially the HLA region, contain lots of variants associated with many traits, and are therefore very useful for prediction.

In LDpred2, we also test more values for  $p$  (17 instead of 7 by default) and for  $h^2$  (3 instead of 1). When testing the grid of hyper-parameters of  $p$  and  $h^2$ , we also allow for testing an option to enable sparse models in LDpred2 (see after). Overall, we test a grid of 102 different values in LDpred2 instead of 7 in LDpred1. We also use 5 times as many iterations in the Gibbs sampler and a larger window size for

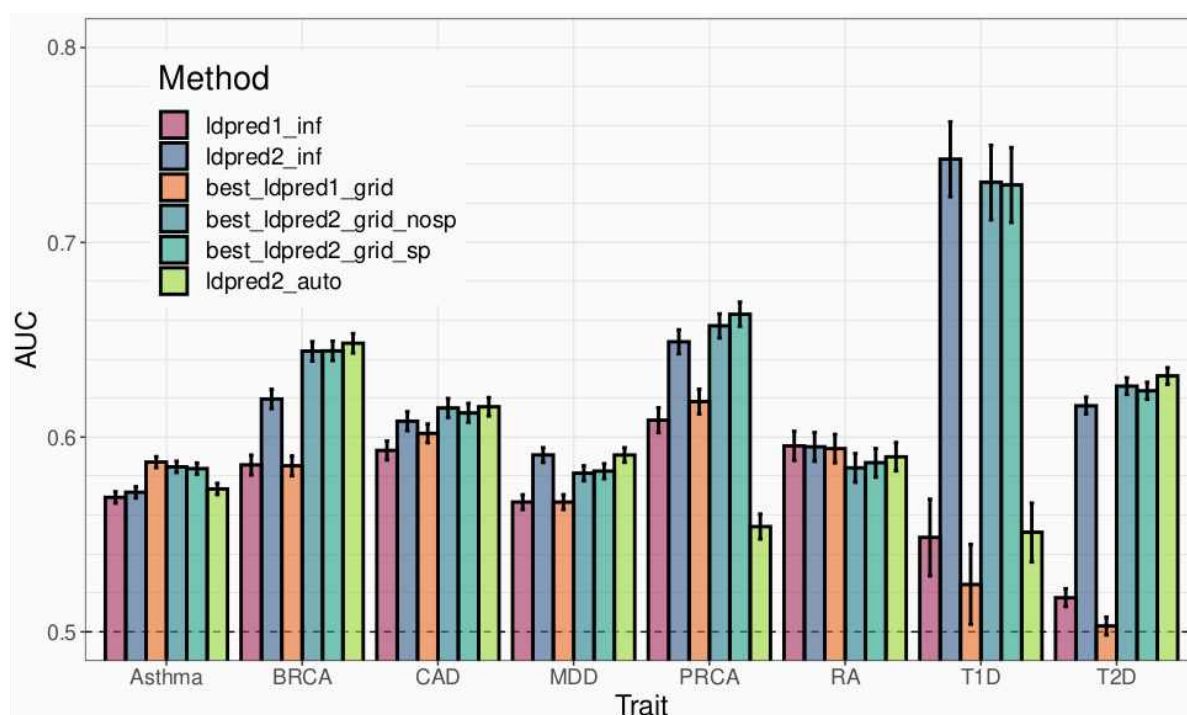


Figure 2: Results of the real data applications using published external summary statistics. Bars present AUC values on the test set of UKBB (mean and 95% CI from 10,000 bootstrap samples).

computing correlations between variants, yet LDpred2 is still as fast as LDpred1. LDpred2 is fast because we provide an efficient parallel implementation in C++. The parallelization is performed on the grid of hyper-parameters, but can also be performed chromosome-wise, as each chromosome is processed independently in LDpred2 (see Methods). As an example, it takes around two hours in total to run LDpred2-inf, LDpred2-grid (108 hyper-parameter values parallelized over 8 cores) and LDpred2-auto for a chromosome with 100,000 variants. It takes 11 minutes to pre-compute the LD matrix for 10,000 individuals and 100,000 variants (using 16 cores).

LDpred2 also comes with two new models. When the sparse option is enabled, it provides models that truly encourage sparsity, as compared to LDpred1 which outputs very small non-zero effect sizes (Janssens and Joyner 2019). It can provide sparsity in effect sizes as much as 98% for small  $p$ , while keeping predictive performance as good as non-sparse models, as opposed to if very small effects were simply discarded (Bolli *et al.* 2019). This sparse model performs equally well as the non-sparse model, therefore we encourage users to test it in the grid and to choose it if it performs better or equally well as the non-sparse model. As for LDpred2-auto, which automatically estimates values for hyper-parameters  $p$  and  $h^2$ , it performs equally well as other LDpred2 models in simulations, but does not perform well for some of the real traits analyzed here. This is expected for T1D because it is mainly composed of large effects in the HLA region (see similar simulation results) and because summary statistics have a small sample size (5913 cases and 8828 controls only). Yet, we do not know yet why LDpred2-auto performs poorly specifically for PRCA. More investigations need to be performed to understand these poor results of LDpred2-auto in these two cases, and to find if this issue can be fixed. Solutions that come to mind include for example running multiple chains with different initial values of  $p$  and keeping the one with the best convergence properties. More stringent

quality control on summary statistics might help as well.

## 4 Methods

### Simulation analyses

We use the UK Biobank data for both real data analyses and simulations (Bycroft *et al.* 2018). We restrict individuals to White-British individuals used in the PC computation of the UK Biobank (i.e. unrelated and quality-controlled). We restrict variants to HapMap3 variants. This results in 337,475 individuals and 1,194,574 variants. We use 10,000 individuals as a validation set for choosing optimal hyper-parameters and for computing correlations between variants (LD matrix  $\mathbf{R}$ ). We use 300,000 other individuals for running logistic GWAS to create summary statistics. We use the remaining 27,475 individuals as test set for evaluating models.

We simulate binary phenotypes with a heritability of  $h^2 = 0.4$  (or 0.3) using a Liability Threshold Model (LTM) with a prevalence of 15% (Falconer 1965). We vary the number of causal variants (300, 3000, 30,000, or 300,000) to match a range of genetic architectures from low to high polygenicity. Liability scores are computed from a model with additive effects only: we compute the liability score of the  $i$ -th individual as  $y_i = \sum_{j \in S_{\text{causal}}} w_j \tilde{G}_{i,j} + \epsilon_i$ , where  $S_{\text{causal}}$  is the set of causal variants,  $w_j$  are weights generated from a Gaussian distribution  $N(0, h^2/|S_{\text{causal}}|)$ ,  $G_{i,j}$  is the allele count of individual  $i$  for variant  $j$ ,  $\tilde{G}_{i,j}$  corresponds to its standardized version (zero mean and unit variance), and  $\epsilon_i$  follows a Gaussian distribution  $N(0, 1 - h^2)$ . Causal variants are chosen randomly anywhere on the genome. We make two other simulation scenarios with 300 or 3000 causal variants randomly chosen in the HLA region (chromosome 6, 25.5-33.5 Mb). Both parts of the  $y_i$ 's are scaled such that the variance of the genetic liability is exactly  $h^2$  and the variance of the total liability is exactly 1. Such simulation of phenotypes based on real genotypes is implemented in function `snpsimuPheno` of R package `bigsnpr`. Each simulation scenario is repeated 10 times and averages of the Area Under the ROC Curve (AUC) are reported.

### Real data analyses

We use the same data as in the simulation analyses. We use the same 10,000 individuals as validation set, and use the remaining 327,475 individuals as test set. We use external published GWAS summary statistics listed in table 1. We defined phenotypes as in Privé *et al.* (2019). For details, please refer to our R code (Software and code availability section).

### From marginal effects to joint effects

In this section, we explain how we can obtain joint effects from summary statistics (marginal effects) and a correlation matrix  $\mathbf{R}$ . Let us denote by  $\mathbf{S}$  the diagonal matrix with standard deviations of the  $m$  variants,  $\mathbf{C}_n = \mathbf{I}_n - \mathbf{1}\mathbf{1}^T/n$  the centering matrix,  $\mathbf{G}$  the genotype matrix of  $n$  individuals and  $m$  variants, and  $\mathbf{y}$  the phenotype vector for  $n$  individuals.



Trait	GWAS citation	GWAS sample size	GWAS #variants
Breast cancer (BRCA)	Michailidou <i>et al.</i> (2017)	137,045 / 119,078	11,792,542
Rheumatoid arthritis (RA)	Okada <i>et al.</i> (2014)	29,880 / 73,758	9,739,303
Type 1 diabetes (T1D)	Censin <i>et al.</i> (2017)	5913 / 8828	8,996,866
Type 2 diabetes (T2D)	Scott <i>et al.</i> (2017)	26,676 / 132,532	12,056,346
Prostate cancer (PRCA)	Schumacher <i>et al.</i> (2018)	79,148 / 61,106	20,370,946
Depression (MDD)	Wray <i>et al.</i> (2018)	59,851 / 113,154	13,554,550
Coronary artery disease (CAD)	Nikpay <i>et al.</i> (2015)	60,801 / 123,504	9,455,778
Asthma	Demenais <i>et al.</i> (2018)	19,954 / 107,715	2,001,280

Table 1: Summary of external GWAS summary statistics used. The GWAS sample size is the number of cases / controls in the GWAS.

When solving a joint model with all variants and an intercept  $\alpha$ , the joint effects  $\gamma_{\text{joint}}$  are obtained by solving

$$\begin{bmatrix} \hat{\alpha} \\ \hat{\gamma}_{\text{joint}} \end{bmatrix} = \left( \begin{bmatrix} \mathbf{1} & \mathbf{G} \end{bmatrix}^T \begin{bmatrix} \mathbf{1} & \mathbf{G} \end{bmatrix} \right)^{-1} \begin{bmatrix} \mathbf{1} & \mathbf{G} \end{bmatrix}^T \mathbf{y}.$$

Using the Woodbury formula, we get

$$\hat{\gamma}_{\text{joint}} = (\mathbf{G}^T \mathbf{C}_n \mathbf{G})^{-1} \mathbf{G}^T \mathbf{C}_n \mathbf{y}.$$

When fitting each variant separately in GWAS, the marginal effects (assuming no covariate) simplify to

$$\hat{\gamma}_{\text{marg}} = \frac{1}{n-1} \mathbf{S}^{-2} \mathbf{G}^T \mathbf{C}_n \mathbf{y}.$$

We further note that the correlation matrix of  $\mathbf{G}$  is  $\mathbf{R} = \frac{1}{n-1} \mathbf{S}^{-1} \mathbf{G}^T \mathbf{C}_n \mathbf{G} \mathbf{S}^{-1}$ . Then we get

$$\hat{\gamma}_{\text{joint}} = \mathbf{S}^{-1} \mathbf{R}^{-1} \mathbf{S} \hat{\gamma}_{\text{marg}}. \quad (1)$$

In practice, the correlation matrix  $\mathbf{R}$  is usually not available but is computed from another dataset. Also note that  $\gamma$  are the effects on the allele scale while we denote by  $\beta = \mathbf{S}\gamma$  the effects of the scaled genotypes.

For the marginal effect  $\hat{\gamma}_j$  of variant  $j$ , let us denote by  $\check{\mathbf{y}}$  and  $\check{\mathbf{G}}_j$  the vectors of phenotypes and genotypes for variant  $j$  residualized from  $K$  covariates, e.g. centering them. Then,

$$(\text{se}(\hat{\gamma}_j))^2 = \frac{(\check{\mathbf{y}} - \hat{\gamma}_j \check{\mathbf{G}}_j)^T (\check{\mathbf{y}} - \hat{\gamma}_j \check{\mathbf{G}}_j)}{(n-K-1) \check{\mathbf{G}}_j^T \check{\mathbf{G}}_j} \approx \frac{\check{\mathbf{y}}^T \check{\mathbf{y}}}{n \check{\mathbf{G}}_j^T \check{\mathbf{G}}_j} \approx \frac{\text{var}(\mathbf{y})}{n \text{var}(\mathbf{g}_j)}.$$

Thus, we can derive  $\text{sd}(\mathbf{g}_j) \approx \frac{\text{sd}(\mathbf{y})}{\text{se}(\hat{\gamma}_j) \sqrt{n}}$  and  $(\text{sd}(\mathbf{g}_j) \hat{\gamma}_j) \approx \frac{\hat{\gamma}_j}{\text{se}(\hat{\gamma}_j)} \frac{\text{sd}(\mathbf{y})}{\sqrt{n}}$ . Let us go back to equation

1. As  $\text{sd}(\mathbf{y})$  is the same for all variants, it is cancelled out by  $\mathbf{S}^{-1}$  and  $\mathbf{S}$ , therefore we can assume that  $\text{var}(\mathbf{y}) = 1$ . Then it justifies the use of the Z-scores  $(\hat{\gamma}_j / \text{se}(\hat{\gamma}_j))$  divided by  $\sqrt{n}$  as input for LDpred (first line of algorithm 1). Then, the effect sizes that LDpred outputs need to be scaled back by multiplying by  $(\text{se}(\hat{\gamma}_j) \sqrt{n})$  (last line of algorithm 1). Note that LDpred1 and other similar methods scale the output dividing by the standard deviation of genotypes. This is correct when  $\text{var}(\mathbf{y}) = 1$  only.

## Overview of LDpred model

LDpred assumes the following model for effect sizes,

$$\beta_j = S_{j,j}\gamma_j \sim \begin{cases} \mathcal{N}\left(0, \frac{h^2}{Mp}\right) & \text{with probability } p, \\ 0 & \text{otherwise,} \end{cases} \quad (2)$$

where  $p$  is the proportion of causal variants,  $M$  the number of variants and  $h^2$  the (SNP) heritability. Vilhjálmsson *et al.* (2015) estimates  $h^2$  using constrained LD score regression (intercept fixed to 1) and recommend testing a grid of hyper-parameter values for  $p$  (1, 0.3, 0.1, 0.03, 0.01, 0.003 and 0.001).

To estimate effect sizes  $\beta_j$ , Vilhjálmsson *et al.* (2015) use a Gibbs sampler, as described in algorithm 1. First, the residualized marginal effect for variant  $j$  is computed as

$$\tilde{\beta}_j = \hat{\beta}_j - \beta_{-j}^T \mathbf{R}_{-j,j} \quad (3)$$

where  $\mathbf{R}_{-j,j}$  is the  $j$ -th column without the  $j$ -th row of the correlation matrix,  $\hat{\beta}$  is the vector of marginal effect sizes, and  $\beta$  is the vector of current effect sizes in the Gibbs sampler. Then, the probability that variant  $j$  is causal is computed as

$$\bar{p}_j = P\left(\beta_j \sim \mathcal{N}(\cdot, \cdot) \mid \tilde{\beta}_j\right) = \frac{\frac{p}{\sqrt{\frac{h^2}{Mp} + \frac{1}{n}}} \exp\left\{-\frac{1}{2} \frac{\tilde{\beta}_j^2}{\frac{h^2}{Mp} + \frac{1}{n}}\right\}}{\frac{p}{\sqrt{\frac{h^2}{Mp} + \frac{1}{n}}} \exp\left\{-\frac{1}{2} \frac{\tilde{\beta}_j^2}{\frac{h^2}{Mp} + \frac{1}{n}}\right\} + \frac{1-p}{\sqrt{\frac{1}{n}}} \exp\left\{-\frac{1}{2} \frac{\tilde{\beta}_j^2}{\frac{1}{n}}\right\}},$$

which we rewrite as

$$\bar{p}_j = \frac{1}{1 + \frac{1-p}{p} \sqrt{1 + \frac{nh^2}{Mp}} \exp\left\{-\frac{1}{2} \frac{n\tilde{\beta}_j^2}{1 + \frac{Mp}{nh^2}}\right\}}. \quad (4)$$

Computing  $\bar{p}_j$  using the second expression is important to avoid numerical issues when  $(n\tilde{\beta}_j^2)$  is large.

Then,  $\beta_j$  is sampled according to

$$\beta_j \mid \tilde{\beta}_j \sim \begin{cases} \mathcal{N}\left(\frac{1}{1 + \frac{Mp}{nh^2}} \tilde{\beta}_j, \frac{1}{1 + \frac{Mp}{nh^2}} \frac{1}{n}\right) & \text{with probability } \bar{p}_j, \\ 0 & \text{otherwise.} \end{cases} \quad (5)$$



Therefore, the posterior mean of  $\beta_j$  is given by

$$\omega_j = \frac{\bar{p}_j \tilde{\beta}_j}{1 + \frac{Mp}{nh^2}}. \quad (6)$$

---

**Algorithm 1** LDpred, with hyper-parameters  $p$  and  $h^2$ , LD matrix  $\mathbf{R}$  and summary statistics  $\hat{\gamma}$ ,  $\text{se}(\hat{\gamma})$  and  $n$

---

```

1:  $\hat{\beta} \leftarrow \frac{\hat{\gamma}}{\text{se}(\hat{\gamma}) \cdot \sqrt{n}}$  ▷ Initialization of scaled marginal effects (see previous section)
2:  $\Omega \leftarrow \mathbf{0}$  ▷ Initialization of posterior means
3: for  $k = 1, \dots, N_{\text{burn-in}} + N_{\text{iter}}$  do ▷ Gibbs iterations
4:   for each variant  $j$  do ▷ All variants
5:     Compute  $\tilde{\beta}_j$  according to (3)
6:     Compute  $\bar{p}_j$  according to (4)
7:     Sample  $\beta_j$  according to (5)
8:     Compute  $\omega_j$  according to (6)
9:   end for
10:  Compute  $\beta^T \mathbf{R} \beta$  (estimate of  $h^2$ ) to check divergence ▷ Return 0s if divergence
11:  if  $k > N_{\text{burn-in}}$  then
12:     $\Omega \leftarrow \Omega + \omega$ 
13:  end if
14: end for
15:  $\Omega \leftarrow \Omega / N_{\text{iter}}$  ▷ Average of all  $\omega$  after burn-in
16: Return  $\Omega \cdot \text{se}(\hat{\gamma}) \cdot \sqrt{n}$  ▷ Return posterior means, scaled back (see previous section)

```

---

## New LDpred2 models

LDpred2 comes with two extensions of the LDpred model.

The first extension consists in estimating  $p$  and  $h^2$  in the model, as opposed to testing several values of  $p$  and estimating  $h^2$  using constrained LD score regression (Bulik-Sullivan *et al.* 2015). This makes LDpred2-auto a method free of hyper-parameters which can therefore be applied directly to data without the need of a validation dataset to choose best-performing hyper-parameters. To estimate  $p$ , we count the number of causal variants (i.e.  $M_c = \sum_j (\beta_j \neq 0)$  in equation 5). We can assume that  $M_c \sim \text{Binom}(M, p)$ , so if we place a prior  $p \sim \text{Beta}(1, 1) \equiv \mathcal{U}(0, 1)$ , we can sample  $p$  from the posterior  $p \sim \text{Beta}(1 + M_c, 1 + M - M_c)$ . Due to complexity reasons, we could not derive a Bayesian estimator of  $h^2$ . Instead, we estimate  $h^2 = \beta^T \mathbf{R} \beta$ , where  $\mathbf{R}$  is the correlation matrix. These parameters  $p$  and  $h^2$  are updated after the inner loop in algorithm 1, then these new values are used in the next iteration of the outer loop.

The second extension, LDpred2-sparse, aims at providing sparse effect size estimates, i.e. some resulting effects are exactly 0. When the sparse solution is sought and when  $\bar{p}_j < p$ , we set  $\beta_j$  and  $\omega_j$  to 0 (lines 6-8 of algorithm 1). Note that LDpred2-auto does not have a sparse option, but it is possible to run LDpred2-sparse

with the estimates of  $p$  and  $h^2$  from LDpred2-auto.

## New strategy for local correlation

LDpred has a window size parameter that needs to be set; for a given variant, correlations with other variants outside of this window are assumed to be 0. The recommended value for this window (in number of variants) is to use the total number of variants divided by 3000, which corresponds to a window radius of around 2 Mb (Vilhjálmsdóttir *et al.* 2015). We have come to the conclusion that this window size is not large enough. Indeed, the human leukocyte antigen (HLA) region of chromosome 6 is 8 Mb long (Price *et al.* 2008). Using a window of 8Mb would be computationally and memory inefficient. Instead, we propose to use genetic distances. Genome-wide, 1 Mb corresponds on average to 1 cM. Yet, the HLA region is only 3 cM long (vs. 8 Mb long). Therefore, genetic distances enable to capture the same LD using a globally smaller window. We provide function `snpr_asGeneticPos` in package `bigsnpr` to easily interpolate physical positions (in bp) to genetic positions (in cM). We recommend to use genetic positions and to use a size parameter of 3 cM when computing the correlation between variants for LDpred2. Note that, in the code, we use `size = 3 / 1000` since parameter `size` is internally multiplied by 1000 in functions of package `bigsnpr`. To prevent storing very small correlations, we also set to 0 all correlations with a p-value larger than `alpha = 0.9` (when testing that the correlation is different from 0).

## New strategy for running LDpred2

We recommend to run LDpred2 per chromosome. Even if it is possible to run LDpred2 genome-wide, this approach has two limitations. First, it can be memory and computationally demanding to do so. For around one million (1M) variants, storing the  $1M \times 1M$  sparse correlation matrix takes more than 32 GB of memory. Doubling to 2M variants would require 128 GB of RAM to store the matrix. Second, as noted in Privé *et al.* (2019), it may be beneficial to assume that architecture of traits may be different for different chromosomes. For example, chromosome 6 clearly encompasses a larger proportion of the heritability of autoimmune diseases compared to other chromosomes (Shi *et al.* 2016). Assuming the same model for genetic effects genome-wide could result in severe model misspecification, which would lead to suboptimal predictive performance. Moreover, since the inverse of a block-diagonal matrix is formed from the inverse of each block, it should be safe to run LDpred2 for each chromosome and then combine the results. Combining results from all chromosomes requires using some appropriate scaling (see tutorial).

## Software and code availability

The newest version of R package bigsnpr can be installed from GitHub (see <https://github.com/privefl/bigsnpr>). A tutorial on the steps to run LDpred2 using some small example data is available at <https://privefl.github.io/bigsnpr/articles/LDpred2.html>.

## Acknowledgements

This research has been conducted using the UK Biobank Resource under Application Number 41181.

F.P. and B.V. are supported by the Danish National Research Foundation (Niels Bohr Professorship to Prof. John McGrath), and also acknowledge the Lundbeck Foundation Initiative for Integrative Psychiatric Research, iPSYCH (R248-2017-2003).

## Declaration of Interests

The authors declare no competing interests.

# References

- Abraham, G., Malik, R., Yonova-Doing, E., Salim, A., Wang, T., Danesh, J., Butterworth, A. S., Howson, J. M., Inouye, M., and Dichgans, M. (2019). Genomic risk score offers predictive performance comparable to clinical risk factors for ischaemic stroke. *Nature Communications*, **10**(1), 1–10.
- Allegrini, A. G., Selzam, S., Rimfeld, K., von Stumm, S., Pingault, J.-B., and Plomin, R. (2019). Genomic prediction of cognitive traits in childhood and adolescence. *Molecular Psychiatry*, **24**(6), 819–827.
- Bolli, A., Di Domenico, P., and Bottà, G. (2019). Software as a service for the genomic prediction of complex diseases. *bioRxiv*, page 763722.
- Bulik-Sullivan, B. K., Loh, P.-R., Finucane, H. K., Ripke, S., Yang, J., Patterson, N., Daly, M. J., Price, A. L., Neale, B. M., of the Psychiatric Genomics Consortium, S. W. G., *et al.* (2015). LD score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nature Genetics*, **47**(3), 291.
- Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L. T., Sharp, K., Motyer, A., Vukcevic, D., Delaneau, O., O’Connell, J., *et al.* (2018). The UK Biobank resource with deep phenotyping and genomic data. *Nature*, **562**(7726), 203–209.
- Censin, J., Nowak, C., Cooper, N., Bergsten, P., Todd, J. A., and Fall, T. (2017). Childhood adiposity and risk of type 1 diabetes: A mendelian randomization study. *PLoS Medicine*, **14**(8), e1002362.
- Choi, S. W. and O’Reilly, P. F. (2019). Prsice-2: Polygenic risk score software for biobank-scale data. *Gigascience*, **8**(7), giz082.
- Demenaïs, F., Margaritte-Jeannin, P., Barnes, K. C., Cookson, W. O., Altmüller, J., Ang, W., Barr, R. G., Beaty, T. H., Becker, A. B., Beilby, J., *et al.* (2018). Multiancestry association study identifies new asthma risk loci that colocalize with immune-cell enhancer marks. *Nature Genetics*, **50**(1), 42.
- Falconer, D. S. (1965). The inheritance of liability to certain diseases, estimated from the incidence among relatives. *Annals of Human Genetics*, **29**(1), 51–76.
- Ge, T., Chen, C.-Y., Ni, Y., Feng, Y.-C. A., and Smoller, J. W. (2019). Polygenic prediction via Bayesian regression and continuous shrinkage priors. *Nature Communications*, **10**(1), 1776.
- Horsdal, H. T., Agerbo, E., McGrath, J. J., Vilhjálmsson, B. J., Antonsen, S., Closter, A. M., Timmermann, A., Grove, J., Mok, P. L., Webb, R. T., *et al.* (2019). Association of childhood exposure to nitrogen dioxide and polygenic risk score for schizophrenia with the risk of developing schizophrenia. *JAMA network open*, **2**(11), e1914401–e1914401.
- Janssens, A. C. J. W. and Joyner, M. J. (2019). Polygenic risk scores that predict common diseases using millions of single nucleotide polymorphisms: is more, better? *Clinical Chemistry*, **65**(5), 609–611.
- Khera, A. V., Chaffin, M., Aragam, K. G., Haas, M. E., Roselli, C., Choi, S. H., Natarajan, P., Lander, E. S., Lubitz, S. A., Ellinor, P. T., *et al.* (2018). Genome-wide polygenic scores for common diseases identify individuals with risk equivalent to monogenic mutations. *Nature Genetics*, **50**(9), 1219–1224.
- Kong, A., Thorleifsson, G., Frigge, M. L., Vilhjálmsson, B. J., Young, A. I., Thorgeirsson, T. E., Benonisdottir, S., Oddsson, A., Halldorsson, B. V., Masson, G., *et al.* (2018). The nature of nurture: Effects of parental genotypes. *Science*, **359**(6374), 424–428.

- Linnér, R. K., Biroli, P., Kong, E., Meddens, S. F. W., Wedow, R., Fontana, M. A., Lebreton, M., Tino, S. P., Abdellaoui, A., Hammerschlag, A. R., *et al.* (2019). Genome-wide association analyses of risk tolerance and risky behaviors in over 1 million individuals identify hundreds of loci and shared genetic influences. *Nature genetics*, **51**(2), 245–257.
- Lloyd-Jones, L. R., Zeng, J., Sidorenko, J., Yengo, L., Moser, G., Kemper, K. E., Wang, H., Zheng, Z., Magi, R., Esko, T., *et al.* (2019). Improved polygenic prediction by Bayesian multiple regression on summary statistics. *Nature Communications*, **10**(1), 1–11.
- Marquez-Luna, C., Gazal, S., Loh, P.-R., Furlotte, N., Auton, A., Price, A. L., 23andMe Research Team, *et al.* (2018). Modeling functional enrichment improves polygenic prediction accuracy in UK Biobank and 23andMe data sets. *bioRxiv*, page 375337.
- Matzaraki, V., Kumar, V., Wijmenga, C., and Zernakova, A. (2017). The MHC locus and genetic susceptibility to autoimmune and infectious diseases. *Genome Biology*, **18**(1), 76.
- Michailidou, K., Lindström, S., Dennis, J., Beesley, J., Hui, S., Kar, S., Lemaçon, A., Soucy, P., Glubb, D., Rostamianfar, A., *et al.* (2017). Association analysis identifies 65 new breast cancer risk loci. *Nature*, **551**(7678), 92.
- Mokhtari, R. and Lachman, H. M. (2016). The major histocompatibility complex (MHC) in schizophrenia: a review. *Journal of Clinical & Cellular Immunology*, **7**(6).
- Nikpay, M., Goel, A., Won, H.-H., Hall, L. M., Willenborg, C., Kanoni, S., Saleheen, D., Kyriakou, T., Nelson, C. P., Hopewell, J. C., *et al.* (2015). A comprehensive 1000 genomes-based genome-wide association meta-analysis of coronary artery disease. *Nature Genetics*, **47**(10), 1121.
- Okada, Y., Wu, D., Trynka, G., Raj, T., Terao, C., Ikari, K., Kochi, Y., Ohmura, K., Suzuki, A., Yoshida, S., *et al.* (2014). Genetics of rheumatoid arthritis contributes to biology and drug discovery. *Nature*, **506**(7488), 376.
- 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.* (2015). Implications of polygenic risk-stratified screening for prostate cancer on overdiagnosis. *Genetics in Medicine*, **17**(10), 789–795.
- Price, A. L., Weale, M. E., Patterson, N., Myers, S. R., Need, A. C., Shianna, K. V., Ge, D., Rotter, J. I., Torres, E., Taylor, K. D., *et al.* (2008). Long-range LD can confound genome scans in admixed populations. *The American Journal of Human Genetics*, **83**(1), 132–135.
- 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 and bigsnpr. *Bioinformatics*, **34**(16), 2781–2787.
- Privé, F., Vilhjálmsson, B. J., Aschard, H., and Blum, M. G. B. (2019). Making the most of clumping and thresholding for polygenic scores. *The American Journal of Human Genetics*, **105**(6), 1213–1221.
- Purcell, S. M., Wray, N. R., Stone, J. L., Visscher, P. M., O'donovan, M. C., Sullivan, P. F., Sklar, P., Ruderfer, D. M., McQuillin, A., Morris, D. W., *et al.* (2009). Common polygenic variation contributes to risk of schizophrenia and bipolar disorder. *Nature*, **460**(7256), 748–752.
- Schumacher, F. R., Al Olama, A. A., Berndt, S. I., Benlloch, S., Ahmed, M., Saunders, E. J., Dadaev, T., Leongamornlert, D., Anokian, E., Cieza-Borrella, C., *et al.* (2018). Association analyses of more than 140,000 men identify 63 new prostate cancer susceptibility loci. *Nature Genetics*, **50**(7), 928.

- Scott, R. A., Scott, L. J., Mägi, R., Marullo, L., Gaulton, K. J., Kaakinen, M., Pervjakova, N., Pers, T. H., Johnson, A. D., Eicher, J. D., *et al.* (2017). An expanded genome-wide association study of type 2 diabetes in Europeans. *Diabetes*, **66**(11), 2888–2902.
- Shi, H., Kichaev, G., and Pasaniuc, B. (2016). Contrasting the genetic architecture of 30 complex traits from summary association data. *The American Journal of Human Genetics*, **99**(1), 139–153.
- Vilhjálmsdóttir, 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). Modeling linkage disequilibrium increases accuracy of polygenic risk scores. *The American Journal of Human Genetics*, **97**(4), 576–592.
- Willoughby, A., Andreassen, P. R., and Toland, A. E. (2019). Genetic testing to guide risk-stratified screens for breast cancer. *Journal of Personalized Medicine*, **9**(1), 15.
- Wray, N. R., Ripke, S., Mattheisen, M., Trzaskowski, M., Byrne, E. M., Abdellaoui, A., Adams, M. J., Agerbo, E., Air, T. M., Andlauer, T. M., *et al.* (2018). Genome-wide association analyses identify 44 risk variants and refine the genetic architecture of major depression. *Nature Genetics*, **50**(5), 668.