

# Package: RapidoPGS (via r-universe)

September 4, 2024

**Title** A Fast and Light Package to Compute Polygenic Risk Scores

**Version** 2.3.0.9001

**Description** Quickly computes polygenic scores from GWAS summary statistics of either case-control or quantitative traits without parameter tuning. Reales,G., Vigorito, E., Kelemen,M., Wallace,C. (2021) <doi:10.1101/2020.07.24.220392> ``RápidoPGS: A rapid polygenic score calculator for summary GWAS data without a test dataset".

**License** GPL-3

**Depends** R (>= 4.3), data.table, RCurl, curl, magrittr

**Imports** dplyr (>= 1.1.3), GenomicRanges (>= 1.52.0), IRanges (>= 2.34.1), bigsnpr (>= 1.12.2), coloc (>= 5.2.3), bigreadr (>= 0.2.5)

**Encoding** UTF-8

**LazyData** true

**Roxygen** list(markdown = TRUE)

**RoxygenNote** 7.2.3

**Suggests** knitr, rmarkdown

**VignetteBuilder** knitr

**Repository** <https://grealesm.r-universe.dev>

**RemoteUrl** <https://github.com/grealesm/rapidopgs>

**RemoteRef** HEAD

**RemoteSha** 05c9cd35c042caf4575ace239a354ee0102f462a

## Contents

create_1000G . . . . .	2
EUR_Id.blocks19 . . . . .	3
EUR_Id.blocks38 . . . . .	4
find_file_in_ftp . . . . .	4
gwascats.download . . . . .	5

logsum . . . . .	6
michailidou19 . . . . .	6
michailidou38 . . . . .	7
rapidopgs_multi . . . . .	8
rapidopgs_single . . . . .	10
sd.prior.est . . . . .	12
sdY.est . . . . .	13
wakefield_pp . . . . .	14
wakefield_pp_quant . . . . .	14

<b>Index</b>	<b>16</b>
--------------	-----------

---

create_1000G	<i>Download 1000 Genomes Phase III panel</i>
--------------	--

---

## Description

create\_1000G downloads and gets 1000 Genomes Phase III panel (hg19) in PLINK format, and apply quality control for being used to compute PGS using rapidopgs\_multi. Given the size of the files, running this function can take long, depending on broadband speed and server status. We also recommend to ensure that there is at least 60GB free space available in disk.

## Usage

```
create_1000G(
  directory = "ref-data",
  remove.related = TRUE,
  qc.maf = 0.01,
  qc.hwe = 1e-10,
  qc.geno = 0,
  autosomes.only = TRUE
)
```

## Arguments

directory	a string indicating the directory to download the panel
remove.related	a logical stating if related individuals should be removed. Default TRUE.
qc.maf	a numeric to set the MAF threshold for variants to be removed. DEFAULT 0.01
qc.hwe	a numeric indicating the threshold for Hardy-Weinberg exact test p-value, below which variants will be removed. DEFAULT 1e-10.
qc.geno	a numeric to set maximum missing call rates for variants. DEFAULT = 0.
autosomes.only	If FALSE, it will include X and Y chromosomes, too.

## Value

bed, fam and bim files for each chromosome in the chosen directory.

**Author(s)**

Guillermo Reales

**Examples**

```
## Not run:  
create_1000G()  
  
## End(Not run)
```

---

EUR_ld.blocks19	<i>LD block architecture for European populations (hg19).</i>
-----------------	---

---

**Description**

A GRanges object containing the LD block for European ancestry, in hg19 build. This dataset was obtained from [Berisa and Pickrell \(2016\)](#), in bed format, then converted to GRanges. See manuscript for more details.

**Usage**

EUR\_ld.blocks19

**Format**

A GRanges object containing 1703 ranges

**seqnames** chromosome

**ranges** start and stop positions for the block

**strand** genomic strand, irrelevant here

**Source**

<https://bitbucket.org/nygcresearch/ldetect-data/src>

---

EUR_ld.blocks38	<i>LD block architecture for European populations (hg38).</i>
-----------------	---

---

### Description

A GRanges object containing the LD block for European ancestry, in hg38 build. This dataset was obtained from [MacDonald et al. \(2022\)](#), in bed format, then transformed to GRanges. See manuscript for more details.

### Usage

```
EUR_ld.blocks38
```

### Format

A GRanges object containing 1361 ranges

**seqnames** chromosome

**ranges** start and stop positions for the block

**strand** genomic strand, irrelevant here

### Source

[https://github.com/jmacdon/LDblocks\\_GRCh38/tree/master/data](https://github.com/jmacdon/LDblocks_GRCh38/tree/master/data)

---

find_file_in_ftp	<i>Finding a file in an FTP directory This is an internal function to help gwascat.download find the right file</i>
------------------	---

---

### Description

Finding a file in an FTP directory This is an internal function to help gwascat.download find the right file

### Usage

```
find_file_in_ftp(ftp_address, acc, hm)
```

### Arguments

ftp_address	a string. An FTP address provided by gwascat.download.
acc	a string containing the accession for the desired study.
hm	a logical. Should it look in the harmonised directory?

**Value**

a data.table containing the dataset.

**Author(s)**

Guillermo Reales

---

gwascat.download	<i>Retrieve GWAS summary datasets from GWAS catalog 'gwascat.download takes a PMID from the user and downloads the associated summary statistics datasets published in GWAS catalog</i>
------------------	---

---

**Description**

This function, takes PUBMED ids as an input, searches at the GWAS catalog for harmonised datasets associated to that, interactively asking the user to choose if there are more than one, and fetches the dataset.

**Usage**

```
gwascat.download(ID, harmonised = TRUE)
```

**Arguments**

ID	a numeric. A PubMed ID (PMID) reference number from a GWAS paper.
harmonised	a logical. Should GWAS catalog harmonised files be pursued? If not available, the function will fall back to non-harmonised

**Details**

If multiple files are available for the same study, R will prompt an interactive dialogue to select a specific file, by number.

**Value**

a character vector containing the url(s) to the dataset(s).

**Author(s)**

Guillermo Reales

---

logsum	<i>Helper function to sum logs without loss of precision</i>
--------	--

---

**Description**

Sums logs without loss of precision This function is verbatim of its namesake in cupcake package ([github.com/ollyburren/cupcake/](https://github.com/ollyburren/cupcake/))

**Usage**

```
logsum(x)
```

**Arguments**

x                    a vector of logs to sum

**Value**

a scalar

**Author(s)**

Chris Wallace

---

michailidou19	<i>Subset of Michailidou BRCA GWAS sumstat dataset.</i>
---------------	---

---

**Description**

A data.table containing a subset of [Michailidou et al., 2017](#) breast cancer summary statistic dataset, in hg19 build. This dataset is freely available in GWAS catalog (see link below). We used "chromosome", "base\_pair\_location columns", removed unnecessary and all-missing columns, and took a random sample of 100,000 SNPs without replacement.

**Usage**

```
michailidou19
```

**Format**

A data.table object containing 100,000 SNPs  
 SNPID, CHR, BP, REF, ALT, ALT\_FREQ, BETA, SE, P  
**SNPID** rsids, or SNP ids  
**CHR** chromosome  
**BP** base position, in hg38  
**REF** reference, or non-effect allele  
**ALT** alternative, or effect allele  
**ALT\_FREQ** effect allele frequency  
**BETA** beta, log(OR), or effect size  
**SE** standard error of beta  
**P** p-value

**Source**

[http://ftp.ebi.ac.uk/pub/databases/gwas/summary\\_statistics/GCST004001-GCST005000/GCST004988/harmonised/29059683-GCST004988-EFO\\_0000305-build37.f.tsv.gz](http://ftp.ebi.ac.uk/pub/databases/gwas/summary_statistics/GCST004001-GCST005000/GCST004988/harmonised/29059683-GCST004988-EFO_0000305-build37.f.tsv.gz)

---

 michailidou38

---

*Subset of Michailidou BRCA GWAS sumstat dataset.*


---

**Description**

A data.table containing a subset of [Michailidou et al., 2017](#) breast cancer summary statistic dataset, in hg38 build. This dataset is freely available in GWAS catalog (see link below). We removed unnecessary and all-missing columns, and rows with missing data at hm\_beta and hm\_effect\_allele\_frequency, and took a random sample of 100,000 SNPs without replacement.

**Usage**

```
michailidou38
```

**Format**

A data.table object containing 100,000 SNPs  
**hm\_rsid** rsids, or SNP ids  
**hm\_chrom** chromosome  
**hm\_pos** base position, in hg38  
**hm\_other\_allele** reference, or non-effect allele  
**hm\_effect\_allele** alternative, or effect allele  
**hm\_beta** beta, log(OR), or effect size  
**hm\_effect\_allele\_frequency** effect allele frequency  
**standard\_error** standard error of beta  
**p\_value** p-value

**Source**

[http://ftp.ebi.ac.uk/pub/databases/gwas/summary\\_statistics/GCST004001-GCST005000/GCST004988/harmonised/29059683-GCST004988-EFO\\_0000305.h.tsv.gz](http://ftp.ebi.ac.uk/pub/databases/gwas/summary_statistics/GCST004001-GCST005000/GCST004988/harmonised/29059683-GCST004988-EFO_0000305.h.tsv.gz)

---

rapidopgs_multi	<i>Compute PGS from GWAS summary statistics using Bayesian sum of single-effect (SuSiE) linear regression using z scores</i>
-----------------	--

---

**Description**

'rapidopgs\_multi' computes PGS from a from GWAS summary statistics using Bayesian sum of single-effect (SuSiE) linear regression using z scores

**Usage**

```
rapidopgs_multi(
  data,
  reference = NULL,
  LDmatrices = NULL,
  N = NULL,
  build = c("hg19", "hg38"),
  trait = c("cc", "quant"),
  ncores = 1,
  alpha.block = 1e-04,
  alpha.snp = 0.01,
  sd.prior = NULL,
  ancestry = "EUR",
  LDblocks = NULL
)
```

**Arguments**

data	a data.table containing GWAS summary statistic dataset with all required information.
reference	a string representing the path to the directory containing the reference panel (eg. "../ref-data").
LDmatrices	a string representing the path to the directory containing the pre-computed LD matrices.
N	a numeric indicating the number of individuals used to generate input GWAS dataset, or a string indicating the column name containing per-SNP sample size.
build	a string indicating the genome build. 'hg19' and 'hg38' are supported. Note that your LD matrices or reference panel should match the build.
trait	a string indicating if trait is a case-control ("cc") or quantitative ("quant").
ncores	a numeric specifying the number of cores (CPUs) to be used. If using pre-computed LD matrices, one core is enough for best performance.



alpha.block	a numeric threshold for minimum P-value in LD blocks. Blocks with minimum P above alpha.block will be skipped. Default: 1e-4.
alpha.snp	a numeric threshold for P-value pruning within LD block. SNPs with P above alpha.snp will be removed. Default: 0.01.
sd.prior	the prior specifies that BETA at causal SNPs follows a centred normal distribution with standard deviation sd.prior. If NULL (default) it will be automatically estimated (recommended).
ancestry	a string indicating the ancestral population (DEFAULT: "EUR", European). If using an alternative population, bear in mind that your LD matrices or reference must be from the same population. You'll also need to provide matching LD.blocks via the LDblocks argument.
LDblocks	a string indicating the path to an alternative LD block file in .RData format. Only required for non-European PGS.

## Details

This function will take a GWAS summary statistic dataset as an input, will assign LD blocks to it, then use user-provided LD matrices or a preset reference panel in Plink format to compute LD matrices for each block. Then SuSiE method will be used to compute posterior probabilities of variants to be causal and generate PGS weights by multiplying those posteriors by effect sizes ( $\beta$ ). Unlike rapidopgs\_single, this approach will assume one or more causal variants.

The GWAS summary statistics file to compute PGS using our method must contain the following minimum columns, with these exact column names:

**CHR** Chromosome  
**BP** Base position (in GRCh37/hg19).  
**REF** Reference, or non-effect allele  
**ALT** Alternative, or effect allele, the one  $\beta$  refers to  
**BETA**  $\beta$  (or log(OR)), or effect sizes  
**SE** standard error of  $\beta$   
**P** P-value for the association test

In addition, quantitative traits must have the following extra column:

**ALT\_FREQ** Minor allele frequency.

Also, for quantitative traits, sample size must be supplied, either as a number, or indicating the column name, for per-SNP sample size datasets (see below). Other columns are allowed, and will be ignored.

Reference panel should be divided by chromosome, in Plink format. Both reference panel and summary statistic dataset should be in GRCh37/hg19. For 1000 Genomes panel, you can use create\_1000G function to set it up automatically.

If prefer to use LD matrices, you must indicate the path to the directory where they are stored. They must be in RDS format, named LD\_chrZ.rds (where Z is the 1-22 chromosome number). If you don't have LD matrices already, we recommend downloading those gently provided by Prive et al., at [https://figshare.com/articles/dataset/European\\_LD\\_reference/13034123](https://figshare.com/articles/dataset/European_LD_reference/13034123). These matrices were computed using for 1,054,330 HapMap3 variants based on 362,320 European individuals of the UK biobank.

**Value**

a data.table containing the sumstats dataset with computed PGS weights.

**Author(s)**

Guillermo Reales, Chris Wallace

**Examples**

```
## Not run:
ss <- data.table(
  CHR=c("4", "20", "14", "2", "4", "6", "6", "21", "13"),
  BP=c(1479959, 13000913, 29107209, 203573414, 57331393, 11003529, 149256398,
    25630085, 79166661),
  REF=c("C", "C", "C", "T", "G", "C", "C", "G", "T"),
  ALT=c("A", "T", "T", "A", "A", "A", "T", "A", "C"),
  BETA=c(0.012, 0.0079, 0.0224, 0.0033, 0.0153, 0.058, 0.0742, 0.001, -0.0131),
  SE=c(0.0099, 0.0066, 0.0203, 0.0171, 0.0063, 0.0255, 0.043, 0.0188, 0.0074),
  P=c(0.2237, 0.2316, 0.2682, 0.8477, 0.01473, 0.02298, 0.08472, 0.9573, 0.07535))
PGS <- rapidopgs_multi(ss, reference = "ref-data/", N = 20000, build = "hg19", trait="cc", ncores=5)

## End(Not run)
```

---

rapidopgs\_single

*Compute PGS from GWAS summary statistics using posteriors from Wakefield's approximate Bayes Factors*

---

**Description**

'rapidopgs\_single computes PGS from a from GWAS summary statistics using posteriors from Wakefield's approximate Bayes Factors

**Usage**

```
rapidopgs_single(
  data,
  N = NULL,
  trait = c("cc", "quant"),
  build = "hg19",
  pi_i = 1e-04,
  sd.prior = if (trait == "quant") {
    0.15
  } else {
    0.2
  },
  filt_threshold = NULL,
  recalc = TRUE,
  reference = NULL
)
```

**Arguments**

<code>data</code>	a data.table containing GWAS summary statistic dataset with all required information.
<code>N</code>	a scalar representing the sample in the study, or a string indicating the column name containing it. Required for quantitative traits only.
<code>trait</code>	a string specifying if the dataset corresponds to a case-control ("cc") or a quantitative trait ("quant") GWAS. If trait = "quant", an ALT_FREQ column is required.
<code>build</code>	a string containing the genome build of the dataset, either "hg19" (for hg19/GRCh37) or "hg38" (hg38/GRCh38). DEFAULT "hg19".
<code>pi_i</code>	a scalar representing the prior probability (DEFAULT: $1 \times 10^{-4}$ ).
<code>sd.prior</code>	the prior specifies that BETA at causal SNPs follows a centred normal distribution with standard deviation sd.prior. Sensible and widely used DEFAULTs are 0.2 for case control traits, and $0.15 * \text{var}(\text{trait})$ for quantitative (selected if trait == "quant").
<code>filt_threshold</code>	a scalar indicating the ppi threshold (if <code>filt_threshold &lt; 1</code> ) or the number of top SNPs by absolute weights (if <code>filt_threshold &gt;= 1</code> ) to filter the dataset after PGS computation. If NULL (DEFAULT), no thresholding will be applied.
<code>recalc</code>	a logical indicating if weights should be recalculated after thresholding. Only relevant if <code>filt_threshold</code> is defined.
<code>reference</code>	a string indicating the path of the reference file SNPs should be filtered and aligned to, see Details.

**Details**

This function will take a GWAS summary statistic dataset as an input, will assign align it to a reference panel file (if provided), then it will assign SNPs to LD blocks and compute Wakefield's ppi by LD block, then will use it to generate PGS weights by multiplying those posteriors by effect sizes ( $\beta$ ). Optionally, it will filter SNPs by a custom filter on ppi and then recalculate weights, to improve accuracy.

Alternatively, if `filt_threshold` is larger than one, RapidoPGS will select the top `filt_threshold` SNPs by absolute weights (note, not ppi but weights).

The GWAS summary statistics file to compute PGS using our method must contain the following minimum columns, with these exact column names:

**CHR** Chromosome

**BP** Base position (in GRCh37/hg19 or GRCh38/hg38). If using hg38, use `build = "hg38"` in parameters

**REF** Reference, or non-effect allele

**ALT** Alternative, or effect allele, the one  $\beta$  refers to

**ALT\_FREQ** Minor/ALT allele frequency in the tested population, or in a close population from a reference panel. Required for Quantitative traits only

**BETA**  $\beta$  (or  $\log(\text{OR})$ ), or effect sizes

**SE** standard error of  $\beta$

If a reference is provided, it should have 5 columns: CHR, BP, SNPID, REF, and ALT. Also, it should be in the same build as the summary statistics. In both files, column order does not matter.

### Value

a data.table containing the formatted sumstats dataset with computed PGS weights.

### Author(s)

Guillermo Reales, Chris Wallace

### Examples

```
sumstats <- data.table(SNPID=c("rs139096444","rs3843766","rs61977545", "rs544733737",
  "rs2177641", "rs183491817", "rs72995775","rs78598863", "rs1411315"),
  CHR=c("4","20","14","2","4","6","6","21","13"),
  BP=c(1479959, 13000913, 29107209, 203573414, 57331393, 11003529, 149256398,
    25630085, 79166661),
  REF=c("C","C","C","T","G","C","C","G","T"),
  ALT=c("A","T","T","A","A","A","T","A","C"),
  BETA=c(0.012,0.0079,0.0224,0.0033,0.0153,0.058,0.0742,0.001,-0.0131),
  SE=c(0.0099,0.0066,0.0203,0.0171,0.0063,0.0255,0.043,0.0188,0.0074))

PGS <- rapidopgs_single(sumstats, trait = "cc")
```

---

sd.prior.est

*Compute Standard deviation prior (SD prior) for quantitative traits using pre-computed heritability.*

---

### Description

sd.prior.est function will take the dataset as an input, a  $h^2$  value obtained from a public repository such as LDhub, (<http://ldsc.broadinstitute.org/ldhub/>), sample size and number of variants, and will provide a sd.prior estimate that can be used to improve prediction performance of RapidoPGS functions on quantitative traits.

### Usage

```
sd.prior.est(data, h2, N, pi_i = 1e-04)
```

### Arguments

data	a data.table containing the GWAS summary statistic input dataset. Must contain SNPID and SE columns.
h2	a numeric. Heritability estimate or $h^2$ (See details).
N	a numeric. Sample size of the GWAS input dataset.
pi_i	a numeric. Prior that a given variant is causal. DEFAULT = 1e-4.

**Author(s)**

Guillermo Reales, Elena Vigorito, Chris Wallace

**Examples**

```
sumstats <- data.table(SNPID=c("4:1479959", "20:13000913", "14:29107209", "2:203573414",
"4:57331393", "6:11003529", "6:149256398", "21:25630085", "13:79166661"),
  REF=c("C", "C", "C", "T", "G", "C", "C", "G", "T"),
  ALT=c("A", "T", "T", "A", "A", "A", "T", "A", "C"),
  ALT_FREQ=c(0.2611, 0.4482, 0.0321, 0.0538, 0.574, 0.0174, 0.0084, 0.0304, 0.7528),
  BETA=c(0.012, 0.0079, 0.0224, 0.0033, 0.0153, 0.058, 0.0742, 0.001, -0.0131),
  SE=c(0.0099, 0.0066, 0.0203, 0.0171, 0.0063, 0.0255, 0.043, 0.0188, 0.0074),
  P=c(0.2237, 0.2316, 0.2682, 0.8477, 0.01473, 0.02298, 0.08472, 0.9573, 0.07535))
sd.prior <- sd.prior.est(sumstats, h2 = 0.2456, N = 45658, pi_i=1e-4)
```

---

sdY.est

*Estimate trait variance, internal function*

---

**Description**

Estimate trait standard deviation given vectors of variance of coefficients, MAF and sample size

**Usage**

```
sdY.est(vbeta, maf, n)
```

**Arguments**

vbeta	vector of variance of coefficients
maf	vector of MAF (same length as vbeta)
n	sample size

**Details**

Estimate is based on  $\text{var}(\hat{\beta}) = \text{var}(Y) / (n * \text{var}(X))$   $\text{var}(X) = 2 * \text{maf} * (1 - \text{maf})$  so we can estimate  $\text{var}(Y)$  by regressing  $n * \text{var}(X)$  against  $1 / \text{var}(\hat{\beta})$  This function is verbatim from its namesake in coloc package ([github.com/chr1swallace/coloc/](https://github.com/chr1swallace/coloc/)), by Chris Wallace

**Value**

estimated standard deviation of Y

**Author(s)**

Chris Wallace

---

wakefield_pp	<i>compute posterior probabilities using Wakefield's approximate Bayes Factors</i> wakefield_pp computes posterior probabilities for a given SNP to be causal for a given SNP under the assumption of a single causal variant.
--------------	--

---

### Description

This function was adapted from its namesake in cupcake package ([github.com/ollyburren/cupcake/](https://github.com/ollyburren/cupcake/)) to no longer require allele frequencies.

### Usage

```
wakefield_pp(beta, se, pi_i = 1e-04, sd.prior = 0.2)
```

### Arguments

beta	a vector of effect sizes ( $\beta$ ) from a quantitative trait GWAS
se	vector of standard errors of effect sizes ( $\beta$ )
pi_i	a scalar representing the prior probability (DEFAULT $1 \times 10^{-4}$ )
sd.prior	a scalar representing our prior expectation of $\beta$ (DEFAULT 0.2). The method assumes a normal prior on the population log relative risk centred at 0 and the DEFAULT value sets the variance of this distribution to 0.04, equivalent to a 95% is in the range of 0.66-1.5 at any causal variant.

### Value

a vector of posterior probabilities.

### Author(s)

Olly Burren, Chris Wallace, Guillermo Reales

---

wakefield_pp_quant	<i>Compute posterior probabilities using Wakefield's approximate Bayes Factors for quantitative traits</i>
--------------------	--

---

### Description

wakefield\_pp\_quant computes posterior probabilities for a given SNP to be causal for a given SNP under the assumption of a single causal variant.

### Usage

```
wakefield_pp_quant(beta, se, sdY, sd.prior = 0.15, pi_i = 1e-04)
```

**Arguments**

beta	a vector of effect sizes ( $\beta$ ) from a quantitative trait GWAS
se	vector of standard errors of effect sizes ( $\beta$ )
sdY	a scalar of the standard deviation given vectors of variance of coefficients, MAF and sample size. Can be calculated using <code>sdY.est</code>
sd.prior	a scalar representing our prior expectation of $\beta$ (DEFAULT 0.15).
pi_i	a scalar representing the prior probability (DEFAULT $1 \times 10^{-4}$ ) The method assumes a normal prior on the population log relative risk centred at 0 and the DEFAULT value sets the variance of this distribution to 0.04, equivalent to a 95% is in the range of 0.66-1.5 at any causal variant.

**Details**

This function was adapted from `wakefield_pp` in `cupcake` package ([github.com/ollyburren/cupcake/](https://github.com/ollyburren/cupcake/))

**Value**

a vector of posterior probabilities.

**Author(s)**

Guillermo Reales, Chris Wallace

# Index

## \* datasets

- EUR\_ld.blocks19, 3
- EUR\_ld.blocks38, 4
- michailidou19, 6
- michailidou38, 7

create\_1000G, 2

EUR\_ld.blocks19, 3

EUR\_ld.blocks38, 4

find\_file\_in\_ftp, 4

gwascat.download, 5

logsum, 6

michailidou19, 6

michailidou38, 7

rapidopgs\_multi, 8

rapidopgs\_single, 10

sd.prior.est, 12

sdY.est, 13

wakefield\_pp, 14

wakefield\_pp\_quant, 14