

Package ‘INPower’

October 18, 2017

Title An R package for computing the number of susceptibility SNPs

Version 1.12.0

Date 2014-03-18

Author Ju-Hyun Park

Description An R package for computing the number of susceptibility SNPs and power of future studies

Maintainer Bill Wheeler <wheelerb@imsweb.com>

Depends R (>= 3.1.0), mvtnorm

Suggests RUnit, BiocGenerics

License GPL-2 + file LICENSE

biocViews SNP

NeedsCompilation no

R topics documented:

| | |
|-------------------|----------|
| INPower | 1 |
| Index | 4 |

| | |
|---------|---|
| INPower | <i>Estimate the number of susceptibility SNPs and power of future studies</i> |
|---------|---|

Description

This function uses the effect sizes for a set of known susceptibility SNPs and the power of detection of these SNPs from the original discovery samples to obtain an estimate of the total number of underlying susceptibility SNP for that trait and the distribution of their effect sizes. The function can further use the estimated number of loci and distribution of effect sizes to evaluate the power for discovery of a future GWAS study (up to three-stage).

Usage

```
INPower(MAFs, betas, pow, sample.size, signif.lvl, k, span=0.5, binary.outcome=TRUE,  
multi.stage.option=NULL, tgv=NULL)
```

Arguments

| | |
|--------------------|---|
| MAFs | Vector of minor allele frequencies associated with the set of known loci |
| betas | Vector of regression effects for the set of known loci under an additive genetic model. For a continuous phenotype analyzed with linear regression model, it is assumed that the outcome has been standardized so that the coefficients correspond to mean change in outcome per unit of s.d. for each copy of the given allele. For a binary outcome analyzed with logistic regression, the regression coefficients should correspond to change in log-odds-ratio per copy of the given allele. |
| pow | A vector representing the powers for the known loci in the original studies that led to their discoveries. Note these power calculations should be carefully done to avoid winner's curse (it is best to obtain effect size estimates from independent replication study) and to take into consideration all complexities of the designs of the original study. If the total SNP set is obtained from a group of studies for a given trait, then the power for an individual marker should reflect the probability of its detection in at least one of the studies. |
| sample.size | Sample size for a future study for which integrated power calculation is desired. For case-control studies, half of the subjects are assumed to be cases and half to be controls. It can take a vector of several sample sizes for the same study as shown in the example below. |
| signif.lvl | The required genome-wide significance level for future study. |
| k | A vector of integer values for which the user would like to calculate probabilities of the type $\Pr(X \geq k)$ to evaluate the probability of detection of at least a specified number of loci in future studies. In addition, the function automatically finds nine values for "k", for which the probabilities are close to 0.1 to 0.9 with an increment of 0.1. |
| span | The parameter which controls the degree of smoothing in loess . It specifies the fraction of SNPs that are used in local linear regression to obtain the estimated number of loci at each effect size. The default is set at 0.5, but we recommend the user to set it at a value depending on the total size of the SNP set so that about 10-20 SNPs are used for local smoothing at each effect size. The total size of the SNP set should be reasonably large (e.g. at least 20 and preferably more) for application of loess . |
| binary.outcome | TRUE/FALSE Is the outcome binary or continuous? |
| multi.stage.option | This option allows to set-up design parameters for the future study if it would be done in multiple stages (up to three). The option has a list of two arguments alpha and pi, where alpha specifies the significance level(s) used for each stage to select markers for the subsequent stage and pi specifies the fraction of subjects who are included in the corresponding stages. The default for the option is NULL, that is, the study is assumed to be single-stage. |
| tgv | An optional argument using which the user can input an estimate of the known total genetic variance (TGV) of the trait that may be available from familial aggregation studies. For a continuous outcome, this could be an estimate of the fraction of the total variance of the trait attributed to heritability. For a binary outcome, this could logarithm of squared sibling-relative-risk that is known to approximate total genetic variance under log-normal model for risk. |

Details

The projections are only shown in the range of effect size for which the original studies had at least 1 percent power. The `loess` fitting procedure, however, may include additional SNPs with smaller effect sizes for local linear smoothing. The user is recommended to remove SNPs that may seem clearly outliers compared to the rest in terms of their effect sizes. By default the program currently removes all SNPs with power less than 0.1 percent from the analysis to avoid undue influence of potentially outlying observations.

Value

A list of two sublists with names `esdist.summary` and `future.study.summary`. The sublist `esdist.summary` contains the estimated number of loci (`t.n.loci`), the genetic variance explained by the estimated number of loci (`gve`), and the estimated number of loci at each different effect size (`es.dist`). Note for linear regression, `gve` is expressed as a percentage of the total variance of the outcome, since it assumed that outcome has been standardized. Further, if an estimate of total genetic variance (TGV) is provided by the user, then the estimate for GVE will be automatically expressed as a percentage of TGV. The sublist `future.study.summary` contains the expected number of loci to be discovered in the future study (`e.discov`), expected genetic variance explained (`e.gve`), and a table of probabilities of discovering at least `k` loci for the different values of `k` (`prob.k`). Note that `e.gve` is defined similarly to `gve`.

References

Park et al. (2010). Estimation of effect size distribution from genome-wide association studies and implications for future discoveries. *Nature Genetics*, 42:570-5.

Examples

```
set.seed(123)
MAFs <- runif(50, min=0.05, max=0.5)
betas <- runif(50, min=-0.5, max=0.5)
pow <- runif(50, min=0.1, max=0.9)
sample.size <- 1000
signif.lvl <- 1e-4
k <- 20

INPower(MAFs, betas, pow, sample.size, signif.lvl, k)
```

Index

*Topic **models**

INPower, 1

INPower, 1

loess, 2, 3