

1

MungeSumstats: A Bioconductor package for the standardisation and quality control of many GWAS summary statistics

Alan E Murphy¹, Nathan G Skene¹

¹ UK Dementia Research Institute at Imperial College London

Abstract

Summary

Genome-wide association studies (GWAS) summary statistics have democratised and accelerated genetic research. However, a lack of standardisation of the file formats used has proven problematic when running secondary analysis tools or performing meta-analysis studies. To address these issues, we have developed MungeSumstats, a Bioconductor R package for the standardisation and quality control of GWAS summary statistics. MungeSumstats can handle the most common summary statistic formats, including variant call format (VCF) producing a reformatted, standardised, tabular summary statistic file.

Contact

Alan Murphy: a.murphy@imperial.ac.uk, Nathan Skene: n.skene@imperial.ac.uk

Availability and implementation

MungeSumstats is available on Bioconductor (v 3.13) and can also be found on Github at:

<https://neurogenomics.github.io/MungeSumstats>

Supplementary information

The analysis deriving the most common summary statistic formats is available at:

<https://al-murphy.github.io/SumstatFormats>

Introduction

Genome-wide association studies (GWAS) summary statistics are used to distribute the most important outputs of GWASs in a manner which does not require the transfer of individual-level personally identifiable information from participants. Summary statistics from past studies tend to become more valuable over time as it becomes possible to meta-analyse them and integrate them with new annotation information through approaches such as Linkage Disequilibrium Score Regression (LDSC)(1), Generalized Gene-Set Analysis of GWAS Data, MAGMA(2), and multi-phenotype investigations(3,4). Summary statistics are also commonly integrated for use in the meta-analysis of GWAS. However, these tools and this integration require a standardised data format which was historically lacking from the field. The diversity of data formats in summary statistics has been a result of the phenotypes in question, for example disease-control or quantitative trait, the software used to perform the analysis, such as PLINK(5) and GCTA(6) or just the preference of the consortium in question.

There have been movements to standardise the summary statistic file format such as the NHGRI-EBI GWAS Catalogue standardised format(7) and the SMR Tool binary format(8). More recently, the variant call format to store GWAS summary statistics (GWAS-VCF)(9) has been developed which has manually converted over 10,000 GWAS to this format. While GWAS-VCF offers a standardised format that future GWAS consortium may adopt, there are still a multitude of past, publicly available GWAS which have not been standardised(10)(11)(12)(13). For instance, although their summary statistics are publicly available, the GWAS for Cerebral small vessel disease(14) is not yet available in VCF format via IEU GWAS. Furthermore, as VCF is not yet the standard for sharing files between geneticists, unpublished GWAS shared internally within genetics consortia or provided by personal genetics companies are still found in a variety of summary statistic formats. As such,

there is a need for tools to move between the various formats in which summary statistics are stored.

The standardisation of GWAS summary statistics also requires quality control to ensure cohesive integration. For example, checking if the non-effect allele from the summary statistics matches the reference sequence from a reference genome to ensure consistent directionality of allelic effects across GWAS. In addition, downstream analysis tools often require a degree of quality control which, in the case of meta-analysis, must be applied across all GWAS. One such example is the removal of all non-biallelic SNPs is a common requirement of all downstream analysis(9).

To address these issues, we introduce MungeSumstats a Bioconductor R package for the rapid standardisation and quality control of many GWAS summary statistics. MungeSumstats can handle the most common summary statistic formats as well as GWAS-VCFs to enable the integrative meta-analysis of diverse GWAS. MungeSumstats also offers a tuneable quality control protocol with defaults for common, best-practice approaches. MungeSumstats capitalises on R's familiar interface, is readily accessible through Bioconductor and utilises an intuitive approach, running with a single line of input code.

Heterogeneity in GWAS Formats

To demonstrate the diversity in summary statistics across GWAS, we analysed a public repository of over 200 publicly available GWAS(15). From this, the most common summary statistics were derived, see Figure 1 for the 12 most common file header formats.

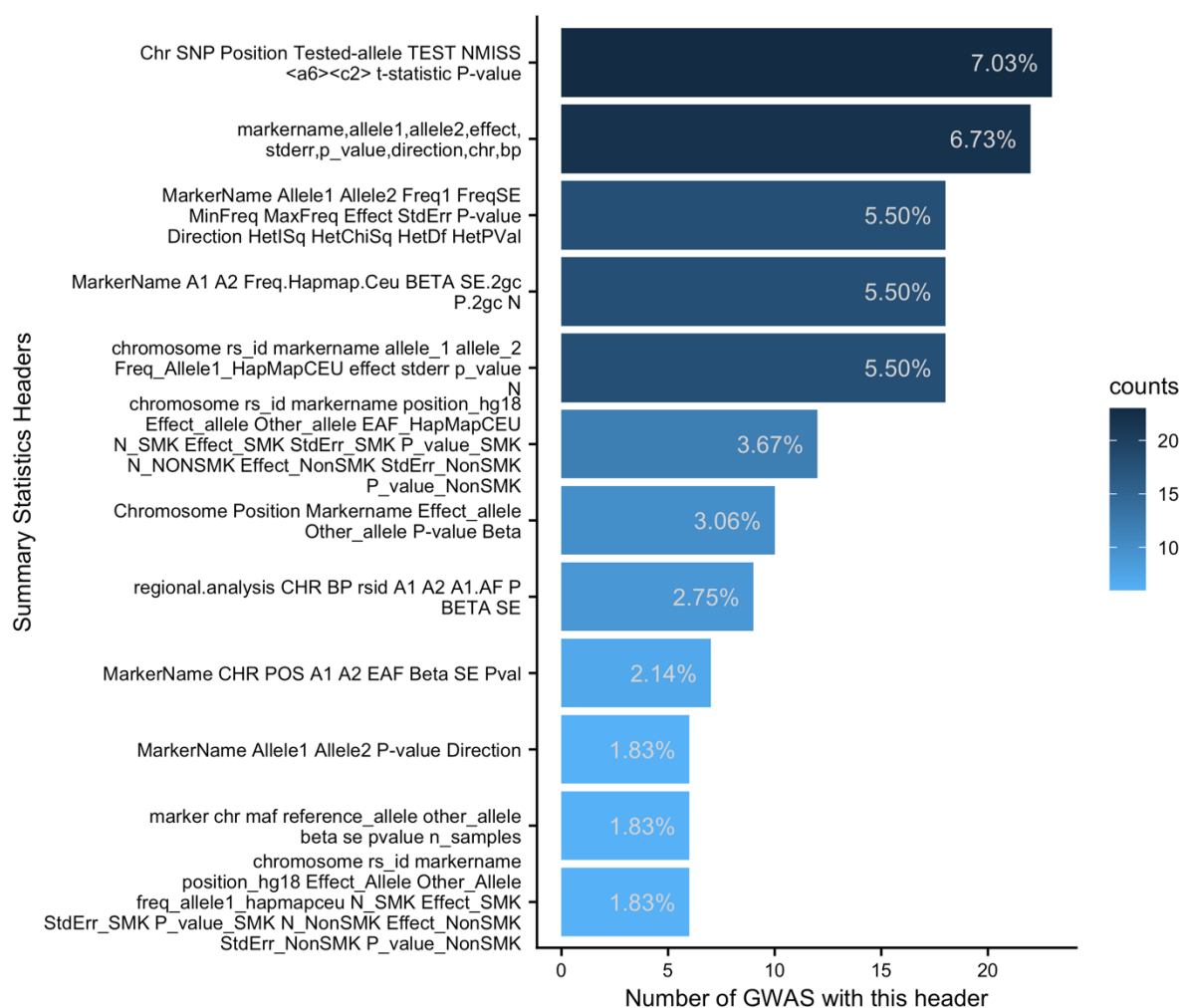


Figure 1: Most common summary statistic formats shows the most common summary statistic formats from a repository of over 200 publicly available GWAS(15). Note that a GWAS can have more than 1 summary statistics file and '<a6><c2>' is the symbol '␣' read into R.

A total of 327 summary statistic files were derived from the analysis which corresponded to 127 unique formats. Thus on average, every 2.5 summary statistic files had a unique format, showing the clear disparity across GWAS. The 12 most common formats, shown in Figure 1, accounted for approximately 47% all summary statistics. MungeSumstats has been tested on these 12 most common formats and is able to standardise their summary statistics.

Implementation

MungeSumstats was implemented using the R programming language (v 4.0) and Bioconductor S4 data infrastructure (v 3.13) enabling the full analysis of summary statistics within the R environment. The package removes the need for external software to perform the standardisation and quality control steps.

MungeSumstats' implementation ensures both memory and speed efficiency through the use of R `data.table` (v.1.14.0)(16), which can take advantage of multi-core parallelization. Moreover, MungeSumstats benefits from Bioconductor's infrastructure for efficient representation of full genomes and their SNPs, using `BSgenome` (v 1.59.2) SNP reference genomes(17). Either Ensembl's GRCh37 or GRCh38 are queried dependent on the build for the particular GWAS. Numerous of MungeSumstats' quality control steps for summary statistics require the use of a reference genome. These steps ensure all SNPs, and their base-pair position, relate to the specified build, ensure consistent directionality of any allelic effect variables and impute any missing, essential information like SNP ID, base-pair position and effect/non-effect allele.

Using these two infrastructures, MungeSumstats conducts more than 25 checks on the inputted summary statistics file, see Table 1 for a description of their use. MungeSumstats is also written to ensure the ease of addition of further checks so if users have summary statistics which can't currently be handled in MungeSumstats, these can be incorporated easily in future releases. Finally, MungeSumstats returns a reformatted, tabular summary statistics file with standardised columns for the information necessary for downstream analysis.

Step	MungeSumstats Check	Description
1	Check VCF format	Check if the input file is in variant call format (VCF) and if so, convert to standard format
2	Check tab, space or comma delimited	If input is space or comma delimited convert to tab delimited.
3	Check for header name synonyms	If any alternative names are found for SNP, BP, CHR, A1, A2, P, Z, OR, BETA, LOG_ODDS, SIGNED_SUMSTAT, N, N_CAS, N_CON, NSTUDY, INFO or FRQ convert them to the standard variable name
4	Check for multiple models or traits in GWAS	If multiple, user must specify one to analyse
5	Check for uniformity in SNP ID	Ensure no mix of RS ID, missing 'rs' prefix and/or CHR:BP
6	Check for CHR:BP:A2:A1 all in one column	Split into separate columns if found
7	Check for CHR:BP in one column	Split into separate columns if found
8	Check for A1/A2 in one column	Split into separate columns if found
9	Check if CHR and/or BP is missing	If so, infer from the chosen reference genome
10	Check if SNP ID is missing	If so, infer from the chosen reference genome
11	Check if A1 and/or A2 are missing	If so, infer from the chosen reference genome
12	Check that vital columns are present	Check for the following necessary columns; SNP, CHR, BP, P, A1, A2
13	Check for one signed/effect column	Check for one of the following necessary columns; Z, OR, BETA, LOG_ODDS, SIGNED_SUMSTAT
14	Check for missing data	If data is missing from any entry, remove the SNP
15	Check for duplicated columns	If there are any remove one
16	Check for small p-values (lower than $5e-324$)	These are not recognised in R and can cause issues with downstream analysis software like LDSC/MAGMA. If any are found they are converted to 0.
17	Check N column	Ensure it is an integer and check if the sample size for any SNP is not greater than mean multiplied by five times the standard deviation. Removes SNPs that have substantial more samples than the rest.
18	Check SNPs are RS ID's	SNP IDs must be RS IDs
19	Check for duplicated rows, based on SNP ID	Duplicates are removed
21	Check for duplicated rows, based on base-pair position	Duplicates are removed
22	Check for SNPs on reference genome	Correct any missing from reference genome using BP and CHR
23	Check INFO score	Remove SNPs with imputation score less than 0.9
24	Check for strand-ambiguous SNPs	Remove strand-ambiguous SNPs if found
25	Check for non-biallelic SNPs (infer from reference genome)	Infer from chosen reference genome and remove any if found
26	Check for allele flipping	The non-effect allele is checked against the reference sequence from the chosen reference genome. The allele columns, along with all effect columns are flipped if allele aligns more closely with the alternative sequence, creating consistent directionality of allelic effects across GWAS.
27	Check for SNPs on chromosome X, Y, and mitochondrial SNPs (MT)	If any are found these are removed.

Table 1: MungeSumstats Implemented Checks lists the quality control and standardisation checks conducted by the package. Here CHR is chromosome, BP is Base-pair position, A1 is the effect allele, A2 is the non-effect allele, N is the sample size, INFO is imputation information score, FRQ is the minor allele frequency (MAF) of the SNP, SNP ID is the single nucleotide polymorphism reference ID and P is the unadjusted p-value. For the effect variables, Z is z-score, OR is odds ratio, LOG_ODDS is the log odds ratio, BETA is the effect size estimate relative to the alternative allele and SIGNED_SUMSTAT is the directional effect size estimate for the summary statistics.

Usage

Once MungeSumstats is installed, usage involves a single line of code or one function call (*format_sumstats*) with the path to the summary statistics file of interest. Then, the path to the reformatted, standardised summary statistic file is returned. MungeSumstats also offers adjustable parameters to manage the quality control steps. These include options to adjust the imputation information score (INFO) cut-off threshold, the number of samples (N) outliers cut-off threshold and whether to remove mitochondrial SNPs or SNPs on the X or Y chromosome (see Table 1). Quality control steps which use a reference genome can also be adjusted such as whether to filter SNPs based on their RS ID's presence on the reference genome, whether to check for allele flipping and whether to remove multi-allelic or strand-ambiguous SNPs. These parameters ensure MungeSumstats can be adjusted to the user's analysis pipelines.

Conclusion

Here, we presented MungeSumstats, a Bioconductor package for the standardisation and quality control of GWAS summary statistics. This package enables integration of summary statistics of vastly different formats, simplifying meta-analysis and summary statistics use in other secondary research applications. The package provides an efficient, user-friendly R-native approach, returning a standardised, tabular format file. This ensures that the summary statistics are accessible to the average user. Moreover, MungeSumstats is written to permit future development of additional standardisation steps if users encounter issues with their specific GWAS.

Acknowledgements

This work was supported by a UKRI Future Leaders Fellowship [grant number MR/T04327X/1] and the UK Dementia Research Institute which receives its funding from UK DRI Ltd, funded by the UK Medical Research Council, Alzheimer's Society and Alzheimer's

Research UK. We thank Alexandru Voda (@alexandruioanvoda) for contributing via github in making an early version of this package work across operating systems.

Data availability

The data to derive the summary statistic formats is open source and collated at:

<https://github.com/mikegloudemans/gwas-download>

Bibliography

1. Bulik-Sullivan B, Finucane HK, Anttila V, Gusev A, Day FR, Loh P-R, et al. An atlas of genetic correlations across human diseases and traits. *Nat Genet.* 2015 Nov;47(11):1236–41.
2. Leeuw CA de, Mooij JM, Heskes T, Posthuma D. MAGMA: Generalized Gene-Set Analysis of GWAS Data. *PLOS Comput Biol.* 2015 Apr 17;11(4):e1004219.
3. Aguirre M, Tanigawa Y, Venkataraman GR, Tibshirani R, Hastie T, Rivas MA. Polygenic risk modeling with latent trait-related genetic components. *Eur J Hum Genet.* 2021 Feb 8;1–11.
4. Tanigawa Y, Li J, Justesen JM, Horn H, Aguirre M, DeBoever C, et al. Components of genetic associations across 2,138 phenotypes in the UK Biobank highlight adipocyte biology. *Nat Commun.* 2019 Sep 6;10(1):4064.
5. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: A Tool Set for Whole-Genome Association and Population-Based Linkage Analyses. *Am J Hum Genet.* 2007 Sep 1;81(3):559–75.
6. Yang J, Lee SH, Goddard ME, Visscher PM. GCTA: A Tool for Genome-wide Complex Trait Analysis. *Am J Hum Genet.* 2011 Jan 7;88(1):76–82.
7. Buniello A, MacArthur JAL, Cerezo M, Harris LW, Hayhurst J, Malangone C, et al. The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.* 2019 Jan 8;47(D1):D1005–12.
8. Zhu Z, Zhang F, Hu H, Bakshi A, Robinson MR, Powell JE, et al. Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. *Nat Genet.* 2016 May;48(5):481–7.
9. Lyon MS, Andrews SJ, Elsworth B, Gaunt TR, Hemani G, Marcora E. The variant call format provides efficient and robust storage of GWAS summary statistics. *Genome Biol.* 2021 Jan 13;22(1):32.
10. Luciano M, Davies G, Summers KM, Hill WD, Hayward C, Liewald DC, et al. The influence of X chromosome variants on trait neuroticism. *Mol Psychiatry.* 2021 Feb;26(2):483–91.
11. Jansen PR, Watanabe K, Stringer S, Skene N, Bryois J, Hammerschlag AR, et al. Genome-wide analysis of insomnia in 1,331,010 individuals identifies new risk loci and functional pathways. *Nat Genet.* 2019 Mar;51(3):394–403.
12. Lin H, van Setten J, Smith AV, Bihlmeyer NA, Warren HR, Brody JA, et al. Common and Rare Coding Genetic Variation Underlying the Electrocardiographic PR Interval. *Circ Genomic Precis Med.* 2018 May;11(5):e002037.
13. McCormack M, Gui H, Ingason A, Speed D, Wright GEB, Zhang EJ, et al. Genetic variation in CFH predicts phenytoin-induced maculopapular exanthema in European-descent patients. *Neurology.* 2018 Jan 23;90(4):e332–41.

14. Sargurupremraj M, Suzuki H, Jian X, Sarnowski C, Evans TE, Bis JC, et al. Cerebral small vessel disease genomics and its implications across the lifespan. *Nat Commun.* 2020 Dec 8;11(1):6285.
15. Gloudemans M. `mikegloudemans/gwas-download` [Internet]. 2021 [cited 2021 Apr 20]. Available from: <https://github.com/mikegloudemans/gwas-download>
16. Dowle M, Srinivasan A. `data.table: Extension of `data.frame`` [Internet]. 2021. Available from: <https://CRAN.R-project.org/package=data.table>
17. Pagès H. `BSgenome: Software infrastructure for efficient representation of full genomes and their SNPs` [Internet]. 2021. Available from: <https://bioconductor.org/packages/BSgenome>