SNVsMutations package to determine the mutation type and counts for a set of single nucleotide variants in a vcf file

Federica Boselli

*federica.boselli@mail.polimi.it

1 June 2023

Abstract

This package takes:a set of mutations (single nucleotide variants, SNVs) in VCF format, the corresponding reference genome (e.g., human genome hg38) and a parameter "context_length" which is a positive, odd integer and determines for each mutation (only SNVs; other mutations like indels can be ignored) the corresponding mutation type. Lastly, the package can provide a graphical output which plots the summarized mutation type counts.

Package

SNVsMutations 0.1.0

Contents

```
1 Using the SNVsMUtations package to determine the mutation type for a set of
single nucleotide variants in a genome
      1.1 Background
            1.1.1 SNVs (Single Nucleotide Variants)
```

1.1.2 VCF file 1.2 Installation and dependencies

1.3 Loading the VCF file

1.4 Getting the mutation types from the vcf file through the SNV_Types() function 1.5 Generating the mutation counts and a final graphical output with the SNV_Counts() function

Using the SNVsMUtations package to determine the mutation type for a set of single nucleotide variants in

Background 1.1 The **SNVsMutations** package is designed to analyze a set of single nucleotide variants (SNVs) in

parameter "context_length." The goal of this package, developed for the Scientific Programming course at Politecnico of Milan as part of the Bioinformatics for Computational Genomics MSc

program, is to determine the mutation types for each SNV and provide a graphical summary of the mutation type counts. To get started with the SNVsMutations package in R, follow this vignette for a guide on how to use it effectively. SNVs (Single Nucleotide Variants) 1.1.1

VCF format, along with the corresponding reference genome (e.g., human genome hg38), and the

Single-nucleotide variant (SNV), also known as single-nucleotide polymorphism (SNP), is the variant of a single nucleotide that occurs at a specific genomic position.

VCF is a text file format (most likely stored in a compressed manner). It contains meta-information

VCF file 1.1.2

Installation and dependencies

To load the packages if they have been already installed:

library(BSgenome.Hsapiens.UCSC.hg38)

library(VariantAnnotation)

vcf2 <- readVcf(file, "hq38")</pre>

characters representing the mutation types.

to the reference genome.

library(GenomicRanges)

Session info

a genome

lines, a header line, and then data lines each containing information about a position in the genome. The format also has the ability to contain genotype information on samples for each position. From: https://samtools.github.io/hts-specs/VCFv4.2.pdf

This package depends from: VariantAnnotation, Biostrings, BSgenome. Hsapiens. UCSC. hg38,

library(BiocStyle)

library(knitr)

1.2

1.4

1.5

SNV_Counts() function

GenomicRanges. So to be able to install this package, these 4 packages must be installed. For BiocManager: example, to manually install package from BiocManager::install("Biostrings") After that, SNVsMutations can be installed from the command line as follows: R CMD install SNVsMutations.tar.gz.

library(ggplot2) library(SNVsMutations) Loading the VCF file 1.3

https://bioconductor.org/packages/release/bioc/vignettes/VariantAnnotation/inst/doc/VariantAnnotation.pdf For example: fl <- system.file("extdata", "chr22.vcf.gz", package="VariantAnnotation")</pre> vcf <- readVcf(fl, "hg38")[1:200]</pre>

In alternative, it is possible to load another vcf file by specifying the path

To load the VCF file it is possible to load the vcf example file from the VariabtAnnotation package

In both cases the vcf file can be read by using the readVcf ()` function. In this case we selected only the first 200 rows to make analysis more easy, smooth and less computationally expensive. In this example vignette we will use chr22.vcf.gz

file <- ('/Users/federicaboselli/Downloads/HG02024_VN049_KHVTrio.chr1.vcf')</pre>

function The SNV_Types() function takes a VCF file, a context length, and a reference genome as input. It processes only the single nucleotide variants (SNVs) present in the VCF file and returns a vector of

Regarding the function inputs: - the vcf file is the one previously loaded by the user (see 'loading the

Getting the mutation types from the vcf file through the SNV_Types()

vcf file') - the context_length parameter should be a positive odd integer. It specifies the number of bases to include both upstream and downstream of each single nucleotide variant (SNV) position in the mutation type string. - the reference genome should be in concordance with the used vcf file and can be loaded by using the BSgenome library. In particular, the BSgenome library should correspond

For example, to load the human reference genome hg38 ref_genome <- Hsapiens</pre> Overall, a possible example on how to use the function can be this one: mut_type <- SNV_Types(vcf, context_length = 3, ref_genome)</pre> mut_type

[1] "C[T>C]C" "C[C>T]G" "C[C>T]C" "T[C>T]G" "C[C>T]G" "G[C>T]G" "G[T>C]G"

[29] "T[C>T]A" "T[T>C]G" "G[C>T]T" "C[C>T]T" "C[C>T]C" "A[C>A]G" "G[C>T]G" [36] "G[C>T]C" "T[C>G]A" "G[C>T]A" "A[C>G]A" "A[T>C]A" "C[T>C]C" "A[C>T]A"

[8] "G[C>T]G" "A[C>T]C" "A[T>C]G" "A[C>T]A" "C[C>T]A" "G[T>C]C" "A[C>T]C" [15] "T[C>T]C" "C[C>A]C" "T[C>T]T" "T[C>T]T" "G[C>T]T" "T[C>T]C" "C[C>T]A" ## [22] "C[C>T]G" "G[T>G]A" "G[C>T]T" "C[T>C]G" "G[C>T]G" "C[C>T]C" "T[C>T]C"

[43] "G[C>T]G" "C[C>T]G" "A[C>T]C" "G[T>C]C" "G[C>T]T" "G[C>T]C" "C[T>A]A" [50] "A[C>T]C" "A[C>T]C" "G[T>C]A" "T[C>T]A" "C[C>T]C" "C[C>T]C" "T[C>G]G" ## [57] "C[T>C]G" "C[C>T]T" "A[C>T]A" "C[C>T]C" "C[T>C]C" "T[C>G]T" "A[C>T]G" [64] "A[T>G]G" "T[C>T]C" "T[C>T]A" "A[C>T]C" "G[T>G]G" "C[T>C]G" "C[C>G]A" ## [71] "C[T>C]T" "C[C>T]T" "G[C>T]G" "G[C>T]C" "G[C>T]T" "G[C>T]A" "C[C>T]A" [78] "T[T>A]C" "T[C>T]A" "C[C>T]A" "G[C>T]T" "G[C>T]G" "G[T>C]T" "A[C>T]A" ## [85] "C[T>C]A" "C[C>T]G" "T[C>T]T" "C[C>T]A" "G[C>T]T" "T[C>T]A" "A[C>T]G" ## [92] "A[T>G]C" "A[C>T]C" "C[C>T]G" "C[C>T]C" "G[C>A]A" "G[C>G]T" "A[T>C]G" [99] "G[C>T]C" "C[C>T]A" "G[C>T]T" "G[C>T]A" "T[C>G]C" "T[C>T]T" "C[C>T]A"

[106] "G[C>T]T" "C[T>A]T" "G[C>A]G" "C[T>G]G" "C[C>G]C" "C[C>T]C" "T[T>C]G"

[120] "T[C>T]G" "C[C>T]A" "C[C>T]G" "A[C>T]A" "C[C>T]C" "C[C>T]G" "T[C>T]A" ## [127] "C[T>C]T" "C[C>T]T" "G[C>G]A" "G[C>T]G" "G[C>A]T" "A[C>T]C" "C[C>T]G" ## [141] "C[C>T]C" "C[C>G]G" "G[C>T]T" "C[C>T]G" "G[C>T]C" "C[T>C]G" "T[C>A]C" ## [148] "A[C>G]G" "G[T>A]G" "G[C>T]G" "G[C>T]G" "G[C>T]G" "A[T>C]A" "G[T>C]A" ## [155] "C[C>T]G" "G[C>T]G" "T[C>T]T" "A[C>G]C" "C[C>T]C" "C[C>T]C" "T[C>G]G" ## [162] "C[C>T]T" "G[T>G]C" "C[C>G]C" "C[C>T]T" "G[T>G]C" "G[C>T]C" "C[C>T]T" ## [169] "A[C>A]C" "G[C>T]A" "C[C>T]A" "A[C>T]A" "C[T>G]T" "G[C>T]C" "T[C>T]A" ## [176] "G[C>G]C" "G[C>T]A" "G[C>T]T" "T[T>A]C" "A[C>T]A" "C[C>T]C" "A[T>C]C" ## [183] "T[C>T]A" "C[C>T]G" "G[T>A]T" "C[T>A]C" "C[T>A]C" "T[C>T]T" "G[C>G]C" ## [190] "C[T>C]G" "T[C>T]C" "G[C>T]T" "T[C>T]C" "T[C>T]C" "A[C>T]A" NB: The mutation "C[G>A]A" corresponds to "T[C>T]G" when considering the reverse strand, indicating a redundancy that can be resolved by converting the mutation types identified for REF bases A and G to their corresponding reverse complements. Therefore, in this function all mutation types are reported with either C or T as the mutated REF base.

Generating the mutation counts and a final graphical output with the

The function SNV_Counts() t takes either a vector of mutation types (mutationstype_vector) or a VCF file, reference genome, and context length as inputs. It counts the occurrences of different mutation types and generates a plot to visualize the counts. If the onlySNV parameter is set to TRUE,

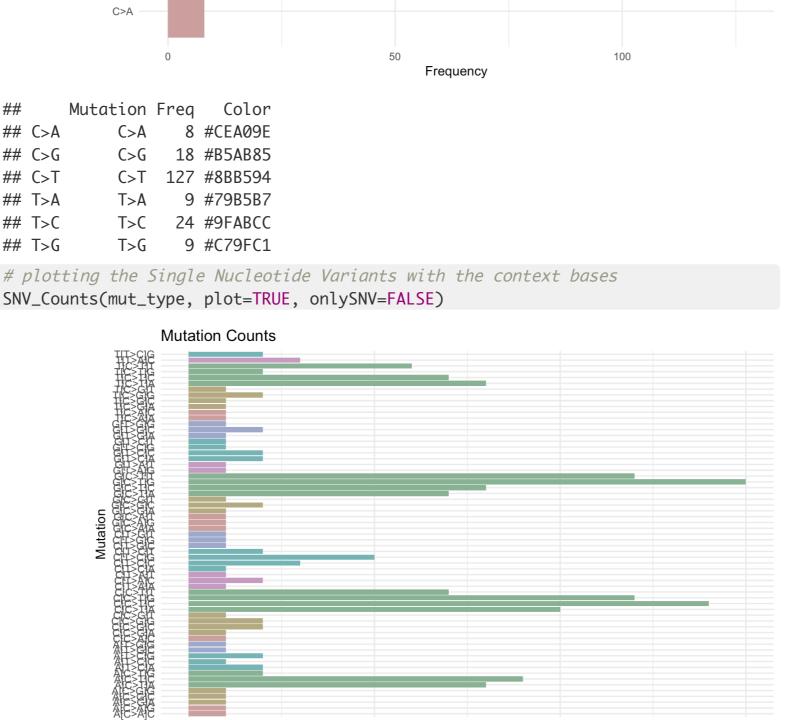
If a mutationstype_vector is provided, it is used directly to count the mutation types. Otherwise, the mutation types are computed using the MutationsType function based on the vcf_file, reference_genome, and context_length parameters. The function also creates a bar plot showing the

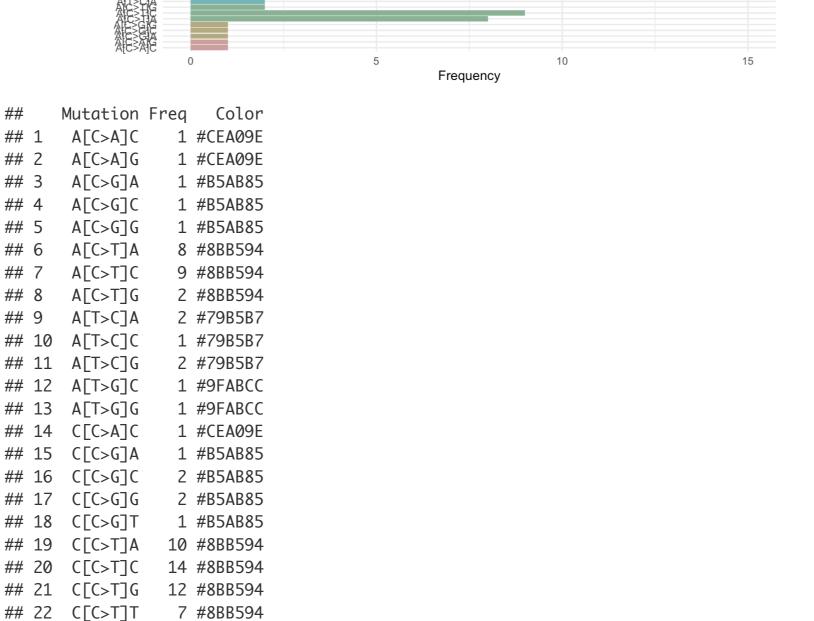
mutation counts with different colors for each mutation type. The plot is displayed if plot is set to TRUE. In particular: If onlySNV is TRUE, the plot shows the mutation types by considering only reference and alternative allele (SNV only). If onlySNV is FALSE, the plot shows the mutation types by also considering the bases upstream and downstream (taking into account the 'context_length'

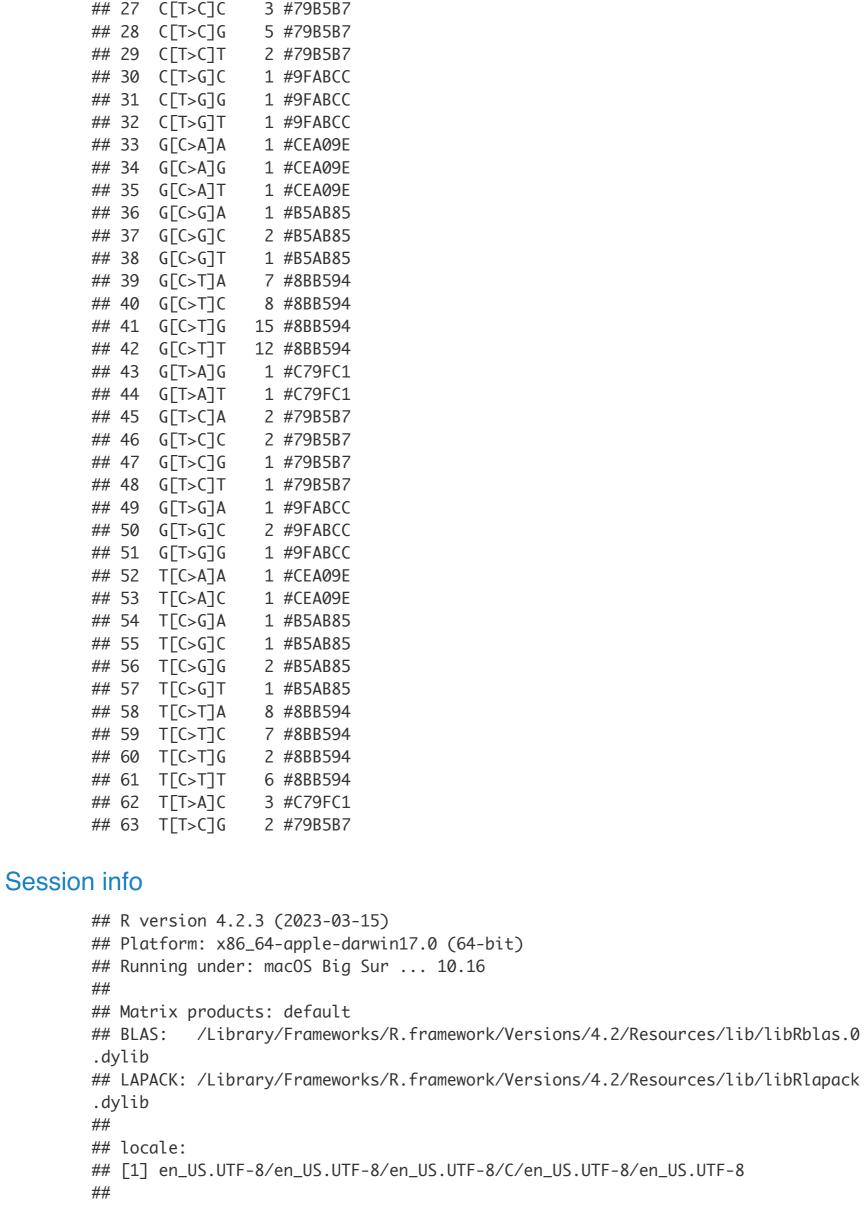
it specifically counts single-nucleotide variant (SNV) mutation types.

variable) of each SNV. Example usage of the SNV_Counts() function using the mut_type object previously generated through the SNV_Types() function: # plotting only the Single Nucleotide Variants

SNV_Counts(mut_type, plot=TRUE, onlySNV=TRUE) **Mutation Counts** T>C







23 C[T>A]A

24 C[T>A]C

25 C[T>A]T

26 C[T>C]A

1 #C79FC1

2 #C79FC1

1 #C79FC1

1 #79B5B7

[15] GenomeInfoDb_1.34.9 IRanges_2.32.0 ## [17] S4Vectors_0.36.2 BiocGenerics_0.44.0 ## [19] knitr_1.43

attached base packages:

other attached packages: ## [1] SNVsMutations_0.1.0

[3] VariantAnnotation_1.44.1

[7] MatrixGenerics_1.10.0

[61] GenomicFeatures_1.50.4

[76] xfun_0.39

[11] rtracklayer_1.58.0

[13] XVector_0.38.0

[5] SummarizedExperiment_1.28.0

stats

[1] stats4

[8] base

BiocStyle_2.26.0 ## loaded via a namespace (and not attached): ## [1] bitops_1.0-7 bit64_4.0.5 filelock_1.0.2 ## [4] progress_1.2.2 httr_1.4.6 tools_4.2.3

[9] BSgenome.Hsapiens.UCSC.hg38_1.4.5 BSgenome_1.66.3

graphics grDevices utils

ggplot2_3.4.2

Biobase_2.58.0

Rsamtools_2.14.0

matrixStats_0.63.0

Biostrings_2.66.0

GenomicRanges_1.50.2

datasets methods

KEGGREST_1.38.0

tibble_3.2.1

memoise_2.0.1

[7] bslib_0.4.2 utf8_1.2.3 R6_2.5.1 ## [10] DBI_1.1.3 colorspace_2.1-0 withr_2.5.0 ## [13] tidyselect_1.2.0 prettyunits_1.1.1 bit_4.0.5 ## [16] curl_5.0.0 compiler_4.2.3 cli_3.6.1

[19] xml2_1.3.4 DelayedArray_0.24.0 labeling_0.4.2 ## [22] bookdown_0.34 sass_0.4.6 scales_1.2.1 ## [25] rappdirs_0.3.3 stringr_1.5.0 digest_0.6.31 ## [28] rmarkdown_2.21 pkgconfig_2.0.3 htmltools_0.5.5

[31] highr_0.10 dbplyr_2.3.2 fastmap_1.1.1 ## [34] rlang_1.1.1 rstudioapi_0.14 RSQLite_2.3.1 ## [37] jquerylib_0.1.4 BiocIO_1.8.0 generics_0.1.3 ## [40] farver_2.1.1 jsonlite_1.8.4 BiocParallel_1.32.6 magrittr_2.0.3 ## [43] dplyr_1.1.2 RCurl_1.98-1.12 ## [46] GenomeInfoDbData_1.2.9 Matrix_1.5-4.1 Rcpp_1.0.10

[49] munsell_0.5.0 fansi_1.0.4 lifecycle_1.0.3 ## [52] stringi_1.7.12 $yaml_2.3.7$ zlibbioc_1.44.0 ## [55] BiocFileCache_2.6.1 grid_4.2.3 blob_1.2.4 ## [58] parallel_4.2.3 lattice_0.21-8 crayon_1.5.2

hms_1.1.3

[64] pillar_1.9.0 codetools_0.2-19 rjson_0.2.21 ## [67] biomaRt_2.54.1 XML_3.99-0.14 glue_1.6.2 ## [70] evaluate_0.21 BiocManager_1.30.20 png_0.1-8 ## [73] vctrs_0.6.2 gtable_0.3.3 cachem_1.0.8

restfulr_0.0.15

[79] GenomicAlignments_1.34.1 AnnotationDbi_1.60.2