

SigProfilerSimulator

Sept 5, 2018

INTRODUCTION

- simulates mutational signatures in cancer to provide a background model for statistical analysis.
- unbiased fashion
- context description
- exome/BED files
- supports GRCh37, GRCh38, mm10
 - See simulating other genomes for info on how to do so

PREREQUISITES

The framework is written in PYTHON, however, it also requires the following additional software with the given versions (or newer) and access to BASH:

| | |
|----------------------------|---|
| PYTHON | version 3.4 or newer |
| FASTRAND | Python module: https://github.com/lemire/fastrand/blob/master/README.md |
| WGET | version 1.19 |
| SigProfilerMatrixGenerator | current version from GitHub: https://github.com/AlexandrovLab/SigProfilerMatrixGenerator |

While the code was developed for application on a local computer, access to a cluster with greater computational power may be required for simulating a large number of mutations.

QUICK START GUIDE

This section will guide you through the minimum steps required to begin simulating mutations:

1. Create a folder with a unique name for your job (ex: BRCA, ncRCC, etc.) within the *vcf_files* folder.
2. Place your vcf/maf/simple text files into this newly created *project* folder. If you are dealing with vcf files, ensure that you have an individual file for each sample of interest.
3. From a command prompt run the mutation_simulation.py script as follows:
`python mutation_simulation.py -g GRCh37 -p BRCA -c 96 INDEL -s 100 **`
4. The script will begin generating all of the required reference files (First time users only).
5. Simulated samples are saved within the *output* folder.

****NOTE:** See commands below for list of all available parameters.

COMMANDS

| | |
|---------------------|--|
| -g or --genome | -> required: Followed by the reference genome (ex: GRCh37, GRCh38, mm10). |
| -p or --project | -> required: Followed by a unique project name (ex: BRCA) |
| -c or --context | -> required: Followed by the desired contexts to simulate. Only one type of single point mutations may be provided (96, 192, 1536, or 3072) along with DINUCs or INDELs (ex: -c 192 DINUC INDEL). |
| -e or --exome | -> optional: Simulates based solely on the exome regions of the genome. |
| -s or --simulations | -> optional: Followed by the desired number of iterated simulations per sample. The default number of iterations is 1 (ex: -s 100 will produce 100 iterations per sample). |
| -u or --update | -> optional: Updates the chromosomes as each mutation is assigned. |
| -b or --bed | -> optional: Followed by the file name of the BED file. Allows for simulations within a set of ranges. (ex: -b regions_for_sim.txt) |
| -i or --indel | -> optional: Simulates INDELs without considering insertions at microhomologies. |
| -S or --Signatures | -> optional: Simulates based upon a matrix that consists of a set of mutational signatures with their respective activities for each sample. This parameter requires that a file is saved within the <i>references/matrix/</i> folder with the following extension: “.mutSignatures” |

FOLDER STRUCTURE

The framework contains three folders (scripts, references, output), the main simulator code, and this readme file [include pictures snapshot]. The scripts folder contains all code required to generate the necessary input files, reference files, and mutational matrix catalogues. The references file contains a *vcf_files* folder, *matrix* folder, and *chromosomes* folder. The *vcf_files* folder is where the user must place their input files (must be in vcf, maf, or simple text file format). The *matrix* folder will contain all of the generated mutational matrix catalogues. The user may upload their own matrix as long as it follows the format given in the *examples* folder, otherwise the simulation code will generate the necessary matrices. The *chromosomes* folder contains a *chrom_string* folder, *context_distributions* folder, *exome* folder, *fasta* folder, *transcripts* folder, and *tsb* folder. NOTE: ALL OF THESE FILES MUST BE PLACED UNDER THE GENOME NAME WITHIN EACH SUBDIRECTORY. The *output* folder is where all simulated data will be stored after completion.

INPUT FILE FORMAT

This tool currently supports *maf*, *vcf*, and *simple text* file formats (See the *examples* folder for examples of each format). The user must provide variant data adhering

to one of these three formats. If the user's files are in *vcf* format, each sample must be saved as a separate file.

DESCRIPTION OF PROVIDED EXAMPLES

An *example_test* folder is provided that contains 10 test samples. This folder can be placed within the *references/vcf_files/* folder and run on the command line as follows:

```
python mutation_simulation.py -g GRCh37 -p example_test -c 192 DINUC -s 10
```

****The SNV context can be changed to 96, 192, 1536, or 3072, and the simulation parameter can be changed to any value.**

The final simulated output will be saved within the *simulation_output/example_test_simulations_GRCh37_[context]/* folder.

SIMULATING ADDITIONAL GENOMES

If the user desires to use a genome other than those currently supported (GRCh37, GRCh38, or mm10), they must:

- 1)Download the individual FASTA files for each chromosome
- 2)Place them into the *references/chromosomes/fasta* folder
- 3)Download the transcriptional data following the format presented in the *transcripts_original* folder.
- 4)Place this file within the *transcripts_original* folder
- 5)Download the exome ranges for the given genome and place in the *references/chromosomes/exome* folder.

ADDITIONAL DATA AND EXAMPLES

COPYRIGHT

This software and its documentation are copyright 2018 as a part of the sigProfiler project. The sigProfilerMatrixGenerator framework is free software and is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details [change to whatever group we should include.

CONTACT INFORMATION

Please address any queries or bug reports to Erik Bergstrom at ebergstr@eng.ucsd.edu