

SomaticSeq Documentation

January 8, 2016

1 Introduction

SomaticSeq is a flexible workflow that has incorporated multiple somatic mutation callers to obtain a combined call set, and then it uses machine learning to distinguish true mutations from false positives from the call set. We have incorporated six somatic mutation caller: MuTect, VarScan2, JointSNVMix, SomaticSniper, VarDict, and MuSE. You may incorporate some or all of those callers into your own pipeline with SomaticSeq.

The manuscript, An ensemble approach to accurately detect somatic mutations using SomaticSeq, is published in Genome Biology 2015, 16:197. The SomaticSeq project is located at <http://bioinform.github.io/somaticseq/>. The data described in the manuscript is also described at <http://bioinform.github.io/somaticseq/data.html>.

SomaticSeq.Wrapper.sh is a bash script that calls a series of scripts to combine the output of the somatic mutation caller(s), after the somatic mutation callers are run. Then, depending on what files are fed to SomaticSeq.Wrapper.sh, it will either train the call set into a classifier, predict high-confidence somatic mutations from the call set, or do nothing.

1.1 Dependencies

- All the .py scripts are written in Python 3 (Yes, Python 3. I deserve brownie points from Python developers). Also required are regex, pysam, numpy, and scipy libraries for Python 3.
- R, plus the ada package in R.
- SnpEff/SnpSift, GATK (for CombineVariants), and BEDTools
- dbSNP in VCF format. COSMIC may be needed in the future.
- At least one of MuTect/Indelocator, VarScan2, JointSNVMix, SomaticSniper, VarDict, and MuSE. Those are the tools we have incorporated in SomaticSeq. If there are other somatic tools that may be good addition to our list, please make the suggestion to us.

2 To use SomaticSeq.Wrapper.sh

The SomaticSeq.Wrapper.sh is a wrapper script that calls a series of programs and procedures. It can be run easily if all the dependencies are met. Some parameters are hard-coded but easily edited. In the next section, we described the workflow in more detail, so you are not dependent on this wrapper script. You can either modify this wrapper script or create your own workflow.

2.1 To train data set into a classifier

To create a trained classifier, ground truth files are required for the data sets. There is also an option to include a list of regions to ignore, where the ground truth is not known in those regions.

```
1 # -M/-I/-V/-v/-J/-S/-D/-U are output VCF files from individual callers.
2 # -i is also optional.
3 SomaticSeq.Wrapper.sh \
4 -M MuTect/variants.snp.vcf \
5 -I Indelocator/variants.indel.vcf \
6 -V VarScan2/variants.snp.vcf \
7 -v VarScan2/variants.indel.vcf \
8 -J JointSNVMix2/variants.snp.vcf \
9 -S SomaticSniper/variants.snp.vcf \
10 -D VarDict/variants.vcf \
11 -U MuSE/variants.snp.vcf \
12 -N matched_normal.bam \
13 -T tumor.bam \
14 -R ada_model_builder.R \
15 -g human_b37.fasta \
16 -c cosmic.b37.v71.vcf \
17 -d dbSNP.b37.v141.vcf \
18 -s $PATH/TO/DIR/snpSift \
19 -G $PATH/TO/GenomeAnalysisTK.jar \
20 -i ignore.bed \
21 -Z truth.snp.vcf \
22 -z truth.indel.vcf \
23 -o $OUTPUT_DIR
```

SomaticSeq.Wrapper.sh supports any combination of the somatic mutation callers we have incorporated into the workflow, so -M/-I/-V/-v/-J/-S/-D/-U are all optional parameters. SomaticSeq will run based on the output VCFs you have provided. It will train SNV and/or indel if you provide the truth.snp.vcf and/or truth.indel.vcf file(s).

2.2 To predict somatic mutation based on trained classifiers

```
1 # The *.RData files are trained classifier from the training mode.
2 SomaticSeq.Wrapper.sh \
3 -M MuTect/variants.snp.vcf \
4 -I Indelocator/variants.indel.vcf \
5 -V VarScan2/variants.snp.vcf \
6 -v VarScan2/variants.indel.vcf \
7 -J JointSNVMix2/variants.snp.vcf \
8 -S SomaticSniper/variants.snp.vcf \
9 -D VarDict/variants.vcf \
10 -U MuSE/variants.snp.vcf \
11 -N matched_normal.bam \
12 -T tumor.bam \
13 -R ada_model_predictor.R \
14 -C sSNV.Classifier.RData \
15 -x sINDEL.Classifier.RData \
16 -g human_b37.fasta \
17 -c cosmic.b37.v71.vcf \
18 -d dbSNP.b37.v141.vcf \
19 -s $PATH/TO/DIR/snpSift \
20 -G $PATH/TO/GenomeAnalysisTK.jar \
21 -o $OUTPUT_DIR
```

2.3 To simply create the tab delimited file without doing any machine learning

Same as the command previously, but not including -R, -C, or -x.

3 The step-by-step SomaticSeq Workflow

We'll describe the workflow here, so you may modify the workflow and/or create your own workflow instead of using the wrapper script described previously.

3.1 Combine the call sets

We use GATK CombineVariants to combine the VCF files from different callers, although it does not matter what tools are used to merge VCF files. To make them compatible with GATK, the VCF files are modified. A somatic call is also tagged with the tool names, so the combined VCF retains those information.

1. Modify MuTect and/or Indelocator output VCF files. Somatic calls will be attached the tag 'CGA' in the INFO. Since MuTect's output VCF do not always put the tumor and normal samples in the same columns, the script uses samtools extract sample name information from the BAM files, and then determine which column belongs to the normal, and which column belongs to the tumor.

```
1 # Modify MuTect and Indelocator's output VCF
2 modify_MuTect.py -infile input.vcf -outfile output.vcf -nbam normal.bam -tbam
   tumor.bam
3
4 # If samtools is not in the PATH:
5 modify_MuTect.py -infile input.vcf -outfile output.vcf -nbam normal.bam -tbam
   tumor.bam -samtools $PATH/TO/samtools
```

Alternatively, you can supply the normal and tumor sample names, instead of supplying the BAM files:

```
1 # Modify MuTect's output VCF
2 # -type snp for MuTect, and -type indel for Indelocator.
3 modify_MuTect.py -type snp -infile input.vcf -outfile output.vcf -nsm
   NormalSampleName -tsm TumorSampleName
```

2. Modify VarScan's output VCF files to be rigorously concordant to VCF format standard, and to attach the tag 'VarScan2' to somatic calls.

```
1 # Do it for both the SNV and indel
2 modify_VJSD.py -method VarScan2 -infile input.vcf -outfile output.vcf
```

3. JointSNVMix2 does not output VCF files. In our own workflow, we have already converted its text file into a basic VCF file with an awk one-liner, which you may see in the Run_5_callers directory, which are:

```

# To avoid text files in the order of terabytes, this awk one-liner keeps
# entries where the reference is not "N", and the somatic probabilities are
# at least 0.95.
2 awk -F "\t" 'NR!=1 && $4!="N" && $10+$11>=0.95'

4 # This awk one-liner converts the text file into a basic VCF file
awk -F "\t" '{print $1 "\t" $2 "\t.\t" $3 "\t" $4 "\t.\t.\tAAAB=" $10 ";AABB="
    $11 "\tRD:AD\t" $5 ":" $6 "\t" $7 ":" $8}'

6

8 ## The actual commands we've used in our workflow:
echo -e '##fileformat=VCFv4.1' > unsorted.vcf
10 echo -e '##INFO=<ID=AAAB,Number=1,Type=Float,Description="Probability of Joint
    Genotype AA in Normal and AB in Tumor">' >> unsorted.vcf
echo -e '##INFO=<ID=AABB,Number=1,Type=Float,Description="Probability of Joint
    Genotype AA in Normal and BB in Tumor">' >> unsorted.vcf
12 echo -e '##FORMAT=<ID=RD,Number=1,Type=Integer,Description="Depth of reference
    -supporting bases (reads1)">' >> unsorted.vcf
echo -e '##FORMAT=<ID=AD,Number=1,Type=Integer,Description="Depth of variant-
    supporting bases (reads2)">' >> unsorted.vcf
14 echo -e '#CHROM\tPOS\tID\tREF\tALT\tQUAL\tFILTER\tINFO\tFORMAT\tNORMAL\tTUMOR'
    >> unsorted.vcf

16 python $PATH/TO/jsm.py classify joint-snv-mix-two genome.GRCh37.fa normal.bam
    tumor.bam trained.parameter.cfg /dev/stdout | awk -F "\t" 'NR!=1 && $4!="N"
    && $10+$11>0.95' | awk -F "\t" '{print $1 "\t" $2 "\t.\t" $3 "\t" $4 "\t.\t.\tAAAB=" $10 ";AABB=" $11 "\tRD:AD\t" $5 ":" $6 "\t" $7 ":" $8}' >>
    unsorted.vcf

```

After that, you'll also want to sort the VCF file. Now, to modify that basic VCF into something that will be compatible with other VCF files under GATK CombineVariants:

```
modify_VJSD.py -method JointSNVMix2 -infile input.vcf -outfile output.vcf
```

4. Modify SomaticSniper's output:

```
modify_VJSD.py -method SomaticSniper -infile input.vcf -outfile output.vcf
```

- VarDict has both SNV and indel, plus some other variants in the same VCF file. Our script will create two files, one for SNV and one for indel, while everything else is ignored for now. By default, LikelySomatic/StrongSomatic and PASS calls will be labeled VarDict. However, in our SomaticSeq paper, based on our experience in DREAM Challenge, we implemented two custom filters to relax the VarDict tagging criteria.

```

1 # Default VarDict tagging criteria, only PASS (and Likely or Strong Somatic):
  modify_VJSD.py -method VarDict -infile input.vcf -outfile output.vcf

3
4 # When running VarDict, if var2vcf-paired.pl is used to generate the VCF file,
  # you may relax the tagging criteria with -filter paired
5 modify_VJSD.py -method VarDict -infile input.vcf -outfile output.vcf -filter
  paired

7 # When running VarDict, if var2vcf-somatic.pl is used to generate the VCF file
  # you may relax the tagging criteria with -filter somatic

```

```
modify_VJSD.py --method VarDict --infile input.vcf --outfile output.vcf --filter
somatic
```

In the SomaticSeq paper, -filter somatic was used because var2vcf_somatic.pl was used to generate VarDict's VCF files. In the SomaticSeq.Wrapper.sh script, however, -filter paired is used because VarDict authors have since recommended var2vcf_paired.pl script to create the VCF files. While there are some differences (different stringencies in some filters) in what VarDict labels as PASS between the somatic.pl and paired.pl scripts, the difference is miniscule after applying our custom filter (which relaxes the filter, resulting in a difference about 5 calls out of 15,000).

The output files will be snp.output.vcf and indel.output.vcf.

6. MuSE was not a part of our analysis in the SomaticSeq paper. We have implemented it later.

```
modify_VJSD.py --method MuSE --infile input.vcf --outfile output.vcf
```

7. Finally, with the VCF files modified, you may combine them with GATK CombineVariants: one for SNV and one for indel separately. There is no particular reason to use GATK CombineVariants. Other combiners should also work. The only useful thing here is to combine the calls, and preserve the tags we have written into each individual VCF file's INFO.

```
1 # Combine the VCF files for SNV. Any or all of the VCF files may be present.
2 # -nt 12 means to use 12 threads in parallel
3 java -jar $PATH/TO/GenomeAnalysisTK.jar -T CombineVariants -R genome.GRCh37.fa
   -nt 12 --setKey null --genotypemergeoption UNSORTED -V mutect.vcf -V
   varscan.snp.vcf -V jointsvnmix.vcf -V snp vardict.vcf -V muse.vcf --out
   CombineVariants.snp.vcf
4 java -jar $PATH/TO/GenomeAnalysisTK.jar -T CombineVariants -R genome.GRCh37.fa
   -nt 12 --setKey null --genotypemergeoption UNSORTED -V indelocator.vcf -V
   varscan.snp.vcf -V indel vardict.vcf --out CombineVariants.indel.vcf
```

8. Use SnpSift to add dbSNP information to the VCF file, since dbSNP information is part of training feature set. Right now, we do not use COSMIC or functional annotation as a part of training feature, but we do have them in the workflow for "future-proofing." We may decide to use those features in the future when we have better data sets for training.

```
1 # SNV
2 java -jar $PATH/TO/SnpSift.jar annotate dbSNP141.vcf CombineVariants.snp.vcf |
   java -jar $PATH/TO/SnpSift.jar annotate COSMIC.vcf - | java -jar $PATH/TO/
   /snpEff.jar GRCh37.75 - > EFF.COSMIC.dbSNP.CombineVariants.snp.vcf
3 # INDEL
4 java -jar $PATH/TO/SnpSift.jar annotate dbSNP141.vcf CombineVariants.indel.vcf
   | java -jar $PATH/TO/SnpSift.jar annotate COSMIC.vcf - | java -jar $PATH/
   TO/snpEff.jar GRCh37.75 - > EFF.COSMIC.dbSNP.CombineVariants.indel.vcf
```

9. This procedure annotates the caller consensus by putting the tool names into the SOURCES in the INFO. You should also use -mincaller 1 to only keep calls where at least one caller has called it somatic. Vast majority of the calls in the previous merged VCF files were REJECT or GERMLINE calls, and may run out of memory in model training if all are included. It also does some rudimentary variant scoring, but the scoring is not utilized in SomaticSeq. You need to run it once for SNV and once for indel.

```

1 # Use -tools to indicate what call sets were combined. CGA=MuTect/Indelocator
2 # In this one, 6 tools were used for SNV
3 score_Somatic.Variants.py -tools CGA VarScan2 JointSNVMix2 SomaticSniper
   VarDict MuSE -infile EFF.COSMIC.dbSNP.CombineVariants.snp.vcf -mincaller 1
   -outfile BINA_somatic.snp.vcf
4
5 # 3 tools for indel
6 score_Somatic.Variants.py -tools CGA VarScan2 VarDict -infile EFF.COSMIC.dbSNP
   .CombineVariants.indel.vcf -mincaller 1 -outfile BINA_somatic.indel.vcf

```

And now, the call sets are combined into one VCF file for SNV and one VCF file for indel.

3.2 For model training: process and annotate the VCF files (union of call sets)

This step may be needed for model training. The workflow in SomaticSeq.Wrapper.sh allows for an exclusion region, e.g., we don't care for anything inside the exclusion region, typically because we don't know if a call made in this region is a false positive or true positive. DREAM Challenge had exclusion regions, e.g., blacklisted regions, etc. Alternatively, you can use an inclusion region instead, to obtain a list of variants that you have done experimental validation, so you know which ones are true mutations and which ones are false positives, to facilitate model training.

```

1 # In the DREAM_Stage_3 directory, we have included an exclusion region BED file as
   an example
2 # This command uses BEDtools to rid of all calls in the exclusion region
3 intersectBed -header -a BINA_somatic.snp.vcf -b ignore -v > somatic.snp.processed.
   vcf
4 intersectBed -header -a BINA_somatic.indel.vcf -b ignore -v > somatic.indel.
   processed.vcf

```

3.3 Convert the VCF file, annotated or otherwise, into a tab separated file

This script works for all VCF files. It extracts information from BAM files as well as VCF files created by the individual callers. If the ground truth VCF file is included, a called variant will be annotated as a true positive, and everything will be annotated as a false positive.

```

1 # -mutect / -sniper / -varscan / -jsm / -vardict / -muse are optional.
2 # -truth is also optional, but it is needed to annotate ground truth information,
   and is required for model training.
3 SSeq-merged.vcf2tsv.py -ref genome.GRCh37.fa -myvcf somatic.snp.processed.vcf -
   truth Ground.truth.snp.vcf -mutect MuTect/variants.snp.vcf.gz -varscan VarScan2
   /variants.snp.vcf -jsm JSM2/variants.vcf -sniper SomaticSniper/variants.vcf -
   vardict VarDict/snp.variants.vcf -tbam tumor.bam -nbam normal.bam -outfile
   Ensemble.sSNV.tsv

```

That was for SNV, and indel is almost the same thing. After version 2.1, we have replaced all information from SAMtools and HaplotypeCaller with information directly from the BAM files. The accuracy differences are negligible with significant improvement in usability and resource requirement.

```
1 # INDEL:
   SSeq_merged.vcf2tsv.py -ref genome.GRCh37.fa -myvcf somatic.indel.processed.vcf -
   truth Ground.truth.indel.vcf -varscan VarScan2/variants.snp.vcf -vardict
   VarDict/indel.variants.vcf -tbam tumor.bam -nbam normal.bam -outfile Ensemble.
   sINDEL.tsv
```

At the end of this, Ensemble.sSNV.tsv and Ensemble.sINDEL.tsv are created.

3.4 Model Training or Mutation Prediction

You can use Ensemble.sSNV.tsv and Ensemble.sINDEL.tsv files either for model training (provided that their mutation status is annotated with 0 or 1) or mutation prediction. This is done with stochastic boosting algorithm we have implemented in R.

Model training:

```
2 # Training:
   R --no-save --args Ensemble.sSNV.tsv < ada_model_builder.R
   R --no-save --args Ensemble.sINDEL.tsv < ada_model_builder.R
```

Ensemble.sSNV.tsv.Classifier.RData and Ensemble.sINDEL.tsv.Classifier.RData will be created from model training.

Mutation prediction:

```
1 # Mutation prediction:
   R --no-save --args Ensemble.sSNV.tsv.Classifier.RData Ensemble.sSNV.tsv Trained.
   sSNV.tsv < ada_model_predictor.R
3 R --no-save --args Ensemble.sINDEL.tsv.Classifier.RData Ensemble.sINDEL.tsv Trained
   .sINDEL.tsv < ada_model_predictor.R
```

After mutation prediction, if you feel like it, you may convert Trained.sSNV.tsv and Trained.sINDEL.tsv into VCF files.

```
1 # If you want probability above 0.7 labeled PASS (-pass 0.7), and between 0.1 and
   0.7 labeled LowQual (-low 0.1):
   # To include calls below 0.1 as REJECTS (-all)# To convert probability values (
   between 0 to 1) into Phred scale in the QUAL column in the VCF file (-phred)
3 SSeq_tsv2vcf.py -tsv Trained.sSNV.tsv -vcf Trained.sSNV.vcf -pass 0.7 -low 0.1 -all
   -phred
   SSeq_tsv2vcf.py -tsv Trained.sINDEL.tsv -vcf Trained.sINDEL.vcf -pass 0.7 -low 0.1
   -all -phred
```

4 Release Notes

Make sure training and prediction use the same version.

4.1 Version 1.0

Version used to generate data in the manuscript and Stage 5 of the ICGC-TCGA DREAM Somatic Mutation Challenge, where SomaticSeq’s results were #1 for INDEL and #2 for SNV.

In the original manuscript, VarDict’s `var2vcf_somatic.pl` script was used to generate VarDict VCFs, and subsequently “-filter somatic” was used for `SSeq_merged.vcf2tsv.py`. Since then (including DREAM Challenge Stage 5), VarDict recommends `var2vcf_paired.pl` over `var2vcf_somatic.pl`, and subsequently “-filter paired” was used for `SSeq_merged.vcf2tsv.py`. The difference in SomaticSeq results, however, is pretty much negligible.

4.2 Version 1.1

Automated the `SomaticSeq.Wrapper.sh` script for both training and prediction mode. No change to any algorithm.

4.3 Version 1.2

Have implemented the following improvement, mostly for indels:

- `SSeq_merged.vcf2tsv.py` can now accept pileup files to extract read depth and DP4 (reference forward, reference reverse, alternate forward, and alternate reverse) information (mainly for indels). Previously, that information can only be extracted from SAMtools VCF. Since the SAMtools or HaplotypeCaller generated VCFs hardly contain any indel information, this option improves the indel model. The `SomaticSeq.Wrapper.sh` script is modified accordingly.
- Extract mapping quality (MQ) from VarDict output if this information cannot be found in SAMtools VCF (also mostly benefits the indel model).
- Indel length now positive for insertions and negative for deletions, instead of using the absolute value previously.

4.4 Version 2.0-rc1

- Removed dependencies for SAMtools and HaplotypeCaller during feature extraction. `SSeq_merged.vcf2tsv.py` extracts those information (plus more) directly from BAM files.
- Allow not only VCF file, but also BED file or a list of chromosome coordinate as input format for `SSeq_merged.vcf2tsv.py`, i.e., use `-mybed` or `-mypos` instead of `-myvcf`.
- Instead of a separate step to annotate ground truth, that can be done directly by `SSeq_merged.vcf2tsv.py`.
- `SSeq_merged.vcf2tsv.py` can annotate dbSNP and COSMIC information directly if BED file or a list of chromosome coordinates are used as input in lieu of an annotated VCF file.
- Consolidated feature sets, e.g., removed some redundant feature sets coming from different resources.

5 To do: planned improvement

- Validate on deep sequencing and low VAF data.
- Dockerize SomaticSeq

6 Known issues

- Some older versions of GATK seem to have problem with FIFO used in the SomaticSeq script (i.e., we had problem with GATK version 2014.1-2.8.1), so we recommend at least version version 2014.4-2 or later.

7 Contact Us

For suggestions, bug reports, or technical support, please email li_tai.fang@bina.roche.com.