

SCclust T10 Tutorial

Alex Krasnitz, Jude Kendall, Junyan Song, Lubomir Chorbadjiev

2018-06-06

Contents

1	Introduction	1
2	Data	1
2.1	Data for the T10 case	1
2.2	Collect the Necessary Data	2
2.3	Explore the Downloaded Data	2
3	Segmentation of Varbin Data	4

1 Introduction

The *SCclust* package implements feature selection based on breakpoints, permutations for FDRs for Fisher test p-values and identification of the clone structure in single cell copy number profiles.

In this tutorial we show how to use *SCclust* package using data, prepared by *sgains* pipeline as described in [Example usage of sGAINS pipeline](#). *SCclust* package is called as the last step in processing data from *sgains* pipeline. In this tutorial we show how *SCclust* package could be used independently from *sgains* pipeline.

We assume that you have an R environment and have installed *SCclust* package as described in the `README.md`.

2 Data

2.1 Data for the T10 case

This tutorial is based on data published in: [Navin N, Kendall J, Troge J, et al. Tumor Evolution Inferred by Single Cell Sequencing. Nature. 2011;472\(7341\):90-94. doi:10.1038/nature09807](#). In particular we will use the data for polygenomic breast tumor T10 case available from SRA. Description of samples for T10 could be found in [Supplementary Table 1 | Summary of 100 Single Cells in the Polygenomic Tumor T10](#)

We are going to run *SCclust* package on prepared by *sgains* pipeline [varbin step](#). You can go through all the step in [sgains T10 tutorial](#) and prepare this data.

For the purposes of this tutorial we recommend you to download already prepared [varbin data](#) from [example data](#). Apart from `varbin` T10 data you will need the binning scheme used in the analysis, that could be found [here](#). And also we will need `cytoBand.txt` for HG19 that you can download it from UCSC Genome Browser.

2.2 Collect the Necessary Data

Let us create a directory, where to store all the data used in this tutorial:

```
mkdir T10data
cd T10data
```

and let us download and extract T10 varbin data:

```
wget -c \
  https://github.com/KrasnitzLab/SCclust/releases/download/v1.0.0RC3/navin_t10_varbin_data.tar.gz
tar zxvf navin_t10_varbin_data.tar.gz
rm navin_t10_varbin_data.tar.gz
```

Let us also download and extract the binning scheme used in preparation of varbin data:

```
wget -c \
  https://github.com/KrasnitzLab/SCclust/releases/download/v1.0.0RC3/hg19_R50_B20k_bins_boundaries.txt.gz
gunzip hg19_R50_B20k_bins_boundaries.txt.gz
```

And finally let us download the cytoBand.txt for Human reference genome *hg19*:

```
wget -c \
  http://hgdownload.cse.ucsc.edu/goldenPath/hg19/database/cytoBand.txt.gz
gunzip cytoBand.txt.gz
```

Our data directory should have following structure:

```
.
|-- T10data
|   |-- cytoBand.txt
|   |-- hg19_R50_B20k_bins_boundaries.txt
|   |-- varbin
|       |-- SRR052047.varbin.20k.txt
|       |-- SRR052148.varbin.20k.txt
|       |-- SRR053437.varbin.20k.txt
|       ...
```

2.3 Explore the Downloaded Data

We are going to use *SCclust* package so let us load it:

```
library("SCclust")
```

2.3.1 Binning schema

```
gc_df <- read.csv("T10data/hg19_R50_B20k_bins_boundaries.txt", header = T, sep='\t')
knitr::kable(head(gc_df))
```

Describe the data.

bin.chrom	bin.start	bin.start.abspos	bin.end	bin.length	mappable.positions	gc.content
chr1	0	0	859077	859077	131390	0.4357746
chr1	859077	859077	999002	139925	131390	0.6280936
chr1	999002	999002	1141973	142971	131391	0.6026537
chr1	1141973	1141973	1280121	138148	131390	0.6284347
chr1	1280121	1280121	1435418	155297	131390	0.5757548

bin.chrom	bin.start	bin.start.abspos	bin.end	bin.length	mappable.positions	gc.content
chr1	1435418	1435418	1603686	168268	131391	0.5690862

2.3.2 Cytobands and Centromeres for HG19

Describe the data.

```
cytobands <- read.csv("T10data/cytoBand.txt", header = F, sep='\t')
knitr::kable(head(cytobands))
```

V1	V2	V3	V4	V5
chr1	0	2300000	p36.33	gneg
chr1	2300000	5400000	p36.32	gpos25
chr1	5400000	7200000	p36.31	gneg
chr1	7200000	9200000	p36.23	gpos25
chr1	9200000	12700000	p36.22	gneg
chr1	12700000	16200000	p36.21	gpos50

The main reason we need `cytoBand.txt` is to get the location of centromeres. Since centromere areas contain a lot of repetitive sequences they are excluded from analysis when segmenting and clustering samples.

To find regions where centromeres are located we are using `calc_centroareas` function:

```
centroareas <- calc_centroareas(cytobands)
knitr::kable(head(centroareas, 5))
```

	chrom	from	to
33	1	120600000	128900000
393	2	83300000	102700000
508	3	87200000	98300000
556	4	48200000	52700000
604	5	46100000	58900000

So, in `centroareas` for each chromosome we have the region where the centromere is located.

2.3.3 Varbin Samples Data

Describe the data.

For each `varbin` sample

```
sample_df <- read.csv("T10data/varbin/SRR052047.varbin.20k.txt", header=T, sep='\t')
knitr::kable(head(sample_df))
```

chrom	chrompos	abspos	bincount	ratio
chr1	0	0	51	0.3327000
chr1	859077	859077	57	0.3718412
chr1	999002	999002	89	0.5805941
chr1	1141973	1141973	53	0.3457471

chrom	chrompos	abspos	bincount	ratio
chr1	1280121	1280121	99	0.6458294
chr1	1435418	1435418	63	0.4109824

```
sample_df <- read.csv("T10data/varbin/SRR052148.varbin.20k.txt", header=T, sep='\t')
knitr::kable(head(sample_df))
```

chrom	chrompos	abspos	bincount	ratio
chr1	0	0	125	0.5530061
chr1	859077	859077	69	0.3052594
chr1	999002	999002	90	0.3981644
chr1	1141973	1141973	57	0.2521708
chr1	1280121	1280121	98	0.4335568
chr1	1435418	1435418	84	0.3716201

3 Segmentation of Varbin Data

```
# centrobins <- calc_regions2bins(gc_df, centroareas)
```