

Longitudinal Analysis of Cancer Evolution with LACE

Daniele Ramazzotti¹, Fabrizio Angaroni², Davide Maspero^{2,3}, Gianluca Ascolani², Isabella Castiglioni^{4,5}, Rocco Piazza¹, Marco Antoniotti², and Alex Graudenzi⁵

¹School of Medicine and Surgery, Univ. of Milan-Bicocca, Monza, Italy.

²Dept. of Informatics, Systems and Communication, Univ. of Milan-Bicocca, Milan, Italy.

³Fondazione IRCCS Istituto Nazionale dei Tumori, Milan, Italy.

⁴Department of Physics "Giuseppe Occhialini", Univ. of Milan-Bicocca, Milan, Italy.

⁵Inst. of Molecular Bioimaging and Physiology, Consiglio Nazionale delle Ricerche (IBFM-CNR), Segrate, Milan, Italy.

April 28, 2020

Overview. LACE (Longitudinal Analysis of Cancer Evolution) is an algorithmic framework that processes single-cell somatic mutation profiles from cancer samples collected at different time points and in distinct experimental settings, in order to describe the evolutionary history of the tumor. The tool can return an high resolution picture of clones' prevalence and their variations, e.g., because of therapies.

LACE can be employed to process single-cell mutational profiles as generated by calling variants from the increasingly available scRNA-seq data, such as the ones obtained by SMARTseq2 protocol.

The output of the method is a longitudinal tree that best fits the input data, modelling both phylogenetic constraints and sc-RNAseq specific noise. Moreover, the package provides a suite of functions to visualize and explore the results.

In this vignette, we give an overview of the package by presenting its main functions.

Contents

1	Using the LACE R package	2
2	<code>sessionInfo()</code>	5

1 Using the LACE R package

We now present an example of longitudinal analysis of cancer evolution with LACE using single-cell data obtained from Rambow, Florian, et al. "Toward minimal residual disease-directed therapy in melanoma." Cell 174.4 (2018): 843-855. The data comprises point mutations for four time points: (1) before treatment, (2) 4 days treatment, (3) 28 days treatment and finally (4) 57 days treatment.

We first load the data.

```
library("LACE")
data(longitudinal_sc_variants)
names(longitudinal_sc_variants)

## [1] "T1_before_treatment" "T2_4_days_treatment" "T3_28_days_treatment"
## [4] "T4_57_days_treatment"
```

The input data (D) can be either a list as shown above or a SummarizedExperiment object. In the latter case, the object needs to include genomic data (0: mutation absent, 1: mutation detected or NA: missing information) in assay field. Following guidelines, the rows of such matrix are variants, while columns are single cells concatenated for all the time points. Furthermore, a rowData field needs to be specified and must be a data.frame with a column named "TimePoint" reporting the label of the experiment for each single cell; we highlight that the ordering of the labels in rowData must match the one of the columns in the assay field. We next show an example.

```
library("SummarizedExperiment", attach.required=FALSE)
T1 = t(longitudinal_sc_variants[[1]])
T2 = t(longitudinal_sc_variants[[2]])
T3 = t(longitudinal_sc_variants[[3]])
T4 = t(longitudinal_sc_variants[[4]])
concat_time_point = cbind(T1,T2,T3,T4)
TimePointLabes = c(rep("T1", ncol(T1)),
                   rep("T2", ncol(T2)),
                   rep("T3", ncol(T3)),
                   rep("T4", ncol(T4)))
longitudinal_SE = SummarizedExperiment(assays = concat_time_point,
                                       colData = data.frame(TimePoint = TimePointLabes))

print(longitudinal_SE)

## class: SummarizedExperiment
## dim: 6 475
## metadata(0):
## assays(1): ''
## rownames(6): ARPC2_2_218249894_C_T CCT8_21_29063389_G_A ...
```

Longitudinal Analysis of Cancer Evolution with LACE

```
## PRAME_22_22551005_T_A RPL5_1_92837514_C_G
## rowData names(0):
## colnames(475): SRR7424153 SRR7424154 ... SRR7424823 SRR7424824
## colData names(1): TimePoint
```

We setup the main parameter in order to perform the inference. First of all, as the three data points may potentially provide sequencing for an unbalanced number of cells, we weight each time point as follows $w_s = (1 - \frac{n_s}{n_T}) / (y - 1)$ in order to account for this. In the formula, e.g., the weight for time point s (w_s) is calculated based on the number of cells observed in the time point (n_s) and the total number of cells in the three time points (n_T). The denominator ($y - 1$, with y being the number of time points, i.e., 3 in our case) aims at normalizing the weights to sum to one.

```
lik_weights = c(0.2308772, 0.2554386, 0.2701754, 0.2435088)
```

The second main parameter to be defined as input is represented by the false positive and false negative error rates, i.e., alpha and beta. We can specify a different rate per time point as a list of rates. When multiple sets of rates are provided, LACE performs a grid search in order to estimate the best set of error rates.

```
alpha = list()
alpha[[1]] = c(0.02, 0.01, 0.01, 0.01)
alpha[[2]] = c(0.10, 0.05, 0.05, 0.05)
beta = list()
beta[[1]] = c(0.10, 0.05, 0.05, 0.05)
beta[[2]] = c(0.10, 0.05, 0.05, 0.05)
head(alpha)

## [[1]]
## [1] 0.02 0.01 0.01 0.01
##
## [[2]]
## [1] 0.10 0.05 0.05 0.05

head(beta)

## [[1]]
## [1] 0.10 0.05 0.05 0.05
##
## [[2]]
## [1] 0.10 0.05 0.05 0.05
```

We can now perform the inference as follows. Notice that `D` can be either the list longitudinal sc variants or the SummarizedExperiment longitudinal SE.

```
inference = LACE(D = longitudinal_sc_variants,
                 lik_w = lik_weights,
                 alpha = alpha,
                 beta = beta,
                 keep_equivalent = FALSE,
                 num_rs = 5,
                 num_iter = 10,
                 n_try_bs = 5,
```

Longitudinal Analysis of Cancer Evolution with LACE

```
num_processes = NA,
seed = 12345,
verbose = FALSE)
```

We notice that the inference resulting on the command above should be considered only as an example; the parameters `num_rs`, `num_iter` and `n_try_bs` representing the number of steps performed during the inference are downscaled to reduce execution time. We refer to the Manual for discussion on default values. We provide within the package results of inferences performed with correct parameters as `RData`.

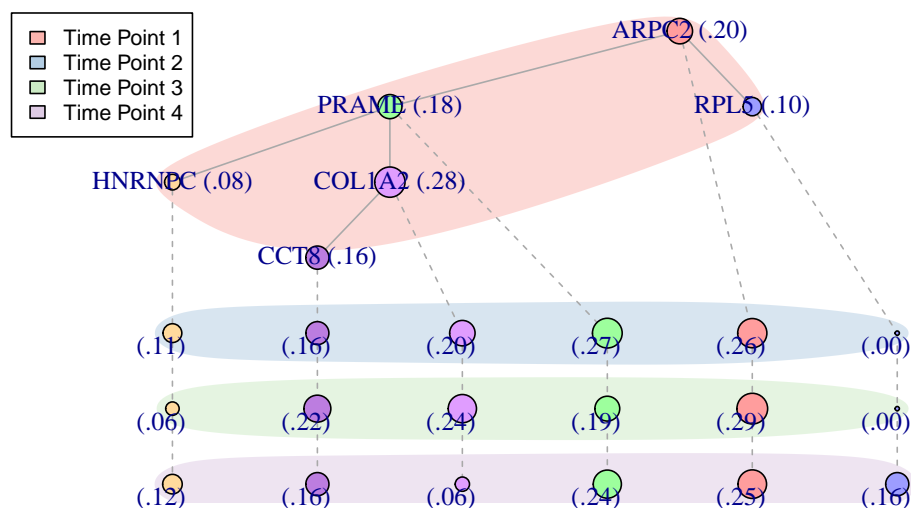
```
data(inference)
print(names(inference))

## [1] "B"                "C"                "corrected_genotypes"
## [4] "clones_prevalence" "relative_likelihoods" "joint_likelihood"
## [7] "clones_summary"   "equivalent_solutions" "error_rates"
```

LACE returns a list of nine elements as results. Namely, `B` and `C` provide respectively the maximum likelihood longitudinal tree and cells attachments; `corrected_genotypes` the corrected genotypes, `clones_prevalence`, the estimated prevalence of any observed clone; `relative_likelihoods` and `joint_likelihood` the estimated likelihoods for each time point and the weighted likelihood; `clones_summary` provide a summary of association of mutations to clones. In equivalent solutions, solutions (`B` and `C`) with likelihood equivalent to the best solution are returned; notice that in the example we disabled this feature by setting equivalent solutions parameter to `FALSE`. Finally, `error_rates` provide the best error rates (alpha and beta) as estimated by the grid search.

We can plot the inferred model using the function `longitudinal.tree.plot`.

```
clone_labels = c("ARPC2", "PRAME", "HNRNPC", "COL1A2", "RPL5", "CCT8")
longitudinal.tree = longitudinal.tree.plot(inference = inference,
                                          labels = "clones",
                                          clone_labels = clone_labels,
                                          legend_position = "topleft")
```



2 sessionInfo()

- R version 4.0.0 (2020-04-24), x86_64-w64-mingw32
- Locale: LC_COLLATE=C, LC_CTYPE=English_United States.1252, LC_MONETARY=English_United States.1252, LC_NUMERIC=C, LC_TIME=English_United States.1252
- Running under: Windows Server 2012 R2 x64 (build 9600)
- Matrix products: default
- Base packages: base, datasets, grDevices, graphics, methods, stats, utils
- Other packages: LACE 1.0.0, SummarizedExperiment 1.18.0, knitr 1.28
- Loaded via a namespace (and not attached): Biobase 2.48.0, BiocGenerics 0.34.0, BiocManager 1.30.10, BiocStyle 2.16.0, DelayedArray 0.14.0, GenomeInfoDb 1.24.0, GenomeInfoDbData 1.2.3, GenomicRanges 1.40.0, IRanges 2.22.0, Matrix 1.2-18, RColorBrewer 1.1-2, RCurl 1.98-1.2, Rcpp 1.0.4.6, RcppZiggurat 0.1.5, Rfast 1.9.9, S4Vectors 0.26.0, XVector 0.28.0, bitops 1.0-6, compiler 4.0.0, digest 0.6.25, evaluate 0.14, grid 4.0.0, highr 0.8, htmltools 0.4.0, igraph 1.2.5, lattice 0.20-41, magrittr 1.5, matrixStats 0.56.0, parallel 4.0.0, pkgconfig 2.0.3, rlang 0.4.5, rmarkdown 2.1, stats4 4.0.0, stringi 1.4.6, stringr 1.4.0, tools 4.0.0, xfun 0.13, yaml 2.2.1, zlibbioc 1.34.0