

SLIDE

SLIDE Pipeline

The SLIDE pipeline is consisted of below two steps. We recommend to run both steps on a computational cluster for optimal computational time.

1. optimizeSLIDE: Calculate and select latent factors (LFs) for multiple input parameter combinations.
2. Reviewing the output of optimizeSLIDE and choose the optimal parameters (delta, lambda and f_size) for rigorous k-fold CV.
3. SLIDEcv: With chosen parameters, perform rigorous k-fold cross-validation.

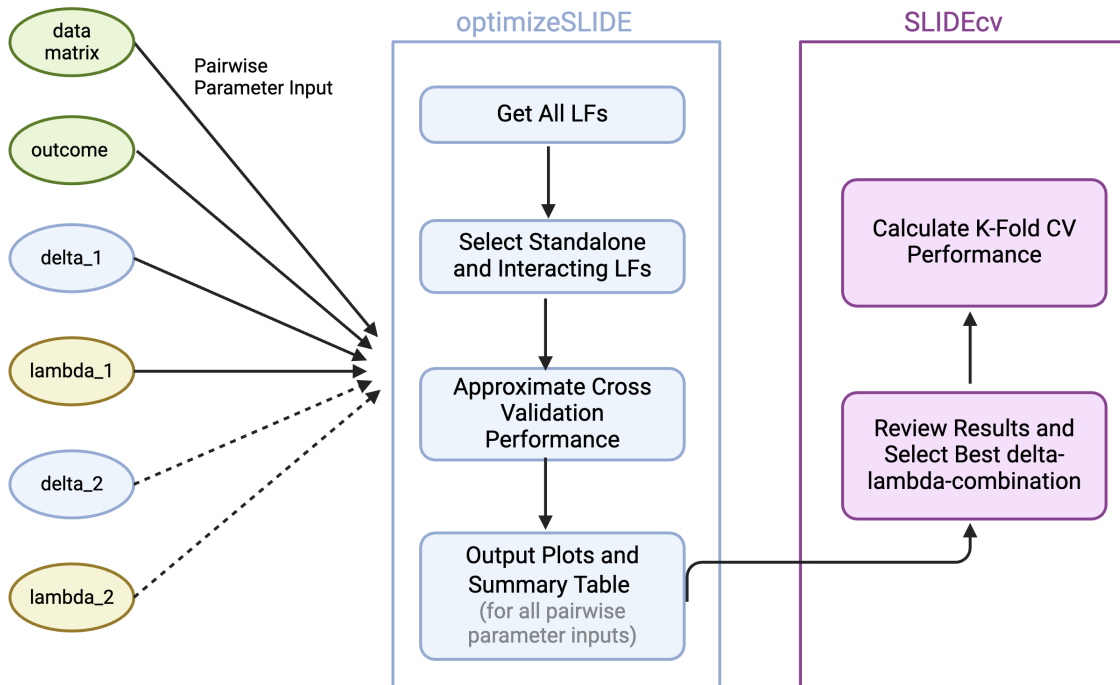


Figure 1: SLIDE workflow

Step 1: Parameter Tuning

1-1 Pre-Processing and Checking Input Data

Due to the assumption-free-nature of SLIDE, there are no data modality limit to the input of SLIDE.

The key input to SLIDE are just two csv files, the data matrix x and the response vector y file post pre-processing such as batch effect correction and/or normalization (for example, for scRNA-seq data, first process with the standard Seurat pipeline).

The x file contains your data in a sample by feature format, such as single cell transcriptomics (**cell by gene**), or spatial proteomics (**region by protein**). The y file contains the responses of the data, such as severity of disease, spatial regions or clonal expansion. Since SLIDE is a regression method, when the response vector has multiple unique values (not just two classes), please make sure there are an ordinal relation between the y values. **Please make sure both csv files have row names and column names.**

In this tutorial, we are going to use the example Systemic Sclerosis dataset we have used in the SLIDE paper. This dataset is a human skin cell scRNA-seq dataset that we have transformed into the pseudo-bulk format.

If you have human single-cell data, our recommended workflow is to pseudobulk your dataset since the cell-to-cell variability is high.

If you have mouse single cell data, with the reduced cell-to-cell variability, you can consider each cell as an sample.

It is extremely important that if the input data is sparse (such as scRNA-seq datasets), please pre-process the data as described in the Preprocessing-and-Filtering vignette.

1-2 Parameters

The input to SLIDE is the path to a YAML file which documents the parameters.

SLIDE accepts many parameters from user input to give user as much freedom to tweak the method to their own data as possible. We have set default values to many parameters that works well with majority of the data modalities. The example yaml files for both continuous/ordinal and binary response can be found in the examples folder. Here, we explain what each of these parameters mean.

```
---
x_path: examples/Ssc_x.csv
y_path: examples/SkinScore.csv
out_path: examples/out
delta:
  - 0.01
  - 0.1
lambda:
  - 0.5
  - 1.0
spec: 0.3
y_factor: FALSE
thresh_fdr: 0.2
y_levels: NULL
eval_type: corr
SLIDE_iter: 500
SLIDE_top_feats: 10
CViter: 10
sampleCV_K: 4
do_interacts: TRUE
```

Figure 2: Example Yaml (Continuous Y)

x_path: a string of the path to the data matrix (x) in the csv format with row names and column names where each row is samples (cells, patients, regions...) and each column is a features (genes, proteins...).

```

---
x_path: path/to/x.csv
y_path: path/to/y.csv
out_path: path/to/output/folder
delta:
  - 0.01
  - 0.1
lambda:
  - 0.5
  - 1
spec: 0.1
y_factor: TRUE
y_levels:
  - 0
  - 1
eval_type: auc
SLIDE_iter: 500
SLIDE_top_feats: 10
CViter: 10
sampleCV_K: 4
do_interacts: TRUE

```

Figure 3: Example Yaml (Binary Y)

y_path: a string of the path of the response vector (y) in the csv format with row names and column names where each row is a sample and the column would be the outcome of interest..

out_path: a string of the path of a folder to store all output files (please see below section to interpret the outputs).

delta: control the number of all latent factors. The higher the delta, the less number of latent factors will be found. Default as 0.01 and 0.1.

lambda: control the sparsity of all latent factors. The higher the lambda, the less number of features will be in a latent factor. Default as 0.5 and 1.

spec: control the number of significant latent factors. The higher the spec, the less number of significant latent factors will be outputted. The desired number of output should be between 5 to 12 LFs. Default as 0.1.

y_factor: set to false if not binary and true if binary.

y_levels: set to null if y is continuous or ordinal. If y is binary, input a list of the correct order relationship such as [0, 1] or [1, 2] depend on the input data.

eval_type: the performance evaluation metric used. corr for continuous Y and auc for binary Y.

SLIDE_iter: the number of times to repeat the SLIDE latent factor selection algorithm. The higher the iteration, the more stable the performance would be. Default as 500.

SLIDE_top_feats: the number of top features to plot from each latent factor. If set as n, a union of the top n weighted features and top n correlated (with y) features will be outputted.

do_interacts (optional): set to false if don't want interacting latent factors. Default as TRUE.

thresh_fdr: set to lower if co-linearity of the features in the data matrix is high. Default as 0.2.

fdr: False discovery control for knockoffs. Default as 0.1.

sampleCV_iter: number of iterations to approximate SLIDE cross-validation results. Default to 500.

sampleCV_K: the number of folds used when approximating cross-validation performance. Default as 4.

2-1 Run optimizeSLIDE

Once your YAML file is ready, we first recommend using a [YAML validator website](#) to ensure your YAML file is correctly formatted.

We can then check if the YAML file is read in correctly by reading it in as an variable.

```
library(SLIDE)
yaml_path = "examples/example_continuous.yaml"
input_params = yaml::read_yaml(yaml_path)

knitr::kable(data.frame(arg = unlist(input_params)))
```

	arg
x_path	examples/Ssc_x.csv
y_path	examples/SkinScore.csv
out_path	examples/out
delta1	0.01
delta2	0.1
lambda1	0.5
lambda2	1
spec	0.3
y_factor	FALSE
thresh_fdr	0.2
eval_type	corr
SLIDE_iter	1000
SLIDE_top_feats	10
sampleCV_iter	500
sampleCV_K	4
do_interacts	TRUE
fdr	0.1

If you have a sparse dataset such as scRNA-seq datasets, we recommend filtering out samples and features that have too many zeros. `zeroFiltering` function will remove samples with more than `g_thresh` number of zeros and features with more than `c_thresh` number of zeros. An appropriate data matrix should at most have around 3-4k features. See [Preprocessing-and-Filtering](#) vignette for more information.

We then check if your data files are formatted correctly.

```
checkDataParams(input_params)
#> Dimension of input X is: 24 by 804 .
#> Dimension of input y is: 24 by 1 .
#> Checking the format and dimensions of input data and response matrices...
#> Checking na values in the input data and response matrices...
#> Checking if yaml file is correct for the input data and response matrices...
```

If everything is formatted correctly, you can now run [Step1](#) of SLIDE. If `sink_file` set to `TRUE`, all print statement will be printed to a txt file.

```
optimizeSLIDE(input_params, sink_file = FALSE)
```

The function above will generate a summary of each `delta` and `lambda` parameter, stored in `summary_table.csv`

```
summary_table = read.csv(paste0(input_params$out_path, "/summary_table.csv"), row.names = 1)

knitr::kable(data.frame(summary_table))
```

delta	lambda	f_size	Num_of_LFs	Num_of_Sig_LFs	Num_of_Interactors	sampleCV_Performance
0.01	0.5	24	172	9	6	0.3713364
0.01	1.0	24	172	NA	NA	NA
0.10	0.5	24	88	4	5	0.7724834
0.10	1.0	24	88	3	3	0.4230769

The `sampleCV_Performance` column will tell us which of our parameters gives latent factors that perform well in short cross validation run. We generally want to pick the parameters that have the highest value for `sampleCV_Performance`. Once we pick parameters, we will run a more rigorous cross validation with more iterations. NA values indicate a model with no significant latent factors.

Finally, we can generate correlation plots for the top features in each significant latent factor.

```
plotCorrelationNetworks(input_params)
```

`optimizeSLIDE` saves various information for each parameter combination, such as latent factor content, model performance against random models, visualization of the LF network and more. [For a more detailed explanation of the outputted files, please see the **Understanding-Outputs** vignette.](#)

3 Run SLIDEcv to do cross-validation on chosen latent factors

After using `optimizeSLIDE` we will have a number of SLIDE models to choose from (and each model will have its own yaml file in its folder, to make running SLIDEcv on a select run easier). We can evaluate the performance for the latent factors in these models by using SLIDEcv. For example, if we thought `delta=0.1` and `lambda=1.0` looked good, we could use the yaml file for that model: To run cross validation use this code:

```
SLIDE::SLIDEcv(yaml_path="examples/out/0.1_1_out/yaml_params.yaml",nrep=20,k=10)
```