

# SCPLearn: Software Manual

Harini Eavani

Harini.Eavani@uphs.upenn.edu

cbica-software@uphs.upenn.edu

September 14, 2015

## 1 Introduction

This software is used to calculate Sparse Connectivity Patterns (SCPs) from resting state fMRI data. SCPs consist of those regions whose connectivity co-varies across subjects. This algorithm was developed as a complementary approach to existing network identification methods. SCPLearn has the following advantages:

1. Does not require thresholding of correlation matrices
2. Allows for both positive and negative correlations
3. Does not constrain the SCPs to have spatial/temporal orthogonality/independence
4. Provides group-common SCPs and subject-specific measures of average correlation within each SCP

## 2 Method

The objective of our method is to find SCPs consisting of functionally synchronous regions, and are smaller than the whole-brain network. The information content within any one of these SCPs is also relatively low, since all the nodes within an SCP are correlated, and express the same information. Hence, if a correlation matrix were constructed for each of these SCPs, it

would show two properties - (1) large number of edges with zero weights, or sparsity and (2) low information content - or rank deficiency.

The input to our method is size  $P \times P$  correlation matrices  $\Sigma_n \succeq 0$ , one for each subject  $n$ ,  $n = 1, 2, \dots, N$ . We would like to find SCPs common to all the subjects, such that a non-negative combination of these SCPs generates the full-correlation matrix  $\Sigma_n$ , for each subject  $n$ . Each of these  $K$  SCPs can be represented by a vector of node-weights  $\mathbf{b}_k$ , where  $-1 \preceq \mathbf{b}_k \preceq 1$ ,  $\mathbf{b}_k \in \mathbf{R}^P$ . Each vector  $\mathbf{b}_k$  reflects the membership of the nodes to the sub-network  $k$ . If  $|\mathbf{b}_k(i)| > 0$ , node  $i$  belongs to the sub-network  $k$ , and if  $\mathbf{b}_k(i) = 0$  it does not. If two nodes in  $\mathbf{b}_k$  have the same sign, then they are positively correlated and opposing sign reflects anti-correlation. Thus, the rank-one matrix  $\mathbf{b}_k \mathbf{b}_k^T$  reflects the correlation behavior of SCP  $k$ . In addition, we constrain these SCPs to be much smaller than the whole-brain network by restricting the  $l_1$ -norm of  $\mathbf{b}_k$  to not exceed a constant value  $\lambda$ .

We would like to approximate each matrix  $\{\Sigma_n\}_{n=1}^N$  by a non-negative combination of SCPs  $\mathbf{B} = [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_K]$ . Thus, we want

$$\begin{aligned} \Sigma_n &\approx \sum_{k=1}^K c_n(k) \mathbf{b}_k \mathbf{b}_k^T + \text{diag}(\mathbf{e}_n) = \mathbf{B} \text{diag}(\mathbf{c}_n) \mathbf{B}^T + \text{diag}(\mathbf{e}_n) \triangleq \hat{\Sigma}_n \\ \|\mathbf{b}_k\|_1 &\leq \lambda, \quad -1 \leq \mathbf{b}_k(i) \leq 1, \quad \mathbf{c}_n \geq 0, \quad \mathbf{e}_n \geq 0 \end{aligned} \quad (1)$$

where  $\text{diag}(\mathbf{c}_n)$  denotes a diagonal matrix with values  $\mathbf{c}_n \in \mathbf{R}_+^K$  along the diagonal. The diagonal matrix  $\text{diag}(\mathbf{e}_n)$  is added to each approximation to make it full-rank. Thus, each subject  $n$  is associated with a vector of  $K$  subject-specific measures  $\mathbf{c}_n$  which are non-negative and reflect the relative contribution of each SCP to the whole-brain functional network in the  $n$ -th subject.

We quantify the approximation between  $\Sigma_n$  and  $\hat{\Sigma}_n$  using the frobenius norm. Note that there is an ambiguity in amplitude between the two factors - if  $\mathbf{b}_k$  and  $c_n(k)$  is a solution,  $\alpha \mathbf{b}_k$  and  $c_n(k)/\alpha^2$  is also a solution for any positive scalar  $\alpha$ . To prevent this, we fix the maximum value in each SCP to unity; i.e.,  $\max_i |\mathbf{b}_k(i)| = 1$ .

Bringing the objective and the constraints together, we have the following optimization problem w.r.t the unknowns  $\mathbf{B}$ ,  $\mathbf{C} = [\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_N]$  and  $\mathbf{E} =$

$[\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n]$ :

$$\begin{aligned}
& \underset{\mathbf{B}, \mathbf{C}}{\text{minimize}} \quad \sum_{n=1}^N \left\| \boldsymbol{\Sigma}_n - \hat{\boldsymbol{\Sigma}}_n \right\|_F^2 \\
& \text{subject to} \quad \left\| \mathbf{b}_k \right\|_1 \leq \lambda, \quad k = 1, \dots, K, \\
& \quad -1 \leq \mathbf{b}_k(i) \leq 1, \quad \max_i |\mathbf{b}_k(i)| = 1, \quad i = 1, \dots, P \\
& \quad \mathbf{c}_n \geq 0, \quad n = 1, \dots, N
\end{aligned} \tag{2}$$

The objective function in the proposed model is non-convex w.r.t both unknown variables  $\mathbf{B}$ ,  $\mathbf{C}$  and  $\mathbf{E}$ . We use the method of alternating minimization to solve for the three unknowns. At each iteration a local minimum is obtained using projected gradient descent. Such a procedure converges to a local minimum. The variable  $\mathbf{B}$  is initialized using k-means run on the time-points;  $\mathbf{C}$  and  $\mathbf{E}$  are initialized to randomly chosen values.

Hierarchical Framework This method can be applied to the data in a hierarchical fashion; first, “primary SCPs” are computed by running the software on the correlation matrices  $\boldsymbol{\Sigma}_n$ . Then, for each primary SCP, multiple smaller “secondary SCPs” can be found by re-applying the method only for those regions that belong to the primary SCP.

Sample weights This software also allows subject-specific weights as input - e.g., in cases when subjects have different numbers of time-points. In this case, the optimization objective is modified to be  $\sum_{n=1}^N \alpha_n \left\| \boldsymbol{\Sigma}_n - \hat{\boldsymbol{\Sigma}}_n \right\|_F^2$ , where  $\alpha_n$  is the subject-specific weight.

## 3 Software

### 3.1 Within main directory

1. Detailed documentation is contained in the latex file [UserManual.tex](#) and the pdf [UserManual.pdf](#) produced from it.
2. Author information is in [AUTHORS.txt](#)
3. All information needed to get started with the software is in [README.md](#)

4. The file `m2cpp.pl` is a helper perl script for doxygen. It is needed to pick up comments from MATLAB script files.
5. The file `INSTALL.txt` has installation and usage information.
6. `Doxyfile` is the configuration file for generating doxygen documentation.
7. `ignore.txt` contains a list of files that doxygen ignores from parsing.
8. The folder `docs/` will contain all the documentation once doxygen is run using `Doxyfile`.
9. The folder `src/` contains source code for SCPLearn.
10. The folder `licenses/` contains license information related to this software.

### 3.2 Within `src/` directory

This software is implemented primarily in MATLAB. For command-line argument parsing, and nifti I/O operations, python code is used.

For ease of use, the software takes as input either rsfMRI voxel time-series data in NIFTI format, or ROI time-series data in mat format. The list of nifti files is provided as a text file.

It also takes as input the node definitions as an atlas in NIFTI space. This atlas must have the ROI regions numbered 1, 2, 3, ... and must be in the same space as the subject data.

1. Main function is `SCPLearn.py`. This function parses the input arguments and calls “`ComputeROIAverages.py`”.
2. `ComputeROIAverages.py` reads the rsfMRI data and the node definitions. It extracts average time-series of each ROI for each subject and saves it in “.mat” files.
3. These mat files are loaded into MATLAB using the function `SCPLearn-FromMatFiles.m`. Then, based on the input parameters, either (1) `SCPLearn.m` is called, which computes only primary SCPs, or (2) `SCPLearn_two_level.m` is called which computes both primary and secondary level SCPs.

4. Both the above functions call [BlockSolverLowRankFrob.m](#), which implements the alternate minimization strategy discussed in the previous section. This function calls (1) [BSolver\\_lowrank\\_frob.m](#) (2) [CSolver\\_lowrank\\_frob.m](#) (3) [ESolver\\_lowrank\\_frob.m](#) until termination criteria are met.
5. [BSolver\\_lowrank\\_frob.m](#) calls two functions (1) [spg\\_04262015.m](#) which implements projected gradient descent. (2) [ProjectionOnUnitBoxSimplex.m](#) implements the projection function for the sparse basis.
6. After the MATLAB code has finished running, [SCPLearn.py](#) calls [replaceLabels\\_nib.py](#). This python code replaces the labels in the atlas with the values of the nodes in each SCP. In this manner, for each SCP, a nifti image is written to disk.
7. The set of python functions in [SCPUtills.py](#) are used for error checking, handling exceptions, file check, etc
8. The MATLAB functions [SCPLearn\\_SplitSample.m](#), [CompareSCPs.m](#), [Hungarian.m](#) implement split-sample SCP reproducibility and error calculation, as described and used in the main paper.
9. The python script [makeAll.py](#) calls mcc to build [SCPLearnFromMatFiles.m](#) as an mcc-executable. It then copies all python files to the build directory.

## 4 Installation

### 4.1 Dependencies

This software has been primarily implemented for Linux operating systems.

- MATLAB Compiler mcc version 5.2 (R2014B)
- MATLAB R2014B
- Python 2.7.9
- Python library numpy 1.7.2
- Python library scipy 0.15.1

- Python library pandas 0.16.2
- Python library nibabel 2.0.1

Make sure all dependencies are met before proceeding with install.

## 4.2 Generating standalone executables

1. Within the `src/` directory, run `makeAll.py` with the location of the install directory as the argument:

```
makeAll.py <installDir>
```

2. Add the install directory to your path by running the following command.

Replace `$installDir` with the location of the install directory from step 1 above.

```
export PATH=$PATH:$installDir
```

## 4.3 Generating documentation

Run `doxygen Doxyfile` from the parent directory. The documentation is built in `docs/`. Open `docs/index.html` in your favorite browser to begin.

# 5 Usage

Usage: `SCPLearn.py [OPTIONS]`

Required Options:

`[-d --data]` Specify the text file with list of nifti files (required)

\*\*\* OR \*\*\*

Specify the text file with list of mat files (required)

Each mat file must contain a variable named 'ts'

'ts' must be a matrix of time-series, size = #ROIs X #timepoints

`[-m --mask]` Specify the nifti ROI/parcel/atlas file (required)

`[-p --prefix]` Specify the prefix of the output file (required)

Options:

`[-s --sparsity]` Specify the sparsity constraint as positive value.  
Default =  $\text{nROIs}/8$

`[-n --numberOfSCPs]` Specify the number of SCPs as a number.  
Default = 10.

`[-r --pruning]` Specify the pruning threshold as a value between  $[0,1]$ .  
Default = 0.7.

`[-l --levels]` Hierarchical learning of SCPs. Not run by default.

`[-o --outputDir]` The output directory to write the results.  
Defaults to the location of the input file

`[-w --workingDir]` Specify a working directory.  
By default a tmp dir is created and used

`[-u --usage | -h --help]` Display this message

`[-v --verbose]` Verbose output

`[-V --Version]` Display version information

Example:

```
SCPLearn.py -d list_of_nifti_files.txt -m Grasp_level5.nii -p SCP_results \
-n 10 -o /sbia/sbiaprfj/BLSA -v
```

Example list\_of\_nifti\_files.txt without sample weights:

```
ProjName_subj_165464.nii.gz
ProjName_subj_26464.nii.gz
ProjName_subj_1054.nii.gz.....
```

Example list\_of\_nifti\_files.txt with sample weights:

```
ProjName_subj_165464.nii.gz,172
ProjName_subj_26464.nii.gz,160
ProjName_subj_1054.nii.gz,120.....
```

To run this software you will need:

1. Text file with list of NIFTI files
2. NIFTI file with node definitions

In the output directory, the software returns:

1. One NIFTI file for each SCP that is generated, with file name `<prefix>_SCP_#.nii.gz`

2. A `<prefix>_SCP_Coeffs.csv` file with the SCP coefficients for all the subjects, indexed by the NIFTI filename that was input
3. A `<prefix>_SCP_basis.csv` file with the SCP basis, in csv format
4. A `<prefix>_SCPs.mat` file with all the outputs in MATLAB “.mat” format
5. A `<prefix>_ts.mat` file with the average time-series from all the subjects in MATLAB “.mat” format
6. A `<subject>_<atlas>.mat` file with the average time-series for each subject

## 6 Testing

To test the software, we have provided synthetic nifti as well as mat data for ten “subjects”. These are located in `src/test`.

To run the software on the synthetic data, run the following command from the `src/` folder :

```
SCPLearn.py -d test/nifti_list.txt -p test_2d -m test/test_2d_mask.nii.gz
-o test/ -n 3 -s 10
```

or

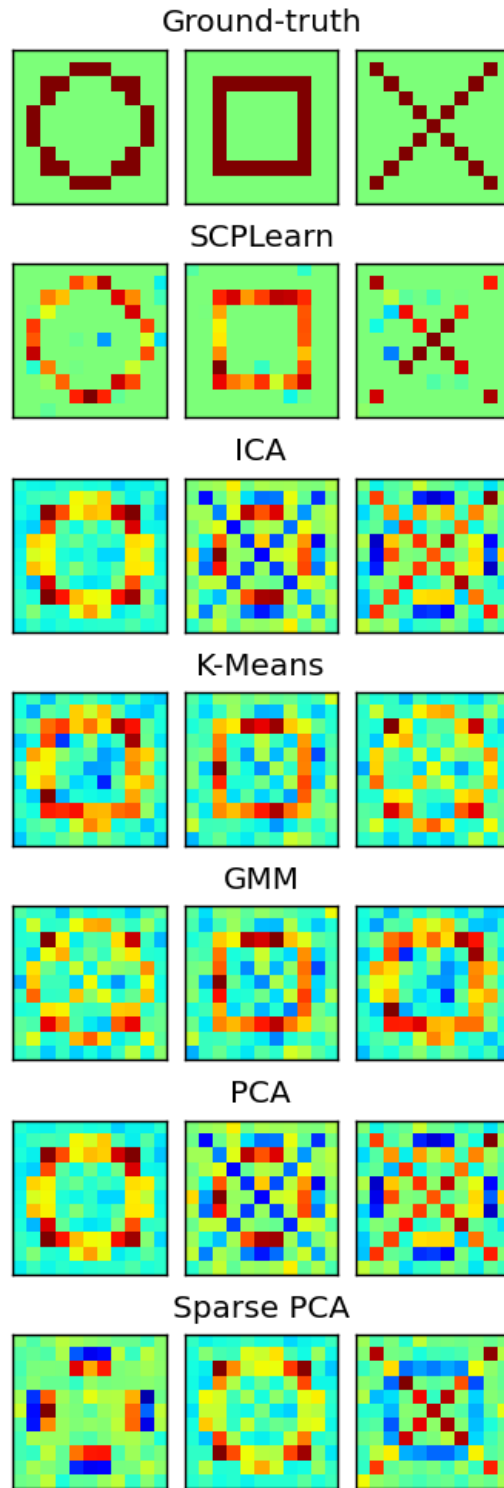
```
SCPLearn.py -d test/mat_list.txt -p test_2d -m test/test_2d_mask.nii.gz
-o test/ -n 3 -s 10
```

This command should return three SCPs, saved as nifti files in `src/test/test_2d_SCP_1,2,3.nii.gz` along with other output.

By design, three patterns were used as “ground-truth” SCPs to generate the synthetic data; a cross, a circle and a square. Was SCPLearn able to clearly separate these three patterns in the output? You can compare your result to the results in `test/Test_2d_all.png`. This image is also shown in Figure 1.

In addition, the NetSim-based synthetic data used in the original paper is provided in `netsim_data_05262014_timeseries.mat`. To test the MATLAB code using this data, run `testSCPLearnCode.m` from the MATLAB command line.





9

Figure 1: Synthetic SCPs used for testing. Results from SCPLearn and other approaches are also shown.

## 7 Citation

If you find this code useful, please cite:

Eavani, H., Satterthwaite, T. D., Filipovych, R., Gur, R. E., Gur, R. C., Davatzikos, C. (2015). Identifying Sparse Connectivity Patterns in the brain using resting-state fMRI. *Neuroimage*, 105, 286-299.