

# c-lasso - a package for constrained sparse and robust regression and classification in Python

09 October 2020

## Summary

We introduce **c-lasso**, a Python package that enables sparse and robust linear regression and classification with linear equality constraints. The underlying statistical forward model is assumed to be of the following form:

$$y = X\beta + \sigma\epsilon \quad \text{subject to} \quad C\beta = 0$$

Here,  $X \in \mathbb{R}^{n \times d}$  is a given design matrix and the vector  $y \in \mathbb{R}^n$  is a continuous or binary response vector. The matrix  $C$  is a general constraint matrix. The vector  $\beta \in \mathbb{R}^d$  contains the unknown coefficients and  $\sigma$  an unknown scale. Prominent use cases are (sparse) log-contrast regression with compositional data  $X$ , requiring the constraint  $1_d^T \beta = 0$  [Aitchison:1984] and the Generalized Lasso which is a *special case* of the described problem (see, e.g., [James:2020], Example 3). The **c-lasso** package provides several estimators for inferring unknown coefficients and scale (i.e., perspective M-estimators [Combettes:2020a]) of the form

$$\min_{\beta \in \mathbb{R}^d, \sigma \in \mathbb{R}_0} f(X\beta - y, \sigma) + \lambda \|\beta\|_1 \quad \text{subject to} \quad C\beta = 0$$

for several convex loss functions  $f(\cdot, \cdot)$ . This includes the constrained Lasso, the constrained scaled Lasso, and sparse Huber M-estimators with linear equality constraints.

## Statement of need

Currently, there is no Python package available that can solve these ubiquitous statistical estimation problems in a fast and efficient manner. **c-lasso** provides algorithmic strategies, including path and proximal splitting algorithms, to solve

the underlying convex optimization problems with provable convergence guarantees. The `c-lasso` package is intended to fill the gap between popular Python tools such as `scikit-learn` which cannot solve these constrained problems and general-purpose optimization solvers such as `cvxpy` that do not scale well for these problems and/or are inaccurate. `c-lasso` can solve the estimation problems at a single regularization level, across an entire regularization path, and includes three model selection strategies for determining the regularization parameter: a theoretically-derived fixed regularization, k-fold cross-validation, and stability selection. We show several use cases of the package, including an application of sparse log-contrast regression tasks for compositional microbiome data, and highlight the seamless integration into R via `reticulate`.

## Functionalities

### Installation and problem instantiation

`c-lasso` is available on pip and can be installed in the shell using

```
pip install c-lasso
```

The central object in the `c-lasso` package is the instantiation of a `c-lasso` problem.

```
# Import the main class of the package
from classo import classo_problem

# Define a c-lasso problem instance with default setting,
# given data X, y, and constraints C.
problem = classo_problem(X,y,C)
```

We next describe what type of problem instances are available and how to solve them.

### Statistical problem formulations

Depending on the type of and the prior assumptions on the data, the noise  $\epsilon$ , and the model parameters, `c-lasso` allows for different estimation problem formulations. More specifically, the package can solve the following four regression-type and two classification-type formulations:

#### **R1 Standard constrained Lasso regression:**

$$\min_{\beta \in \mathbb{R}^d} \|X\beta - y\|^2 + \lambda \|\beta\|_1 \quad \text{subject to} \quad C\beta = 0$$

This is the standard Lasso problem with linear equality constraints on the  $\beta$  vector. The objective function combines Least-Squares (LS) for model fitting with the  $L_1$ -norm penalty for sparsity.

```
# Formulation R1
problem.formulation.huber = False
problem.formulation.concomitant = False
problem.formulation.classification = False
```

#### **R2 Constrained sparse Huber regression:**

$$\min_{\beta \in \mathbb{R}^d} h_\rho(X\beta - y) + \lambda \|\beta\|_1 \quad \text{subject to} \quad C\beta = 0$$

This regression problem uses the Huber loss  $h_\rho$  as objective function for robust model fitting with an  $L_1$  penalty and linear equality constraints on the  $\beta$  vector. The default parameter  $\rho$  is set to 1.345 [Huber:1981].

```
# Formulation R2
problem.formulation.huber = True
problem.formulation.concomitant = False
problem.formulation.classification = False
```

#### **R3 Constrained scaled Lasso regression:**

$$\min_{\beta \in \mathbb{R}^d, \sigma \in \mathbb{R}_0} \frac{\|X\beta - y\|^2}{\sigma} + \frac{n}{2}\sigma + \lambda \|\beta\|_1 \quad \text{subject to} \quad C\beta = 0$$

This formulation is the default problem formulation in `c-lasso`. It is similar to *R1* but allows for joint estimation of the (constrained)  $\beta$  vector and the standard deviation  $\sigma$  in a concomitant fashion [Combettes:2020a; Combettes:2020b].

```
# Formulation R3
problem.formulation.huber = False
problem.formulation.concomitant = True
problem.formulation.classification = False
```

#### **R4 Constrained sparse Huber regression with concomitant scale estimation:**

$$\min_{\beta \in \mathbb{R}^d, \sigma \in \mathbb{R}_0} \left( h_\rho \left( \frac{X\beta - y}{\sigma} \right) + n \right) \sigma + \lambda \|\beta\|_1 \quad \text{subject to} \quad C\beta = 0$$

This formulation combines *R2* and *R3* allowing robust joint estimation of the (constrained)  $\beta$  vector and the scale  $\sigma$  in a concomitant fashion [Combettes:2020a; Combettes:2020b].

```

# Formulation R4
problem.formulation.huber = True
problem.formulation.concomitant = True
problem.formulation.classification = False

```

**C1** Contrained sparse classification with Square Hinge loss:

$$\min_{\beta \in \mathbb{R}^d} \sum_{i=1}^n l(y_i x_i^\top \beta) + \lambda \|\beta\|_1 \quad \text{subject to} \quad C\beta = 0$$

where  $x_i$  denotes the  $i^{th}$  row of  $X$ ,  $y_i \in \{-1, 1\}$ , and  $l(\cdot)$  is defined for  $r \in \mathbb{R}$  as:

$$l(r) = \begin{cases} (1-r)^2 & \text{if } r \leq 1 \\ 0 & \text{if } r \geq 1 \end{cases}$$

This formulation is similar to *R1* but adapted for classification tasks using the Square Hinge loss with (constrained) sparse  $\beta$  vector estimation [Lee:2013].

```

# Formulation C1
problem.formulation.huber = False
problem.formulation.concomitant = False
problem.formulation.classification = True

```

**C2** Contrained sparse classification with Huberized Square Hinge loss:

$$\min_{\beta \in \mathbb{R}^d} \sum_{i=1}^n l_\rho(y_i x_i^\top \beta) + \lambda \|\beta\|_1 \quad \text{subject to} \quad C\beta = 0.$$

This formulation is similar to *C1* but uses the Huberized Square Hinge loss  $l_\rho$  for robust classification with (constrained) sparse  $\beta$  vector estimation [Rosset:2007]:

$$l_\rho(r) = \begin{cases} (1-r)^2 & \text{if } \rho \leq r \leq 1 \\ (1-\rho)(1+\rho-2r) & \text{if } r \leq \rho \\ 0 & \text{if } r \geq 1 \end{cases}$$

This formulation can be selected in *c-lasso* as follows:

```

# Formulation C2
problem.formulation.huber = True
problem.formulation.concomitant = False
problem.formulation.classification = True

```

## Optimization schemes

The problem formulations *R1-C2* require different algorithmic strategies for efficiently solving the underlying optimization problems. The **c-lasso** package implements four published algorithms with provable convergence guarantees. The package also includes novel algorithmic extensions to solve Huber-type problems using the mean-shift formulation [Mishra:2019]. The following algorithmic schemes are implemented:

- Path algorithms (*Path-Alg*): This algorithm follows the proposal in [Gaines:2018; Jeon:2020] and uses the fact that the solution path along  $\lambda$  is piecewise-affine [Rosset:2007]. We also provide a novel efficient procedure that allows to derive the solution for the concomitant problem *R3* along the path with little computational overhead.
- Douglas-Rachford-type splitting method (*DR*): This algorithm can solve all regression problems *R1-R4*. It is based on Douglas-Rachford splitting in a higher-dimensional product space and makes use of the proximity operators of the perspective of the LS objective [Combettes:2020a; Combettes:2020b]. The Huber problem with concomitant scale *R4* is reformulated as scaled Lasso problem with mean shift vector [Mishra:2019] and thus solved in  $(n + d)$  dimensions.
- Projected primal-dual splitting method (*P-PDS*): This algorithm is derived from [Briceno:2020] and belongs to the class of proximal splitting algorithms, extending the classical Forward-Backward (FB) (aka proximal gradient descent) algorithm to handle an additional linear equality constraint via projection. In the absence of a linear constraint, the method reduces to FB.
- Projection-free primal-dual splitting method (*PF-PDS*): This algorithm is a special case of an algorithm proposed in [Combettes:2012] (Eq. 4.5) and also belongs to the class of proximal splitting algorithms. The algorithm does not require projection operators which may be beneficial when  $C$  has a more complex structure. In the absence of a linear constraint, the method reduces to the Forward-Backward-Forward scheme.

The following table summarizes the available algorithms and their recommended use for each problem:

|           | <i>Path-Alg</i>                              | <i>DR</i>               | <i>P-PDS</i> | <i>PF-PDS</i>               |
|-----------|--|-------------------------|--------------|-----------------------------|
| <i>R1</i> | use for large $\lambda$ and path computation | use for small $\lambda$ | possible     | use for complex constraints |
| <i>R2</i> | use for large $\lambda$ and path computation | use for small $\lambda$ | possible     | use for complex constraints |

|           | <i>Path-Alg</i>                                    | <i>DR</i>               | <i>P-PDS</i> | <i>PF-PDS</i> |
|-----------|--|-------------------------|--------------|---------------|
| <i>R3</i> | use for large $\lambda$<br>and path<br>computation | use for small $\lambda$ | -            | -             |
| <i>R4</i> | -  | only option             | -            | -             |
| <i>C1</i> | only option  | -                       | -            | -             |
| <i>C2</i> | only option  | -                       | -            | -             |

The following Python snippet shows how to select a specific algorithm:

```
problem.numerical_method = "Path-Alg"
# Alternative options: "DR", "P-PDS", and "PF-PDS"
```

## Computation modes and model selection

The `c-lasso` package provides several computation modes and model selection schemes for tuning the regularization parameter.

- *Fixed Lambda*: This setting lets the user choose a fixed parameter  $\lambda$  or a proportion  $l \in [0, 1]$  such that  $\lambda = l \times \lambda_{\max}$ . The default value is a scale-dependent tuning parameter that has been derived in [Shi:2016] and applied in [Combettes:2020b].
- *Path Computation*: This setting allows the computation of a solution path for  $\lambda$  parameters in an interval  $[\lambda_{\min}, \lambda_{\max}]$ . The solution path is computed via the *Path-Alg* scheme or via warm-starts for other optimization schemes.
- *Cross Validation*: This setting allows the selection of the regularization parameter  $\lambda$  via k-fold cross validation for  $\lambda \in [\lambda_{\min}, \lambda_{\max}]$ . Both the Minimum Mean Squared Error (or Deviance) (MSE) and the “One-Standard-Error rule” (1SE) are available [Hastie:2009].
- *Stability Selection*: This setting allows the selection of the  $\lambda$  via stability selection [Meinshausen:2010; Lin:2014; Combettes:2020b]. Three modes are available: selection at a fixed  $\lambda$  [Combettes:2020b], selection of the  $q$  first variables entering the path (default setting), and of the  $q$  largest coefficients (in absolute value) across the path [Meinshausen:2010].

The Python syntax to use a specific computation mode and model selection is exemplified below:

```
# Example how to perform path computation and cross-validation:
problem.model_selection.LAMfixed = False
problem.model_selection.PATH = True
problem.model_selection.CV = True
problem.model_selection.StabSel = False
```

```
# Example how to add stability selection to the problem instance
problem.model_selection.StabSel = True
```

Each model selection procedure has additional meta-parameters that are described in the Documentation.

## Computational examples

### Toy example using synthetic data

We illustrate the workflow of the `c-lasso` package on synthetic data using the built-in routine `random_data` which enables the generation of test problem instances with normally distributed data  $X$ , sparse coefficient vectors  $\beta$ , and constraints  $C \in \mathbb{R}^{k \times d}$ .

Here, we use a problem instance with  $n = 100$ ,  $d = 100$ , a  $\beta$  with five non-zero components,  $\sigma = 0.5$ , and a zero-sum constraint.

```
from classo import classo_problem, random_data

n,d,d_nonzero,k,sigma =100,100,5,1,0.5
(X,C,y),sol = random_data(n,d,d_nonzero,k,sigma,zerosum=True, seed = 123 )
print("Relevant variables : {}".format(list(numpy.nonzero(sol)) ) )

problem = classo_problem(X,y,C)

problem.formulation.huber = True
problem.formulation.concomitant = False
problem.formulation.rho = 1.5

problem.model_selection.LAMfixed = True
problem.model_selection.PATH = True
problem.model_selection.LAMfixedparameters.rescaled_lam = True
problem.model_selection.LAMfixedparameters.lam = 0.1

problem.solve()

print(problem.solution)
```

We use formulation *R2* with  $\rho = 1.5$ , computation mode and model selections *Fixed Lambda* with  $\lambda = 0.1\lambda_{\max}$ , *Path computation*, and *Stability Selection* (as per default).

The corresponding output reads:

```
Relevant variables : [43 47 74 79 84]
```

```

LAMBDA FIXED :
    Selected variables :  43    47    74    79    84
    Running time :  0.294s

PATH COMPUTATION :
    Running time :  0.566s

STABILITY SELECTION :
    Selected variables :  43    47    74    79    84
    Running time :  5.3s

```

c-lasso allows standard visualization of the computed solutions, e.g., coefficient plots at fixed  $\lambda$ , the solution path, the stability selection profile at the selected  $\lambda$ , and the stability selection profile across the entire path.

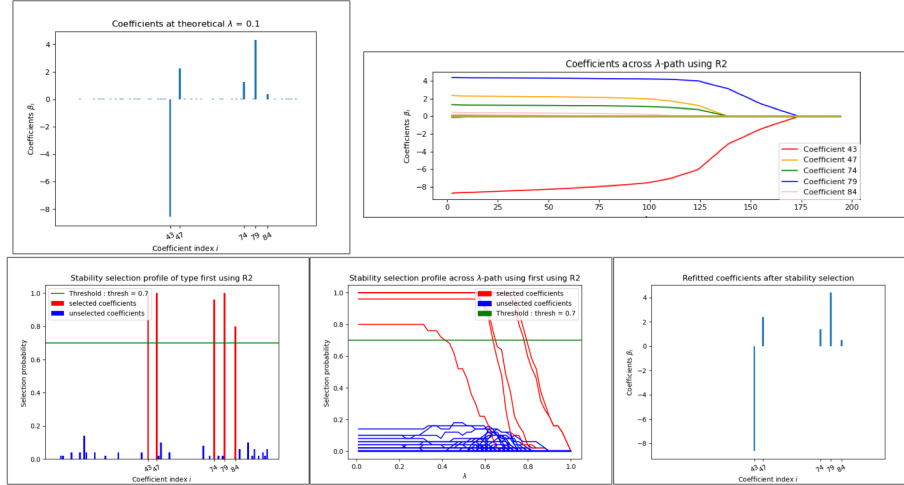


Figure 1: Visualizations after calling problem.solution

For this tuned example, the solutions at the fixed lambda and with stability selection recover the oracle solution. The solution vectors are stored in `problem.solution` and can be directly accessed for each mode/model selection.

*# Access to the estimated coefficient vector at a fixed lambda*  
`problem.solution.LAMfixed.beta`

Note that the run time for this  $d = 100$ -dimensional example for a single path computation is about 0.5 seconds on a standard Laptop.



## Log-contrast regression on gut microbiome data

We next illustrate the application of `c-lasso` on the COMBO microbiome dataset [Lin:2014;Shi:2016;Combettes:2020b]. Here, the task is to predict the Body Mass Index (BMI) of  $n = 96$  participants from  $d = 45$  relative abundances of bacterial genera, and absolute calorie and fat intake measurements. The code snippet for this example is available in the `README.md` and the example notebook.

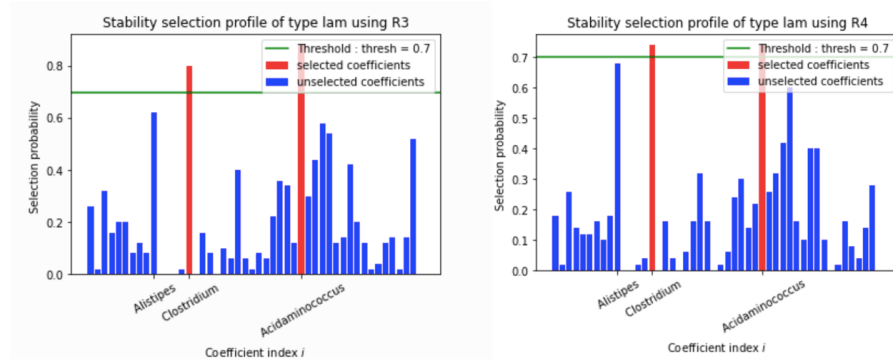


Figure 2: Stability selection profiles of problems R3/R4 on the COMBO data

Stability selection profiles using formulation  $R3$  (left) and  $R4$  (right) on the COMBO dataset, reproducing Figure 5a in [Combettes:2020b].

## Calling `c-lasso` in R

The `c-lasso` package also integrates with R via the R package `reticulate`. We refer to `reticulate`'s manual for technical details about connecting python environments and R. A successful use case of `c-lasso` is available in the R package `trac` [Bien:2020], enabling tree-structured aggregation of predictors when features are rare.

The code snippet below shows how `c-lasso` is called in R to perform regression at a fixed  $\lambda = 0.1\lambda_{\max}$ . In R, `X` and `C` need to be of `matrix` type, and `y` of `array` type.

```
problem <- classo$classo_problem(X=X,C=C,y=y)
problem$model_selection$LAMfixed <- TRUE
problem$model_selection$StabSel <- FALSE
problem$model_selection$LAMfixedparameters$rescaled_lam <- TRUE
problem$model_selection$LAMfixedparameters$lam <- 0.1
problem$solve()

# Extract coefficient vector with tidy-verse
beta <- as.matrix(map_df(problem$solution$LAMfixed$beta, as.numeric))
```

## Acknowledgements

The work of LS was conducted at and financially supported by the Center for Computational Mathematics (CCM), Flatiron Institute, New York, and the Institute of Computational Biology, Helmholtz Zentrum München. We thank Dr. Leslie Greengard (CCM and Courant Institute, NYU) for facilitating the initial contact between LS and CLM. The work of PLC was supported by the National Science Foundation under grant DMS-1818946.

## References