

# Sparse Logistic Regression ToolBox ver1.2alpha

(Updated on 2009/07/28 by Okito Yamashita)

The Sparse Logistic Regression toolbox (**SLR toolbox** hereafter) is a suite of MATLAB functions for solving classification problems. It provides a solution for binary or multi-class classification problem. The unique feature is the weight parameters of the classifier are learned in a sparse way. Thus the algorithm estimates the weight parameters while it automatically finds important features. Due to this unique feature, SLR is applicable to a high-dimensional classification problem where the number of samples is much less than the number of features without suffering from 'the overfitting problem' to some extent. In addition, SLR releases users from the feature selection task that is usually very time-consuming and requires experience. The algorithms are designed for problems of which feature size is up to several thousand and of which sample size is up to several thousand.

~~~~~

## Installation

~~~~~

SLR toolbox is a suite of MATLAB functions and scripts. MATLAB, a commercial engineering mathematics package, is required to use SLR toolbox. A couple of functions require the optimization toolbox (see below). Codes in the toolbox were written for MATLAB ver7.0.1 or later under UNIX. This toolbox has originally been developed by Okito Yamashita in ATR Computational Neuroscience laboratories for personal use.

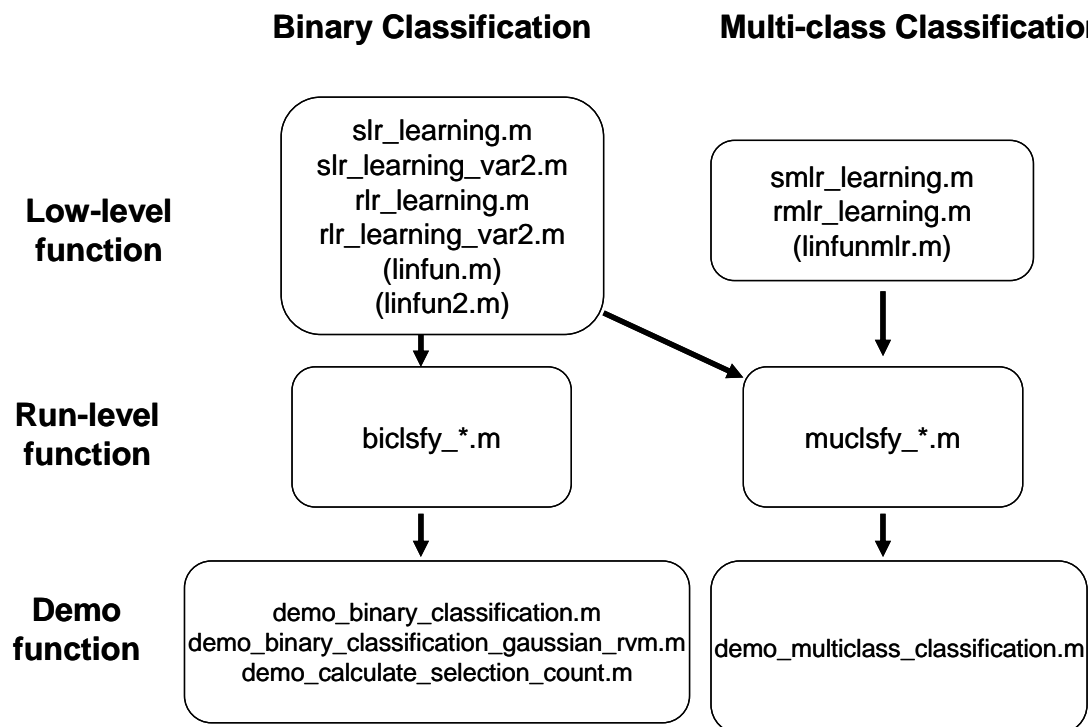
To get installed the toolbox, you just download and unzip the file (SLR1.2alpha.zip) wherever you like. [Please start from demo functions \('demo\\_\\*.m'\) to learn how the functions in SLR toolbox work.](#)

In two demo functions ('demo\_{binary, multiclass}\_classification.m'), you can work with simulated data sets as well as datasets acquired from two real experiments by changing a variable 'data' on the top of each code. You can download these two experimental data. Please download 'TESTDATA.zip' and unzip it in the parallel level where you unzip SLR toolbox.

## Functions in the toolbox

Functions in SLR toolbox can be separated into 5 classes; “low-level functions”, “run-level functions”, “demo functions”, “common functions”, “functions from others” (see figure 1).

- “Low-level functions” are the functions that implement parameter learning in the classifier. Most functions take a pair of label and a feature matrix as inputs and output estimated weight parameters.
- “Run-level functions” are the functions that implement whole procedure to solve a classification problem (i.e. normalizing features → learning parameters → testing learning parameters). Most functions take pairs of label and a feature matrix for learning and testing as inputs and output estimated parameters as well as some performance measure. There are seven functions for the binary classification problem and six functions for the multi-class classification problem (see below). [It should be sufficient for most of users to know how to use these functions.](#)
- “Demo functions” are for demonstration purpose. You can learn how functions in the toolbox work using some examples. [Please run these functions at first.](#)
- “Common functions” are the functions that are commonly used in toolbox.
- “Functions from others” are the functions that are borrowed from EEGLAB toolbox and MATLAB File exchange.



Common  
functions

calc\_SCval.m calc\_label.m calc\_percor.m gen\_gm\_data.m  
gen\_gmgm\_data.m gen\_simudata.m gen\_sin\_data.m  
label2num.m normalize\_feature.m  
randmn.m separate\_train\_test.m slr\_\*.m

Functions  
from others

finputcheck.m xyrefline.m

Figure 1: Functions in SLR toolbox

### List of “Run-level functions”

In “run-level functions”, each of seven functions of which name starts from ‘biclsfy\_’ implements five binary classifiers and each of six functions of which name starts from ‘muclsfy\_’ implements six multi-class classifiers (see Table 1). Classifier in the row of ‘Sparse’ has feature selection property while classifiers in the row of ‘Others’ do not (for comparison purpose).

Table 1: Supported classifiers in SLR toolbox (Run-level functions)

	Binary classification	Multi-class classification
<b>Sparse</b>	SLR-LAP (biclsfy_slrlap.m) SLR-VAR (biclsfy_slrvar.m) L1-SLR-LAP (biclsfy_l1slrlap.m) L1-SLR-C (biclsfy_l1slrc.m)	SMLR (muclsfy_smlr.m) SLR-LAP-1vsR (muclsfy_slrlapovrm.m) SLR-VAR-1vsR (muclsfy_slrvarovrm.m) SLR-VAR-1vs1 (muclsfy_slrvarovo.m)
<b>Others</b>	RLR-LAP (biclsfy_rlrslap.m) RLR-VAR (biclsfy_rlrvar.m) RVM (biclsfy_rvm.m)	RMLR (muclsfy_rmlr.m) RLR-VAR-1vsR (muclsfy_rlrvarovrm.m)

SLR = Sparse Logistic Regression

L1-SLR = L1-norm Sparse Logistic Regression

RLR = Regularized Logistic Regression

RVM = Relevance Vector Machine

SMLR = Sparse Multinomial Logistic Regression

RMLR = Regularized Multinomial Logistic Regression

LAP = with Laplace approximation

VAR = with variational approximation

C = with component-wise update

1vsR = One-Versus-Rest

1vs1 = One-Versus-One

### **Binary classifiers**

Two sparse classifiers and three non-sparse classifiers are supported so far.

- **SLR-LAP**: SLR with Laplace approximation. The marginal posterior-distribution of weight parameters is approximated by multivariate Gaussian distribution (see ref.[1] for details). This was developed for the research in ref.[1]. **The optimization toolbox is required ('fminunc.m').**
- **SLR-VAR**: SLR with variational approximation. The logistic function is approximated by Gaussian distribution using a variational parameter (see ref.[2.3] or section 10.6 of ref.[4]). The faster and the less memory.
- **RLR-LAP**: Regularized logistic regression with Laplace approximation. This is not sparse algorithm. The regularization parameter is automatically determined by the algorithm. **The optimization toolbox is required ('fminunc.m').**
- **RLR-VAR**: Regularized logistic regression with variational approximation. This is not sparse algorithm. The regularization parameter is automatically determined by the algorithm.
- **RVM**: Relevance Vector Machine as proposed by Tipping (see ref.[5]). Bayesian version of Support Vector Machine (SVM). The linear and Gaussian kernels are supported. The Gaussian kernel RVM is only non-linear classifier supported in this toolbox.
- **L1-SLR-LAP**: L1-norm based sparse logistic regression (see ref.[7]). This classifier also provides sparse solution. The extent of sparsity depends on one tunable parameter. **This program has not been debugged carefully yet. The optimization toolbox is required ('fminunc.m').**
- **L1-SLR-C**: L1-norm based sparse logistic regression (see ref.[6]). This is fast and requires less memory. **This program has not been debugged carefully yet.**

Two sparse classifiers, SLR-LAP and SLR-VAR, are derived from the identical probabilistic model (see ref[1]) but different approximation to the posterior distribution. The difference between RLR-LAP and RLR-VAR is as in the same way. **Among two sparse classification methods, I recommend that you start using SLR-VAR because it is faster and require less memory.**

### **Multi-class classifier**

Three sparse classifiers and two non-sparse classifiers are supported so far.

- **SMLR** : Sparse Multinomial Logistic Regression (see ref.[1]). The multinomial distribution is used for observation. In general, memory and time required are huge. **The optimization toolbox is required ('fminunc.m').**
- **SLR-LAP-1vsR** : Combination of SLR-LAP classifiers. One-versus-the rest scheme is used. **The optimization toolbox is required ('fminunc.m').**
- **SLR-VAR-1vsR** : Combination of SLR-VAR classifiers. One-versus-the rest scheme is used. The faster computation and the less memory.
- **RMLR** : Regularized Multinomial Logistic Regression. One regularization parameter common to all the classes is automatically estimated. **The optimization toolbox is required ('fminunc.m').**
- **RLR-VAR-1vsR** : Combination of SLR-VAR classifiers. One-versus-the rest schemes is used. One regularization parameter per each class is automatically estimated.
- **SLR-VAR-1vs1**: Combination of SLR-VAR classifiers. One-versus-the one scheme is used. The faster computation and the less memory.

SMLR is a true multinomial classifier that uses multinomial distribution for observation. SLR-LAP-1vsR and SLR-VAR-1vsR are consisting of combination of sparse binary classifiers. Especially one-versus-the-rest scheme is employed (see chapter4 of ref[4] for example). In this release, one-versus-one scheme is also supported for SLR-VAR. In theory, SMLR is the best classifier for multi-class problem since the model learning take into account all the information among classes. But in my experience SLR-LAP-1vsR or SLR-VAR-1vsR, SLR-VAR-1vs1 do perform as well as SMLR probably due to a small number of training samples. Thus I recommend to start from SLR-VAR-1vs1 that is fast and requires less memory.

~~~~~

### How to use “Run-level functions”

~~~~~

The details of “[biclfsfy\\_slrvar.m](#)” are explained as an example of binary classifiers. You can easily understand the structure and options of other run-level functions since they share the same structure and the similar but slightly different options with “[biclfsfy\\_slrvar.m](#)”.

The function has following format;

- Syntax

[ww, ix\_eff\_all, errTable\_tr, errTable\_te, parm, AXall,Ptr,Pte] =...

biclsfy\_slrvar(x\_train, t\_train, x\_test, t\_test, varargin).

- Example of Usage

```
[ww, ix_eff, errTable_tr, errTable_te] = biclsfy_slrvar(xtr, ttr, xte, tte, 'nlearn', 300, 'mean',  
'displaytext', 0);
```

- Input variables

Name	Brief explanation	Type	Details
x_train	Training feature matrix	N x D matrix	Each row represents a feature vector for a specific sample.
t_train	Training label vector	N x 1 vector	The $n$ th element represents the task label for the $n$ th row of training feature matrix (i.e. feature vector of the $n$ th sample). Elements of this vector must consist of any two integer values (e.g. {0,1} or {1,2}). .
x_test	Test feature matrix	M x D matrix	Same form as x_train
t_test	Test label vector	M x 1 vector	Same form as t_train

\* The variables N and M represents the number of training and test samples, respectively.

\* The variable D represents the number of features.

- Optional Input variables

Name	Brief explanation	Type	Details
scale_mode	Mode of scaling when normalizing the feature matrices	String ['each']	Each row represents a feature vector for a specific sample.
mean_mode	Mode of mean correction when normalizing the feature matrices	String ['each']	The $n$ th element represents the task label for the $n$ th row of training feature matrix (i.e. feature vector of the $n$ th sample). Elements of this vector must consist of any two integer values (e.g. {0,1} or {1,2}). .
ax0	Initial relevance parameter	Scalar [1]	This value is set to all the relevance parameters.
nlearn	The number of iterations for learning a classifier	Scalar[1000]	

nstep	The number of iteration for displaying an intermediate result during learning iterations	Scalar[100:	Invalid when 'displaytext' option is 0.
amax	Maximum relevance parameter	Scalar [10^8]	Features of which relevance parameters exceed this value are removed from a classifier. This threshold is important to avoid computational instability.
usebias	Flag for adding a constant term or not	Boolean [1]	If 1, the constant vector ones(N,1) is added to the last column of feature matrices.
norm_sep	Flag for applying normalization to training and test data using different scaling and mean correction factors.	Boolean [0]	
displaytext	Flag for displaying a result during learning iterations	Boolean [1]	
invhessian	Flag for using inverse of Hessian matrix during learning a classifier	Boolean [1]	If 1, an inverse of Hessian matrix (size D x D) is used in every learning iteration. If 0, equivalent matrix manipulations (size N x N) is used.

\* The default values usually work reasonably well.

#### - Output variables

Name	Brief explanation	Type	Details
ww	Weight vector,	D (or D+1) x 1 matrix	The last element corresponds to a constant vector if usebias = 1 or the size of a vector is D+1.
ix_eff_all	Indices of features selected	1 x 1 cell	
errTable_tr	Confusion matrix for training data	2 x 2 matrix	The 1 <sup>st</sup> row represents the number of samples classified as 1 or 2 when the true label is 1. The 2 <sup>nd</sup> row represents the number of samples classified as 1 or 2 when the true label is 2.

errTable_te	Confusion matrix for test data	2 x 2 matrix	Same form as errTable_tr
parm	Parameter struct	struct	
Axall	History of relevance vector updates	? x D vector	
Ptr	Probabilistic output of training data	N x 2 vector	The 1 <sup>st</sup> and 2 <sup>nd</sup> column represent the probability of training samples classified as label 1 and 2, respectively.
Pte	Probabilistic output of test data	M x 2 vector	Same form as Ptr

\* The variables N and M represents the number of training and test samples, respectively.

\* The variable D represents the number of features.

In the function, the following processing is successively executed;

1. check optional variables
2. normalize the feature matrices
3. add a bias regressor to the feature matrices depending on 'parm.usebias'
4. learning weight parameters (boundary parameters) of a classifier
5. evaluating percent correct for training and test data.

The step 4 is the main step where the low-level function "slr\_learning\_var2" is called.

If you know how to call this low-level function, you can modify this run-level function for your purpose easily.

## FAQ

Q1. What does the normalization step exactly do ?

A1. The normalization is simply to apply scaling and mean-correction to each element in the feature matrix so that all the elements have appropriate range (basically -1 to 1). This step is important to avoid an ill-conditioned problem in computation.

If 'mean\_mode' = 'each' (default), the mean is calculated in a feature-wise way (i.e. each feature has different mean). If 'mean\_mode' = 'all', the mean is calculated as average over all the elements (i.e all the feature shares one common mean value). In the same way, 'scale\_mode' defines the way to compute scaling factors.

If 'norm\_sep' = 0 (default), the mean and scaling is computed using only training data and then they are applied to both training and test feature matrices. On the other hand, if



'norm\_sep' = 1, the mean and scaling is separately computed for training and test data (this is not recommended).

If both 'mean\_mode' and 'scale\_mode' are each, the following affine transformation is applied,

$$\tilde{\mathbf{x}}_d = \frac{\mathbf{x}_d - m_d}{a_d} \quad \begin{cases} m_d = \text{mean}(\mathbf{x}_d) \\ a_d = \max_n \{ |x_{dn} - m_d| \} \end{cases}$$

where  $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_D] : N \times D$  is a feature matrix.

Q2. How do I compute percent correct from 'errTable\_tr' or 'errTable\_te'?

A2. Use the "calc\_percor.m" function in the toolbox like

```
> percent_tr = calc_percor(errTable_tr)
```

Q3. Why is the size of the output variable 'ww' the number of feature plus one?

A3. If 'use\_bias' = 1 (default), the bias-term is automatically concatenated to the last row of the input feature matrix. The last row or value of 'ww' is the weight of this bias regressor.

Q4. Is it possible to speed up the learning step?

A4. If you use the 'biclsfy\_slrvar.m' function, changing optional input 'invhessian' may improve the computation speed dramatically without affecting classification performance. If the dimension of your feature matrix is larger than the number of training samples, setting 'invhessian' to 0 is fast. Otherwise setting 'invhessian' to 1 is fast.

## ~~~~~ **Mathematical Basics** ~~~~~

Logistic regression (LR) is a well-known classifier originally developed in statistics. SLR is a Bayesian extension of LR in which a sparseness prior is imposed on LR. In the literature, two kinds of the sparsenss prior has been suggested; Automatic Relevance Determination (ARD) prior ([8,9]) and Laplace prior ([5]). In this toolbox the ARD prior is employed. Please see the appendix of [1] for equations of the model and derivation of the algorithm. It should be noted that for SLR having feature selection property, the boundary must be linear (but not using linear kernel). This is because in this case each feature has its own weight parameter and thereby sparse estimation of weight parameters can be interpreted as removing irrelevant features.

~~~~~

## Referencing the toolbox

~~~~~

When using this tool for a paper please refer to the following paper:

**[Yamashita O, Sato MA, Yoshioka T, Tong F, Kamitani Y \(2008\).](#)**

**Sparse estimation automatically selects voxels relevant for the decoding of fMRI activity patterns. *Neuroimage*. Oct 1;42(4):1414-29.**

The above manuscript contains basics of SLR (SLR-LAP) and applications to fMRI decoding.

~~~~~

## Feedback & Bug report

~~~~~

Any feedback and bug report are welcome. Please keep contact with me ([oyamashi@atr.jp](mailto:oyamashi@atr.jp)). I would like to respond as quickly as possible.

~~~~~

## Licencing & Copy Right

~~~~~

SLR toolbox is free but copyright software, distributed under the terms of the GNU General Public Licence as published by the Free Software Foundation. Further details on "copyleft" can be found at <http://www.gnu.org/copyleft/>. No formal support or maintenance is provided or implied.

-----

## Acknowledgements

-----

This toolbox is brought to you by ATR Computational Neuroscience laboratories in Kyoto. This research was supported in part by the NICT, Honda Research Institute, the SCOPE, SOUMU, the Nissan Science Foundation, and grants from the National Eye Institute to FT (R01 EY017082 and R01 EY14202).

-----

## References

-----  
[1] Yamashita O, Sato MA, Yoshioka T, Tong F, Kamitani Y (2008).

Sparse estimation automatically selects voxels relevant for the decoding of fMRI activity patterns. *Neuroimage*. Oct 1;42(4):1414-29.

[2] Jaakkola TS, Jordan MI (2000), Bayesian parameter estimation via variational methods, *Statistics and Computing*, 10, pp.25--37

[3] Bishop C, Tipping ME (2000), Variational relevance vector machines, *Proceedings of the 16th Conference in Uncertainty in Artificial Intelligence*, pp.46-53

[4] Bishop C (2006), Pattern recognition and machine learning, Springer, New York

[5] Tipping ME (2001), Sparse Bayesian Learning and the Relevance Vector Machine, *J Machine Learning Research*, 1, pp.211-244

[6] Krishnapuram B, Carin L, Figueiredo MAT, and Hartemink AJ (2005), Sparse Multinomial Logistic Regression Fast Algorithms and Generalization Bounds, *IEEE Trans. Pattern Analysis and Machine Intelligence*, 27, pp.957-968

[7] Krishnapuram B, Hartemink AJ, Carin L and Figueiredo MAT (2004), A Bayesian Approach to Joint Feature Selection and Classifier Design, *IEEE Trans. Pattern Analysis and Machine Intelligence*, 26, pp.1105-1111

[8] MacKay D (1992), Bayesian interpolation, *Neural Computation*, 4, pp.415-447

[9] Neal RM (1996), Bayesian Learning for Neural Networks, *Lecture Notes in Statistics* 118,. Springer