

LIBOL: A Library for Online Learning Algorithms

Version 0.3.0

Steven C.H. Hoi

Other contributors: Jialei Wang, Peilin Zhao, Ji Wan

School of Computer Engineering

Nanyang Technological University

Singapore 639798

E-mail: chhoi@ntu.edu.sg

<http://LIBOL.stevenhoi.org/>

December 12, 2013

Abstract

LIBOL is an open-source library for large-scale online learning, which consists of a large family of efficient and scalable online learning algorithms for large-scale online classification tasks. We have offered easy-to-use command-line tools and examples for users and developers. We also have made comprehensive documents available for both beginners and advanced users. LIBOL is not only a machine learning tool, but also a comprehensive experimental platform for conducting online learning research.

Contents

1	Introduction	3
1.1	Suitable Tasks	3
1.1.1	Binary Classification	3
1.1.2	Multiclass Classification	5
1.2	Main Features	6
1.3	Summary of Main Algorithms	6
1.4	How to Use the LIBOL Package	7
1.4.1	Data Formats and Folders	7
1.4.2	Examples	8
1.4.3	Parameter Settings	9
1.4.4	How to Add New Algorithms	10
1.5	Documentation	11
1.6	Design	11
2	Online Learning Algorithms for Binary Classification	12
2.1	First-Order Algorithms	12
2.1.1	Perceptron	12
2.1.2	Approximate Large Margin Algorithm (ALMA)	12
2.1.3	Relaxed Online Maximum Margin Algorithm (ROMMA)	13
2.1.4	Online Gradient Descent (OGD) Algorithm	13
2.1.5	Passive-Aggressive (PA) learning Algorithms	14
2.2	Second-Order Algorithms	16
2.2.1	Second-Order Perceptron (SOP)	16
2.2.2	Confidence-weighted learning (CW)	16
2.2.3	Improved Ellipsoid Method (IELLIP)	17
2.2.4	Adaptive Regularization Of Weights (AROW)	17
2.2.5	New Adaptive Regularization of Weights (NAROW)	19
2.2.6	Normal HERD (NHERD)	19
2.2.7	Soft Confidence-weighted learning (SCW)	21
3	Online Learning Algorithms for Multiclass Classification	23
3.1	First-Order Algorithms	23
3.1.1	Multiclass Perceptron	23
3.1.2	Multiclass Relaxed Online Maximum Margin Algorithm	24
3.1.3	Multiclass Online Gradient Descent Algorithm	24
3.1.4	Multiclass Passive Aggressive Algorithms	25
3.2	Second-Order Algorithms	27
3.2.1	Multiclass Confidence Weighted Learning	27
3.2.2	Multiclass Adaptive Regularization Of Weights	28
3.2.3	Multiclass Soft Confidence-Weighted Learning	28
4	Conclusions, Revision, and Citation	31

1 Introduction

Online learning represents an important family of efficient and scalable machine learning algorithms for large-scale applications. In general, online learning algorithms are fast, simple, and often make few statistical assumptions, making them applicable to a wide range of applications. Online learning has been actively studied in several communities, including machine learning, statistics, and artificial intelligence. Over the past years, a variety of online learning algorithms have been proposed, but so far there is very few comprehensive library which includes most of the state-of-the-art algorithms for researchers to make easy side-by-side comparisons and for developers to explore their various applications.

In this work, we develop LIBOL as an easy-to-use online learning tool that consists a large family of existing and recent state-of-the-art online learning algorithms for large-scale online classification tasks. In contrast to many existing software for large-scale data classification, LIBOL enjoys significant advantages for massive-scale classification in the era of big data nowadays, especially in efficiency, scalability, parallelization, and adaptability. The software is available at <http://libol.stevenhoi.org/>.

1.1 Suitable Tasks

Currently the implemented toolbox is suitable for online binary classification and multi-class classification tasks, which are discussed with more details as follows.

1.1.1 Binary Classification

Online binary classification operates on a sequence of data examples with time stamps. At each step t , the learner receives an incoming example $\mathbf{x}_t \in \mathcal{X}$ in a d -dimensional vector space, i.e., $\mathcal{X} = \mathbb{R}^d$. It first attempts to predict the class label of the incoming instance,

$$\hat{y}_t = \text{sgn}(f(\mathbf{x}_t; \mathbf{w}_t)) = \text{sgn}(\mathbf{w}_t \cdot \mathbf{x}_t) \in \mathcal{Y}$$

and $\mathcal{Y} = \{-1, +1\}$ for binary classification tasks. After making the prediction, the true label $y_t \in \mathcal{Y}$ is revealed, and the learners then computes the loss $\ell(y_t, \hat{y}_t)$ based on some criterion to measure the difference between the learner’s prediction and the revealed true label y_t . Based on the result of the loss, the learner finally decides when and how to update the classification model at the end of each learning step. The following algorithmic framework gives an overview of most online learning algorithms ¹ for linear (binary) classification tasks, where $\Delta(\mathbf{w}_t; (\mathbf{x}_t, y_t))$ denotes the update function for upgrading the classification models. Different online learning algorithms in general are distinguished in terms of different definitions and designs of the loss function $\ell(\cdot)$ and their various update functions $\Delta(\cdot)$, although they might be founded on different theories and principles.

¹Second-order online learning algorithms follow a slightly different procedure.

Algorithm 1: LIBOL: A Framework for Online Linear Binary Classification.

```
1 Initialize:  $\mathbf{w}_1 = 0$ 
2 for  $t = 1, 2, \dots, T$  do
3   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
4   The learner predicts the class label:  $\hat{y}_t = \text{sgn}(f(\mathbf{x}_t; \mathbf{w}_t))$ ;
5   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
6   The learner calculates the suffered loss:  $\ell(\mathbf{w}_t; (\mathbf{x}_t, y_t))$ ;
7   if  $\ell(\mathbf{w}_t; (\mathbf{x}_t, y_t)) > 0$  then
8     The learner updates the classification model:
9      $\mathbf{w}_{t+1} \leftarrow \Delta(\mathbf{w}_t; (\mathbf{x}_t, y_t))$ ;
10  end
11 end
```

In particular, this software consists of 16 different online algorithms and their variants for binary classification, and 13 online algorithms and variants for multiclass classification. In general, they can be grouped into two major categories: (i) first-order learning [13, 2], and (ii) second-order learning [4, 14, 15]. Examples algorithms in the first-order learning category include the following list of classical and popular algorithms:

- Perceptron: the classical online learning algorithm [13];
- ALMA: Approximate Maximal Margin Classification Algorithm [10];
- ROMMA: the Relaxed Online Maximum Margin algorithms [11];
- OGD: the Online Gradient Descent (OGD) algorithms [18];
- PA: Passive Aggressive (PA) algorithms [2];

In recent years, to improve the efficacy of first-order learning methods, the second-order online learning algorithms have been actively explored. One major family of second order learning algorithms typically assume the weight vector follows a Gaussian distribution $\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ with mean vector $\boldsymbol{\mu} \in \mathbb{R}^d$ and covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$. The model parameters, including both the mean vector and the covariance matrix are updated in the online learning process. Example of the second-order online learning algorithms include the following:

- SOP: the Second-Order Perceptron (SOP) algorithm [1];
- CW: the Confidence-Weighted (CW) learning algorithm [4];
- IELLIP: online learning algorithms by improved ellipsoid method [15];
- AROW: the Adaptive Regularization of Weight Vectors [5];
- NAROW: New variant of Adaptive Regularization [12];
- NHERD: the Normal Herding method via Gaussian Herding [7]
- SCW: the recently proposed Soft Confidence Weighted algorithms [14].

For the details of each of the above algorithms, please refer to the detailed descriptions in Section 2.

1.1.2 Multiclass Classification

Similar to online binary classification, online *multiclass* learning is performed over a sequence of training examples $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$. Unlike binary classification where $y_t \in \{-1, +1\}$, in multiclass learning, each class assignment $y_t \in \mathcal{Y} = \{1, \dots, k\}$, making it a more challenging problem. We use \hat{y}_t to represent the class label predicted by the online learning algorithm. Online multiclass classification algorithms [8, 9] learn multiple hypotheses/classifiers, one classifier for each class in \mathcal{Y} , leading to a total of k classifiers that are trained for the classification task. The predicted label is the one, which is associated with the largest prediction value, i.e.

$$\hat{y}_t = \arg \max_{i \in \{1, \dots, k\}} \mathbf{w}_{t,i} \cdot \mathbf{x}_t.$$

After the prediction, the true label $y_t \in \mathcal{Y}$ will be disclosed, the learner then compute the loss based on some criterion to measure the difference between the prediction and the true label. Based on the results of the loss, the learner finally decides when and how to update the k classifiers at the end of each learning step.

In particular, this software package consists of 13 different online multiclass classification algorithms and their variants. In general, they can be grouped into two major categories: (i) first-order learning [8, 2], and (ii) second-order learning [3, 6].

Specifically, the first-order learning algorithms only keep updating k classification functions (for the k different labels), which only utilize the first-order information of the received instances. The main first-order algorithms implemented in this toolbox include:

- Perceptron: the classical multiclass online learning algorithm [8];
- ROMMA: the Multiclass Relaxed Online Maximum Margin algorithms [11];
- OGD: the Multiclass Online Gradient Descent (OGD) algorithms [18];
- PA: Multiclass Passive Aggressive (PA) algorithms [2];

Unlike the first-order learning algorithms, the second-order online learning aim to better exploit the underlying structures between features. Specifically, the second-order online learning algorithms typically assume the weight vector follows Gaussian distributions $\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ with mean vector $\boldsymbol{\mu} \in \mathbb{R}^{kd}$ and covariance matrix $\Sigma \in \mathbb{R}^{kd \times kd}$. The main second-order algorithms implemented in this toolbox include:

- CW: the Multiclass Confidence-Weighted (CW) learning algorithm [3];
- AROW: the Multiclass Adaptive Regularization of Weight Vectors [6];
- SCW: the recently proposed Soft Confidence Weighted algorithms [14].

For the details of all the algorithms, please refer to the section 2.

1.2 Main Features

- **Comprehensiveness:** A large family of existing online binary and multiclass classification algorithms have been implemented in this software package;
- **Extendibility:** One can easily implement a new algorithm by only implementing the key updating module;
- **Usability:** It is easy to use the main functions to evaluate one algorithm by comparing with all the algorithms;

1.3 Summary of Main Algorithms

Table 1 gives a summary of the family of implemented algorithms in this software package. For simplicity, we exclude some incremental variants of the algorithms in this list.

Table 1: Summary of the Implemented Algorithms.

Problem Type	Methodology	Algorithm	Description	Section
Binary Classification	First-Order	Perceptron	The perceptron algorithm [13]	2.1.1
		ALMA	Approximate Large Margin Algorithm [10]	2.1.2
		ROMMA	Relaxed Online Maximum Margin Algorithms [11]	2.1.3
		OGD	Online Gradient Descent [18]	2.1.4
		PA	Passive Aggressive (PA) algorithms [2]	2.1.5
		SOP	Second-Order Perceptron [1]	2.2.1
	Second-Order	CW	Confidence-Weighted (CW) learning [4]	2.2.2
		IELLIP	Improved Ellipsoid method [15]	2.2.3
		AROW	Adaptive Regularization of Weight Vectors [5]	2.2.4
		NAROW	New variant of Adaptive Regularization [12]	2.2.5
		NHERD	Normal Herding method via Gaussian Herding [7]	2.2.6
		SCW	Soft Confidence Weighted algorithms [14]	2.2.7
Multiclass Classification	First-Order	Perceptron	The perceptron algorithm [13]	3.1.1
		ROMMA	Approximate Large Margin Algorithm [10]	3.1.2
		OGD	Online Gradient Descent [18]	3.1.3
		PA	Passive Aggressive (PA) algorithms [2]	3.1.4
	Second-Order	CW	Confidence-Weighted (CW) learning [4]	3.2.1
		AROW	Adaptive Regularization of Weight Vectors [5]	3.2.2
		SCW	Soft Confidence Weighted algorithms [14]	3.2.3

1.4 How to Use the LIBOL Package

The current version of LIBOL package includes a MATLAB library, a C library and command-line tools for the learning task. The data formats used by this software package are compatible with some of the very popular machine learning and data mining packages, such as LIBSVM, SVM-light, and WEKA, etc.

To evaluate an algorithm on a dataset with a specific format, one can call the function:

```
$ demo(TaskType, Algorithm, Dataset, DataFormat)
```

where the details of the input parameters are listed as in the Table 2.

To easily compare all the online learning algorithms for one problem type on one dataset, you can call the functions as

```
$ run_experiment(TaskType, Dataset, DataFormat)
```

where the details of the input parameters are listed as in the Table 2.

Table 2: The input parameters for function “demo.m”.

Parameter	Description
TaskType	The type of the online learning problem. Two optional values: ‘bc’ and ‘mc’. ‘bc’ denotes binary classification. ‘mc’ denotes multiclass classification.
Algorithm	The name of the evaluated algorithm, which is in the folder - algorithms. For example, ‘Perceptron’ and ‘Perceptron_c’ for the ‘bc’ problem. ‘Perceptron’ and ‘Perceptron_c’ are Matlab and C versions, respectively. For example, ‘PerceptronM’ for the ‘mc’ problem.
Dataset	The name of the dataset for evaluation For example: svmguide3. For example: ‘ionosphere.arff’.
DataFormat	The format of the dataset. By default, use ‘mat’ version ‘libsvm’ and ‘arff’ denote the format for LIBSVM and WEKA

1.4.1 Data Formats and Folders

The types of data formats supported by this software package include (i) the matlab data format (denoted as “mat” for short), (ii) the “libsvm” data format (commonly used in LIBSVM and SVM-light), and the “arff” data format used in WEKA. Specifically, for the “mat” data format, it must contain a $n \times d$ data matrix (“data”), in which the first column stores the true class labels of n data points, and the rest columns store the d -dimensional features of the n data points in the data set. More specifically, given the data matrix “data”, we can retrieve feature matrix and the label vector by $x_t = \text{data}(:, 2 : d)$; and $y = \text{data}(:, 1)$; respectively.

All the datasets to be evaluated on should be placed into the folder “\LIBOL\data”. The dataset can be in formats of “mat” (MATLAB), LIBSVM, and WEKA versions.

The Matlab and C versions of the algorithms listed in subsection 1.1 are located at the folders “\LIBOL\algorithms” and “\LIBOL\algorithms_mex”, respectively. The main functions “demo” and “run_experiment” are located at the folder “\LIBOL”. To use these functions, you can use the command “addpath(genpath(['root\LIBOL']));”, where “root” is the path to the folder “LIBOL”.

1.4.2 Examples

To illustrate the online learning procedure, we take two data sets from the LIBSVM website, including one small data set “svmguide3” with 1243 instances and one large data set “ijcnn1” with 141,691 instances. In the following example, we use the default “Perceptron” algorithm to demo the usage of LIBOL for a binary classification (‘bc’) task:

```
$ demo('bc', 'Perceptron', 'svmguide3')
```

The results output by the above command are summarized as follows:

Algorithm:	mistake rate	nb of updates	cpu time (seconds)
Perceptron	0.3318 +/- 0.0118	412.45 +/- 14.66	0.0516 +/- 0.0008

To ease researchers to run a full set of experiments for side-by-side comparisons of different algorithms, we offer a very easy-to-use example program as follows:

```
$ run_experiment('bc', 'svmguide3')
```

The above command will run side-by-side comparison of varied online learning algorithms on the given data set fully automatically, including all the parameter settings and selection. The full set of experimental results will be generated by the library automatically, as shown in Table 1 and Figure 1.

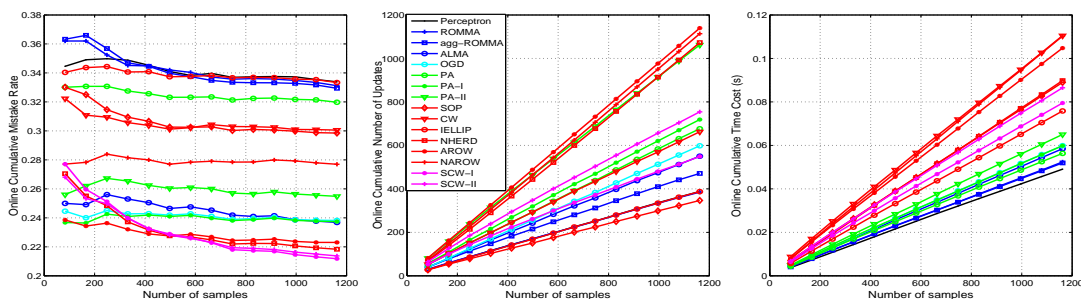


Figure 1: Comparison of a variety of online learning algorithms on dataset “svmguide3”.

Table 3: Comparison of a variety of online learning algorithms on two data sets.

Dataset:	svmguid3 (#samples=1243,#dimensions=36)			ijcnn1 (#samples=141,691,#dimensions=22)		
Algorithm	mistake	# updates	time (s)	mistake	# updates	time (s)
Perceptron	0.332 ± 0.012	412.4 ± 14.7	0.052 ± 0.002	0.106 ± 0.000	15059.9 ± 65.1	5.668 ± 0.064
ROMMA	0.329 ± 0.019	409.3 ± 23.2	0.056 ± 0.001	0.101 ± 0.001	14284.2 ± 89.8	5.823 ± 0.111
aROMMA	0.328 ± 0.018	500.1 ± 29.1	0.055 ± 0.001	0.101 ± 0.001	14776.0 ± 101.8	5.665 ± 0.062
ALMA	0.230 ± 0.006	592.9 ± 6.6	0.062 ± 0.001	0.071 ± 0.000	21474.0 ± 80.4	6.662 ± 0.127
OGD	0.237 ± 0.003	636.5 ± 4.2	0.064 ± 0.002	0.095 ± 0.000	27465.8 ± 31.1	6.369 ± 0.107
PA	0.318 ± 0.013	721.1 ± 18.3	0.060 ± 0.001	0.102 ± 0.001	33847.9 ± 135.2	5.955 ± 0.091
PA1	0.236 ± 0.002	763.4 ± 11.5	0.064 ± 0.001	0.077 ± 0.000	28376.3 ± 84.2	6.352 ± 0.109
PA2	0.253 ± 0.007	1131.5 ± 15.6	0.069 ± 0.001	0.081 ± 0.000	61093.8 ± 199.3	6.876 ± 0.114
SOP	0.297 ± 0.012	369.1 ± 14.8	0.095 ± 0.002	0.102 ± 0.001	14470.7 ± 81.3	10.616 ± 0.096
IELLIP	0.332 ± 0.013	412.7 ± 16.1	0.081 ± 0.002	0.119 ± 0.001	16876.8 ± 72.9	8.079 ± 0.082
CW	0.299 ± 0.011	704.6 ± 19.1	0.118 ± 0.002	0.093 ± 0.001	30648.3 ± 166.2	9.499 ± 0.110
NHERD	0.217 ± 0.007	1150.5 ± 27.7	0.096 ± 0.002	0.084 ± 0.001	86660.7 ± 2692.6	9.735 ± 0.133
AROW	0.222 ± 0.004	1219.5 ± 8.1	0.112 ± 0.002	0.082 ± 0.000	74247.1 ± 846.3	10.164 ± 0.069
NAROW	0.276 ± 0.042	1193.8 ± 23.2	0.118 ± 0.002	0.095 ± 0.008	103843.6 ± 8841.3	12.027 ± 0.467
SCW	0.206 ± 0.004	593.4 ± 13.9	0.085 ± 0.002	0.060 ± 0.002	11077.3 ± 678.3	7.921 ± 0.144
SCW2	0.212 ± 0.009	802.0 ± 73.2	0.092 ± 0.002	0.070 ± 0.001	30833.8 ± 2116.8	8.681 ± 0.150

1.4.3 Parameter Settings

Setting proper parameters plays a nontrivial role in affecting the empirical performance of different online learning algorithms. Some online learning algorithms (e.g., Perceptron) are parameter-free, but some may have multiple parameters. Table 4 gives a summary of parameters and their default settings by different online learning algorithms. We note that different algorithms may use different notations for the same kind of parameter. To be consistent, we use the same notation for the parameter with the same/similar meaning, and will indicate their original notation in appropriate place. Below we discuss a few commonly used parameters.

The first common parameter is C , which is typically used to trade off between a regularization term and a loss term. For example, in passive-aggressive (PA) learning algorithms, parameter C is to trade off between the passiveness (regularization) and aggressiveness (loss). One exceptional case is OGD where we use C as a constant for defining the learning rate, $\frac{C}{\sqrt{t}}$, where $C > 0$ and t is the iteration no. The second parameter is η , which is commonly used in some second order algorithms. For example, for the family of confidence-weighted learning algorithms, η is a parameter used in defining a key parameter Φ of the loss function, i.e., $\phi = \Phi^{-1}(\eta)$, where Φ is the cumulative function of the normal distribution. The last parameter is a , typically used for initializing the covariance matrix in the second order algorithms, i.e, $\Sigma = a * I$, where I is an identity matrix. For most cases, parameter a is not too sensitive and typically fixed to 1.

Finally, to enable fair side-by-side comparisons between different algorithms, we choose the best parameters for each algorithm automatically via a parameter validation scheme (“CV_algorithm.m”). In particular, we randomly draw a permutation of the whole dataset as the input training data to determine the best parameters for different algorithms. For the above commonly used parameters, we search for the best parameter C in the range of $[2^{-4}, 2^{-3}, \dots, 2^4]$, and the best parameter η in the range of $[0.55, 0.60, \dots, 0.95]$. More detailed settings can be found in (“CV_algorithm.m”).

Table 4: Summary of default parameter settings by different algorithms.

Function	Algorithm	C	a	eta	Other parameters/Remark	Section
Perceptron()	Perceptron	NA	NA	NA	parameter-free	2.1.1
ALMA()	ALMA	$C=\sqrt{2}$	NA	NA	$\alpha = 0.9, B=\frac{1}{\alpha}, \text{norm } p=2$	2.1.2
ROMMA()	ROMMA	NA	NA	NA	parameter-free	2.1.3
AROMMA()	aggressive ROMMA	NA	NA	NA	parameter-free	2.1.3
OGD()	OGD	$C=1$	NA	NA	$\eta_t = \frac{C}{\sqrt{t}}, \text{loss_type} = 1$	2.1.4
PA()	PA	NA	NA	NA	parameter-free	2.1.5
PA1()	PA-I	$C=1$	NA	NA	NA	2.1.5
PA2()	PA-II	$C=1$	NA	NA	NA	2.1.5
SOP()	Second-order Perceptron	NA	$a=1$	NA	parameter-free	2.2.1
CW()	CW	NA	$a=1$	eta=0.70	$\Sigma = a * I$	2.2.2
IELLIP()	Improved Ellipsoid	NA	$a=1$	NA	$b=0.3, c=0.1$	2.2.3
AROW()	AROW	$C=1$	$a=1$	NA	$r=C, \Sigma = a * I$	2.2.4
NAROW()	New AROW variant	$C=1$	$a=1$	$\Sigma = a * I$	$b = C$	2.2.5
NHERD()	Normal Herd	$C=1$	$a=1$	NA	$\gamma = \frac{1}{C}, \Sigma = a * I$	2.2.6
SCW()	SCW-I	$C=1$	$a=1$	eta=0.75	$\Sigma = a * I$	2.2.7
SCW2()	SCW-II	$C=1$	$a=1$	eta=0.90	$\Sigma = a * I$	2.2.7
M_PerceptronM()	multi-class Perceptron	NA	NA	NA	parameter-free	2.1.1
M_PerceptronS()	multi-class Perceptron	NA	NA	NA	parameter-free	2.1.1
M_PerceptronU()	multi-class Perceptron	NA	NA	NA	parameter-free	2.1.1
M_ROMMA()	multi-class ROMMA	NA	NA	NA	parameter-free	2.1.3
M_AROMMA()	multi-class aROMMA	NA	NA	NA	parameter-free	2.1.3
M_OGD()	multi-class OGD	NA	NA	NA	$\eta_t = \frac{C}{\sqrt{t}}, \text{loss_type} = 1$	2.1.4
M_PA()	multi-class PA	NA	NA	NA	parameter-free	2.1.5
M_PA1()	multi-class PA-I	$C=1$	NA	NA	NA	2.1.5
M_PA2()	multi-class PA-II	$C=1$	NA	NA	NA	2.1.5
M_CW()	multi-class CW	NA	$a=1$	eta=0.75	$\Sigma = a * I$	2.2.2
M_AROW()	multi-class AROW	$C=1$	$a=1$	NA	$r=C, \Sigma = a * I$	2.2.4
M_SCW()	multi-class SCW-I	$C=1$	$a=1$	eta=0.75	$\Sigma = a * I$	2.2.7
M_SCW2()	multi-class SCW-II	$C=1$	$a=1$	eta=0.75	$\Sigma = a * I$	2.2.7

1.4.4 How to Add New Algorithms

A salient property of this library is that it provides a fairly easy-to-use testbed to facilitate online learning researchers to develop their new algorithms and conduct side-by-side comparisons with the state-of-the-art algorithms with the minimal efforts. To facilitate this for beginners, we have provided a template named “NEW_ALGORITHM.m” in which a user can quickly implement their new online algorithm with the step by step instructions. More specifically, adding a new algorithm has to address two major issues:

- What is the condition for making an update? This is usually equivalent to defining a proper loss function (e.g., a hinge loss $l_t = \max(0, 1 - y_t * f_t)$) such that an update occurs wherever the loss is nonzero, i.e., ($l_t > 0$).
- How to perform the update on the classifier (i.e., the weight vector \mathbf{w}) whenever

the condition is satisfied? For example, Perceptron updates $w = w + y_t * x_t$;

- Are there some parameters in your new algorithm? If so, you need to do some initializations, including (i) modify the "init_options.m" file by adding the initialization (under the line of "case 'NEW_ALGORITHM'"); (ii) modify the "CV_algorithm.m" by adding the auto parameter selection process (under the line of "case 'NEW_ALGORITHM'").

1.5 Documentation

The LIBOL package comes with comprehensive documentation. The README file describes the setup and usage. Users can read the "Quick Start" section to begin shortly. All the functions and related data structures are explained in detail. If the README file does not give the information users want, they can also check the online FAQ. In addition to software documentation, theoretical properties of some algorithms and comparisons can be found in [14]. The authors are also willing to answer any further questions.

1.6 Design

The design principle is to keep the package simple, easy to read and extend. All codes follow the MATLAB (and Octave-compatible) standard and need no external libraries (except for the support of popular data formats, such as LIBSVM and WEKA datasets for which existing libraries are included). The reason to choose MATLAB/Octave is to facilitate machine learning researchers who can quickly implement a new algorithm with a new idea, and compare it with a large family of existing algorithms without spending much time and efforts in programming parts (C/C++ or Java). We also understand that purely MATLAB and Octave implementations are not efficient for large-scale applications. To address this limitation, we have also provided C/C++ implementations for the core learning functions. Both MATLAB and C implementations produce the identical results except different execution times.

When implementing a new algorithm, a user can first implement a MATLAB version quickly, and then test the efficacy of their MATLAB algorithm by comparing with the family of existing algorithms on smaller data sets. After that, the user can then decide if it is necessary to do C/C++ implementation for their new algorithm. For the library design, all the online learning algorithms can be called via the uniform "ol_train(Y,X,options)" function by setting proper options, where Y denotes the true label vector, X denotes the feature matrix, and "options" is a "struct" with a set of predefined parameters. More details about this function can be found in Appendix A of this manual.

In general, LIBOL is written in a modular way, in which one can easily develop a new algorithm and make side-by-side comparisons with the existing ones in the package. Thus, we hope that LIBOL is not only a machine learning tool, but also a comprehensive experimental platform for conducting online learning research.

2 Online Learning Algorithms for Binary Classification

In a regular binary classification task, the goal of online learning for classification is to minimize the cumulative mistake number over the entire sequence of data examples. There are some other atypical settings where other learning objectives are preferred. In literature, many linear algorithms have been proposed for online classification. They can be generally grouped into two major categories: (i) first-order online learning algorithms, and (ii) second-order online learning algorithms; Below we review the basics of these algorithms in detail.

2.1 First-Order Algorithms

In this section, we provide the details of the first-order learning algorithms and the corresponding functions (located at /algorithms/Binary-Classification). The first-order learning algorithms only keep updating one classification function, which only utilizes the first-order information of the examples.

Specifically, the struct “model” in the codes contains a field “ \mathbf{w} ”, which is the linear classifier updated round by round. In addition, “model” also contains various parameters for different algorithms, which will be explained in every algorithm.

2.1.1 Perceptron

The Perceptron [13] algorithm is the earliest and simplest approach for online learning.

Algorithm 2: Perceptron.

```
1 Initialize:  $\mathbf{w}_1 = 0$ 
2 for  $t = 1, 2, \dots, T$  do
3   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
4   The learner predicts the class label:  $\hat{y}_t = \text{sign}(\mathbf{w}^\top \mathbf{x}_t)$ ;
5   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
6   The learner calculates the suffered loss:  $l_t = \mathbb{I}(\hat{y}_t \neq y_t)$ ;
7   if  $l_t > 0$  then
8     The learner updates the classification model:
9      $\mathbf{w} = \mathbf{w} + y_t \mathbf{x}_t$ ;
10  end
11 end
```

Perceptron does not contain any parameters.

2.1.2 Approximate Large Margin Algorithm (ALMA)

The Approximate Large Margin Algorithm (ALMA) [10] aims to learn an approximate large margin classifier.

Algorithm 3: ALMA: Approximate Large Margin Algorithm .

```
1 Parameters:  $p \in [2, \dots, 10]$ (default=2),  $\alpha \in (0, 1]$ ,  $C > 0$  (default= $\sqrt{2}$ ),  $B = \frac{1}{\alpha}$ 
2 Initialize:  $\mathbf{w}_1 = 0$ 
3 for  $t = 1, 2, \dots, T$  do
4   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
5   The learner predicts the class label:  $\hat{y}_t = \text{sign}(\mathbf{w}^\top \mathbf{x}_t)$ ;
6   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
7   The learner calculates the suffered loss:  $l_t = \max(0, (1 - \alpha)B \frac{\sqrt{p-1}}{\sqrt{k}} - \frac{y_t \mathbf{w}_t \cdot \mathbf{x}_t}{\|\mathbf{x}_t\|})$ ;
8   if  $l_t > 0$  then
9     The learner updates the classification model:
10     $\mathbf{w} = (\mathbf{w} + \frac{C}{\sqrt{p-1}\sqrt{k}} y_t \frac{\mathbf{x}_t}{\|\mathbf{x}_t\|}) / \max(1, \|\mathbf{w} + \frac{C}{\sqrt{p-1}\sqrt{k}} y_t \frac{\mathbf{x}_t}{\|\mathbf{x}_t\|}\|)$ 
11     $k = k + l_t$ ;
12  end
13 end
```

The implemented algorithm is $\text{ALMA}_p(\alpha)$, with the norm parameter $p = 2$. The “model” contains the parameter:

- $\alpha \in (0, 1]$: which determines the final margin $(1 - \alpha)\gamma_*$, where γ_* is the best margin
- C : typically fixed to the optimal value $\sqrt{2}$ as suggested by the theorem in the paper;
- B : typically set to the optimal value $\frac{1}{\alpha}$ as suggested by the theorem in the paper.

2.1.3 Relaxed Online Maximum Margin Algorithm (ROMMA)

The Relaxed Online Maximum Margin Algorithm(ROMMA) [11] and its aggressive version agg-ROMMA. The ROMMA and agg-ROMMA algorithm do not contain parameters.

2.1.4 Online Gradient Descent (OGD) Algorithm

The online gradient descent algorithm in [18] exploits the gradient descent updating approach for optimizing the objective function defined by different type of loss functions. In our implementation, we set the learning rate η_t to $\frac{C}{\sqrt{t}}$, where $C > 0$ is a constant learning rate parameter and t is the number of learning rounds. We have implemented several different types of losses, and a user can choose different types of loss functions by setting the parameter `loss_type`, e.g., setting `loss_type = 0` for 0-1 loss, `loss_type = 1` for hinge loss, `loss_type = 2` for logistic loss, and `loss_type = 3` for square loss.

Algorithm 4: ROMMA: Relaxed Online Maximum Margin Algorithm.

```

1 Initialize:  $\mathbf{w}_1 = 0$ 
2 for  $t = 1, 2, \dots, T$  do
3   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
4   The learner predicts the class label:  $\hat{y}_t = \text{sign}(\mathbf{w}^\top \mathbf{x}_t)$ ;
5   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
6   The learner calculates the suffered loss:
      
$$l_t = \begin{cases} \mathbb{I}(\hat{y}_t \neq y_t) & \text{ROMMA ,} \\ \max(0, 1 - y_t(\mathbf{w}_t \cdot \mathbf{x}_t)) & \text{agg-ROMMA .} \end{cases}$$

7   if  $l_t > 0$  then
8     The learner updates the classification model:
      
$$\mathbf{w}_{t+1} = \left( \frac{\|\mathbf{x}_t\|^2 \|\mathbf{w}_t\|^2 - y_t \mathbf{w}_t \cdot \mathbf{x}_t}{\|\mathbf{x}_t\|^2 \|\mathbf{w}_t\|^2 - (\mathbf{w}_t \cdot \mathbf{x}_t)^2} \right) \mathbf{w}_t + \left( \frac{\|\mathbf{w}_t\|^2 (y_t - \mathbf{w}_t \cdot \mathbf{x}_t)}{\|\mathbf{x}_t\|^2 \|\mathbf{w}_t\|^2 - (\mathbf{w}_t \cdot \mathbf{x}_t)^2} \right) \mathbf{x}_t$$

9   end
10 end
```

2.1.5 Passive-Aggressive (PA) learning Algorithms

The family of online Passive-Aggressive (PA) learning algorithms [2] is formulated to trade off the objective of minimizing the distance between the learnt classifier and the previous classifier, and the objective of minimizing the loss of the learnt classifier suffered on the current instance.

Formally, the PA algorithm is formulated as the following online optimization problem

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_t\|^2, \quad \text{s.t. } \ell(\mathbf{w}; (\mathbf{x}_t, y_t)) = \max(0, 1 - y_t(\mathbf{w} \cdot \mathbf{x}_t)) = 0$$

Furthermore, PA is extended to PA-I and PA-II algorithms, which can better handle the non-separable and noisy case. PA-I is formulated as

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_t\|^2 + C \ell(\mathbf{w}; (\mathbf{x}_t, y_t))$$

where $\ell(\mathbf{w}; (\mathbf{x}_t, y_t)) = \max(0, 1 - y_t(\mathbf{w} \cdot \mathbf{x}_t))$ and $C > 0$, which is used to trade of the passiveness and aggressiveness.

Finally, PA-II is formulated as

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_t\|^2 + C (\ell(\mathbf{w}; (\mathbf{x}_t, y_t)))^2$$

where $\ell(\mathbf{w}; (\mathbf{x}_t, y_t)) = \max(0, 1 - y_t(\mathbf{w} \cdot \mathbf{x}_t))$ and $C > 0$, which is used to trade of the passiveness and aggressiveness.

The struct ‘‘model’’ contains the following parameter

- $C > 0$, which trades off the passiveness and aggressiveness

Algorithm 5: OGD: Online Gradient Descent.

1 Parameters: loss_type (default=1), $C > 0$ (default=1);
2 Initialize: $\mathbf{w}_1 = 0$;
3 **for** $t = 1, 2, \dots, T$ **do**
4 The learner receives an incoming instance: $\mathbf{x}_t \in \mathcal{X}$;
5 The learner predicts the class label: $\hat{y}_t = \text{sign}(\mathbf{w}^\top \mathbf{x}_t)$;
6 The true class label is revealed from the environment: $y_t \in \mathcal{Y}$;
7 The learner calculates the suffered loss:

$$l_t = \begin{cases} \mathbb{I}(y_t \neq \hat{y}_t) & \text{loss_type=0 (0-1 loss),} \\ \max(0, 1 - y_t(\mathbf{w}_t \cdot \mathbf{x}_t)) & \text{loss_type=1 (hinge loss),} \\ \log(1 + \exp(-y_t(\mathbf{w}_t \cdot \mathbf{x}_t))) & \text{loss_type=2 (logistic loss),} \\ 0.5 * (\mathbf{w}_t \cdot \mathbf{x}_t - y_t)^2 & \text{loss_type=3 (square loss),} \end{cases}$$

8 **if** $l_t > 0$ **then**
9 Set learning rate: $\eta_t = \frac{C}{\sqrt{t}}$
9 The learner updates the classification model:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \begin{cases} \eta_t y_t \mathbf{x}_t & \text{loss_type=0 (0-1 loss),} \\ \eta_t y_t \mathbf{x}_t & \text{loss_type=1 (hinge loss),} \\ \eta_t y_t \mathbf{x}_t * \frac{1}{1 + \exp(y_t(\mathbf{w}_t \cdot \mathbf{x}_t))} & \text{loss_type=2 (logistic loss),} \\ \eta_t (y_t - \mathbf{w}_t \cdot \mathbf{x}_t) \mathbf{x}_t & \text{loss_type=3 (square loss),} \end{cases}$$

10 **end**
11 **end**

Algorithm 6: PA: Passive-Aggressive learning algorithms.

```
1 Parameter:  $C > 0$  (default=1)
2 Initialize:  $\mathbf{w}_1 = 0$ 
3 for  $t = 1, 2, \dots, T$  do
4   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
5   The learner predicts the class label:  $\hat{y}_t = \text{sign}(\mathbf{w}^\top \mathbf{x}_t)$ ;
6   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
7   The learner calculates the suffered loss:  $l_t = \max(0, 1 - y_t(\mathbf{w}_t \cdot \mathbf{x}_t))$ ;
8   if  $l_t > 0$  then
9     The learner updates the classification model:
10      
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \begin{cases} \frac{\ell(\mathbf{w}_t; (\mathbf{x}_t, y_t))}{\|\mathbf{x}_t\|^2} y_t \mathbf{x}_t & \text{PA ,} \\ \min\{C, \frac{\ell(\mathbf{w}_t; (\mathbf{x}_t, y_t))}{\|\mathbf{x}_t\|^2}\} y_t \mathbf{x}_t & \text{PA-I ,} \\ \frac{\ell(\mathbf{w}_t; (\mathbf{x}_t, y_t))}{\|\mathbf{x}_t\|^2 + \frac{1}{2C}} y_t \mathbf{x}_t & \text{PA-II.} \end{cases}$$

11   end
12 end
```

2.2 Second-Order Algorithms

To better exploring the underlying structure between features, the second-order online learning algorithms typically assume the weight vector follows a Gaussian distribution $\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ with mean vector $\boldsymbol{\mu} \in \mathbb{R}^d$ and covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$.

Specifically, the struct “model” in the codes contains two fields “ $\boldsymbol{\mu}$ ” and “ Σ ”, which are the mean vector and covariance matrix for the learnt Gaussian distribution. In addition, “model” also contains various parameters for different algorithms, which will be explained in every algorithm.

2.2.1 Second-Order Perceptron (SOP)

The Second-Order Perceptron(SOP) [1] could be considered as the second-order counterpart of Perceptron by maintaining a covariance matrix Σ_t .

The struct “model” contains the following parameter:

- $a > 0$, which is used to initialize the $\Sigma_1 = a\mathbf{I}$, where \mathbf{I} is identity matrix.

2.2.2 Confidence-weighted learning (CW)

In confidence-weighted (CW) learning [4], the weight distribution $\mathcal{N}(\boldsymbol{\mu}, \Sigma)$ is updated by minimizing the Kullback-Leibler divergence between the new weight distribution and the old one while ensuring that the probability of correct classification is greater than a

Algorithm 7: SOP: Second-Order Perceptron.

```
1 Parameter:  $a > 0$  (default=1)
2 Initialize:  $\mathbf{w}_1 = 0, \Sigma_1 = a\mathbf{I}$ 
3 for  $t = 1, 2, \dots, T$  do
4   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
5   The learner predicts the class label:  $\hat{y}_t = \text{sign}(\mathbf{w}_t^\top \mathbf{x}_t)$ , where
    $\mathbf{w}_t = (\Sigma_t + \mathbf{x}_t \mathbf{x}_t^\top)^{-1} \boldsymbol{\mu}_t$ ;
6   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
7   The learner calculates the suffered loss:  $l_t = \mathbb{I}(\hat{y}_t \neq y_t)$ ;
8   if  $l_t > 0$  then
9     The learner updates the classification model:
10     $\boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t + y_t \mathbf{x}_t, \quad \Sigma_{t+1} = \Sigma_t + \mathbf{x}_t \mathbf{x}_t^\top$ ;
11  end
12 end
```

threshold as follows:

$$(\boldsymbol{\mu}_{t+1}, \Sigma_{t+1}) = \arg \min_{\boldsymbol{\mu}, \Sigma} D_{KL}(\mathcal{N}(\boldsymbol{\mu}, \Sigma), \mathcal{N}(\boldsymbol{\mu}_t, \Sigma_t))$$
$$s.t. \quad Pr_{\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)}[y_t(\mathbf{w} \cdot \mathbf{x}_t) \geq 0] \geq \eta$$

The struct “model” contains the following parameters:

- $a > 0$, which is used to initialize the $\Sigma_1 = a\mathbf{I}$, where \mathbf{I} is identity matrix;
- $\eta \in (0.5, 1]$, which is the probability required for the updated distribution on the current instance

2.2.3 Improved Ellipsoid Method (IELLIP)

The improved ellipsoid method [15] is an improved version of the classical ellipsoid method for online learning, which is modified so that it is able to address the inseparable case.

The struct “model” contains the following parameters:

- $c \in (0, 1], b \in (0, 1)$: parameters controlling the memory of online learning

2.2.4 Adaptive Regularization Of Weights (AROW)

Unlike the original CW learning algorithm, the Adaptive Regularization Of Weights (AROW) learning introduces the adaptive regularization of the prediction function when processing each new instance in each learning step, making it more robust than CW to sudden changes of label noise in the learning tasks. In particular, the optimization is formulated as follows:

$$(\boldsymbol{\mu}_{t+1}, \Sigma_{t+1}) = \arg \min_{\boldsymbol{\mu}, \Sigma} D_{KL}(\mathcal{N}(\boldsymbol{\mu}, \Sigma), \mathcal{N}(\boldsymbol{\mu}_t, \Sigma_t)) + \frac{1}{2\gamma} \ell^2(\boldsymbol{\mu}; (\mathbf{x}_t, y_t)) + \frac{1}{2\gamma} \mathbf{x}_t^\top \Sigma_t \mathbf{x}_t$$

Algorithm 8: CW: Confidence Weighted algorithm.

```
1 Parameter:  $a > 0$ (default=1)
2 Initialize:  $\mathbf{w}_1 = 0$ ,  $\Sigma_1 = a\mathbf{I}$ , ( $\Phi$  is the cumulative function of the normal
   distribution),  $\psi = 1 + \frac{\phi^2}{2}$ , and  $\zeta = 1 + \phi^2$ 
3 for  $t = 1, 2, \dots, T$  do
4   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
5   The learner predicts the class label:  $\hat{y}_t = \text{sign}(\mathbf{w}_t^\top \mathbf{x}_t)$ , where  $\mathbf{w}_t = \boldsymbol{\mu}_t$ ;
6   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
7   The learner computes the updating coefficients:
8    $u_t = \frac{1}{4}(-\alpha_t v_t \phi + \sqrt{\alpha_t^2 v_t^2 \phi^2 + 4v_t})^2$ ,  $v_t = \mathbf{x}_t^\top \Sigma_t \mathbf{x}_t$ ,  $m_t = y_t(\boldsymbol{\mu}_t \cdot \mathbf{x}_t)$ 
9    $\alpha_t = \max\{0, \frac{1}{v_t \zeta}(-m_t \psi + \sqrt{m_t^2 \frac{\phi^4}{4} + v_t \phi^2 \zeta})\}$ ,  $\beta_t = \frac{\alpha_t \phi}{\sqrt{u_t + v_t \alpha_t \phi}}$ ;
10  The learner calculates the suffered loss:  $l_t = \mathbb{I}(\alpha_t > 0)$ ;
11  if  $l_t > 0$  then
12    The learner updates the classification model:
13     $\boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t + \alpha_t y_t \Sigma_t \mathbf{x}_t$ ,  $\Sigma_{t+1} = \Sigma_t - \beta_t \Sigma_t \mathbf{x}_t \mathbf{x}_t^\top \Sigma_t$ ;
14  end
15 end
```

Algorithm 9: IELLIP: Improved Ellipsoid Method.

```
1 Parameter:  $a > 0$ (default=1),  $0 < c < 1$ (default=0.1),  $0 < b < 1$ (default=0.3)
2 Initialize:  $\mathbf{w}_1 = 0$ ,  $\Sigma_1 = a\mathbf{I}$ 
3 for  $t = 1, 2, \dots, T$  do
4   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
5   The learner predicts the class label:  $\hat{y}_t = \text{sign}(\mathbf{w}_t^\top \mathbf{x}_t)$ , where  $\mathbf{w}_t = \boldsymbol{\mu}_t$ ;
6   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
7   The learner calculates the suffered loss:  $l_t = \mathbb{I}(\hat{y}_t \neq y_t)$ ;
8   if  $l_t > 0$  then
9     The learner updates the classification model:
10     $\boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t + \alpha_t \Sigma_t \mathbf{g}_t$ ,  $\Sigma_{t+1} = \frac{1}{1-c_t}(\Sigma_t - c_t \Sigma_t \mathbf{g}_t \mathbf{g}_t^\top \Sigma_t)$ 
11     $\alpha_t = \frac{\alpha \gamma - y_t \boldsymbol{\mu}_t^\top \mathbf{x}_t}{\sqrt{\mathbf{x}_t^\top \Sigma_t \mathbf{x}_t}}$ ,  $\mathbf{g}_t = \frac{y_t \mathbf{x}_t}{\sqrt{\mathbf{x}_t^\top \Sigma_t \mathbf{x}_t}}$ ,  $c_t = cb^t$ ,  $0 < c \leq 1$ ,  $0 < b < 1$ ;
12  end
13 end
```

where $\ell^2(\boldsymbol{\mu}; (\mathbf{x}_t, y_t)) = (\max\{0, 1 - y_t(\boldsymbol{\mu} \cdot \mathbf{x}_t)\})^2$ and γ is a regularization parameter.

Algorithm 10: AROW: Adaptive Regularization Of Weights.

```

1 Parameters:  $a > 0$ (default=1),  $r > 0$ (default=1)
2 Initialize:  $\mathbf{w}_1 = 0$ ,  $\Sigma_1 = a\mathbf{I}$ , parameter  $r$ 
3 for  $t = 1, 2, \dots, T$  do
4   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
5   The learner predicts the class label:  $\hat{y}_t = \text{sign}(\mathbf{w}_t^\top \mathbf{x}_t)$ , where  $\mathbf{w}_t = \boldsymbol{\mu}_t$ ;
6   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
7   The learner calculates the suffered loss:  $l_t = \max(0, 1 - y_t \mathbf{w}_t^\top \mathbf{x}_t)$ ;
8   if  $l_t > 0$  then
9     The learner updates the classification model:
10     $\boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t + \alpha_t \Sigma_t y_t \mathbf{x}_t$ ,  $\Sigma_{t+1} = \Sigma_t - \beta_t \Sigma_t \mathbf{x}_t \mathbf{x}_t^\top \Sigma_t$ 
11     $\alpha_t = \ell(\boldsymbol{\mu}_t; (\mathbf{x}_t, y_t)) \beta_t$ ,  $\beta_t = \frac{1}{\mathbf{x}_t^\top \Sigma_t \mathbf{x}_t + r}$ ;
12  end
13 end
```

The struct “model” contains the following parameters:

- $r > 0$: the trade-off between the regularization and the loss

2.2.5 New Adaptive Regularization of Weights (NAROW)

The New Adaptive Regularize Of Weights (NAROW) is an algorithm that interpolates between a second-order algorithm with adaptive-second-order-information, like AROW, and one with fixed-second-order-information. Even the bound is in between these two worlds, the matrix Σ_t is updated only less frequently than AROW, preventing its eigenvalues from growing too much.

The struct “model” contains the following parameters:

- $b > 0$: a threshold for deciding the adaptive update

2.2.6 Normal HERD (NHERD)

The normal herd algorithm is introduced based upon constraining the velocity flow over a distribution of weight vectors. In particular, it is designed to effectively herd a Gaussian weight vector distribution by trading off velocity constraints with a loss function. Formally, the updated is performed as follows:

$$\boldsymbol{\mu}_{t+1} = A_t \boldsymbol{\mu}_t + \mathbf{b}_t, \quad \Sigma_{t+1} = A_t \Sigma_t A_t^\top, \quad (A_t, \mathbf{b}_t) = \arg \min_{A, \mathbf{b}} \mathbb{E}_{\mathbf{w}_t \sim \mathcal{N}(\boldsymbol{\mu}_t, \Sigma_t)} C_t(A_t \mathbf{w}_t + \mathbf{b}_t),$$

where

$$C_t(\mathbf{w}) = \left[\frac{1}{2} (\mathbf{w} - \mathbf{w}_t)^\top \Sigma_t^{-1} (\mathbf{w} - \mathbf{w}_t) + C \ell(\mathbf{w}; (\mathbf{x}_t, y_t)) \right],$$

Algorithm 11: NAROW: New Adaptive Regularization of Weights.

```
1 Parameters:  $a > 0$ (default=1),  $b > 0$ (default=1)
2 Initialize:  $\mathbf{w}_1 = 0$ ,  $\Sigma_1 = a\mathbf{I}$ 
3 for  $t = 1, 2, \dots, T$  do
4   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
5   The learner predicts the class label:  $\hat{y}_t = \text{sign}(\mathbf{w}_t^\top \mathbf{x}_t)$ , where  $\mathbf{w}_t = \Sigma_t \boldsymbol{\mu}_t$ ;
6   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
7   The learner calculates the suffered loss:  $l_t = \max(0, 1 - y_t \mathbf{w}_t^\top \mathbf{x}_t)$ ;
8   if  $l_t > 0$  then
9     The learner updates the classification model:
10     $\boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t + y_t \mathbf{x}_t$ ,  $\Sigma_{t+1} = \Sigma_t - \beta_t \Sigma_t \mathbf{x}_t \mathbf{x}_t^\top \Sigma_t$ 
11     $\beta_t = \frac{1}{\mathbf{x}_t^\top \Sigma_t \mathbf{x}_t + \gamma_t}$ ,  $\gamma_t = \frac{\mathbf{x}_t^\top \Sigma_t \mathbf{x}_t}{b \mathbf{x}_t^\top \Sigma_t \mathbf{x}_t - 1}$ , when  $\mathbf{x}_t^\top \Sigma_t \mathbf{x}_t \geq 1/b$  and  $\gamma_t = +\infty$ ;
12  end
13 end
```

and $\ell(\mathbf{w}; (\mathbf{x}_t, y_t)) = (\max\{0, 1 - y_t \mathbf{w}_t^\top \mathbf{x}_t\})^2$.

To analytically solve the update problem, the above optimization is further uniformly relaxed as the following objective

$$\begin{aligned} (\boldsymbol{\mu}_{t+1}, A_t) &= \arg \min_{\boldsymbol{\mu}, A} \frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_t)^\top \Sigma_t^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_t) + C \ell(\boldsymbol{\mu}; (\mathbf{x}_t, y_t)) \\ &\quad + \frac{1}{2} \text{Tr}((A - I)^\top \Sigma_t^{-1} (A - I) \Sigma_t) + \frac{C}{2} \mathbf{x}_t^\top A \Sigma_t A^\top \mathbf{x}_t \end{aligned}$$

which enjoy a closed-form solution.

Algorithm 12: NHERD: Normal HERD Algorithm.

```
1 Parameters:  $a > 0$ (default=1),  $C > 0$ (default=1)
2 Initialize:  $\mathbf{w}_1 = 0$ ,  $\Sigma_1 = a\mathbf{I}$ 
3 for  $t = 1, 2, \dots, T$  do
4   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
5   The learner predicts the class label:  $\hat{y}_t = \text{sign}(\mathbf{w}_t^\top \mathbf{x}_t)$ , where  $\mathbf{w}_t = \boldsymbol{\mu}_t$ ;
6   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
7   The learner calculates the suffered loss:  $l_t = \mathbb{I}((\max\{0, 1 - y_t \mathbf{w}_t^\top \mathbf{x}_t\})^2 > 0)$ ;
8   if  $l_t > 0$  then
9     The learner updates the classification model:
10     $\boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t + \alpha_t \Sigma_t y_t \mathbf{x}_t$ ,  $\Sigma_{t+1} = \Sigma_t - \beta_t \Sigma_t \mathbf{x}_t \mathbf{x}_t^\top \Sigma_t$ 
11     $\alpha_t = \frac{\max\{0, 1 - y_t \boldsymbol{\mu}_t^\top \mathbf{x}_t\}}{\mathbf{x}_t^\top \Sigma_t \mathbf{x}_t + 1/C}$ ,  $\beta_t = \frac{C^2 \mathbf{x}_t \Sigma_t \mathbf{x}_t^\top + 2C}{(1 + C \mathbf{x}_t \Sigma_t \mathbf{x}_t^\top)^2}$ ;
12  end
13 end
```

The struct “model” contains the following parameters:

- $C > 0$: a trade-off between regularization and the loss

2.2.7 Soft Confidence-weighted learning (SCW)

This section presents the soft confidence-weighted (SCW) learning [14], which aims to address the limitation of the CW and AROW learning. After the introduction of the following loss function:

$$\ell^\phi(\mathcal{N}(\boldsymbol{\mu}, \Sigma); (\mathbf{x}_t, y_t)) = \max\left(0, \phi\sqrt{\mathbf{x}_t^\top \Sigma \mathbf{x}_t} - y_t \boldsymbol{\mu} \cdot \mathbf{x}_t\right),$$

where $\phi = \Phi^{-1}(\eta)$, it is easy to verify that the optimization problem of the original CW can be re-written as follows

$$\begin{aligned} (\boldsymbol{\mu}_{t+1}, \Sigma_{t+1}) &= \arg \min_{\boldsymbol{\mu}, \Sigma} D_{KL}(\mathcal{N}(\boldsymbol{\mu}, \Sigma) \| \mathcal{N}(\boldsymbol{\mu}_t, \Sigma_t)) \\ &s.t. \ell^\phi(\mathcal{N}(\boldsymbol{\mu}, \Sigma); (\mathbf{x}_t, y_t)) = 0, \phi > 0 \end{aligned}$$

The original CW learning method employs a very aggressive updating strategy by changing the distribution as much as necessary to satisfy the constraint imposed by the current example. Although it results in the rapid learning effect, it could force to wrongly change the parameters of the distribution dramatically when handling a mislabeled instance. Such undesirable property makes the original CW algorithm performs poorly in many real-world applications with relatively large noise.

To overcome the above limitation of the CW learning problem, Soft Confidence-Weighted (SCW) learning method is proposed as follows:

$$(\boldsymbol{\mu}_{t+1}, \Sigma_{t+1}) = \arg \min_{\boldsymbol{\mu}, \Sigma} D_{KL}(\mathcal{N}(\boldsymbol{\mu}, \Sigma) \| \mathcal{N}(\boldsymbol{\mu}_t, \Sigma_t)) + C \ell^\phi(\mathcal{N}(\boldsymbol{\mu}, \Sigma); (\mathbf{x}_t, y_t))$$

where C is a parameter to tradeoff the passiveness and aggressiveness. The above formulation of the Soft Confidence-Weighted algorithm is denoted as ‘‘SCW-I’’ for short.

Furthermore, employing a squared penalty leads to the second formulation of SCW learning (denoted as ‘‘SCW-II’’ for short):

$$(\boldsymbol{\mu}_{t+1}, \Sigma_{t+1}) = \arg \min_{\boldsymbol{\mu}, \Sigma} D_{KL}(\mathcal{N}(\boldsymbol{\mu}, \Sigma) \| \mathcal{N}(\boldsymbol{\mu}_t, \Sigma_t)) + C \ell^\phi(\mathcal{N}(\boldsymbol{\mu}, \Sigma); (\mathbf{x}_t, y_t))^2$$

The struct ‘‘model’’ contains the following parameters:

- $a > 0$, which is used to initialize the $\Sigma_1 = a\mathbf{I}$, where \mathbf{I} is identity matrix;
- $\phi > 0$, which is $\phi = \Phi^{-1}(\eta)$, where η is the probability required for the updated distribution on the current instance (default=0.75 for SCW-I and 0.90 for SCW-II);
- $C > 0$, which is used to trade-off between keeping the previous information and minimizing the current loss

Algorithm 13: SCW: Soft Confidence Weighted algorithms.

```

1 Parameters:  $a > 0$ (default=1),  $C > 0$ (default=1),  $\eta > 0$ 
2 Initialize:  $\mathbf{w}_1 = 0$ ,  $\Sigma_1 = a\mathbf{I}$ ,  $\phi = \Phi^{-1}(\eta)$ ,  $\psi = 1 + \frac{\phi^2}{2}$  and  $\zeta = 1 + \phi^2$ 
3 for  $t = 1, 2, \dots, T$  do
4   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
5   The learner predicts the class label:  $\hat{y}_t = \text{sign}(\mathbf{w}_t^\top \mathbf{x}_t)$ , where  $\mathbf{w}_t = \boldsymbol{\mu}_t$ ;
6   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
7   The learner calculates the suffered loss:  $l_t = \mathbb{I}(\max\{0, 1 - y_t \mathbf{w}_t^\top \mathbf{x}_t\} > 0)$ ;
8   if  $l_t > 0$  then
9     The learner updates the classification model:
10     $\boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t + \alpha_t y_t \Sigma_t \mathbf{x}_t$ ,  $\Sigma_{t+1} = \Sigma_t - \beta_t \Sigma_t \mathbf{x}_t^\top \mathbf{x}_t \Sigma_t$ 
11
12    
$$\alpha_t = \begin{cases} \min\{C, \max\{0, \frac{1}{v_t \zeta}(-m_t \psi + \sqrt{m_t^2 \frac{\phi^4}{4} + v_t \phi^2 \zeta})\}\} & \text{SCW-I,} \\ \max\{0, \frac{-(2m_t n_t + \phi^2 m_t v_t) + \gamma_t}{2(n_t^2 + n_t v_t \phi^2)}\} & \text{SCW-II.} \end{cases}$$

13
14    
$$\beta_t = \frac{\alpha_t \phi}{\sqrt{u_t + v_t \alpha_t \phi}}, u_t = \frac{1}{4}(-\alpha_t v_t \phi + \sqrt{\alpha_t^2 v_t^2 \phi^2 + 4v_t})^2,$$

15
16    
$$v_t = \mathbf{x}_t^\top \Sigma_t \mathbf{x}_t, m_t = y_t (\boldsymbol{\mu}_t \cdot \mathbf{x}_t), \gamma_t = \phi \sqrt{\phi^2 m_t^2 v_t^2 + 4n_t v_t (n_t + v_t \phi^2)}, \text{ and}$$

17
18    
$$n_t = v_t + \frac{1}{2C}.$$

19   end
20 end

```

3 Online Learning Algorithms for Multiclass Classification

3.1 First-Order Algorithms

In this section, we provide the details of the first-order learning algorithms and the corresponding functions (located at /algorithms/Multiclass-Classification). The first-order learning algorithms only keep updating k classification functions (for the k different labels), which only utilize the first-order information of the examples.

Specifically, the struct “model” in the codes contains a field “W”, which is a $k \times d$ matrix, whose i -th row is the linear classifier for the i -th label. In addition, “model” also contains various different parameters for different algorithms, which will be listed in every algorithm.

3.1.1 Multiclass Perceptron

Crammer and Singer extended the binary Perceptron algorithm to a family of multiclass Perceptron algorithms [8], by allocating different weights on the support vectors added to the classifiers corresponding to those error-set $E = \{i \neq y : W_i \cdot \mathbf{x} \geq W_y \cdot \mathbf{x}\}$. According to different allocation strategies, three variants of multiclass Perceptron algorithms are proposed: max-score multiclass Perceptron, uniform multiclass Perceptron, and proportion multiclass Perceptron.

Algorithm 14: M-Perceptron: multiclass Perceptron algorithms.

```
1 Initialize:  $\mathbf{W}_1 = 0$ 
2 for  $t = 1, 2, \dots, T$  do
3   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
4   The learner predicts the class label:  $\hat{y}_t = \arg \max_{i=1}^k (W_{t,i} \cdot \mathbf{x}_t)$ ;
5   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
6   The learner calculates the suffered loss:  $l_t = \mathbb{I}(\hat{y}_t \neq y_t)$ .;
7   if  $l_t > 0$  then
8     The learner updates the classification model:
9      $W_{t+1,i} = W_{t,i} + \alpha_{t,i} \mathbf{x}_t$ 
10  end
11 end
```

For max-score multiclass update,

$$\alpha_{t,i} = \begin{cases} -1 & i = \arg \max_{j=1}^k W_{t,j} \cdot \mathbf{x}_t \\ 1 & i = y_t \\ 0 & \text{otherwise} \end{cases}$$

For uniform multiclass update,

$$\alpha_{t,i} = \begin{cases} -1/|E_t| & i \in E_t, E_t = \{i \neq y_t : W_{t,i} \cdot \mathbf{x}_t \geq W_{t,y_t} \cdot \mathbf{x}_t\} \\ 1 & i = y_t \\ 0 & \text{otherwise} \end{cases}$$

For promotion multiclass update,

$$\alpha_{t,i} = \begin{cases} -\frac{[W_{t,i} \cdot \mathbf{x}_t - W_{t,y_t} \cdot \mathbf{x}_t]_+}{\sum_{j=1}^k [W_{t,j} \cdot \mathbf{x}_t - W_{t,y_t} \cdot \mathbf{x}_t]_+} & i \neq y_t \\ 1 & i = y_t \end{cases}$$

where $[z]_+ = \max(z, 0)$.

3.1.2 Multiclass Relaxed Online Maximum Margin Algorithm

Algorithm 15: M-ROMMA: Multiclass Relaxed Online Maximum Margin Algorithm algorithms.

```

1 Initialize:  $\mathbf{W}_1 = 0$ 
2 for  $t = 1, 2, \dots, T$  do
3   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
4   The learner predicts the class label:  $\hat{y}_t = \arg \max_{i=1}^k (W_{t,i} \cdot \mathbf{x}_t)$ ;
5   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
6   The learner calculates the suffered loss:


$$l_t = \begin{cases} \mathbb{I}(\hat{y}_t \neq y_t) & \text{Multiclass ROMMA} \\ \mathbb{I}(\max(0, 1 - (W_{t,y_t} \cdot \mathbf{x}_t - \max_{j \neq y_t} W_{t,j} \cdot \mathbf{x}_t)) > 0) & \text{Multiclass agg-ROMMA} \end{cases}$$


7   if  $l_t > 0$  then
8     The learner updates the classification model:
9      $W_{t+1,i} = c_t W_{t,i} + \alpha_{t,i} d_t \mathbf{x}_t$ 
10     $c_t = \frac{2\|\mathbf{x}_t\|^2 \|W\|^2 - (W_{t,y_t} \cdot \mathbf{x}_t - \max_{j \neq y_t} W_{t,j} \cdot \mathbf{x}_t)}{2\|\mathbf{x}_t\|^2 \|W\|^2 - (W_{t,y_t} \cdot \mathbf{x}_t - \max_{j \neq y_t} W_{t,j} \cdot \mathbf{x}_t)^2}$ 
11     $d_t = \frac{\|W\|^2 (1 - (W_{t,y_t} \cdot \mathbf{x}_t - \max_{j \neq y_t} W_{t,j} \cdot \mathbf{x}_t))}{2\|\mathbf{x}_t\|^2 \|W\|^2 - (W_{t,y_t} \cdot \mathbf{x}_t - \max_{j \neq y_t} W_{t,j} \cdot \mathbf{x}_t)^2}$ 
12    
$$\alpha_{t,i} = \begin{cases} -1 & i = \arg \max_{j=1}^k W_{t,j} \cdot \mathbf{x}_t \\ 1 & i = y_t \\ 0 & \text{otherwise} \end{cases}$$

13  end
14 end

```

The ROMMA and agg-ROMMA algorithm do not contains parameters.

3.1.3 Multiclass Online Gradient Descent Algorithm

The online gradient descent algorithm [18] applies the gradient descent updating approach together with the following hinge loss:

$$\ell(W; (\mathbf{x}_t, y_t)) = \max \left(0, 1 - (W_{y_t} \cdot \mathbf{x}_t - \max_{j \neq y_t} W_j \cdot \mathbf{x}_t) \right)$$

Algorithm 16: M-OGD: Multiclass Online Gradient Descent Algorithm.

```

1 Initialize:  $\mathbf{W}_1 = 0$ 
2 for  $t = 1, 2, \dots, T$  do
3   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
4   The learner predicts the class label:  $\hat{y}_t = \arg \max_{i=1}^k (W_{t,i} \cdot \mathbf{x}_t)$ ;
5   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
6   The learner calculates the suffered loss:  $l_t = \mathbb{I}(\ell(W_t; (\mathbf{x}_t, y_t)) > 0)$ ;
7   if  $l_t > 0$  then
8     The learner updates the classification model:
9      $W_{t+1,i} = W_{t,i} + \alpha_{t,i} \mathbf{x}_t$ 
10
11     
$$\alpha_{t,i} = \begin{cases} -1/\sqrt{t} & i = \arg \max_{j=1}^k W_{t,j} \cdot \mathbf{x}_t \\ 1/\sqrt{t} & i = y_t \\ 0 & \text{otherwise} \end{cases}$$

12   end
13 end

```

In our implementation, the algorithm does not contains any parameters.

3.1.4 Multiclass Passive Aggressive Algorithms

The family of online multiclass Passive-Aggressive (PA) learning [2] is extended from its binary version, which is formulated to trade off between passiveness (minimizing the distance between the learnt classifier and the previous classifier) and aggressiveness (minimizing the loss of the learnt classier suffered on the current instance).

Formally, the PA algorithm is formulated as the following online optimization problem:

$$W_{t+1} = \arg \min_W \frac{1}{2} \|W - W_t\|^2, \quad \text{s.t. } \ell(W; (\mathbf{x}_t, y_t)) = \max \left(0, 1 - (W_{y_t} \cdot \mathbf{x}_t - \max_{j \neq y_t} W_j \cdot \mathbf{x}_t) \right) = 0$$

Similar to binary cases, the multiclass Passive Aggressive algorithm can be extended to multiclass Passive Aggressive I and II algorithms, which can better handle the non-separable case. Specifically, Multiclass Passive Aggressive I is formulated as

$$W_{t+1} = \arg \min_W \frac{1}{2} \|W - W_t\|^2 + C \ell(W; (\mathbf{x}_t, y_t))$$

Multiclass Passive Aggressive II is formulated as

$$W_{t+1} = \arg \min_W \frac{1}{2} \|W - W_t\|^2 + C \ell(W; (\mathbf{x}_t, y_t))^2$$

Algorithm 17: M-PA: Multiclass Passive Aggressive Algorithms.

```

1 Initialize:  $\mathbf{W}_1 = 0, C > 0$ 
2 for  $t = 1, 2, \dots, T$  do
3   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
4   The learner predicts the class label:  $\hat{y}_t = \arg \max_{i=1}^k (W_{t,i} \cdot \mathbf{x}_t)$ ;
5   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
6   The learner calculates the suffered loss:  $l_t = \mathbb{I}(\ell(W_t; (\mathbf{x}_t, y_t)) > 0)$ ;
7   if  $l_t > 0$  then
8     The learner updates the classification model:
9      $W_{t+1,i} = W_{t,i} + \alpha_{t,i} \mathbf{x}_t$ ;
10  end
11 end

```

For multiclass Passive Aggressive algorithm, the closed-form solution to $\alpha_{t,i}$ is

$$\alpha_{t,i} = \begin{cases} -\frac{\ell(W_t; (\mathbf{x}_t, y_t))}{2\|\mathbf{x}_t\|^2} & i = \arg \max_{j=1}^k W_{t,j} \cdot \mathbf{x}_t \\ \frac{\ell(W_t; (\mathbf{x}_t, y_t))}{2\|\mathbf{x}_t\|^2} & i = y_t \\ 0 & \text{otherwise} \end{cases}$$

For multiclass Passive Aggressive I algorithm, the closed-form solution to $\alpha_{t,i}$ is

$$\alpha_{t,i} = \begin{cases} -\min(C, \frac{\ell(W_t; (\mathbf{x}_t, y_t))}{2\|\mathbf{x}_t\|^2}) & i = \arg \max_{j=1}^k W_{t,j} \cdot \mathbf{x}_t \\ \min(C, \frac{\ell(W_t; (\mathbf{x}_t, y_t))}{2\|\mathbf{x}_t\|^2}) & i = y_t \\ 0 & \text{otherwise} \end{cases}$$

For multiclass Passive Aggressive II algorithm, the closed-form solution to $\alpha_{t,i}$ is

$$\alpha_{t,i} = \begin{cases} -\frac{\ell(W_t; (\mathbf{x}_t, y_t))}{2\|\mathbf{x}_t\|^2 + \frac{1}{2C}} & i = \arg \max_{j=1}^k W_{t,j} \cdot \mathbf{x}_t \\ \frac{\ell(W_t; (\mathbf{x}_t, y_t))}{2\|\mathbf{x}_t\|^2 + \frac{1}{2C}} & i = y_t \\ 0 & \text{otherwise} \end{cases}$$

3.2 Second-Order Algorithms

To better exploring the underlying structure between features, the second-order online learning algorithms typically assume the weight vector follows Gaussian distributions $\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ with mean vector $\boldsymbol{\mu} \in \mathbb{R}^{kd}$ and covariance matrix $\Sigma \in \mathbb{R}^{kd \times kd}$.

Specifically, the struct “model” in the codes contains two fields “ $\boldsymbol{\mu} \in \mathbb{R}^{kd}$ ” and “ $\Sigma \in \mathbb{R}^{kd \times kd}$ ”, which are the mean vector and covariance matrix for the learnt Gaussian distribution. In addition, “model” also contains various parameters for different algorithms, which will be explained in every algorithm.

3.2.1 Multiclass Confidence Weighted Learning

In multiclass confidence-weighted learning [3], the weight distribution $\mathcal{N}(\boldsymbol{\mu}, \Sigma)$ is updated with one strategy similar with the binary case. To make the proposed algorithms more efficient, the authors proposed three variant of multiclass confidence weighted algorithms. Here, we implemented the most efficient one.

Specifically, we introduced a new label-dependent feature, $\psi(\mathbf{x}, i) = [\mathbf{0}^\top, \dots, \mathbf{x}^\top, \dots, \mathbf{0}^\top]^\top$, where $\mathbf{0}, \mathbf{x} \in \mathbb{R}^d$, only the i -th position of $\psi(\mathbf{x}, i)$ is \mathbf{x} and the others are $\mathbf{0}$.

Algorithm 18: M-CW: Multiclass Confidence Weighted Algorithms.

```

1 Initialize:  $\boldsymbol{\mu}_1 = \mathbf{0}$ ” and “ $\Sigma_1 = \mathbf{I}$ ,  $\phi = \Phi^{-1}(\eta)$  and  $\eta \in (0.5, 1]$ 
2 for  $t = 1, 2, \dots, T$  do
3   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
4   The learner predicts the class label:  $\hat{y}_t = \arg \max_{i=1}^k (\mathbf{w}_t \cdot \psi(\mathbf{x}_t, i))$ , where
    $\mathbf{w}_t = \boldsymbol{\mu}_t$ ;
5   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
6   The learner computes all the updating coefficients:
7    $m_t = \boldsymbol{\mu}_t^\top \Delta \psi_t$ ,  $v_t = \Delta \psi_t^\top \Sigma_t \Delta \psi_t$ 
8    $\alpha_t = \max \left\{ 0, \frac{-(1+2\phi m_t) + \sqrt{(1+2\phi m_t)^2 - 8\phi(m_t - \phi v_t)}}{4\phi v_t} \right\}$ ,  $\beta_t = \frac{1}{1/(2\alpha_t \phi) + v_t}$ ;
9   The learner calculates the suffered loss:  $l_t = \mathbb{I}(\alpha_t > 0)$ ;
10  if  $l_t > 0$  then
11    The learner updates the classification model:
12     $\boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t + \alpha_t \Sigma_t \Delta \psi_t$ ,  $\Sigma_{t+1} = \Sigma_t - \beta_t \Sigma_t \Delta \psi_t \Delta \psi_t^\top \Sigma_t$ 
13    where  $\Delta \psi_t = \psi(\mathbf{x}_t, y_t) - \psi(\mathbf{x}_t, \tilde{y}_t)$ ,  $\tilde{y}_t = \arg \max_{i=1, i \neq y_t}^k (\boldsymbol{\mu}_i \cdot \mathbf{x}_t)$ ;
14  end
15 end
```

The struct “model” contains the following parameters:

- $\eta \in (0.5, 1]$, which is the probability required for the updated distribution on the current instance

3.2.2 Multiclass Adaptive Regularization Of Weights

The multiclass AROW [6] is a natural extension of of binary AROW algorithm. To extend the binary algorithm, the authors proposed two kinds of updating methods. Here, the implemented version is the most efficient one.

Algorithm 19: M-AROW: Multiclass Adaptive Regularization Of Weights.

```

1 Initialize:  $\boldsymbol{\mu}_1 = \mathbf{0}$  and " $\Sigma_1 = \mathbf{I}$ ,  $\phi = \Phi^{-1}(\eta)$ ,  $\eta \in (0.5, 1]$ , and  $r > 0$ 
2 for  $t = 1, 2, \dots, T$  do
3   The learner receives an incoming instance:  $\mathbf{x}_t \in \mathcal{X}$ ;
4   The learner predicts the class label:  $\hat{y}_t = \arg \max_{i=1}^k (\mathbf{w}_t \cdot \psi(\mathbf{x}_t, i))$ , where
    $\mathbf{w}_t = \boldsymbol{\mu}_t$ ;
5   The true class label is revealed from the environment:  $y_t \in \mathcal{Y}$ ;
6   The learner computes all the updating coefficients:
7    $m_t = \boldsymbol{\mu}_t^\top \Delta \psi_t$ ,  $v_t = \Delta \psi_t^\top \Sigma_t \Delta \psi_t$ ,  $\alpha_t = \max \{0, 1 - m_t\} \beta_t$ ,  $\beta_t = \frac{1}{r + v_t}$ ;
8   The learner calculates the suffered loss:  $l_t = \mathbb{I}(m_t < 1)$ ;
9   if  $l_t > 0$  then
10    The learner updates the classification model:
11     $\boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t + \alpha_t \Sigma_t \Delta \psi_t$ ,  $\Sigma_{t+1} = \Sigma_t - \beta_t \Sigma_t \Delta \psi_t \Delta \psi_t^\top \Sigma_t$ 
12    where  $\Delta \psi_t = \psi(\mathbf{x}_t, y_t) - \psi(\mathbf{x}_t, \tilde{y}_t)$ ,  $\tilde{y}_t = \arg \max_{i=1, i \neq y_t}^k (\boldsymbol{\mu}_i \cdot \mathbf{x}_t)$ ;
13  end
14 end
```

The struct “model” contains the following parameters:

- $r \in (0, +\infty)$,

3.2.3 Multiclass Soft Confidence-Weighted Learning

In Multiclass Soft Confidence-Weighted learning, we assume the each prototype vector \mathbf{w}_i follows the Gaussian distribution with the mean vector $\boldsymbol{\mu}_i$ and the covariance matrix Σ_i . For simplicity, we assume each prototype $i \in [1, k]$ share the same covariance matrix Σ . In multiclass learning setting, we want to make that the probability of the score of the class of the smallest score among all relevant classes to be higher than the score of the class of the highest score among all irrelevant classes is larger than a threshold η . In mathematical form, the constrain can be expressed as follows:

$$Pr_{\mathbf{w}_{t,r_t} \sim \mathcal{N}(\boldsymbol{\mu}_{t,r_t}, \Sigma_t), \mathbf{w}_{t,s_t} \sim \mathcal{N}(\boldsymbol{\mu}_{t,s_t}, \Sigma_t)} [(\mathbf{w}_{t,r_t} \cdot \mathbf{x}_t) \geq (\mathbf{w}_{t,s_t} \cdot \mathbf{x}_t)] \geq \eta \quad (1)$$

where

$$r_t = \arg \min_{r \in Y_t} \boldsymbol{\mu}_{t,r} \cdot \mathbf{x}_t, \quad s_t = \arg \max_{s \notin Y_t} \boldsymbol{\mu}_{t,s} \cdot \mathbf{x}_t \quad (2)$$

To overcome the above limitation of the CW learning problem, we propose a Soft Confidence-Weighted (SCW) learning method, which aims to soften the aggressiveness

of the CW updating strategy. The idea of the SCW learning is inspired by the variants of PA algorithms (PA-I and PA-II) and the adaptive margin. In particular, we formulate the optimization of SCW for learning the soft-margin classifiers as follows:

$$(\boldsymbol{\mu}_{t+1,r_t}, \boldsymbol{\mu}_{t+1,s_t}, \Sigma_{t+1}) = \arg \min_{\boldsymbol{\mu}_r, \boldsymbol{\mu}_s, \Sigma} D_{KL}(\mathcal{N}(\boldsymbol{\mu}_r, \Sigma) \| \mathcal{N}(\boldsymbol{\mu}_{t,r_t}, \Sigma_t)) \quad (3)$$

$$+ D_{KL}(\mathcal{N}(\boldsymbol{\mu}_s, \Sigma) \| \mathcal{N}(\boldsymbol{\mu}_{t,s_t}, \Sigma_t)) \quad (4)$$

$$+ C\ell^\phi(\mathcal{N}(\boldsymbol{\mu}_r, \boldsymbol{\mu}_s, \Sigma_t); (\mathbf{x}_t, y_t)) \quad (5)$$

where C is a parameter to tradeoff the passiveness and aggressiveness. We denoted the above formulation of the Multiclass Soft Confidence-Weighted algorithm, as ‘‘M-SCW1’’ for short.

The closed-form solution of the optimization (3) is expressed as:

$$\begin{aligned} \boldsymbol{\mu}_{t+1,r_t} &= \boldsymbol{\mu}_{t,r_t} + \alpha_t y_t \Sigma_t \mathbf{x}_t, & \boldsymbol{\mu}_{t+1,s_t} &= \boldsymbol{\mu}_{t,s_t} - \alpha_t y_t \Sigma_t \mathbf{x}_t, \\ \Sigma_{t+1} &= \Sigma_t - \beta_t \Sigma_t \mathbf{x}_t^T \mathbf{x}_t \Sigma_t \end{aligned}$$

where the updating coefficients are as follows:

$$\alpha_t = \min\{C, \max\{0, \frac{1}{2v_t\psi}(-m_t\psi + \sqrt{m_t^2\psi^2 - m_t^2\psi + 2\psi\phi^2v_t})\}\}$$

$$\beta_t = \frac{\alpha_t\phi}{\sqrt{2u_t} + v_t\alpha_t\phi}$$

where $u_t = \frac{1}{8}(-\alpha_t v_t \phi + \sqrt{\alpha_t^2 v_t^2 \phi^2 + 8v_t})^2$, $v_t = \mathbf{x}_t^T \Sigma_t \mathbf{x}_t$, $m_t = \boldsymbol{\mu}_{t,r_t} \cdot \mathbf{x}_t - \boldsymbol{\mu}_{t,s_t} \cdot \mathbf{x}_t$, $\phi = \Phi^{-1}(\eta)$, $\psi = 1 + \frac{\phi^2}{2}$.

Similar to the variant of PA, we can also modify the above formulation by employing a squared penalty, leading to the second formulation of Multiclass SCW learning (denoted as ‘‘M-SCW2’’ for short):

$$(\boldsymbol{\mu}_{t+1,r_t}, \boldsymbol{\mu}_{t+1,s_t}, \Sigma_{t+1}) = \arg \min_{\boldsymbol{\mu}_r, \boldsymbol{\mu}_s, \Sigma} D_{KL}(\mathcal{N}(\boldsymbol{\mu}_r, \Sigma) \| \mathcal{N}(\boldsymbol{\mu}_{t,r_t}, \Sigma_t)) \quad (6)$$

$$+ D_{KL}(\mathcal{N}(\boldsymbol{\mu}_s, \Sigma) \| \mathcal{N}(\boldsymbol{\mu}_{t,s_t}, \Sigma_t)) \quad (7)$$

$$+ C\ell^\phi(\mathcal{N}(\boldsymbol{\mu}_r, \boldsymbol{\mu}_s, \Sigma_t); (\mathbf{x}_t, y_t))^2 \quad (8)$$

The closed-form solution of the optimization (6) is:

$$\begin{aligned} \boldsymbol{\mu}_{t+1,r_t} &= \boldsymbol{\mu}_{t,r_t} + \alpha_t y_t \Sigma_t \mathbf{x}_t, & \boldsymbol{\mu}_{t+1,s_t} &= \boldsymbol{\mu}_{t,s_t} - \alpha_t y_t \Sigma_t \mathbf{x}_t, \\ \Sigma_{t+1} &= \Sigma_t - \beta_t \Sigma_t \mathbf{x}_t^T \mathbf{x}_t \Sigma_t \end{aligned}$$

The updating coefficients are as follows:

$$\alpha_t = \max\{0, \frac{-(2m_t\rho_t + \phi^2 m_t v_t) + \gamma_t}{2(\rho_t^2 + \rho_t v_t \phi^2)}\}, \quad \beta_t = \frac{\alpha_t\phi}{\sqrt{2u_t} + v_t\alpha_t\phi}$$

where $\gamma_t = \phi\sqrt{\phi^2 m_t^2 v_t^2 + 8\rho_t v_t(\rho_t + v_t\phi^2)}$, and $\rho_t = 2v_t + \frac{1}{2C}$.

In the implementation of MSCW1 an MSCW2 functions, you can specify the struct ‘‘model’’ with the following parameters:

- $a > 0$, which is used to initialize the $\Sigma_1 = a\mathbf{I}$, where \mathbf{I} is identity matrix;
- $\phi > 0$, which is $\phi = \Phi^{-1}(\eta)$, where η is the probability required for the updated distribution on the current instance
- $C > 0$, which is used to trade-off between keeping the previous information and minimizing the current loss

Algorithm 20: Multiclass SCW learning algorithms (M-SCW)

```

1 INPUT: parameters  $C > 0, \eta > 0$ .
2 INITIALIZATION:  $\boldsymbol{\mu}_{1,1}, \dots, \boldsymbol{\mu}_{1,k} = (0, \dots, 0)^\top, \Sigma_1 = I$ . for  $t = 1, \dots, T$  do
3   Receive an example  $\mathbf{x}_t \in \mathbb{R}^d$ ;
4   Make prediction:  $\hat{Y}_t = \arg \max_r (\boldsymbol{\mu}_{t,r} \cdot \mathbf{x}_t), r \in [1, k]$ ;
5   Receive true label  $Y_t$ ;
6   Suffer loss  $\ell^\phi(\mathcal{N}(\boldsymbol{\mu}_{t,r_t}, \boldsymbol{\mu}_{t,s_t}, \Sigma_t); (\mathbf{x}_t, Y_t))$ ;
7   if  $\ell^\phi(\mathcal{N}(\boldsymbol{\mu}_{t,r_t}, \boldsymbol{\mu}_{t,s_t}, \Sigma_t); (\mathbf{x}_t, Y_t)) > 0$  then
8      $\boldsymbol{\mu}_{t+1,r_t} = \boldsymbol{\mu}_{t,r_t} + \alpha_t y_t \Sigma_t \mathbf{x}_t$ ,
9      $\boldsymbol{\mu}_{t+1,s_t} = \boldsymbol{\mu}_{t,s_t} - \alpha_t y_t \Sigma_t \mathbf{x}_t$ ,
10     $\Sigma_{t+1} = \Sigma_t - \beta_t \Sigma_t \mathbf{x}_t \mathbf{x}_t^\top \Sigma_t$ 
11    where  $\alpha_t$  and  $\beta_t$  are computed as follows:

    MSCW1 =  $\begin{cases} \alpha_t = \min\{C, \max\{0, \frac{1}{2v_t\psi}(-m_t\psi + \sqrt{m_t^2\psi^2 - m_t^2\psi + 2\psi\phi^2v_t})\}\} \\ \beta_t = \frac{\alpha_t\phi}{\sqrt{2u_t+v_t\alpha_t\phi}} \end{cases}$ 

    where  $u_t = \frac{1}{8}(-\alpha_tv_t\phi + \sqrt{\alpha_t^2v_t^2\phi^2 + 8v_t})^2, v_t = \mathbf{x}_t^\top \Sigma_t \mathbf{x}_t, m_t =$ 
 $\boldsymbol{\mu}_{t,r_t} \cdot \mathbf{x}_t - \boldsymbol{\mu}_{t,s_t} \cdot \mathbf{x}_t, \phi = \Phi^{-1}(\eta), \psi = 1 + \frac{\phi^2}{2}$ .

    MSCW2 =  $\begin{cases} \alpha_t = \max\{0, \frac{-(2m_t\rho_t + \phi^2m_tv_t) + \gamma_t}{2(\rho_t^2 + \rho_tv_t\phi^2)}\} \\ \beta_t = \frac{\alpha_t\phi}{\sqrt{2u_t+v_t\alpha_t\phi}} \end{cases}$ 

    where  $\gamma_t = \phi\sqrt{\phi^2m_tv_t^2 + 8\rho_tv_t(\rho_t + v_t\phi^2)}$ , and  $\rho_t = 2v_t + \frac{1}{2C}$ .
12   end
13 end

```

4 Conclusions, Revision, and Citation

LIBOL is an easy-to-use open source package for efficient and scalable on-line linear classification. It is currently one of comprehensive online learning software packages that include the largest number of diverse online learning algorithms for online classification. LIBOL is still being improved by improvements from practical users and new research results [16, 17]. The ultimate goal is to make easy learning with massive-scale data streams to tackle the emerging grand challenge of big data mining.

Revision History

- Version 0.1.0 was released on 28 December 2012. This version mainly includes MATLAB implementation for binary classification. In addition to Dr Hoi, other contributors include Jialei Wang and Peilin Zhao.
- Version 0.2.0 was released on 27 July 2013. This version includes several major improvements: (i) C++ implementations for the main functions; (ii) online learning algorithms for multi-class classification; (iii) new design and framework. In addition to Dr Hoi, other contributors include Jialei Wang and Peilin Zhao.
- Version 0.2.3 was released on 21 September 2013. This version was mad compatible with Octave. Some bugs in the C++ implementations were fixed. In addition to Dr Hoi, Ji Wan helped fix bugs in this version.
- Version 0.3.0 was released on 12 December 2013. This version was improved by offering a template for adding new algorithms and providing clear guide; detailed documentation was provided for parameter settings; some bugs in the implementations were also fixed. In addition to Dr Hoi, Ji Wan helped test this version.

More details of the revision history can be found in the LIBOL website:

<http://LIBOL.stevenhoi.org/changelog.html>

Welcome to send us your suggestions/corrections for LIBOL by the following email:

chhoi@ntu.edu.sg

Citation

In citing LIBOL in your papers, please use the following reference:

S.C.H. Hoi, J. Wang, P. Zhao. LIBOL: A Library for Online Learning Algorithms. Nanyang Technological University, 2012. <http://LIBOL.stevenhoi.org>.

Appendix

Appendix A: Framework Design and Main Functions

In this appendix, we give detailed documentations for the main functions in the software. The following figure shows an overall framework of the main function call for learning a model from input data. We explain each of these functions below.

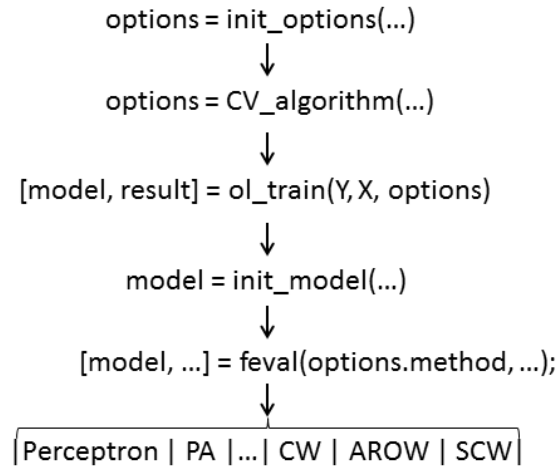


Figure 2: The framework of the function calls in the learning process.

```
function [model, result] = ol_train(Y, X, options)
%This is the main function to call an online algorithm for training a model from data.
INPUT:
-Y: the label vector, e.g., Y(t) is the label of t-th instance;
-X: training data matrix, e.g., X(t,:) denotes the features of t-th instance;
-options: a struct of predefined parameters and training settings
OUTPUT:
-model: a struct of the weight vector (w) and the SV indexes
-result: a struct of storing training results

function [options] = init_options(method, n, task_type)
%This is to create the "options" structure given a specific method and a specific task
INPUT:
-method: method name
-n: number of training instances in the database
-task_type: type of task (bc or mc)
OUTPUT:
-options: the generated options
```



```

function [options] = CV_algorithm(Y, X, options)
%This is to choose the best parameters by running the validation automatically.
INPUT:
-Y: the label vector, e.g., Y(t) is the label of t-th instance;
-X: training data, e.g., X(t,:) denotes features of t-th instance;
-options: a struct of predefined parameter settings;
OUTPUT:
-options: an output struct with the chosen parameters

```

```

function [model] = init_model(options, d, nb_class)
%This is to initialize a model prior to invoking a specific online algorithm.
INPUT:
-options: the method name and parameter setting
-d: feature dimensionality of input training data
-nb_class: number of classes
OUTPUT:
-model: a struct of the weight vector (w) and the SV indexes

```

```

function [model, hat_y_t, l_t] = feval(options.method, y_t, x_t, model)
%This is to invoke an online algorithm for learning a model from an instance (x_t, y_t).
For example, if "option.method='Perceptron'", this is equivalent to invoking
[model, hat_y_t, l_t] = Perceptron(y_t, x_t, model)
INPUT:
-x_t: feature of the instance
-y_t: class label of the instance
-model: a struct of the weight vector (w) and the SV indexes
OUTPUT:
-model: a struct of the weight vector (w) and the SV indexes
-hat_y_t: the predicted class label of the instance
-l_t: the loss suffered by the learner

```

References

- [1] N. Cesa-Bianchi, A. Conconi, and C. Gentile. A second-order perceptron algorithm. *SIAM J. Comput.*, 34(3):640–668, 2005.
- [2] K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, and Y. Singer. Online passive-aggressive algorithms. *Journal of Machine Learning Research*, 7:551–585, 2006.

- [3] K. Crammer, M. Dredze, and A. Kulesza. Multi-class confidence weighted algorithms. In *EMNLP*, pages 496–504, 2009.
- [4] K. Crammer, M. Dredze, and F. Pereira. Exact convex confidence-weighted learning. In *NIPS*, pages 345–352, 2008.
- [5] K. Crammer, A. Kulesza, and M. Dredze. Adaptive regularization of weight vectors. In *NIPS*, pages 414–422, 2009.
- [6] K. Crammer, A. Kulesza, and M. Dredze. Adaptive regularization of weight vectors. *Machine Learning*, 91(2):155–187, 2013.
- [7] K. Crammer and D. D. Lee. Learning via gaussian herding. In *NIPS*, pages 451–459, 2010.
- [8] K. Crammer and Y. Singer. Ultraconservative online algorithms for multiclass problems. *Journal of Machine Learning Research*, 3:951–991, 2003.
- [9] M. Fink, S. Shalev-Shwartz, Y. Singer, and S. Ullman. Online multiclass learning by interclass hypothesis sharing. In *Proceedings of the 25th International Conference on Machine learning (ICML’06)*, pages 313–320, 2006.
- [10] C. Gentile. A new approximate maximal margin classification algorithm. *Journal of Machine Learning Research*, 2:213–242, 2001.
- [11] Y. Li and P. M. Long. The relaxed online maximum margin algorithm. *Machine Learning*, 46(1-3):361–387, 2002.
- [12] F. Orabona and K. Crammer. New adaptive algorithms for online classification. In *NIPS*, pages 1840–1848, 2010.
- [13] F. Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psych. Rev.*, 7:551–585, 1958.
- [14] J. Wang, P. Zhao, and S. C. H. Hoi. Exact soft confidence-weighted learning. In *ICML*, 2012.
- [15] L. Yang, R. Jin, and J. Ye. Online learning by ellipsoid method. In *ICML*, page 145, 2009.
- [16] P. Zhao, S. C. H. Hoi, and R. Jin. Double updating online learning. *Journal of Machine Learning Research*, 12:1587–1615, 2011.
- [17] P. Zhao, S. C. H. Hoi, R. Jin, and T. Yang. Online auc maximization. In *ICML*, pages 233–240, 2011.
- [18] M. Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *ICML*, pages 928–936, 2003.