

ClassificalO: machine learning for classification graphical user interface

Raeuf Roushangar^{1,2}, George I. Mias^{1,2,*}

¹Department of Biochemistry and Molecular Biology,

²Institute for Quantitative Health Science and Engineering,

Michigan State University, East Lansing MI 48824, USA

Abstracts

Machine learning methods are being used routinely by scientists in many research areas, typically requiring significant statistical and programming expertise. Here we present ClassificalO, an open-source Python graphical user interface for machine learning classification for the scikit-learn Python library. ClassificalO provides an interactive way to train, validate, and test data on a range of classification algorithms. The software enables fast comparisons within and across classifiers, and facilitates uploading and exporting of trained models, and both validation, and testing data results. ClassificalO enables biomedical and other researchers to apply machine learning algorithms to their research in an interactive point-and-click way.

SUPPLEMENTARY CONTENTS

TABLE 1: Classification Algorithms.....3

ClassicalIO User Manual.....4

SUPPLEMENTARY FILES (S1-7 available with manuscript).

	<i>File Name</i>	<i>Description</i>
1.	<i>S1_Iris_Dependent_DataSet.csv</i>	Iris data set (150 data points)
2.	<i>S2_Iris_Target.csv</i>	Iris Target data set (150 labels)
3.	<i>S3_Iris_Testing_DataSet.csv</i>	Iris Testing data set (150 data points)
4.	<i>S4_Iris_FeatureNames.csv</i>	Example Iris features (2 features: sepal length and petal width)
5.	<i>S5_LogisticRegression_IrisTrainedModel.pkl</i>	Example ClassicalIO trained model using logistic regression
6.	<i>S6_IrisTrainValidationResult.csv</i>	Example ClassicalIO testing result using logistic regression
7.	<i>S7_IrisTestingResult.csv</i>	Example ClassicalIO validation result using logistic regression

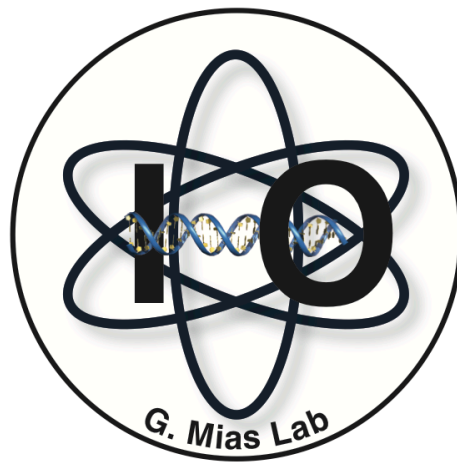
All Supplementary Files below (S8-S14) are referenced in the user manual, and are available online at GitHub (due to journal file size restrictions):

(<https://github.com/gmiaslab/manuals/tree/master/ClassicalIO/Supplementary%20Files>)

8.	<i>S8_GPL570_GSE99039_Dependent_DataSet.csv</i>
9.	<i>S9_GPL570_GSE99039_Target.csv</i>
10.	<i>S10_Y_ChromosomeGenes_FeatureNames.csv</i>
11.	<i>S11_GPL570_GSE18781_Testing_DataSet.csv</i>
12.	<i>S12_LinearSVC_GPL570_GSE99039TrainedModel.pkl</i>
13.	<i>S13_GPL570_GSE99039TrainValidationResult.csv</i>
14.	<i>S14_GPL570_GSE18781TestingResult.csv</i>

CLASSIFIER	Scikit-learn FUNCTION USED	IMMUTABLE PARAMETERS
1: Logistic regression	LogisticRegression	class_weight = None
2: Passive Aggressive	PassiveAggressiveClassifier	class_weight = None n_iter= None
3: Perceptron	Perceptron	class_weight = None
4: Classifier using Ridge regression	RidgeClassifier	class_weight = None
5: Stochastic Gradient Descent - SGD	SGDClassifier	—
6: Linear Discriminant Analysis	LinearDiscriminantAnalysis	shrinkage= None priors = None
7: Quadratic Discriminant Analysis	QuadraticDiscriminantAnalysis	store_covariances = None priors = None
8: Linear Support Vector	LinearSVC	class_weight = None
9: Nu-Support Vector	NuSVC	class_weight = None
10: C-Support Vector	SVC	class_weight = None
11: k-Nearest Neighbors	KNeighborsClassifier	metric_params = None
12: Nearest centroid	NearestCentroid	—
13: Radius Nearest Neighbors	RadiusNeighborsClassifier	metric_params = None
14: Gaussian Process Classification (GPC)	GaussianProcessClassifier	kernel = None
15: Naive Bayes for Multivariate Bernoulli Models	BernoulliNB	class_prior = None
16: Gaussian Naive Bayes	GaussianNB	class_prior = None
17: Naive Bayes for Multinomial Models	MultinomialNB	class_prior = None
18: Decision Tree	DecisionTreeClassifier	class_weight = None
19: Extremely Randomized Tree	ExtraTreeClassifier	min_impurity_split = None class_weight = None
20: AdaBoost	AdaBoostClassifier	base_estimator = None
21: Bagging	BaggingClassifier	base_estimator = None
22: Extra Trees	ExtraTreesClassifier	class_weight = None
23: Random Forest	RandomForestClassifier	class_weight = None
24: Label Propagation	LabelPropagation	—
25: Neural network Multi-layer Perceptron	MLPClassifier	—

Table 1. A list of all 25 classification algorithms, their corresponding scikit-learn functions, and immutable (unchangeable) parameters with their default values.



ClassificaIO

Machine Learning for Classification Graphical
User Interface User Manual 1.1.3 (07/2018)

Summary:

ClassificalO is an open-source Python graphical user interface (GUI) for supervised machine learning classification for the scikit-learn module [1]. ClassificalO aims to provide an easy-to-use interactive way to train, validate, and test data on a range of classification algorithms. The GUI enables fast comparisons within and across classifiers, and facilitates uploading and exporting of trained models, and both validated, and tested data results.

Dependencies:

ClassificalO is a Python package with the following external dependencies:

- Tkinter ≥ 8.6.7, Pillow ≥ 5.1.0, pandas == 0.22.0, numpy ≥ 1.14.3, scikit-learn ≥ 0.19.1, and scipy ≥ 1.1.0

Prerequisites:

ClassificalO requires Python version 3.6 or higher and can be used on Mac OS X High Sierra, Linux (tested on Ubuntu), and Windows 10 operating systems. To avoid any system errors, crashes, and crude fonts, we recommend to not install ClassificalO using integrated environment package installers – i.e. native installation of ClassificalO is highly encouraged using pip. In case you do not have pip installed, you must install it first.

Installation Instructions:

1. Mac or Windows

To install the current release use pip in the terminal:

```
$ pip install ClassificalO
```

Alternatively, you can install directly from GitHub using:

```
$ pip install git+https://github.com/gmiaslab/ClassificalO/
```

2. Linux

First install the current release of tkinter and pip:

```
$ sudo apt-get install python3-tk  
$ sudo apt-get install python3-pip
```

To install the current ClassificalO release use pip:

```
$ pip3 install ClassificalO
```

Alternatively, you can install directly from GitHub using:

```
$ pip install git+https://github.com/gmiaslab/ClassificalO/
```

Getting started with ClassificalO:

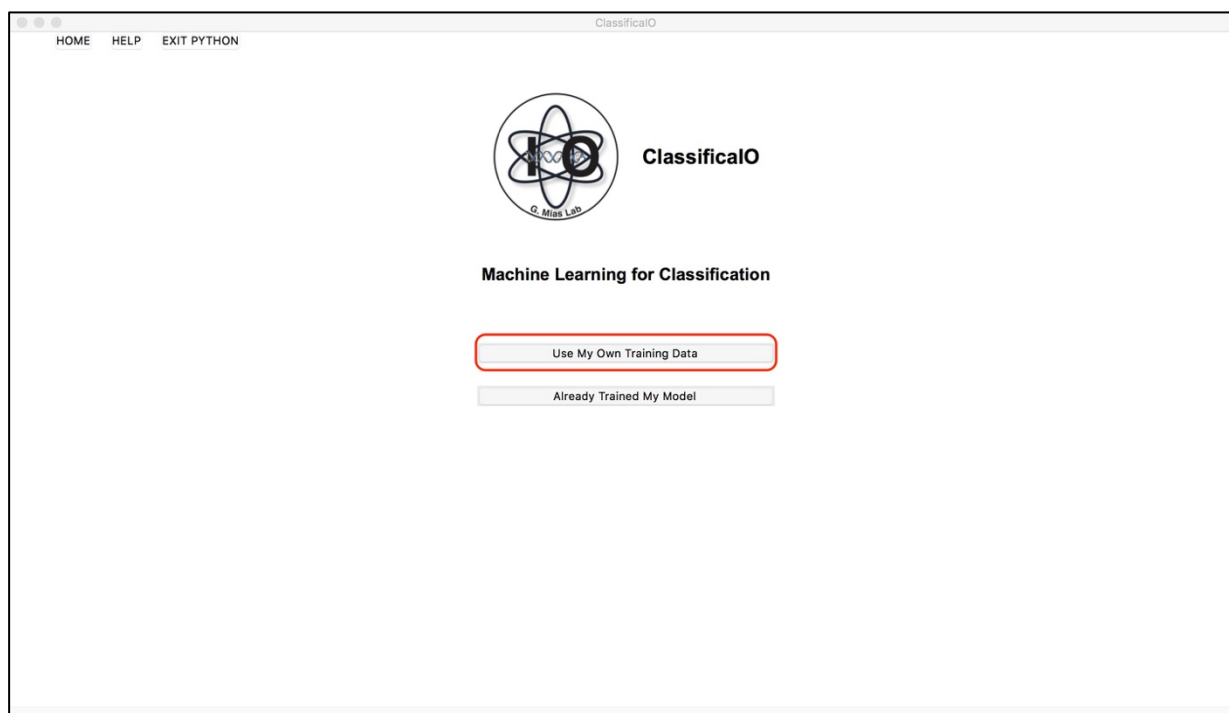
Please Note:

- ClassificalO supports comma-separated values (CSV) input files only.
- In this document we use the machine learning Iris dataset [2, 3] (150 data points) as an example to illustrate the interface, to demonstrate model training, validation, and testing, as well as the data formats that ClassificalO relies on.
- In the classification example below, we use 70% of the Iris dataset (105 data points) for model training and 30% (45 data points) for model testing.
- You can download the Iris and the processed gene expression datasets, as well as all files used in this manual from:

<https://github.com/gmiaslab/manuals/tree/master/ClassificalO/Supplementary%20Files>

After installing ClassificalO, please run it from the terminal using Python:

```
$ python3
>>> from ClassificalO import ClassificalO
>>> ClassificalO.gui()
```

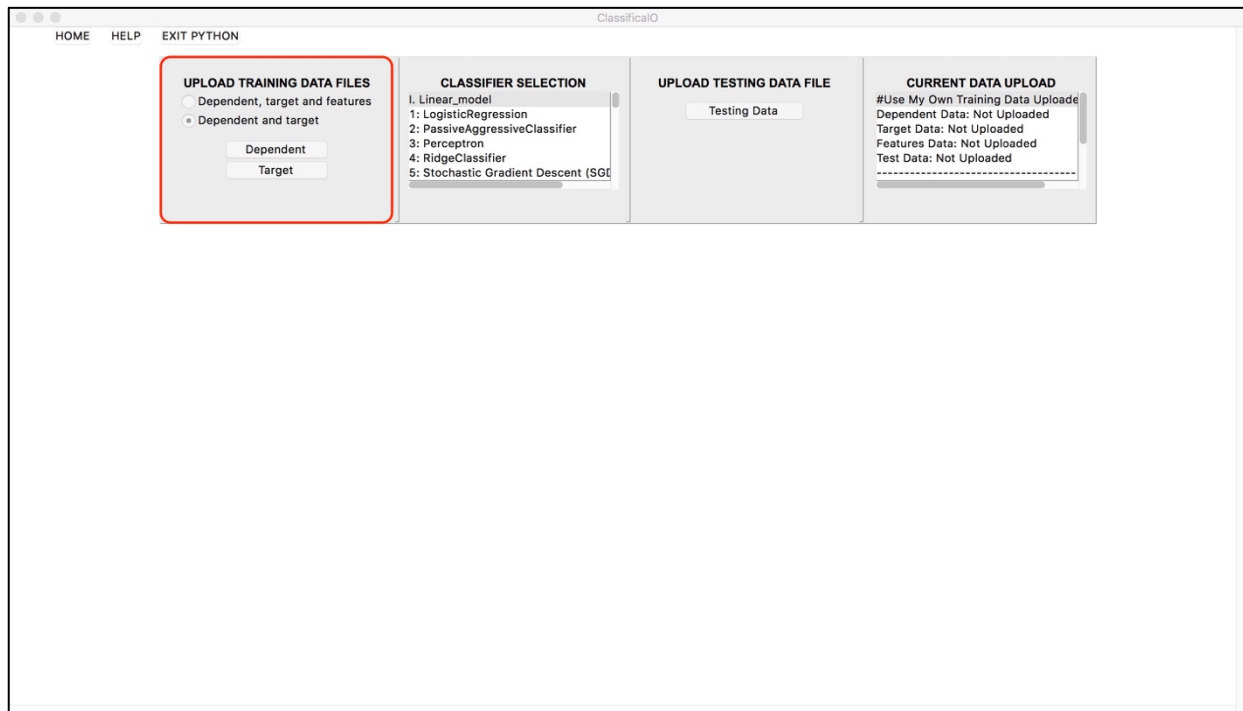


We note here the name is case sensitive (i.e. the 'IO' is capitalized). Once ClassificalO's main window appears on your screen, you can click on 'Use My Own Training Data' button and start your supervised machine learning classification project.

Iris dataset prediction using a logistic regression classifier:

Training data input:

You first need to make a selection (either 'Dependent and Target' or 'Dependent, Target and Features') from the 'UPLOAD TRAINING DATA FILES' panel to upload training data files. For this example, we select the 'Dependent and Target' button.



To begin uploading data files, click the corresponding buttons in the 'UPLOAD TRAINING DATA FILES' panel: a file selector directs you to upload both, dependent data file (e.g. **Fig. 1(a)**) and target data file (e.g. **Fig. 1(b)**). Once a file is uploaded to ClassificalIO, the file name and directory are automatically saved in the 'CURRENT DATA UPLOAD' panel (**Fig. 2**). This updatable log allows for tracking current data files in-use, and maintains a history of all files uploaded to the software.

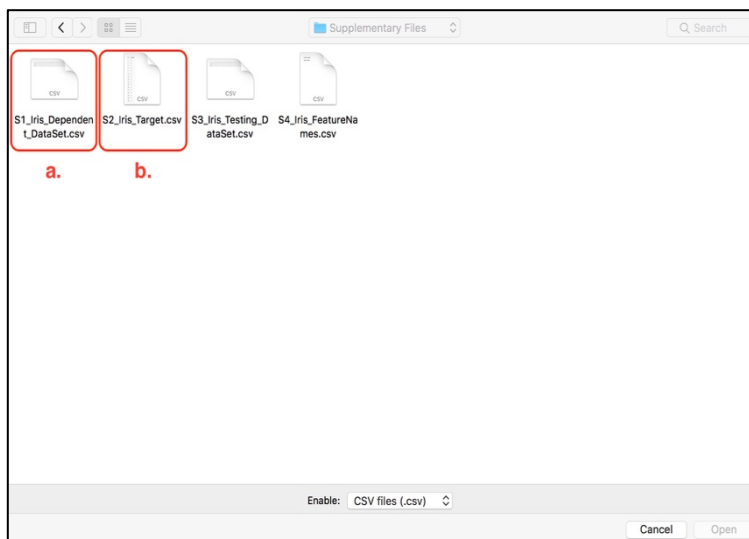


Fig. 1. Graphical Control Element Dialog Box. **(a)** Dependent data file selected for upload. **(b)** Selected target data file to upload. N.B. each file selection has to be done one at a time.

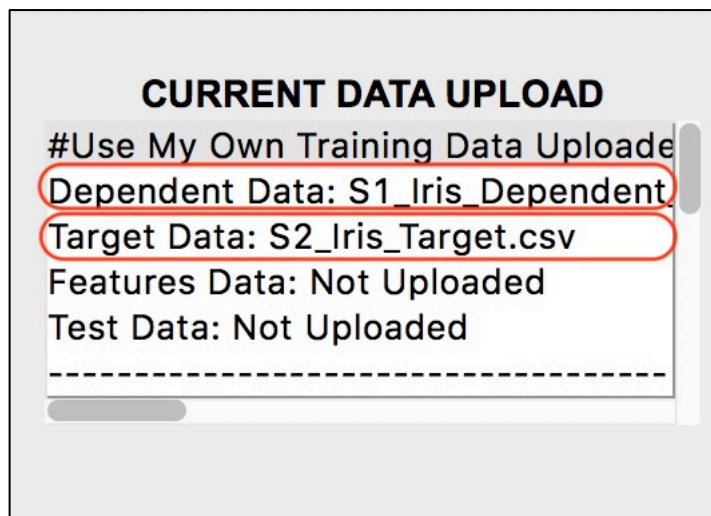


Fig. 2. Current Data Upload Panel. Both dependent and target data file names shown (red boxes). Scroll down for uploaded data files directories.

Data format:

Data formats are shown in **Fig. 3(a)** for dependent data and **Fig. 3(b)** for target data. The dependent data represent the data on which the model will depend on for learning and the target data is the annotation, i.e. what is going to be predicted. In this example, the dependent data have 4 rows and 105 columns. For the dependent data, each row is an attribute (also known as feature) and each column is an object (also known as an observation or a sample). Thus, the header row enumerates the objects, and the header column names the attributes. The values in the file represent the measurement made for each of the objects(columns) for each of the attributes (rows).

		Objects																							
		A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X
Attributes	1		1	2	5	6	7	11	12	14	16	18	20	21	22	23	24	25	26	27	28	29	30	31	32
	2	sepal length	5.1	4.9	5	5.4	4.6	5.4	4.8	4.3	5.7	5.1	5.1	5.4	5.1	4.6	5.1	5.1	5.1	5.1	5.1	5.1	5.1	5.1	5.1
	3	sepal width	3.5	3	3.6	3.9	3.4	3.7	3.4	3	4.4	3.5	3.8	3.4	3.7	3.6	3.3	3.3	3.3	3.3	3.3	3.3	3.3	3.3	3.3
	4	petal length	1.4	1.4	1.4	1.7	1.4	1.5	1.6	1.1	1.5	1.4	1.5	1.7	1.5	1	1.7	1.7	1.7	1.7	1.7	1.7	1.7	1.7	1.7
	5	petal width	0.2	0.2	0.2	0.4	0.3	0.2	0.2	0.1	0.4	0.3	0.3	0.2	0.4	0.2	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5

Fig. 3(a) Dependent Data. Example of partial dependent data file format. Testing data (not shown) uses the same format.

For the target data, we have 105 rows and 2 columns (note: for the target data, the rows correspond to the objects and column is the class per object). The values in the “ids” column in the target data must match the Objects header row in the dependent data, the columns headers must match, otherwise an error will occur. Hence, the number of columns (i.e. objects) in the dependent data must also match the number of rows in the target data (i.e. each object has a unique “id” and must be assigned a target class for training). Finally, the “target” column in the target data must be numerically-valued.

Objects	Attributes	
	A	B
	id	target
	2	1
	3	2
	4	5
	5	6
	6	7
	7	11
	8	12
	9	14
	10	16
	11	18
	12	20
	13	21
	14	22
	15	23
	16	24
		0

Fig. 3(b) Target Data. Example of partial target data file format where the targets correspond to setosa = 0, versicolor = 1, and virginica = 2. Versicolor and virginica are not visible in this screenshot.

Classifier selection:

Once you have uploaded all required training data files, you can select between 25 different machine learning classification algorithms in the 'CLASSIFIER SELECTION' panel.

The screenshot shows the ClassifaiO web interface. At the top, there are navigation links for 'HOME' and 'HELP'. The main content area is divided into four panels: 'UPLOAD TRAINING DATA FILES', 'CLASSIFIER SELECTION', 'UPLOAD TESTING DATA FILE', and 'CURRENT DATA UPLOAD'. The 'CLASSIFIER SELECTION' panel is highlighted with a red box and contains a list of classifiers: 1. Linear_model, 2. LogisticRegression (selected), 3. PassiveAggressiveClassifier, 4. Perceptron, 5. RidgeClassifier, and 6. Stochastic Gradient Descent (SGD). Below this list, there is a description of Logistic Regression and a 'Submit' button. The bottom section of the interface contains three panels: 'CONFUSION MATRIX, MODEL ACCURACY & ERROR', 'TRAINING RESULT: ID — ACTUAL — PREDICTION', and 'TESTING RESULT: ID — PREDICTION', each with an 'Export' button.

Here are all classification algorithms in order of appearance in the 'CLASSIFIER SELECTION' panel. Immutable (unchangeable) parameters with their default values are also listed for each classifier in the parentheses:

I. Linear_model

- 1: LogisticRegression. (class_weight = None)
- 2: PassiveAggressiveClassifier. (class_weight = None, n_iter= None)
- 3: Perceptron. (class_weight = None)
- 4: RidgeClassifier. (class_weight = None)
- 5: Stochastic Gradient Descent (SGDClassifier).

II. Discriminant_analysis

- 6: LinearDiscriminantAnalysis. (shrinkage= None, priors = None)
- 7: QuadraticDiscriminantAnalysis. (store_covariances = None, priors = None)

III. Support vector machines (SVMs)

- 8: LinearSVC. (class_weight = None)
- 9: NuSVC. (class_weight = None)
- 10: SVC. (class_weight = None)

IV. Neighbors

- 11: KNeighborsClassifier. (metric_params = None)
- 12: NearestCentroid.
- 13: RadiusNeighborsClassifier. (metric_params = None)

V. Gaussian_process

- 14: GaussianProcessClassifier. (kernel = None)

VI. Naive_bayes

- 15: BernoulliNB. (class_prior = None)
- 16: GaussianNB. (class_prior = None)
- 17: MultinomialNB. (class_prior = None)

VII. Trees

- 18: DecisionTreeClassifier. (class_weight = None)
- 19: ExtraTreeClassifier. (min_impurity_split = None, class_weight = None)

VIII. Ensemble

- 20: AdaBoostClassifier. (base_estimator = None)
- 21: BaggingClassifier. (base_estimator = None)
- 22: ExtraTreesClassifier. (class_weight = None)
- 23: RandomForestClassifier. (class_weight = None)

IX. Semi_supervised

- 24: LabelPropagation.

X. Neural_network

- 25: MLPClassifier.

The following will populate once you make a classifier selection:

- **Fig. 4(a):** The classifier definition with a clickable underlined link “learn more” in blue, which, when clicked opens an external web-browser to the scikit-learn documentation for the selected classifier.
- **Fig. 4(b):** Interactive way to select between train-validate split and cross-validation methods (radio buttons), which are necessary to prevent/minimize training model overfitting.
- **Fig. 4(c):** classifier parameters, to provide you with a point-and-click interface to set, modify, and test the influence of each parameter on your data

The screenshot displays the ClassifierIO web interface. At the top, there are navigation links: HOME, HELP, and EXIT PYTHON. The main content area is divided into four panels: UPLOAD TRAINING DATA FILES, CLASSIFIER SELECTION, UPLOAD TESTING DATA FILE, and CURRENT DATA UPLOAD. The CLASSIFIER SELECTION panel shows a list of classifiers: 1. Linear_model, 1: LogisticRegression (selected), 2. PassiveAggressiveClassifier, 3. Perceptron, 4. RidgeClassifier, and 5. Stochastic Gradient Descent (SGD). Below this, a red box labeled 'a.' highlights the classifier definition for LogisticRegression, which includes a clickable link "Learn more." and a description: "Logistic regression, despite its name, is a linear model for classification rather than regression. Logistic regression is also known in the literature as logit regression, maximum-entropy classification (MaxEnt) or the log-linear classifier." Below this, a red box labeled 'b.' highlights the training methods section, which includes a slider for Train Sample Size (%) set to 75, and radio buttons for K-fold Cross-Validation (set to 10) and Train Sample Size (%). A red box labeled 'c.' highlights the classifier parameters section, which includes various settings such as random_state, multi_class, intercept_scaling, verbose, fit_intercept, dual, warm_start, integer, penalty, max_iter, n_jobs, solver, tol, and C, all set to their default values. A Submit button is located below the parameters section. At the bottom, there are three panels: CONFUSION MATRIX, MODEL ACCURACY & ERROR, TRAINING RESULT: ID — ACTUAL — PREDICTION, and TESTING RESULT: ID — PREDICTION, each with an Export button.

Fig. 4. Selected Logistic Regression Classifier. The interface for each selected classifier, has uniform features. **(a)** Classifier definition is displayed, together with an underlined clickable link that reads “Learn more” next to the classifier name. **(b)** Training methods with ‘Train Sample Size (%)’ method selected. **(c)** The classifier parameters set to their default values.

Model training, evaluation, validation and result output:

You can now click 'submit' to train your classifier using the uploaded training data files 'Dependent and Target' in this example, and evaluate your result. Or, alternatively you can upload testing data first, and then click 'submit' to train and test a classifier on your uploaded data at the same time. For this example, **first**: we train a selected classifier, 'LogisticRegression', using its default parameters, and default train-validate split method 'Train Sample Size (%)', and then, **second**: we upload testing data to test the trained model.

After clicking 'submit', our selected classifier, 'LogisticRegression' for this example, is trained using the loaded training data, 'Dependent and Target' for this example.

Notes

ClassificalO always shuffles your training data before splitting to eliminate mini batch effects.

Internally, when 'Train Sample Size (%)' method is selected, ClassificalO uses the scikit-learn *train_test_split* method, to allow for fast training data split into training and validation subsets. With this method the parameter set is *train_size*, which takes the train sample size set by you (e.g. Train Sample Size (%): set to 75% means *train_size* = 0.75 and *test_size* = 0.25).

If the 'K-fold Cross-Validation' method is selected instead, ClassificalO uses the scikit-learn *cross_val_predict* method where the training data is split into k-sets. The model is trained on k-1 of the folds followed by a validation step on the remaining part of the data. This will be repeated for each of the k-folds.

After training is completed, the confusion matrix, classifier accuracy and error are displayed in the 'CONFUSION MATRIX, MODEL ACCURACY & ERROR' panel (**Fig. 5(a)**). Model validation data results are displayed in the 'TRAINING RESULT: ID – ACTUAL – PREDICTION' panel (**Fig. 5(b)**) with each data point ID is the 1st value, actual target value is displayed 2nd, and predicted target value 3rd, where the predictions correspond to the iris flower species, with 0=setosa, 1=versicolor, and 2=virginica.

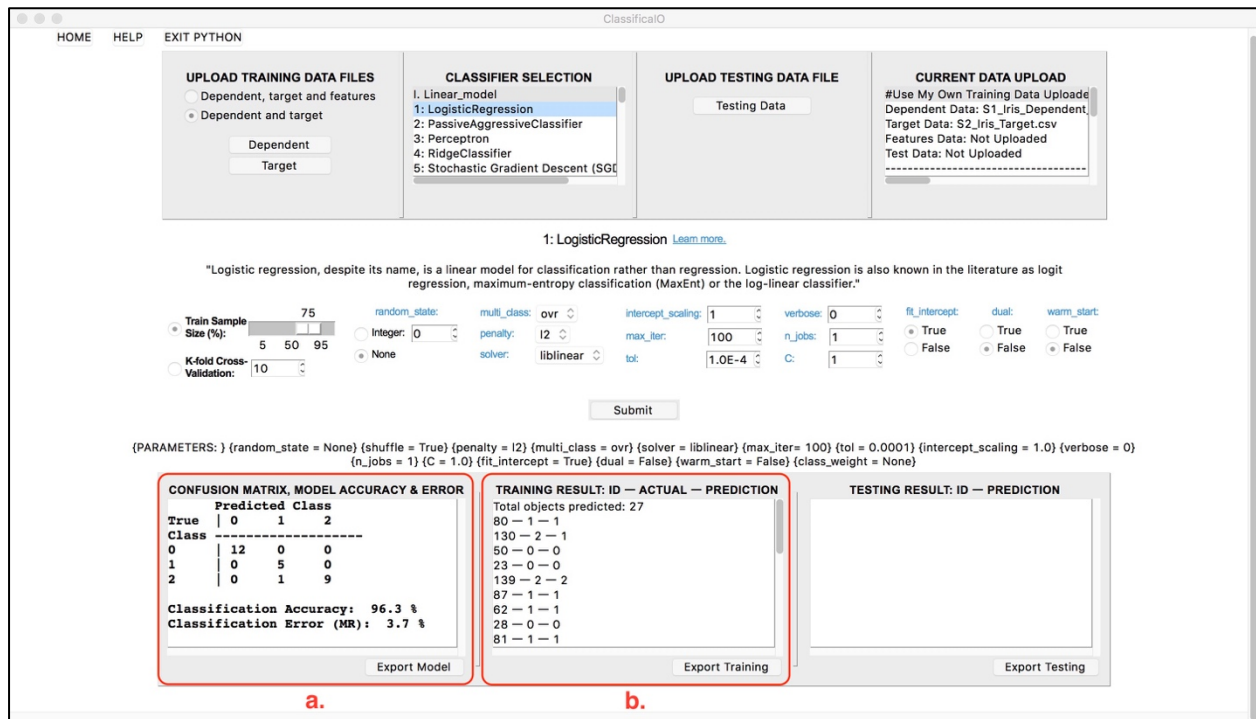


Fig. 5. Trained Logistic Regression Classifier. **(a)** Trained model using 78 data points (75% of 105 data points), classifier evaluation (confusion matrix, model accuracy and error). **(b)** Model validated using 27 data points (25% of 105 data points).

Testing data input and result output:

To test your trained model, first upload the testing data file by clicking the 'Testing Data' button in the 'UPLOAD TESTING DATA FILE' panel (**Fig. 6(a)**). Once clicked, a file selector directs you to upload the testing data file, see **Fig. 1** Once testing data is uploaded, the file name is automatically saved in the 'CURRENT DATA UPLOAD' panel (outlined in the red box in the figure) to indicate that your file has been uploaded. The Testing Data file format is the same as for the dependent data file, see **Fig. 3(a)**.

After clicking 'Submit', testing results are displayed in the 'TESTING RESULT: ID – PREDICTION' panel (**Fig. 6(b)**) with each data point ID shown 1st, and the corresponding predicted target value displayed after it 2nd, separated by a hyphen.

The screenshot shows the ClassifierIO web interface. At the top, there are navigation links: HOME, HELP, EXIT PYTHON. The main interface is divided into several panels:

- UPLOAD TRAINING DATA FILES:** Includes radio buttons for 'Dependent, target and features' and 'Dependent and target', with buttons for 'Dependent' and 'Target'.
- CLASSIFIER SELECTION:** A list of classifiers: 1: Linear_model, 2: LogisticRegression (selected), 3: PassiveAggressiveClassifier, 4: Perceptron, 5: RidgeClassifier, 6: Stochastic Gradient Descent (SGD).
- UPLOAD TESTING DATA FILE:** A button labeled 'Testing Data' is highlighted with a red box and labeled 'a.'.
- CURRENT DATA UPLOAD:** A panel outlined in red and labeled 'b.' showing the following text:

```
#Use My Own Training Data Upload
Dependent Data: S1_Iris_Dependent
Target Data: S2_Iris_Target.csv
Features Data: Not Uploaded
Test Data: S3_Iris_Testing_DataSet.csv
```

Below these panels, there is a description of Logistic Regression and a 'Submit' button. The parameters section shows the following values:

- Train Sample Size (%): 75
- random_state: Integer: 0
- multi_class: ovr
- penalty: l2
- solver: liblinear
- intercept_scaling: 1
- max_iter: 100
- tol: 1.0E-4
- verbose: 0
- n_jobs: 1
- C: 1
- fit_intercept: True
- dual: False
- warm_start: False

The parameters are displayed as: {PARAMETERS: } {random_state = None} {shuffle = True} {penalty = l2} {multi_class = ovr} {solver = liblinear} {max_iter= 100} {tol = 0.0001} {intercept_scaling = 1.0} {verbose = 0} {n_jobs = 1} {C = 1.0} {fit_intercept = True} {dual = False} {warm_start = False} {class_weight = None}

Below the parameters, there are three panels:

- CONFUSION MATRIX, MODEL ACCURACY & ERROR:** Shows a confusion matrix and classification accuracy of 96.3% and error of 3.7%.
- TRAINING RESULT: ID – ACTUAL – PREDICTION:** Shows a list of training results for 27 objects.
- TESTING RESULT: ID – PREDICTION:** A panel outlined in red showing testing results for 45 objects, with each object ID followed by its predicted value (e.g., 3 - 0, 4 - 0, etc.).

Fig. 6. Tested Logistic Regression Classifier. **(a)** Upload testing data panel. **(b)** Model tested using 45 data points.

Result export:

Now you are ready to export your trained model to preserve it for future use without having to retrain. Simply, click the 'Export Model' button (**Fig. 5(a)**) and save your model. Your exported ClassificalO model can then be used for future testing on new data in the 'Already Trained My Model' window in ClassificalO, shown below.

The screenshot shows the ClassificalO web application interface. At the top, there is a navigation bar with links for HOME, HELP, and EXIT PYTHON. The main content area is divided into several sections. On the left, there is a section titled 'UPLOAD TRAINING MODEL FILE' with a 'Model File' button. In the center, there is a section titled 'UPLOAD TESTING DATA FILE' with a 'Testing Data' button. On the right, there is a section titled 'CURRENT DATA UPLOAD' which contains a list of items: '#Already Trained My Model Uploade', 'Model: Not Uploaded', 'Test Data: Not Uploaded', and '#Upload History'. Below these sections, there is a 'Submit' button. At the bottom, there is a section titled 'CONFUSION MATRIX, MODEL ACCURACY & ERROR' and 'TESTING RESULT: ID — PREDICTION'. This section contains two large empty boxes for displaying results. An 'Export Testing' button is located at the bottom right of this section.

ClassificalO model input:

You will need to upload ClassificalO model by clicking the 'Model File' button in the 'UPLOAD TRAINING MODEL FILE' panel (**Fig. 7(a)**). Once clicked, a file selector directs you to upload a ClassificalO trained model. Also, you will need to upload a testing data file (the testing data file format is the same as explained above), by clicking the 'Testing Data' button in the "UPLOAD TESTING DATA FILE" panel (**Fig. 7(b)**). Once a ClassificalO model and testing data files are uploaded, files names are automatically displayed in the 'CURRENT DATA UPLOAD' panel (**Fig. 7(c)**).

After clicking 'submit', the uploaded model preset parameters will populate (**Fig. 7(d)**) to show the classifier used to originally train the uploaded model. The confusion matrix, classifier accuracy and error of trained model are then displayed in the 'CONFUSION MATRIX, MODEL ACCURACY & ERROR' panel (**Fig. 7(e)**). Testing data results are displayed in the 'Testing RESULT: ID – PREDICTION' panel (**Fig. 7(f)**) with the data point ID shown 1st, followed by a hyphen and the predicted value displayed right after it.

The screenshot shows the ClassificalO web interface with the following panels:

- (a) UPLOAD TRAINING MODEL FILE:** Contains a 'Model File' button.
- (b) UPLOAD TESTING DATA FILE:** Contains a 'Testing Data' button.
- (c) CURRENT DATA UPLOAD:** Displays the names of the uploaded files: '#Already Trained My Model Upload', 'Model: S5_LogisticRegression_IrisTr', and 'Test Data: S3_Iris_Testing_DataSet.c'. It also has an '#Upload History' button.
- (d) CLASSIFIER:** Displays the model's parameters: 'PARAMETERS: ', 'random_state = None', 'shuffle = True', 'penalty = l2', 'multi_class = ovr', 'solver = liblinear', 'max_iter = 100', 'tol = 0.0001', 'intercept_scaling = 1.0', 'verbose = 0', 'n_jobs = 1', 'C = 1.0', 'fit_intercept = True', 'dual = False', 'warm_start = False', 'class_weight = None'.
- (e) CONFUSION MATRIX, MODEL ACCURACY & ERROR:** Displays a confusion matrix and model evaluation metrics.

	Predicted Class		
True	0	1	2
Class			
0	6	0	0
1	0	11	1
2	0	0	9

Classification Accuracy: 96.3 %
Classification Error (MR): 3.7 %
- (f) TESTING RESULT: ID – PREDICTION:** Displays the testing results for 10 data points.

ID	Prediction
48	0
55	1
56	1
57	2
67	2
69	1
72	1
76	1
79	1
84	2

Fig. 7. 'Already Trained My Model' window. **(a)** Upload ClassificalO trained model panel. **(b)** Upload testing data panel. **(c)** Current data upload panel with both model and testing data files names shown (red boxes). **(d)** Model preset parameters. **(e)** Trained model result and model evaluation (confusion matrix, model accuracy and error). **(f)** Model testing result.

Results Export:

Full results (trained models, and both trained/validated, and tested data) for both windows (**'Use My Own Training Data'** and **'Already Trained My Model'**) can be exported as CSV files for later use, e.g. further analysis, publication, sharing, etc. (for an example on the export data file formats, see the **Supplementary files S6 and S7**).

Additional Examples:

Ex1. Iris dataset prediction using k-nearest neighbors classifier:

Here we follow the same process as above but select k-nearest neighbors classification algorithm 'KNeighborsClassifier' rather than logistic regression. The Classification Accuracy shown here to be 92.59 % with low classification error. When the train Sample Size (%) change to 25% means $train_size = 0.25$ and $test_size = 0.75$, metric set to Chebyshev for distance metric, and n-neighbors to 10 for number of neighbors, the Classification Accuracy goes down to 63.29 % with much higher classification error. Please note that you may get different results when running the program, given the random selection of samples used for training/validation.

a.

UPLOAD TRAINING DATA FILES
☐ Dependent, target and features
☒ Dependent and target
 Dependent:
 Target:

CLASSIFIER SELECTION
 8: LinearSVC
 9: NuSVC
 10: SVC
 11: KNeighborsClassifier
 12: NearestCentroid

UPLOAD TESTING DATA FILE

CURRENT DATA UPLOAD
☒ #Use My Own Training Data Upload
 Dependent Data: S1_Iris_Dependent
 Target Data: S2_Iris_Target.csv
 Features Data: Not Uploaded
 Test Data: S3_Iris_Testing_DataSet.csv

11: KNeighborsClassifier [Learn more](#)

Neighbors-based classification is a non-generalizing learning: it does not attempt to construct a general internal model, but simply stores instances of the training data. KNeighborsClassifier implements learning based on the k nearest neighbors of each query point, where k is an integer value specified by the user.

Train Sample Size (%): 75
 random_state:
 Integer: 0
 None
 algorithm: auto
 weights: uniform
 metric: minkowski
 n_neighbors: 5
 leaf_size: 30
 n_jobs: 1
 p: 2
 K-fold Cross-Validation: 10

(PARAMETERS:) {random_state = None} {shuffle = True} {metric = minkowski} {weights = uniform} {algorithm = auto} {n_neighbors = 5} {leaf_size = 30} {n_jobs = 1} {p = 2} {metric_params = None}

CONFUSION MATRIX, MODEL ACCURACY & ERROR

True \ Predicted Class	Predicted Class		
	0	1	2
Class			
0	12	0	0
1	0	4	2
2	0	0	9

 Classification Accuracy: 92.59 %
 Classification Error (MR): 7.41 %

TRAINING RESULT: ID — ACTUAL — PREDICTION
 Total objects predicted: 27
 36 — 0 — 0
 73 — 1 — 2
 23 — 0 — 0
 85 — 1 — 1
 22 — 0 — 0
 64 — 1 — 2
 25 — 0 — 0
 43 — 0 — 0
 149 — 2 — 2

TESTING RESULT: ID — PREDICTION
 Total objects tested: 45
 3 — 0
 4 — 0
 8 — 0
 9 — 0
 10 — 0
 13 — 0
 15 — 0
 17 — 0
 19 — 0

b.

UPLOAD TRAINING DATA FILES
☐ Dependent, target and features
☒ Dependent and target
 Dependent:
 Target:

CLASSIFIER SELECTION
 8: LinearSVC
 9: NuSVC
 10: SVC
 11: KNeighborsClassifier
 12: NearestCentroid

UPLOAD TESTING DATA FILE

CURRENT DATA UPLOAD
☒ #Use My Own Training Data Upload
 Dependent Data: S1_Iris_Dependent
 Target Data: S2_Iris_Target.csv
 Features Data: Not Uploaded
 Test Data: S3_Iris_Testing_DataSet.csv

11: KNeighborsClassifier [Learn more](#)

Neighbors-based classification is a non-generalizing learning: it does not attempt to construct a general internal model, but simply stores instances of the training data. KNeighborsClassifier implements learning based on the k nearest neighbors of each query point, where k is an integer value specified by the user.

Train Sample Size (%): 25
 random_state:
 Integer: 0
 None
 algorithm: auto
 weights: uniform
 metric: chebyshev
 n_neighbors: 10
 leaf_size: 30
 n_jobs: 1
 p: 2
 K-fold Cross-Validation: 10

(PARAMETERS:) {random_state = None} {shuffle = True} {metric = chebyshev} {weights = uniform} {algorithm = auto} {n_neighbors = 10} {leaf_size = 30} {n_jobs = 1} {p = 2} {metric_params = None}

CONFUSION MATRIX, MODEL ACCURACY & ERROR

True \ Predicted Class	Predicted Class		
	0	1	2
Class			
0	25	0	0
1	0	25	0
2	0	29	0

 Classification Accuracy: 63.29 %
 Classification Error (MR): 36.71 %

TRAINING RESULT: ID — ACTUAL — PREDICTION
 Total objects predicted: 79
 24 — 0 — 0
 140 — 2 — 1
 123 — 2 — 1
 54 — 1 — 1
 73 — 1 — 1
 23 — 0 — 0
 40 — 0 — 0
 70 — 1 — 1
 127 — 2 — 1

TESTING RESULT: ID — PREDICTION
 Total objects tested: 45
 3 — 0
 4 — 0
 8 — 0
 9 — 0
 10 — 0
 13 — 0
 15 — 0
 17 — 0
 19 — 0

Fig. 8. (a) selected k-nearest neighbors classifier with trained and tested the data using the default parameters values, **(b)** Same classifier selected with trained and tested data but using different parameters values.

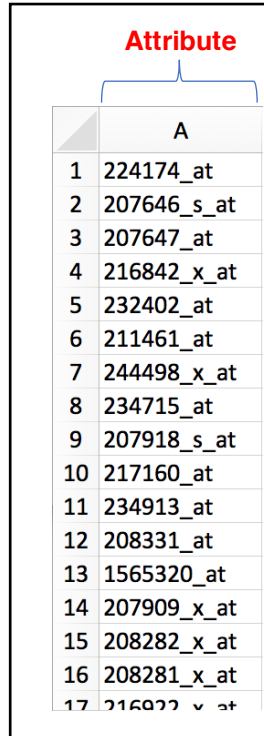
Ex2. Gene expression sex prediction using linear support vector classifier:

In this example we use gene expression data from microarrays to predict the sex of the donor who provided the experimental sample. This is useful when the sex of the donor is unknown, and the information is needed for metadata analysis. We use two raw microarray gene expression datasets obtained from the Gene Expression Omnibus (GEO), GSE99039[4], 121 samples for training and GSE18781[5], 25 samples for testing. The data were generated using Affymetrix Human Genome U133 Plus 2.0 Array, (GPL570). Gene expression arrays need to first be processed and normalized: We first preprocessed all datasets (microarray raw CEL files) separately using the RMA (Robust Multi- Array Average) algorithm [6] implemented in R. Preprocessing consisted of performing background noise correction, quantile normalization of the intensities to obtain the same distribution of intensities across arrays, and summarization using median polish procedure to summarize the data in each probe set to obtain a single value corresponding to the gene expression of a probe set. The expression data were logarithmically transformed (base 2 log), and exported as CSV files.

We select 'Dependent, target and features' button from the 'UPLOAD TRAINING DATA FILES' (**Fig. 9(a)**) and select 'LinearSVC' as classifier (**Fig. 9(b)**). We use the buttons to upload the data: GSE32140 data were used for the Dependent Data button (Supplementary File: S8), with Target data using values as annotated on GEO (Supplementary File: S9), Y chromosomes genes were used for the Features (Supplementary File: S10), and GSE18781 was used for Testing Data (Supplementary File: S11) as an example. The file names are displayed in the 'CURRENT DATA UPLOAD' panel (**Fig. 9(c)**) after selection has been made.

For the classifier, we use the default parameters, and 'K-fold Cross-Validation' set to 10 (**Fig. 9(d)**) to train our model. After training is completed, the confusion matrix, classifier accuracy and error are displayed in the 'CONFUSION MATRIX, MODEL ACCURACY & ERROR' panel (**Fig. 9(e)**). Model validation data results are displayed in the 'TRAINING RESULT: ID – ACTUAL – PREDICTION' panel (**Fig. 9(f)**) with each data point ID is the 1st value, actual target value is displayed 2nd, and predicted target value 3rd where Male = 0 and Female = 1. Testing results are displayed in the 'TESTING RESULT: ID – PREDICTION' panel (**Fig. 9(g)**) with each data point ID shown 1st, and the corresponding predicted target value displayed after it 2nd, separated by a hyphen. Full results (trained models, and both trained/validated, and tested data) were then exported (for more details on the export data file, see Supplementary Files: S13-15).

Note on Features annotation Fig. 10: To identify the Y chromosomes genes, we used the biomaRt [7] which mapped the Y probes on Affymetrix Human Genome U133 Plus 2.0 Array, GPL570, to the human reference genome (for complete list of all 96 Y chromosome gene ids, see file S12).



Attribute	
	A
1	224174_at
2	207646_s_at
3	207647_at
4	216842_x_at
5	232402_at
6	211461_at
7	244498_x_at
8	234715_at
9	207918_s_at
10	217160_at
11	234913_at
12	208331_at
13	1565320_at
14	207909_x_at
15	208282_x_at
16	208281_x_at
17	216922_x_at

Fig. 10. Features Data. Example of partial features data file format where each Affymetrix probe id correspond to a Y chromosome gene.

References

1. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V *et al*: **Scikit-learn: Machine Learning in Python**. *J Mach Learn Res* 2011, **12**:2825-2830.
2. Fisher RA: **The use of multiple measurements in taxonomic problems**. *Ann Eugenetic* 1936, **7**:179-188.
3. Anderson E: **The Irises of the Gaspé peninsula**. *Bulletin of American Iris Society* 1935, **59**:2-5.
4. Shamir R, Klein C, Amar D, Vollstedt EJ, Bonin M, Usenovic M, Wong YC, Maver A, Poths S, Safer H *et al*: **Analysis of blood-based gene expression in idiopathic Parkinson disease**. *Neurology* 2017, **89**(16):1676-1683.
5. Sharma SM, Choi D, Planck SR, Harrington CA, Austin CR, Lewis JA, Diebel TN, Martin TM, Smith JR, Rosenbaum JT: **Insights in to the pathogenesis of axial spondyloarthritis based on gene expression profiles**. *Arthritis Res Ther* 2009, **11**(6):R168.
6. Bolstad BM, Irizarry RA, Astrand M, Speed TP: **A comparison of normalization methods for high density oligonucleotide array data based on variance and bias**. *Bioinformatics* 2003, **19**(2):185-193.
7. Durinck S, Moreau Y, Kasprzyk A, Davis S, De Moor B, Brazma A, Huber W: **BioMart and Bioconductor: a powerful link between biological databases and microarray data analysis**. *Bioinformatics* 2005, **21**(16):3439-3440.