

# **ClassificalO: machine learning for classification graphical user interface**

Raeuf Roushangar<sup>1,2</sup>, George I. Mias<sup>1,2,\*</sup>

<sup>1</sup>Department of Biochemistry and Molecular Biology,

<sup>2</sup>Institute for Quantitative Health Science and Engineering,

Michigan State University, East Lansing MI 48824, USA

## **Abstracts**

**Background:** Machine learning methods and algorithms have been used routinely by scientists in many research areas. However, significant statistical and programming expertise are typically required to use these methods and algorithms. Thus, an easy-to-use graphical user interface software that enables biomedical researchers wanting to apply machine learning methods in their research in biology is essential, and can facilitate further use of such methodology.

**Results:** Here we present ClassificalO, an open-source Python graphical user interface for machine learning classification for the scikit-learn Python module. ClassificalO aims to provide an easy-to-use interactive way to train, validate, and test data on a range of state-of-the-art classification algorithms. The software enables fast comparisons within and across classifiers, and facilitates uploading and exporting of trained models, and both validated, and tested data results. ClassificalO is implemented as a Python package and is available for download and installation through the Python Package Index (PyPI) (<http://pypi.python.org/pypi/ClassificalO>) and it can be deployed using the “import” function once installed. The software is distributed under an MIT license and source code is available for download (for Mac OS X, Linux and Microsoft Windows) through PyPI and GitHub (<http://github.com/gmiaslab/ClassificalO>), and at <https://doi.org/10.5281/zenodo.1133266>.

**Conclusions:** ClassificalO facilitates the use of machine learning algorithms through a graphical user interface (GUI), and can help biomedical and other researchers with broad machine learning

background to use machine learning, and apply it to their research in a simple and interactive point-and-click way.

## SUPPLEMENTARY CONTENTS

**TABLE S1: Classification Algorithms ..... 3**

**SUPPLEMENTARY INFORMATION: ClassificalO Software Manual..... 4**

**ATTACHMENTS: SUPPLEMENTARY FILES (Files S1-7 available with manuscript).**

	<b><i>File Name</i></b>	<b><i>Description</i></b>
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 ClassificalO trained model using logistic regression
6.	<i>S6_IrisTrainValidationResult.csv</i>	Example ClassificalO testing result using logistic regression
7.	<i>S7_IrisTestingResult.csv</i>	Example ClassificalO validation result using logistic regression

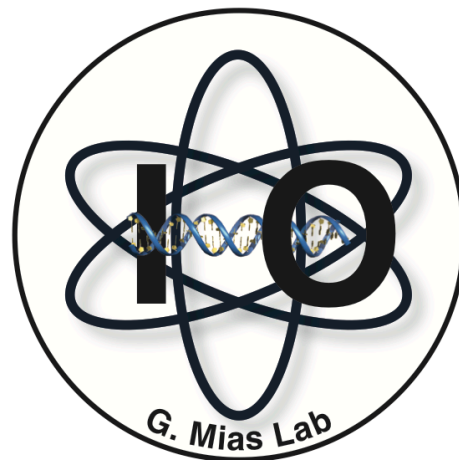
**All files below (S8-S14) are referenced in the manual, and are available online at GitHub (due to journal file size restrictions):**

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

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

<b>CLASSIFIER</b>	<b>Scikit-learn FUNCTION USED</b>	<b>IMMUTABLE PARAMETERS</b>
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 S1.** A list of all 25 classification algorithms, their corresponding scikit-learn functions, and immutable (unchangeable) parameters with their default values.



## **ClassificalO**

Machine Learning for Classification Graphical  
User Interface User Manual 1.1.2 (05/2018)

## **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  $\geq$  8.6.7, Pillow  $\geq$  5.1.0, pandas  $\geq$  0.22.0, numpy  $\geq$  1.14.3, scikit-learn  $\geq$  0.19.1, and scipy  $\geq$  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/
```

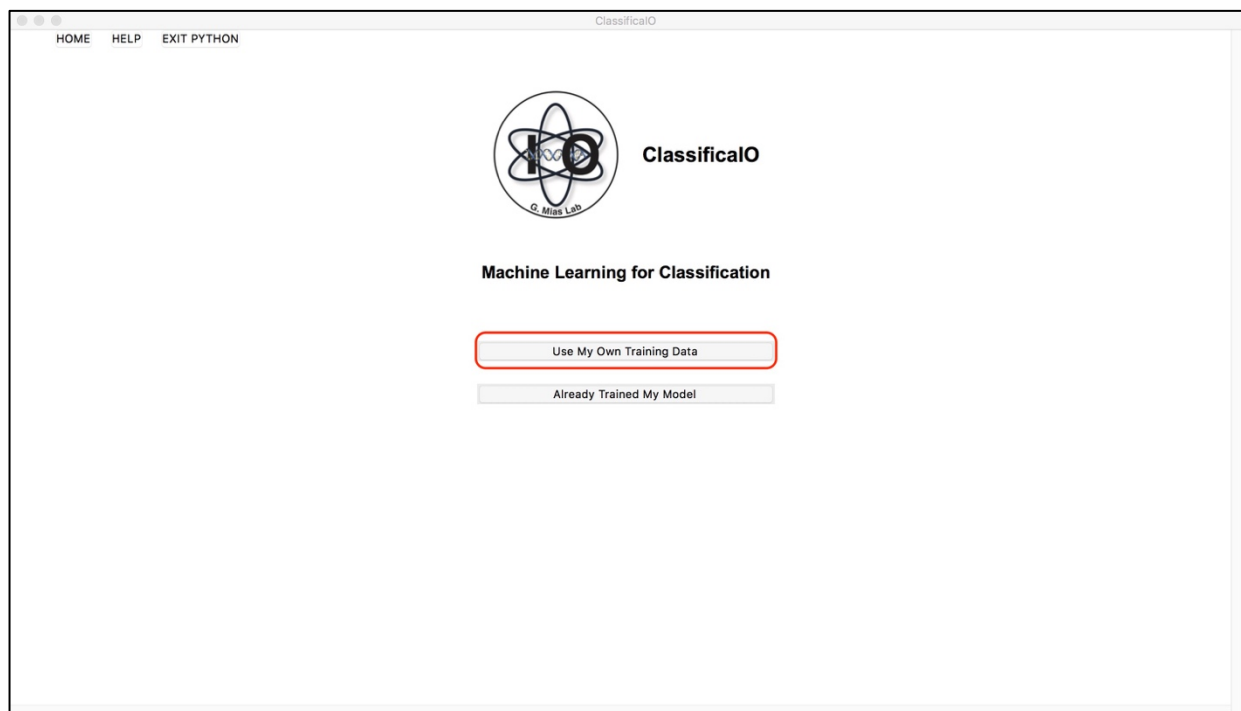
## Getting started:

### **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()
```

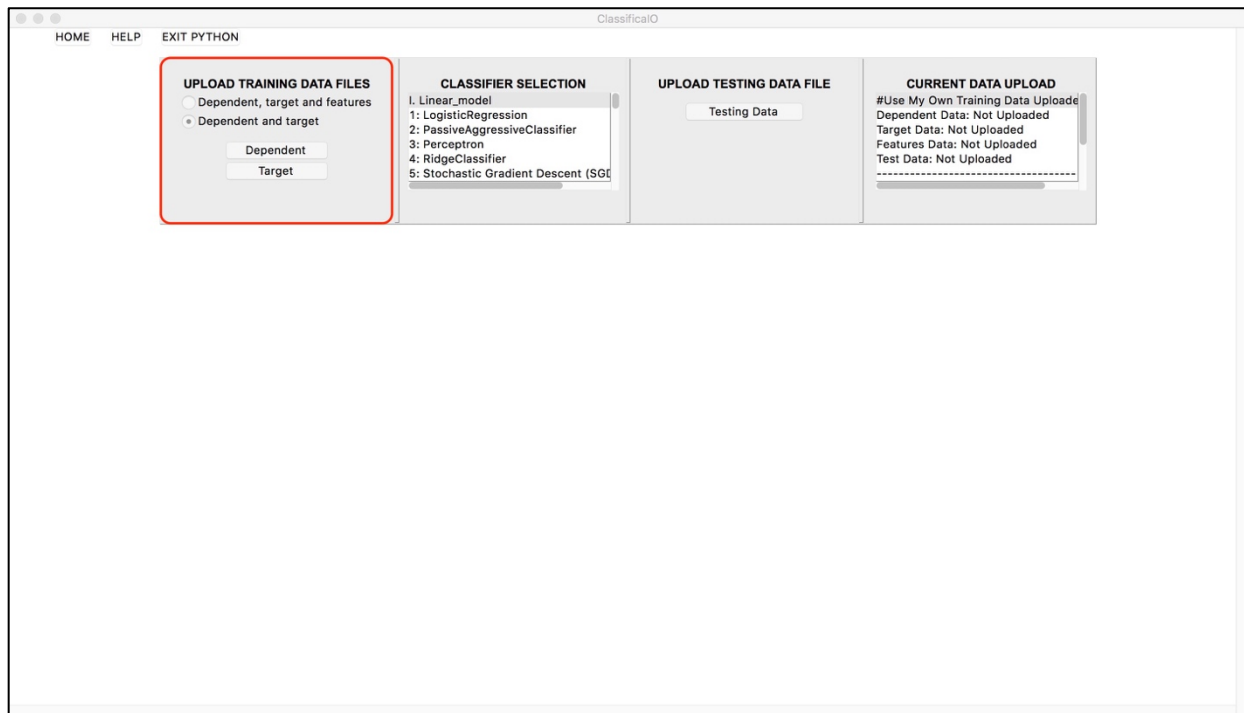


We note here the name is case sensitive (i.e. the 'IO' is capitalized). Once ClassificalO's main window appears on your screen, you can click on 'Use My Own Training Data' button and start your supervised machine learning classification project.

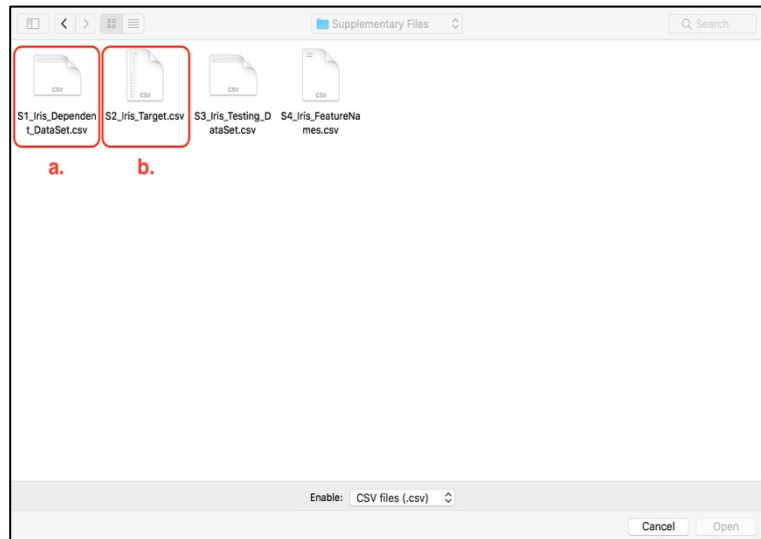
## Iris dataset prediction using a logistic regression classifier:

### 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.

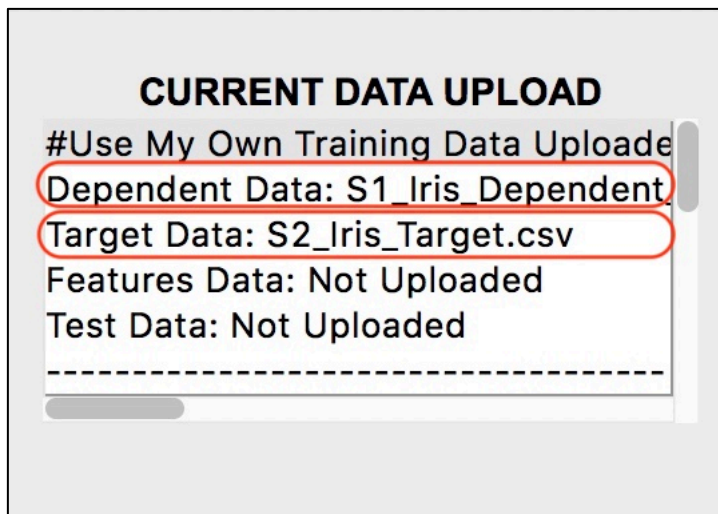


To begin uploading data files, click the corresponding buttons in the 'UPLOAD TRAINING DATA FILES' panel: a file selector directs you to upload both, dependent data file (**Supplementary Figure 1.a**) and target data file (**Supplementary Figure 1.b**). Once a file is uploaded to ClassificalO, the file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



**Supplementary Figure 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.

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







## 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 displays the ClassifaiO web application 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 five classifiers: 1: LogisticRegression, 2: PassiveAggressiveClassifier, 3: Perceptron, 4: RidgeClassifier, and 5: Stochastic Gradient Descent (SGD). Below this list, the '1: LogisticRegression' classifier is selected, and its parameters are displayed. These parameters include 'Train Sample Size (%)' (75), 'K-fold Cross-Validation' (10), 'random\_state' (Integer: 0), 'penalty' (l2), 'max\_iter' (100), 'verbose' (0), 'fit\_intercept' (True), 'dual' (True), 'warm\_start' (True), 'multi\_class' (ovr), 'tol' (1.0E-4), 'n\_jobs' (1), 'C' (1), 'solver' (liblinear), and 'intercept\_scaling' (1). A 'Submit' button is located below the parameters. At the bottom of the interface, there are three panels: 'CONFUSION MATRIX, MODEL ACCURACY & ERROR', 'TRAINING RESULT: ID — ACTUAL — PREDICTION', and 'TESTING RESULT: ID — PREDICTION'. Each panel has an 'Export' button (Export Model, Export Training, Export Testing).

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:

- **Supplementary Figure 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.
- **Supplementary Figure 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.
- **Supplementary Figure 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 ClassificaiO web interface. At the top, there are navigation links: HOME, HELP, and EXIT PYTHON. The main content area is divided into several sections:

- 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 (highlighted), 3: PassiveAggressiveClassifier, 4: Perceptron, 5: RidgeClassifier, and 6: Stochastic Gradient Descent (SGD). A red box highlights this section.
- UPLOAD TESTING DATA FILE:** Includes a "Testing Data" button.
- CURRENT DATA UPLOAD:** A section for uploading training and testing data files, with fields for "Dependent Data", "Target Data", "Features Data", and "Test Data".

Below the classifier selection, a red box labeled 'a.' highlights the definition of Logistic Regression: "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." A blue underlined link "Learn more." is visible next to the classifier name.

Below the definition, a red box labeled 'b.' highlights the training methods section. It includes a "Train Sample Size (%)" slider set to 75, and a "K-fold Cross-Validation" dropdown set to 10. A red box labeled 'c.' highlights the classifier parameters section, which includes various settings like "random\_state", "multi\_class", "intercept\_scaling", "verbose", "fit\_intercept", "dual", "warm\_start", "penalty", "max\_iter", "n\_jobs", "solver", "tol", and "C".

At the bottom, there are three sections for results: "CONFUSION MATRIX, MODEL ACCURACY & ERROR", "TRAINING RESULT: ID — ACTUAL — PREDICTION", and "TESTING RESULT: ID — PREDICTION". Each section has an "Export" button.

**Supplementary Figure 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 (**Supplementary Figure 5.a**). Model validation data results are displayed in the 'TRAINING RESULT: ID – ACTUAL – PREDICTION' panel (**Supplementary Figure 5.b**) with each data point ID is the 1<sup>st</sup> value, actual target value is displayed 2<sup>nd</sup>, and predicted target value 3<sup>rd</sup>, where the predictions correspond to the iris flower species, with 0=setosa, 1=versicolor, and 2=virginica.



## 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 (**Supplementary Figure 6.a**). Once clicked, a file selector directs you to upload the testing data file, see **Supplementary Figure 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 **Supplementary Figure 3.a**.

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

The screenshot displays the ClassificationIO web interface. Panel (a) shows the 'UPLOAD TESTING DATA FILE' section with a 'Testing Data' button. The 'CURRENT DATA UPLOAD' section on the right shows the uploaded file 'S3\_Iris\_Testing\_DataSet.csv'. Below this, the '1: LogisticRegression' model is selected, and its parameters are configured. Panel (b) shows the 'TESTING RESULT: ID – PREDICTION' section, which displays the results for 45 data points. The results are as follows:

ID	Prediction
148	2
61	1
115	2
126	2
133	2
130	2
58	1
90	1
52	1

The 'TESTING RESULT: ID – PREDICTION' section also shows the total objects tested: 45.

**Supplementary Figure 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 (**Supplementary Figure 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 displays 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 three columns: 'UPLOAD TRAINING MODEL FILE' with a 'Model File' button, 'UPLOAD TESTING DATA FILE' with a 'Testing Data' button, and 'CURRENT DATA UPLOAD' which contains a text area with pre-defined upload status messages and an '#Upload History' link. Below these columns is a 'Submit' button. At the bottom, there are two large empty boxes labeled 'CONFUSION MATRIX, MODEL ACCURACY & ERROR' and 'TESTING RESULT: ID — PREDICTION', with an 'Export Testing' button located at the bottom right of the interface.



## ClassificalO model input:

You will need to upload ClassificalO model by clicking the 'Model File' button in the 'UPLOAD TRAINING MODEL FILE' panel (**Supplementary Figure 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 (**Supplementary Figure 7.b**). Once a ClassificalO model and testing data files are uploaded, files names are automatically displayed in the 'CURRENT DATA UPLOAD' panel (**Supplementary Figure 7.c**).

After clicking 'submit', the uploaded model preset parameters will populate (**Supplementary Figure 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 (**Supplementary Figure 7.e**). Testing data results are displayed in the 'Testing RESULT: ID – PREDICTION' panel (**Supplementary Figure 7.f**) with the data point ID shown 1<sup>st</sup>, followed by a hyphen and the predicted value displayed right after it.

The screenshot displays the ClassificalO web interface with the following components:

- Navigation Bar:** HOME, HELP, EXIT PYTHON
- Panel a. UPLOAD TRAINING MODEL FILE:** Contains a 'Model File' button.
- Panel b. UPLOAD TESTING DATA FILE:** Contains a 'Testing Data' button.
- Panel c. CURRENT DATA UPLOAD:** Displays the upload status: '#Already Trained My Model Uploade', 'Model: S5\_LogisticRegression\_IrisTr', 'Test Data: S3\_Iris\_Testing\_DataSet.c', and '#Upload History'.
- Submit Button:** A central button to initiate the process.
- Panel d. CLASSIFIER PARAMETERS:** Shows the preset parameters for the classifier: ('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').
- Panel e. CONFUSION MATRIX, MODEL ACCURACY & ERROR:** Displays a confusion matrix and model performance 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 (NR): 3.7 %
- Panel f. TESTING RESULT: ID – PREDICTION:** Displays the results for 16 data points.

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

Export Testing

**Supplementary Figure 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.

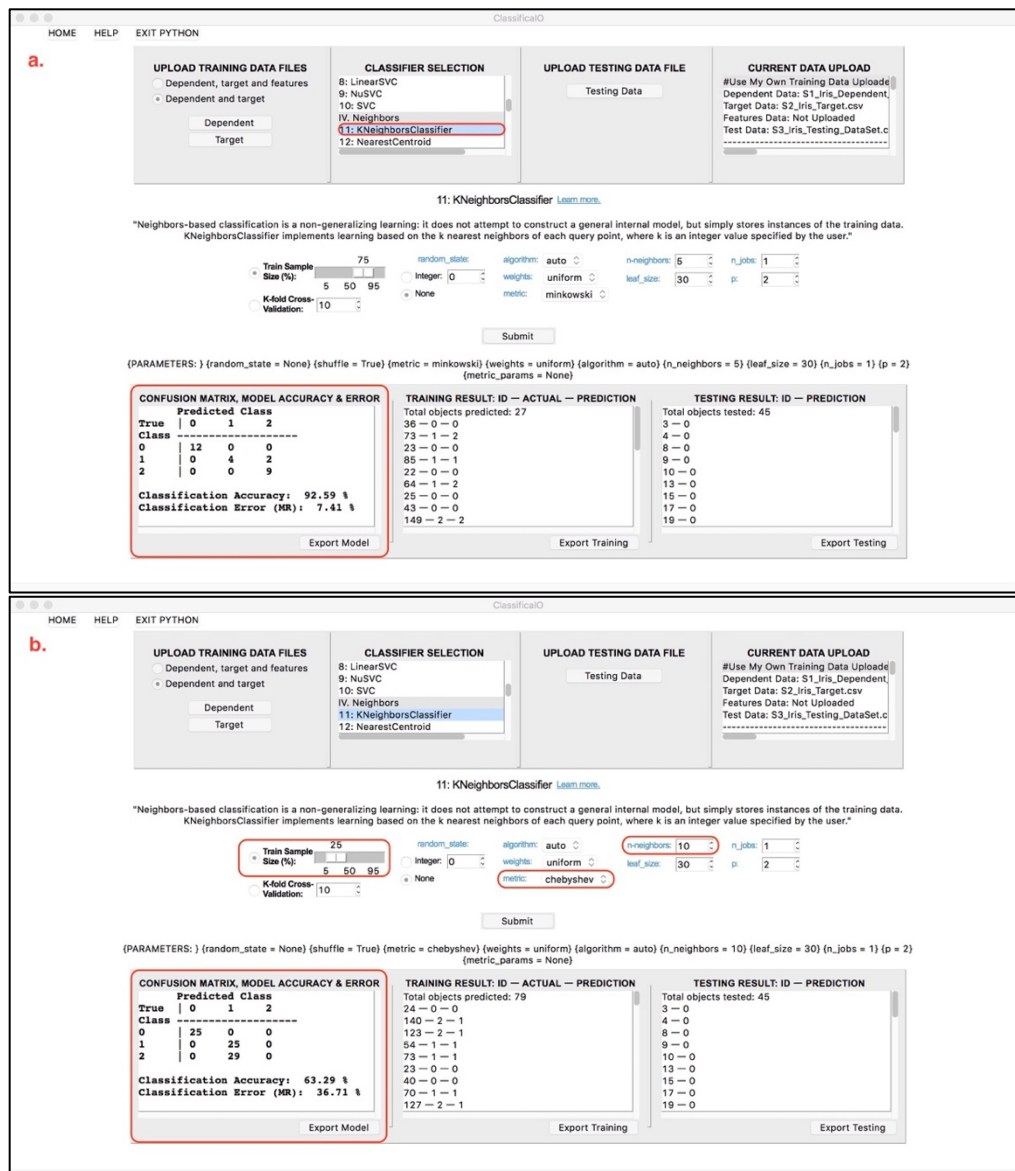
## **Results Export:**

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 **Additional 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.



**Supplementary Figure 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.

## **Ex. 2. 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' (**Supplementary Figure 9.a**) and select 'LinearSVC' as classifier (**Supplementary Figure 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 (**Supplementary Figure 9.c**) after selection has been made.

For the classifier, we use the default parameters, and 'K-fold Cross-Validation' set to 10 (**Supplementary Figure 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 (Supplementary Figure 9.e). Model validation data results are displayed in the 'TRAINING RESULT: ID – ACTUAL – PREDICTION' panel (Supplementary Figure 9.f) with each data point ID is the 1<sup>st</sup> value, actual target value is displayed 2<sup>nd</sup>, and predicted target value 3<sup>rd</sup> where Male = 0 and Female = 1. Testing results are displayed in the 'TESTING RESULT: ID – PREDICTION' panel (Supplementary Figure 9.g) with each data point ID shown 1<sup>st</sup>, and the corresponding predicted target value displayed after it 2<sup>nd</sup>, 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 formats, see the Supplementary files S13-15).



**Note on Features annotation:** 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 Supplementary file S12).



	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	216822_x_at

**Supplementary Figure 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.