CSCI 578 Machine Learning Competition 2022: Butterfly Classification

Instructor: Murat Dundar

Due date: 4/24/2022, 11:59pm

Please type in your assignment and submit it online through Canvas. No hard copies will be accepted.

Task 3: Multiple-instance, Open-set Butterfly Classification using DiNo Features

Dataset: All the data and splits are the same as in Task 2. The content of *whole.mat* and *parts.mat* files are described in Tables 1 and 2. 2D Part IDs are explained in Figure 1.

Goal: The goal of the third and final task is to predict species labels for testing samples in an open-set setting. For every sample id in the $test_sampleid$ variable in the whole.mat file you will still submit one prediction. However, these predictions may not have to be one of the species names included in the training set. There are dozens of unseen classes in the test set. Your classifier should also detect samples of these unseen classes. These samples should be labeled as unseen in the spreadsheet.

Approach: You are allowed to use any classifier algorithms covered in the class. SVM, Logistic Regression, and Student-t classifier are some possible alternatives. Implementing extensions of these models for multiple instance may help improve their prediction accuracy on this problem. You are allowed to use any normalization and dimensionality reduction techniques.

Submission: You are going to submit a CSV spreadsheet with predicted labels. The order of predictions should match the order of the samples in the *test_sampleid* variable in the *whole.mat* file. You can include up to three columns in this spreadsheet, one column for each of the different approaches you have tried. You will also send a one-page PDF document explaining each approach you used in detail. This document will not be graded but if not submitted may disqualify your submission. Your submissions should be reproducible. Your training and testing scripts may be requested if needed.

Evaluation: Your submission will be evaluated by the harmonic mean of mean class accuracy and detection F1 score. The detection F1 score will be computed by $\frac{2TP}{2TP+FP+FN}$, where TP denotes the number of unseen class test samples identified as unseen, FP denotes the number of seen class test samples identified as unseen, and FN

Variable name	Size	Description	Type
train_classid	7849x1	Species names as class labels	String
train_feats	7849x384	DiNo Features	Float
train_imgid	7849x1	Image IDs	Integer
train_sampleid	7849x1	Specimen IDs	Integer
val_classid	1379x1	Species names as class labels	String
val_feats	1379x384	DiNo Features	Float
val_imgid	1379x1	Image IDs	Integer
val_sampleid	1379x1	Specimen IDs	Integer
test_feats	2454x384	DiNo Features	Float
test_imgid	2454x1	Image IDs	Integer
test_sampleid	2454x1	Specimen IDs	Integer

Table 1: Variables in the whole.mat file.

Variable name	Size	Description	Type
train_classid	70641x1	Species names as class labels	String
train_feats	70641x384	DiNo Features	Float
train_imgid	70641x1	Image IDs	Integer
train_sampleid	70641x1	Specimen IDs	Integer
train_tileid	70641x2	2D Tile IDs	Integer
val_classid	12411x1	Species names as class labels	String
val_feats	12411x384	DiNo Features	Float
val_imgid	12411x1	Image IDs	Integer
val_sampleid	12411x1	Specimen IDs	Integer
val_tileid	12411x2	2D IDs for each image part	integer
test_feats	22086x384	DiNo Features	Float
test_imgid	22086x1	Image IDs	Integer
test_sampleid	22086x1	Specimen IDs	Integer
test_tileid	22086x2	2D IDs for each image part	Integer

Table 2: Variables in the parts.mat file.

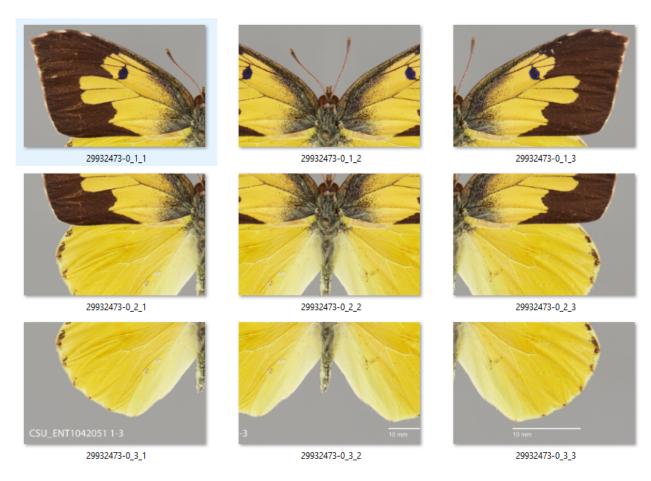


Figure 1: Image is split into nine parts by a 3×3 grid with 50% overlap between parts. Part IDs are the last two numbers in the image file name in the figure.

denotes the number of unseen class test samples classified into one of the seen classes. Top performer will receive the top score of 150 points. The rest of the submissions will be scored relative to the top performer.