# Machine Learning Project Report
# Praveen Mareedu
# N16407897

## What do we have ?

We have an input file '**cs6923Project.mat**' which internally consists of train_label.csv , train.csv , test.csv files .

**train.csv**– This file consists of 50000 (rows) x 77 (columns) matrix of data inside (X), which can be used to train the algorithm , In the below approach , I have used 90% of the data to preprocess and 10% of the data to cross validate the results . Although this ration can be varied , 90%:10% gave the best accuracy scores on the test data.

**train_label.csv** : This file consist of 50000 (rows) x 1 (column) matrix of data inside (Y), which gives the label information for each row in train file .

**test.csv** : This file consists of 50000 [rows] x 77 [columns] same as train file but for which we need to predict the label information .

**test_label.csv :** File generated with predicted label information for the data in test.csv file .

## Algorithms to consider :

This is a supervised learning problem . We need to use classification Machine Learning algorithms to predict output labels. Some of the best classification algorithms that can be used are **Logistic Regression , SVM (Support Vector Machines), Discriminant Analysis , K-Nearest neighbors etc**.

## Prerequisites :

This program is developed in Python using popular open source machine learning library 'sklearn' .
-Need **SKLEARN** packages installed on the machine .
-Need Python – **Scipy , Numpy** installed on the machine.
Please check out the end of the document to find the installation links.

# Algorithms used to test :

In this project I tried to use two famous techniques Logistic Regression , SVM to train the classification algorithm and predict labels for the test data . Logistic regression being one of the easy approach and SVM because it produces high accuracy scores in the test runs. Algorithm speeds and memory consumptions are referred in the URL mentioned in the end of the document.

# Parameters that can be tweaked and Results :

## SVM :

**Kernel :** rbf (Gaussian) , linear , poly , sigmoid.

**cache_size:** To specify the kernel cache size.

**C :** Penalty Parameter  for the error term.

**gamma :** Kernel coefficient for 'rbf', 'poly' and 'sigmoid'

## Logistic:

**C :** Regularization term

# Results :

# Cache Size 1000 and kernel RBF
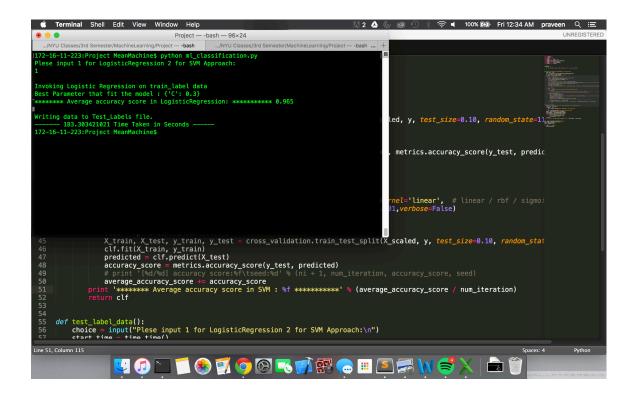
## Cache Size 1000 and Kernel Function Sigmoid:



## Cache Size 200 and Linear Kernel :

## Logistic Regression :



We can notice the time taken to execute the program varies in the each run with cache size . As the cache size increases the time taken to computer the program decreases and also with kernel RBF the accuracy levels are high with multiple runs .

## References :

**Classification Algorithms and their speed , memory consumption :**

http://www.mathworks.com/help/stats/supervised-learning-machine-learning-workflow-and-algorithms.html?refresh=true

**Scipy , Numpy :**

http://www.scipy.org/install.html

**Scikit Learn Kit :**

http://scikit-learn.org/stable/install.html