

# Implementation of Adaboost-M1 Classifier Based on Linear SVM

Sami Alperen Akgun  
sami.alperen.akgun@gmail.com

## A. Visualization of Data

The visualization of given dataset containing class A and class B data can be seen in Figure 1.

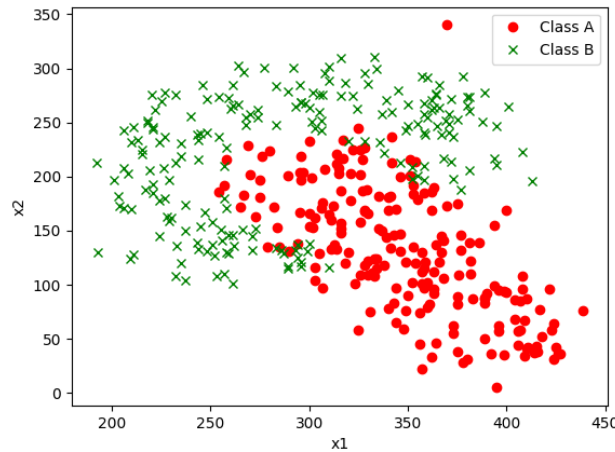


Fig. 1. Given data for Question 2

## B. Training SVMs and Check Accuracy via 10-times-10-fold Cross Validation

Linear SVM classifiers were trained for different penalization factors ( $c$  values). Decision boundaries were plotted for these values as it can be seen in Figure 2. Decision boundaries with margins can be seen in Figure 3. 10-times-10-fold cross validation were also applied on these linear SVM classifiers with different  $c$  values. Obtained results can be seen below:

```
1 C value: 0.1
2 Mean of the accuracy: 0.7975152439024389 --> 79.75%
3 Var of the accuracy: 1.1149851435049755e-06
4 Best accuracy: 0.90625 --> 90.62%
5
6 C value: 1
7 Mean of the accuracy: 0.7980564024390244
8 Var of the accuracy: 4.437610043916163e-06
9 Best accuracy: 0.9375
10
11 C value: 10
12 Mean of the accuracy: 0.7964634146341464
13 Var of the accuracy: 2.014219261754062e-06
14 Best accuracy: 0.926829268292683
15
16 C value: 100
17 Mean of the accuracy: 0.7988948170731708
18 Var of the accuracy: 1.4532859742308036e-06
19 Best accuracy: 0.926829268292683
```

## C. Select SVM with best Penalization Factor ( $c$ ) and Implement Adaboost-M1 based on that

After analyzing obtained results in Part 2, one can clearly see that mean accuracy corresponding to different  $c$  values are almost the same for all cases. However,  $c=1$  case has the best accuracy value (not mean accuracy) among them. Therefore, a linear SVM classifier with  $c$  value 1 was selected as a weak classifier for Adaboost approach. Nonetheless, selecting  $c$  value as 0.1 might be another good solution since actually we are interested in mean accuracies after 10-times-10-fold cross validation and they are almost the same for all cases. Higher values of  $c$  were not selected since they don't offer any better

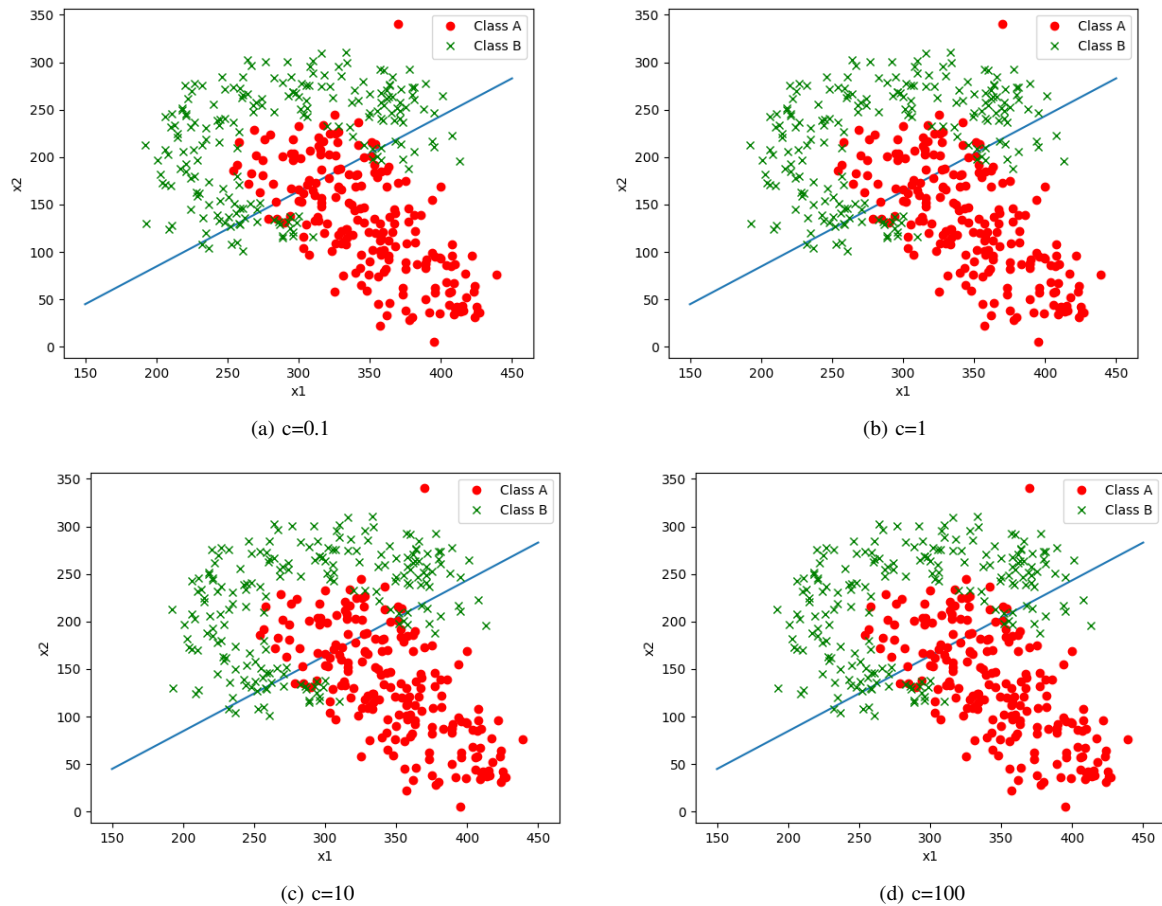
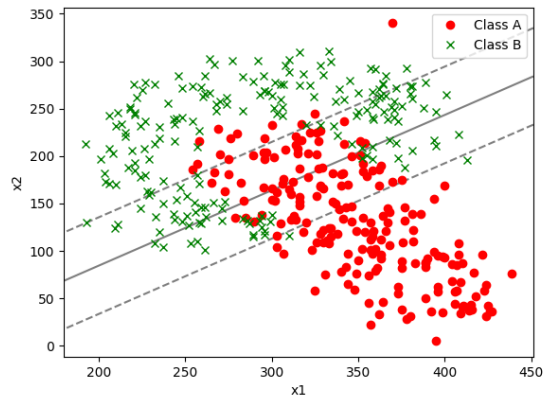


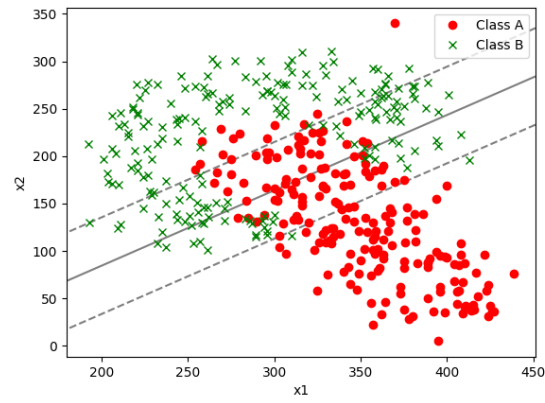
Fig. 2. Different decision boundaries of linear SVM classifiers

classification performance, but their convergence time is way larger than lower values, i.e. they are computationally expensive. The implemented function for Adaboost-M1 approach can be seen below:

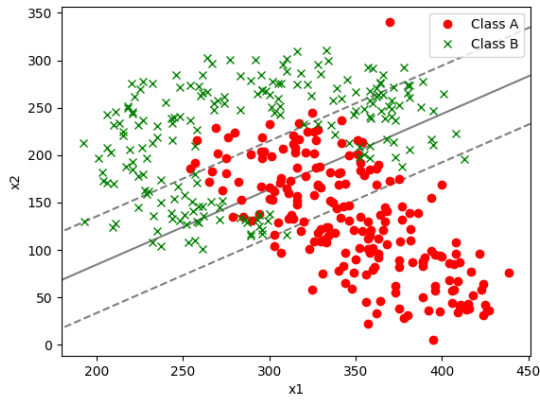
```
1 def adaboost_M1(X_data,Y_data,c_value,T):
2     """
3     This function applies 10 fold Cross validation on the given input data
4     Inputs: X_data --> input data contains features and their values
5            Y_data --> input data contains output labels
6            c_value --> Penalization factor for linear SVM classifier
7            T --> Max number of weak learners, i.e. max number of iterations
8            for adaboost mainloop
9     Output: h --> fitted linear SVM model
10    """
11
12    m = X_data.shape[0] # Total number of samples
13    n = X_data.shape[1] # Total number of features
14    h_list = [] # The list to hold all hypothesis, i.e. learned models
15    D = [] # List to hold all weights for each hypothesis
16    beta_list = [] # List to hold all beta values
17    # Initialize weight list with 1/m for all samples
18    D.append( np.ones(m) * (1/m) )# --> shape = (m,)
19
20    t = 0
21    while t<T:
22        # Pick 100 examples as a training set out of m data with probability
23        # distribution D[t]
24
25        training_indices = np.random.choice(D[t].shape[0],100,replace=False,p=D[t])
26        training_X = X_data[training_indices,:]
27        training_Y = Y_data[training_indices]
28        h = svm.SVC(kernel='linear', C = c_value) # current hypothesis
```



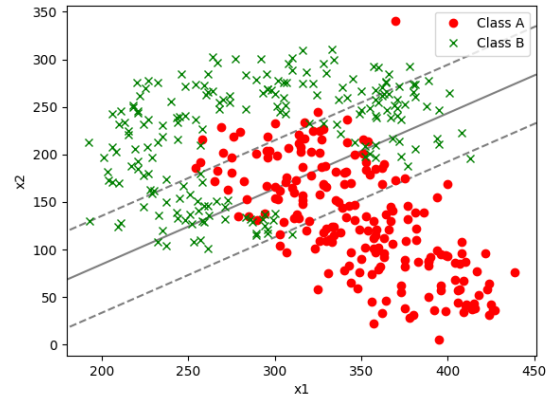
(a)  $c=0.1$



(b)  $c=1$



(c)  $c=10$



(d)  $c=100$

Fig. 3. Different decision boundaries of linear SVM classifiers with margins

```

29 h.fit(training_X,training_Y) # train using training set
30 Y_predicted = h.predict(X_data) # predict using all samples
31
32 epsilon = h_error(Y_predicted,Y_data,D[t])
33
34 # if error is more than 50%, start the loop again
35 # This means we are discarding the classifier corresponding error
36 # rate more than 50% and we are trying to resample another training
37 # set --> Therefore, in the end we will have 50 trained classifiers
38 if epsilon >= 0.5:
39     # Not update t, so that we can iterate with the same t value
40     continue
41
42 beta = epsilon / (1-epsilon)
43 D.append( update_weights(Y_predicted,Y_data,beta,D[t]) ) #--> D[t+1]
44 h_list.append(h)
45 beta_list.append(beta)
46 t+=1 # While loop!
47
48 return h_list, beta_list, D

```

#### D. Check Accuracy of Adaboost-M1 classifier via 10-times-10-fold Cross Validation

Obtained mean and variance of accuracy for 10-times-10-fold cross validation approach on adaboost-M1 classifier (with weak classifier linear SVM and  $c=1$ ) can be seen below:

```

1 Mean of the accuracy: 0.8372256097560975 --> 83.72%
2 Var of the accuracy: 7.904603999270937e-07
3 Best accuracy: 0.915609756097561 --> 91.56%

```

#### E. Plot Decision Boundary of Adaboost M1 Classifier and Compare Accuracy of it with Linear SVM

Decision boundary of the ensemble model can be seen in Figure 4. If one compare this decision boundary with the decision boundary of weak classifier (linear SVM with  $c$  value 1) in Figure 2, s(he) can clearly see that thanks to boosting approach, the trained classifier were improved.

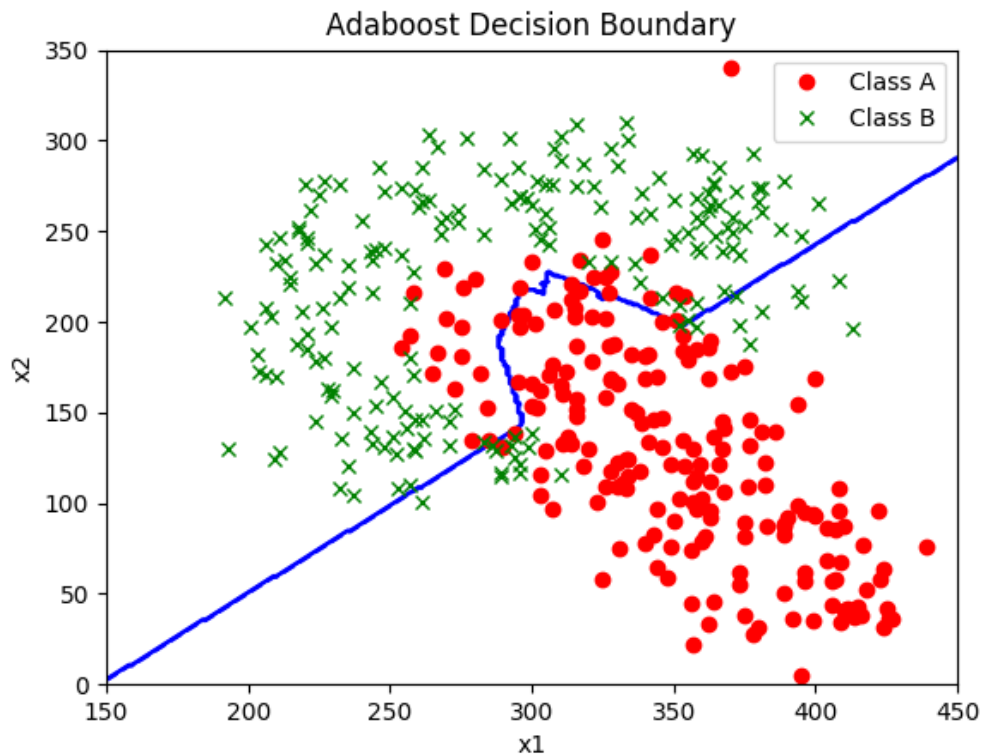


Fig. 4. Decision Boundary of the Ensemble Model