

ECE 5424/CS 5824 - Advanced Machine Learning - Problem Set 2

Question 1.1:

Question 1 :

$$\min_{\lambda, f} \|\lambda\|^2 + C \sum_{i=1}^n f_i$$

$$\text{s.t. } \left. \begin{aligned} y_i (\lambda x_i + \lambda_0) &\geq 1 - f_i \\ f_i &\geq 0 \end{aligned} \right\} \text{ for } i=1, \dots, n$$

checking the constraints in above optimization problem.

$$f_i \geq 1 - y_i (\lambda x_i + \lambda_0) \rightarrow \text{lower bound on } f_i$$

$$f_i \geq 0$$

hence  $f_i = (1 - y_i (\lambda x_i + \lambda_0))_+ = \max(0, 1 - y_i (\lambda x_i + \lambda_0))$

$$f_i = \begin{cases} 1 - y_i (\lambda x_i + \lambda_0) & \text{if } 1 - y_i (\lambda x_i + \lambda_0) > 0 \\ 0 & \text{otherwise} \end{cases}$$

say  $t = y_i (\lambda x_i + \lambda_0)$ ,  $\ell_{\text{hinge}}(t)$  is

$$\ell_{\text{hinge}}(y_i (\lambda x_i + \lambda_0)) = \begin{cases} 1 - y_i (\lambda x_i + \lambda_0) & \text{if } y_i (\lambda x_i + \lambda_0) < 1 \\ 0 & \text{otherwise} \end{cases}$$

Hence,  $f_i = \ell_{\text{hinge}}(y_i (\lambda x_i + \lambda_0))$  can be written,

and it automatically includes the constraints

$$y_i (\lambda x_i + \lambda_0) \geq 1 - f_i \quad \& \quad f_i \geq 0 \quad \text{for } i=1, \dots, n$$

in the primal problem. We can directly plug in this term to objective function and drop the constraints, here we end up with;

$$\min_{\lambda} \|\lambda\|^2 + C \sum_{i=1}^n \ell_{\text{hinge}}(y_i (\lambda x_i + \lambda_0))$$

Question 1.2:

Dual formulation for soft-margin SVM :

$$\min_{\lambda, \lambda_0} \frac{1}{2} \|\lambda\|_2^2 + C \sum_{i=1}^n f_i$$

$$\text{s.t.} \quad \left. \begin{aligned} 1 - y_i(\lambda^T x_i + \lambda_0) - f_i &< 0 \\ -f_i &< 0 \end{aligned} \right\} \text{ for } i=1, \dots, n$$

we know  $L(\alpha, x) = f(x) + \sum_{i=1}^n \alpha_i g_i(x)$ , hence

$$L(\alpha, \beta, [\lambda, \lambda_0]) = \frac{1}{2} \|\lambda\|_2^2 + C \sum_{i=1}^n f_i + \sum_{i=1}^n \alpha_i (1 - y_i(\lambda^T x_i + \lambda_0) - f_i) + \sum_{i=1}^n \beta_i (-f_i), \quad \alpha_i, \beta_i \geq 0$$

↓  
additional  
dual  
variable

$$\Theta_D(\alpha, \beta) = \min_{\lambda, \lambda_0} L([\lambda, \lambda_0], \alpha, \beta)$$

to find min of Lagrangian we can set the gradient to 0.

$$\frac{\nabla L}{\nabla \lambda} = \lambda + \sum_{i=1}^n -\alpha_i y_i x_i = 0 \Rightarrow \lambda = \sum_{i=1}^n \alpha_i y_i x_i$$

$$\frac{\partial L}{\partial \lambda_0} = \sum_{i=1}^n \alpha_i (-y_i) = 0 \Rightarrow \sum_{i=1}^n \alpha_i y_i = 0$$

what if we evaluate  $\alpha$  such that  $\sum_{i=1}^n \alpha_i y_i \neq 0$

if  $\sum_{i=1}^n \alpha_i y_i > 0$ , we can take  $\lambda_0 \rightarrow \infty$  hence  $\Theta_D \rightarrow -\infty$

since the goal of  $\alpha$  is to maximize dual objective we will make sure the condition  $\sum_{i=1}^n \alpha_i y_i = 0$  is met.

↓  
 $\lambda, \lambda_0$  optimal values  
 defining  $\Theta_D$ .

$$\text{for } \lambda^* = \sum_{i=1}^n \alpha_i y_i x_i, \quad \lambda_0^* = 0$$

$$\Theta_D(\alpha) = \mathcal{L}([\lambda^*, \lambda_0^*], \alpha, \beta)$$

$$= \frac{1}{2} \|\lambda^*\|_2^2 + C \sum_{i=1}^n \beta_i + \sum_{i=1}^n \alpha_i (1 - y_i (\lambda^{*T} x_i) - \beta_i) - \sum_{i=1}^n \beta_i \beta_i$$

$$= \frac{1}{2} \|\lambda^*\|_2^2 + \lambda^{*T} \sum_{i=1}^n (-\alpha_i y_i x_i) + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \beta_i - \sum_{i=1}^n \beta_i \beta_i$$

$$= -\frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \alpha_i \alpha_k y_i y_k x_i^T x_k + \sum_{i=1}^n \alpha_i + \sum_{i=1}^n \alpha_i - \beta_i - \sum_{i=1}^n \beta_i \beta_i$$

$$\max_{\alpha} \Theta_D(\alpha) \quad \text{s.t.} \quad \alpha \geq 0, \quad \sum_{i=1}^n \alpha_i y_i = 0 \quad \sum_{i=1}^n (-\alpha_i - \beta_i) \beta_i$$

$$\beta \geq 0$$

$$\& \quad 0 \leq \alpha_i \leq C \\ 0 \leq \beta_i \leq C$$

↓  
since the problem is maximization  
&  $-\beta_i \beta_i$  will be negative  
it is best to choose  $\beta_i = 0$ .

here we can write:

$$\max_{\alpha} \Theta_D(\alpha) = \max_{\alpha} \left[ -\frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \alpha_i \alpha_k y_i y_k x_i^T x_k + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i \beta_i \right]$$

$$\text{where} \quad 0 \leq \alpha_i \leq C \quad \& \quad \sum_{i=1}^n \alpha_i y_i = 0$$



## Question 2: Programming Assignment

### Task 1:

```
## STUDENT: YOUR CODE STARTS HERE
# Task: Append '1' to the beginning of each vector.
# Hint: You can use data_features.toarray() to transform data_features into a numpy array
# The output should be a numpy array named data_mat
data_features = data_features.toarray()
n = data_features.shape[0] #n=3000
m = data_features.shape[1] #m=4500
data_mat = np.zeros((n, m+1))

for i in range(n):
    # copy the elements of the existing array into the data_mat's row, starting at index 1
    data_mat[i][1:] = data_features[i]

    # insert the value 1 at the beginning of the data_mat's first index
    data_mat[i][0] = 1

## STUDENT: CODE ENDS
print ('The updated size: ',data_mat.shape)
```

The updated size: (3000, 4501)

The size of data\_mat is (3000,4501)

Task 2:

Derivation of gradient of loss function with respect to weight array:

$$\tilde{\Theta}_t = \underset{\tilde{\Theta}}{\operatorname{argmin}} L(\tilde{\Theta}) = \underset{\tilde{\Theta}}{\operatorname{argmin}} \sum_{i=1}^n \ln(1 + e^{-y_i \tilde{\Theta}_t^T \tilde{x}_i})$$

↓ we want to find  $\tilde{\Theta}$  such that this summation is minimized.

$$\frac{\nabla L}{\nabla \tilde{\Theta}} = \frac{\nabla}{\nabla \tilde{\Theta}} \left( \sum_{i=1}^n \ln(1 + e^{-y_i \tilde{\Theta}_t^T \tilde{x}_i}) \right)$$

$$\frac{\nabla L}{\nabla \tilde{\Theta}} = \sum_{i=1}^n \left( \frac{\nabla}{\nabla \tilde{\Theta}} (\ln(1 + e^{-y_i \tilde{\Theta}_t^T \tilde{x}_i})) \right)$$

$$\text{first, } \frac{\partial L}{\partial \theta_0} = \sum_{i=1}^n \frac{\nabla}{\nabla \theta_0} (\ln(1 + e^{-y_i \tilde{\Theta}_t^T \tilde{x}_i}))$$

$$\text{say } u = 1 + e^{-y_i \tilde{\Theta}_t^T \tilde{x}_i}$$

$$\frac{\partial (\ln(1 + e^{-y_i \tilde{\Theta}_t^T \tilde{x}_i}))}{\partial \theta_0} = \frac{\partial \ln(u)}{\partial u} \cdot \frac{\partial u}{\partial \theta_0} \quad \text{by chain rule}$$

$$\text{we know } \frac{\partial \ln(u)}{\partial u} = \frac{1}{u},$$

$$\frac{\partial u}{\partial \theta_0} = \frac{\partial (1 + e^{-y_i \tilde{\Theta}_t^T \tilde{x}_i})}{\partial \theta_0} = \frac{\partial (e^{-y_i \tilde{\Theta}_t^T \tilde{x}_i})}{\partial \theta_0}, \quad \text{say } -y_i \tilde{\Theta}_t^T \tilde{x}_i = z$$

$$\frac{\partial (e^{-y_i \tilde{\Theta}_t^T \tilde{x}_i})}{\partial \theta_0} = \frac{\partial e^{(z)}}{\partial z} \cdot \frac{\partial z}{\partial \theta_0}, \quad \text{where } \frac{\partial e^{(z)}}{\partial z} = e^z$$

$$\frac{\partial z}{\partial \theta_0} = \frac{\partial (-y_i \tilde{\theta}_t^T \tilde{x}_i)}{\partial \theta_0}$$

since  $-y_i \tilde{\theta}_t^T \tilde{x}_i$  has the form of:  $y_i (\theta_0 \cdot 1 + \theta_1 x_1 + \dots + \theta_d x_d)$

where  $\tilde{x} = [1 \ x_1 \ x_2 \ \dots \ x_d]^T$ ,  $\tilde{\theta} = [\theta_0 \ \theta_1 \ \dots \ \theta_d]^T$

$$\frac{\partial (-y_i \tilde{\theta}_t^T \tilde{x}_i)}{\partial \theta_0} = -y_i, \text{ finally substituting all findings:}$$

$$\frac{\partial L}{\partial \theta_0} = \sum_{i=1}^n \left( \frac{\partial \ln(u)}{\partial u} \cdot \frac{\partial u}{\partial \theta_0} \right) \quad \frac{\partial u}{\partial \theta_0} = \frac{\partial (1 + e^{-y_i \tilde{\theta}_t^T \tilde{x}_i})}{\partial \theta_0}$$

$$= \frac{1}{1 + e^{-y_i \tilde{\theta}_t^T \tilde{x}_i}} \cdot \frac{\partial u}{\partial \theta_0}, \quad \frac{\partial u}{\partial \theta_0} = \frac{\partial (1 + e^{-y_i \tilde{\theta}_t^T \tilde{x}_i})}{\partial \theta_0} = \frac{\partial (e^{-y_i \tilde{\theta}_t^T \tilde{x}_i})}{\partial \theta_0}$$

$$= \frac{1}{1 + e^{-y_i \tilde{\theta}_t^T \tilde{x}_i}} \cdot \frac{\partial (e^{-y_i \tilde{\theta}_t^T \tilde{x}_i})}{\partial \theta_0}, \quad \frac{\partial (e^{-y_i \tilde{\theta}_t^T \tilde{x}_i})}{\partial \theta_0} = e^z \cdot \frac{\partial z}{\partial \theta_0}$$

where  $z = -y_i \tilde{\theta}_t^T \tilde{x}_i$

$$= \frac{1}{1 + e^{-y_i \tilde{\theta}_t^T \tilde{x}_i}} \cdot e^{(-y_i \tilde{\theta}_t^T \tilde{x}_i)} \cdot \frac{\partial z}{\partial \theta_0}, \quad \text{where } \frac{\partial z}{\partial \theta_0} = -y_i$$

hence

$$\frac{\partial L}{\partial \theta_0} = \sum_{i=1}^n \left( \frac{-y_i \cdot e^{(-y_i \tilde{\theta}_t^T \tilde{x}_i)}}{1 + e^{-y_i \tilde{\theta}_t^T \tilde{x}_i}} \right), \quad \text{for other } \theta_t \text{'s in vector } \theta = [\theta_0 \ \theta_1 \ \dots \ \theta_d]$$

only  $\frac{\partial z}{\partial \theta_t}$  term will change and be equal to  $-y_i \cdot x_t^{(i)}$

i.e.,

$$\frac{\partial L}{\partial \theta_t} = \sum_{i=1}^n \left( \frac{-y_i \cdot x_t^{(i)} \cdot e^{(-y_i \tilde{\theta}_t^T \tilde{x}_i)}}{1 + e^{-y_i \tilde{\theta}_t^T \tilde{x}_i}} \right)$$

(2)

Therefore for more compact expression with vectors we can write

$$\frac{\nabla L}{\nabla \tilde{\theta}_t} = \sum_{i=1}^n \left( \frac{-y_i \tilde{x}_i e^{-y_i \tilde{\theta}_t^T \tilde{x}_i}}{1 + e^{-y_i \tilde{\theta}_t^T \tilde{x}_i}} \right), \text{ where } x^{(i)} = \begin{bmatrix} 1 & x_1 & x_2 & \dots & x_d \end{bmatrix}$$

$\downarrow$   
 $(x_0)$

$$\left[ \frac{\nabla L}{\nabla \theta_0}, \frac{\nabla L}{\nabla \theta_1}, \dots, \frac{\nabla L}{\nabla \theta_d} \right]$$

as the final form.

code:

```
# STUDENT: PRINT THE OUTPUT AND COPY IT TO THE SOLUTION FILE
my_weights = np.ones(data_mat.shape[1]) # a weight of all 1s
derivative = weight_derivative(my_weights, train_data, train_labels)
#len(my_weights)

#print(derivative[0][4500])
print(derivative[:10])
```

```
[ 1.23415752e+03 -4.13993755e-08  1.00000000e+00  9.99993856e-01
  1.99987630e+00  9.99859072e-01  9.52574127e-01  3.59772256e+01
  2.99996572e+00 -1.38879439e-11]
```

First 10 elements of the derivative array



Task 3:

initial weights = [1 ... 1], all 1 vector

step size = 0.01

tolerance = 5

```
Here are the final weights after convergence:  
[-1.57105795  1.06197822  0.44064367 ... -2.00719313  1.  
 0.52556552]
```

Task 4:

```
## STUDENT: CODE STARTS HERE  
## Pull out the parameters (theta_0, theta) of the logistic regression model  
theta0 = final_weights[0]  
#theta = [ $\theta_1$   $\theta_2$  ...  $\theta_d$ ]  
#theta[1] =  $\theta_2$ , theta[2] =  $\theta_3$ ,  
theta = final_weights[1:4501]  
## STUDENT: CODE ENDS HERE  
print ('y intercept: ',theta0)  
print ('theta1 and theta2: ',theta[1],theta[2])  
  
y intercept: -1.5710579458497245  
theta1 and theta2: 0.4406436676371895 0.3361974703293962
```

In this question y intercept is -1.57 is theta0, and theta array is in form = [ $\theta_1$   $\theta_2$   $\theta_3$  ...  $\theta_{4500}$ ], so theta array consists of the weights, not bias. Hence when the code returns theta[1] it returns  $\theta_2$  and for theta[2] it returns  $\theta_3$ . If question wanted to it to print  $\theta_1, \theta_2$  then the print statement should be updated as print(theta[0],theta[1]) but since the template is given such I haven't changed it.



#### Task 5:

The prediction is done such that given data point if it's probability of having label +1 is bigger than 0.5, we predict that points label as +1 otherwise we say it is -1.

```
# STUDENT: copy the output of this section to the solution file

## Get predictions on training and test data
preds_train = model_predict(train_data, final_weights)
preds_test = model_predict(test_data, final_weights)

## Compute errors
errs_train = np.sum((preds_train > 0.0) != (train_labels > 0.0))
errs_test = np.sum((preds_test > 0.0) != (test_labels > 0.0))

print ("Training error: ", float(errs_train)/len(train_labels))
print ("Test error: ", float(errs_test)/len(test_labels))

Training error:  0.0076
Test error:  0.18
```

#### Task 6:

code:

```
def margin_counts(feature_matrix, weights, gamma):
    ## Return number of points for which  $\Pr(y=1)$  lies in  $[0, 0.5 - \gamma)$  or  $(0.5 + \gamma, 1]$ 

    # Input:
    # feature_matrix: numpy array of size n by d+1, where n is the number of data points, and d+1 is the
    #               note we have included the dummy feature as the first column of the feature_matrix
    # weights: weight vector to start with, a numpy vector of dimension d+1
    # gamma: the margin value
    # Output:
    # number of points for which  $\Pr(y=1)$  lies in  $[0, 0.5 - \gamma)$  or  $(0.5 + \gamma, 1]$ 

    ## STUDENT: YOUR CODE HERE

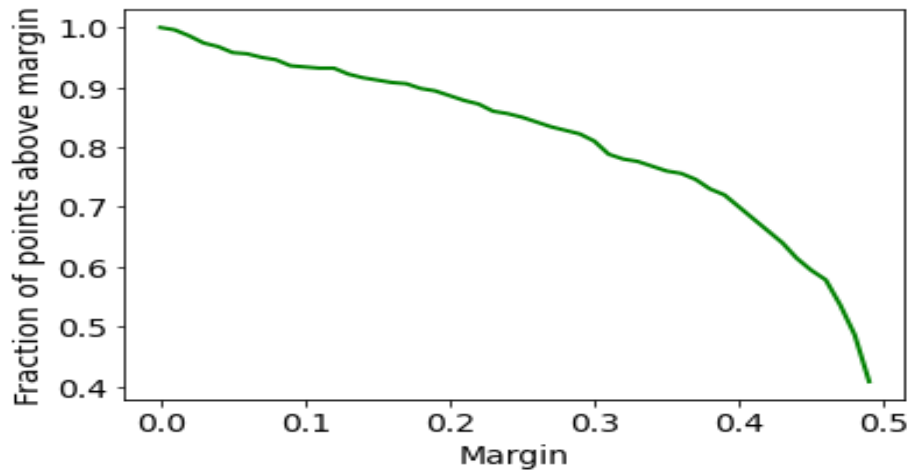
    #y_pred = np.zeros(feature_matrix.shape[0])
    j = 0
    for i in range(feature_matrix.shape[0]):
        z = np.dot(weights, feature_matrix[i])
        prob_1 = 1/(1+np.exp(-z))
        if prob_1 >= 0.5 + gamma :
            j = j + 1
        if prob_1 < 0.5 - gamma:
            j = j + 1

    return j

## STUDENT: CODE ENDS
```

given the data points  $x_i$ , total number of points that lie in interval  $P(Y=1|X=x_i) \geq 0.5 + \gamma$  and  $P(Y=1|X=x_i) < 0.5 - \gamma$  are returned.

plot:



Task 7:

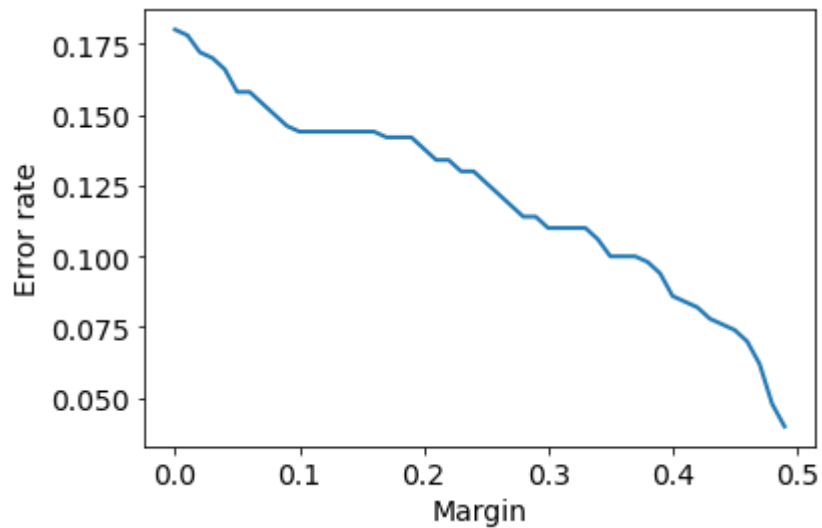
Copy the code and the output plot (i.e., visualization of the relationship between margin and error rate) to the solution file. What do you observe from the plot?

code:

```
def margin_errors(feature_matrix, labels, weights, gamma):  
    ## Return error of predictions that lie in intervals  $[0, 0.5 - \gamma)$  and  $(0.5 + \gamma, 1]$   
  
    ## STUDENT: YOUR CODE HERE  
    y_pred = np.zeros(feature_matrix.shape[0])  
    for i in range(feature_matrix.shape[0]):  
        z = np.dot(weights, feature_matrix[i])  
        prob_1 = 1/(1+np.exp(-z))  
        if prob_1 >= 0.5:  
            y_pred[i] = 1  
        else:  
            y_pred[i] = -1  
  
    #return y_pred  
  
    error_counter = 0  
  
    for i in range(feature_matrix.shape[0]):  
        z = np.dot(weights, feature_matrix[i])  
        prob_1 = 1/(1+np.exp(-z))  
        if prob_1 >= 0.5 + gamma or prob_1 <= 0.5 - gamma:  
            if y_pred[i] != labels[i]:  
                error_counter = error_counter + 1  
  
    #no_1 = # of points misclassified when  $y = 1$  for the points  $P(Y=1|x) \Rightarrow 0.5 + \gamma$  and  $y = 0$  for the  
    #no_1 = np.count_nonzero(error_margin)  
  
    return error_counter/(feature_matrix.shape[0])  
  
    ## STUDENT: YOUR CODE ENDS
```

In the function first prediction array is constructed with  $\gamma = 0$ . Then these results compared with the classifications done when the  $\gamma$  in model varies and the error is recorded. Plot indicates that best results are observed when  $\gamma$  is 0. Hence the largest margin/buffer is most probably achieved when  $\gamma$  is set to 0.

plot:



Task 8:

```
top 10 positive world  
awful  
deliciously  
nicest  
greatest  
loved  
perfected  
lovely  
beautifully  
fantasy  
world
```

```
top 10 negative world  
poorly  
disapppointment  
badly  
wasted  
worth  
disappointment  
dog  
maker  
die  
notable
```