

CSCE 5063-001: Assignment 2

Due 11:59pm Wednesday, October 6, 2021

1 Implementation of SVM via Gradient Descent

You will use the SVM to predict if a person doesn't have diabetes (negative class) or has diabetes (positive class), where the data is from National Institute of Diabetes and Digestive and Kidney Diseases. A total of 8 features will be used, including **Preg** (number of times pregnant), **Pres** (Diastolic blood pressure), **Mass** (body mass index), etc. There are a total number of 760 patients in the data.

The dataset is given in file `data.txt`, which contains the data description, and a matrix with 9 columns and 760 rows. Columns 1-8 are input features, and column 9 is the class label.

Note: For simplicity, in this assignment, you don't need to do the normalization and training/testing splitting, and you don't need to shuffle the data either as it is already shuffled.

You will implement the soft margin SVM using different gradient descent methods as we learned in class. To recap, to estimate the \mathbf{w}, b of the soft margin SVM, we can minimize the loss function:

$$J(\mathbf{w}, b) = \frac{1}{2} \sum_{j=1}^n (w_j)^2 + C \sum_{i=1}^m \max \left\{ 0, 1 - y^{(i)} \left(\sum_{j=1}^n w_j x_j^{(i)} + b \right) \right\}. \quad (1)$$

In order to minimize the function, we first obtain the gradient with respect to w_j , the j th item in the vector \mathbf{w} , as follows:

$$\nabla_{w_j} J(\mathbf{w}, b) = \frac{\partial J(\mathbf{w}, b)}{\partial w_j} = w_j + C \sum_{i=1}^m \frac{\partial L(x^{(i)}, y^{(i)})}{\partial w_j}, \quad (2)$$

where:

$$\frac{\partial L(x^{(i)}, y^{(i)})}{\partial w_j} = \begin{cases} 0 & \text{if } y^{(i)}(\mathbf{x}^{(i)}\mathbf{w} + b) \geq 1 \\ -y^{(i)}x_j^{(i)} & \text{otherwise.} \end{cases}$$

We should also compute the gradient with respect to b , which needs to be derived by yourself. Now, we will implement and compare the following gradient descent techniques.

For all methods, use $C = 10$.

Algorithm 1: Batch Gradient Descent: Iterate through the entire dataset and update the parameters as follows:

```
 $k = 0;$ 
while convergence criteria not reached do
  for  $j = 1, \dots, n$  do
    Update  $w_j \leftarrow w_j - \eta \nabla_{w_j} J(\mathbf{w}, b);$ 
  Update  $b \leftarrow b - \eta \nabla_b J(\mathbf{w}, b);$ 
  Update  $k \leftarrow k + 1;$ 
```

1. Batch Gradient Descent

where,

k is the number of iterations,

n is the dimensions of \mathbf{w} ,

η is the learning rate of the gradient descent, and

$\nabla_{w_j} J(\mathbf{w}, b)$ is the value computed from computing Eq. (2) above and $\nabla_b J(\mathbf{w}, b)$ is the value computed from your answer in question (a) below.

The *convergence criteria* for the above algorithm is $\Delta_{\%cost} < \epsilon$, where

$$\Delta_{\%cost} = \frac{|J_{k-1}(\mathbf{w}, b) - J_k(\mathbf{w}, b)| \times 100}{J_{k-1}(\mathbf{w}, b)}, \quad (3)$$

where,

$J_k(\mathbf{w}, b)$ is the value of Eq. (1) at k th iteration, and

$\Delta_{\%cost}$ is computed at the end of each iteration of the while loop.

Initialize $\mathbf{w} = \mathbf{0}, b = 0$ and compute $J_0(\mathbf{w}, b)$ with these values.

For this method, it is recommended to use $\eta = 0.0000000001$ and $\epsilon = 0.03$.

2. Stochastic Gradient Descent

where,

k is the number of iterations,

m is the number of examples in the training dataset,

n is the dimensions of \mathbf{w} ,

i is the training example currently used for computing gradient,

η is the learning rate of the gradient descent, and

$\nabla_{w_j} J^{(i)}(\mathbf{w}, b)$ is defined for a single training example as follows:

$$\nabla_{w_j} J^{(i)}(\mathbf{w}, b) = \frac{\partial J^{(i)}(\mathbf{w}, b)}{\partial w_j} = w_j + C \frac{\partial L(x^{(i)}, y^{(i)})}{\partial w_j}.$$

Note that you will also have to derive $\nabla_b J^{(i)}(\mathbf{w}, b)$, but it should be similar to your solution to question (a) below.

Algorithm 2: Stochastic Gradient Descent: Go through the dataset and update the parameters, one training example at a time, as follows:

```

Randomly shuffle the training data;
 $i = 0, k = 0$ ;
while convergence criteria not reached do
    for  $i = 0, \dots, m - 1$  do
        Take training example  $(i)$ ;
        for  $j = 1, \dots, n$  do
            Update  $w_j \leftarrow w_j - \eta \nabla_{w_j} J^{(i)}(\mathbf{w}, b)$ ;
        Update  $b \leftarrow b - \eta \nabla_b J^{(i)}(\mathbf{w}, b)$ ;
        Update  $k \leftarrow k + 1$ ;

```

The *convergence criteria* is $\Delta_{\%cost} < \epsilon$ similar to that of BGD (Eq. (3)). The difference is that, it is not computed at the end of each iteration, but is computed at the end of the while loop, i.e., compute the relative difference of loss values of two successive while loops.

For this method, it is recommended to use $\eta = 0.00000001, \epsilon = 0.4$. Note that the data is already shuffled so you don't have to.

3. Mini Batch Gradient Descent where,

Algorithm 3: Mini Batch Gradient Descent: Go through the dataset in batches of predetermined size and update the parameters as follows:

```

Randomly shuffle the training data;
 $i = 0, k = 0$ ;
while convergence criteria not reached do
    for  $i = 0, batch\_size, 2batch\_size, \dots, m - batch\_size$  do
        Take training examples  $(i), (i + 1), \dots, (i + batch\_size - 1)$  as a batch  $I$ ;
        for  $j = 1, \dots, n$  do
            Update  $w_j \leftarrow w_j - \eta \nabla_{w_j} J^I(\mathbf{w}, b)$ ;
        Update  $b \leftarrow b - \eta \nabla_b J^I(\mathbf{w}, b)$ ;
        Update  $k \leftarrow k + 1$ ;

```

k is the number of iterations,
 m is the number of examples in the training dataset,
 n is the dimensions of \mathbf{w} ,
 $batch_size$ is the number of training examples in each batch,
 I is the batch currently used for computing gradient,
 η is the learning rate of the gradient descent, and

$\nabla_{w_j} J^I(\mathbf{w}, b)$ is defined for a batch of training example as follows:

$$\nabla_{w_j} J^I(\mathbf{w}, b) = \frac{\partial J^I(\mathbf{w}, b)}{\partial w_j} = w_j + C \sum_{i \in I} \frac{\partial L(x^{(i)}, y^{(i)})}{\partial w_j}.$$

The *convergence criteria* is $\Delta_{\%cost} < \epsilon$ similar to that of SGD, i.e., it is computed at the end of the while loop.

For this method, it is recommended to use $\eta = 0.00000001$, $\epsilon = 0.5$, *batch_size* = 4. Note that the data is already shuffled so you don't have to.

Question (a)

Notice that we have not given you the equation for $\nabla_b J(\mathbf{w}, b)$.

Task: What is $\nabla_b J(\mathbf{w}, b)$ used for the Batch Gradient Descent Algorithm?

(Hint: It should be very similar to $\nabla_{w_j} J(\mathbf{w}, b)$.)

Question (b)

Task: Implement the SVM algorithm for the above mentioned gradient descent techniques.

Use $C = 10$ for all the techniques. **Note: update w in iteration $i + 1$ using the values computed in iteration i .** Do not update using values computed in the current iteration!

Run your implementation on the given dataset.

Task: (1) Plot the value of the cost function $J_k(\mathbf{w}, b)$ vs. the number of iterations (k). The diagram should have graphs from all the three techniques on the same plot. (2) Report the total time taken for convergence by each of the gradient descent techniques. You may notice that the total time taken by SGD is even larger than that of BGD. This is because we compute the cost function $J_k(\mathbf{w}, b)$ in each iteration of SGD for plotting, which is not necessary if we don't plot the diagram. Remove the code for computing the cost function $J_k(\mathbf{w}, b)$ in each iteration of SGD and MBGD and measure the time again.

For debugging, BGD should converge in about 60 iterations, SGD in about 4,500 iterations, and MBGD in about 1,000 iterations. All three gradient descent algorithms should converge with cost values around 6065.

What to submit

1. Equation for $\nabla_b J(\mathbf{w}, b)$.
2. Plot of $J_k(\mathbf{w}, b)$ vs. the number of iterations (k).
3. The convergence times in seconds.
4. The source code.