

# (PSL) Coding Assignment 5

## Contents

Introduction . . . . .	1
The SGD Algorithm . . . . .	1
Pegasos Algorithm . . . . .	1
Test Your Function . . . . .	2

## Introduction

In this assignment, you will be implementing a linear Support Vector Machine (SVM) classifier from scratch using stochastic gradient descent (SGD).

As we discussed in class, traditional SVMs often solve the dual problem, which involves a quadratic objective function subject to linear constraints. While this approach can be efficient for small-scale tasks, it becomes less practical for large-scale problems. In such cases, we can leverage the benefits of SGD to directly solve the primal problem.

## The SGD Algorithm

The SGD algorithm works as follows:

1. Start by choosing a random initial value of parameters
2. Loop Over Epochs:
  - Within each epoch, iterate over each data point in the training dataset. An epoch is a complete pass through all the training data.
3. Update the Gradient:
  - For each data point, calculate the gradient of the loss function with respect to the current parameter values. This gradient represents the direction of steepest ascent.
4. Update Parameters:
  - For each parameter, calculate the step size as :  $\text{step size} = \text{gradient} * \text{learning rate}$ .
  - Update new parameters as :  $\text{new params} = \text{old params} - \text{step size}$
5. Repeat Until Convergence:
  - Repeat steps 3 to 4 for each data point in the dataset. Continue this process for a fixed number of epochs or until convergence criteria are met.

## Pegasos Algorithm

The Pegasos (Primal Estimated sub-GrAdient SOLver for SVM) algorithm, proposed by Shalev-Shwartz et al. (2011) [Paper Link], is an application of SGD.

Recall that the primal problem of linear SVM can be expressed in a “Loss + Penalty” format without constraints:

$$\frac{1}{n} \sum_{i=1}^n [1 - y_i(x_i^t \beta + \alpha)]_+ + \frac{\lambda}{2} \|\beta\|^2$$

where  $\alpha$  is the intercept and  $\beta$  is the  $p$ -dimensional coefficient vector.

Here  $\lambda$  is not the Lagrange multiplier in the dual problem; rather, it is a tuning parameter associated with  $\|\beta\|^2$ , consistent with the notation used in the Pegasos paper. This  $\frac{\lambda}{2}$  corresponds to  $\nu$  in Eq (6) on page 10 of lec\_W11.2\_SVM [Slides Link].

The **Pegasos Algorithm** can be summarized as follows:

1. initialize  $\beta = 0_{p \times 1}$ ,  $\alpha_1 = 0$ , and  $t = 0$
2. for  $epoch = 1, 2, \dots, T$  do
  - for  $i = 1, 2, \dots, n$  do
    - $t = t + 1$ ,  $\eta_t = \frac{1}{t\lambda}$
    - update  $\beta_{t+1} \leftarrow \beta_t - \eta_t \Delta_t$
    - update  $\alpha_{t+1} \leftarrow \alpha_t - \eta_t \delta_t$

Here  $\eta_t$  is the learning rate, and  $\Delta_t$  and  $\delta_t$  are the (sub)gradient of  $J_i(\beta, \alpha)$  when  $\beta = \beta_t$  and  $\alpha = \alpha_t$ :

$$J_i(\beta, \alpha) = \frac{\lambda}{2} \|\beta\|^2 + [1 - y_i(x_i^t \beta + \alpha)]_+$$

$$\Delta_t = \begin{cases} \lambda \beta_t - y_i x_i & \text{if } y_i(x_i^t \beta_t + \alpha_t) < 1 \\ \lambda \beta_t & \text{otherwise} \end{cases}$$

$$\delta_t = \begin{cases} -y_i & \text{if } y_i(x_i^t \beta_t + \alpha_t) < 1 \\ 0 & \text{otherwise} \end{cases}$$

## Test Your Function

- Implement the **Pegasos Algorithm**.
  - Use a fixed number of epochs, e.g.,  $T = 20$ .
  - In each epoch, before going through the dataset, consider randomizing the order of the data points. To achieve this, you should set random seeds for shuffling. For this assignment, the seeds used for shuffling **do not need to be associated with your UIN**.
- Test your code with the provided training (200 samples) and test (600 samples) datasets, which are subsets of the MNIST data. Each dataset consists of 257 columns, with the first 256 columns representing the features, and the last column indicating the label (either 5 or 6).
  - `coding5_train.csv` [Link]
  - `coding5_test.csv` [Link]
- If you set  $y_i = 1$  for label 5 and  $y_i = -1$  for label 6, then after obtaining  $\beta$  and  $\alpha$ , for a test point  $x_*$ , you'll predict:
  - $y_* = 5$ , if  $x_*^t \beta + \alpha > 0$ ; otherwise, 6.
- Report **confusion tables** on the **training and test** datasets.