

## Problem Set 7, Oct 27, 2020 (SVM)

**Goals.** The goal of this exercise is to

- Implement and debug Support Vector Machine (SVM) using SGD and coordinate descent.
- Derive updates for the coordinate descent algorithm for the dual optimization problem for SVM.
- Implement and debug the coordinate descent algorithm.
- Compare it to the primal solution.

**Setup, data and sample code.** Obtain the folder `labs/ex07` of the course github repository

[github.com/epfml/ML\\_course](https://github.com/epfml/ML_course)

We will finally depart from using the height-weight dataset and instead use a toy dataset from scikit-learn. We have provided sample code templates that already contain useful snippets of code required for this exercise.

### 1 Support Vector Machines using SGD

Until now we have implemented linear and logistic regression to do classification. In this exercise we will use the Support Vector Machine (SVM) for classification. As we have seen in the lecture notes, the original optimization problem for the Support Vector Machine (SVM) is given by

$$\min_{\mathbf{w} \in \mathbb{R}^D} \sum_{n=1}^N \ell(y_n \mathbf{x}_n^\top \mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad (1)$$

where  $\ell : \mathbb{R} \rightarrow \mathbb{R}$ ,  $\ell(z) := \max\{0, 1 - z\}$  is the *hinge loss* function. Here for any  $n$ ,  $1 \leq n \leq N$ , the vector  $\mathbf{x}_n \in \mathbb{R}^D$  is the  $n^{th}$  data example, and  $y_n \in \{\pm 1\}$  is the corresponding label.

#### Problem 1 (SGD for SVM):

Implement stochastic gradient descent (SGD) for the original SVM formulation (1). That is in every iteration, pick one data example  $n \in [N]$  uniformly at random, and perform an update on  $\mathbf{w}$  based on the (sub-)gradient of the  $n^{th}$  summand of the objective (1). Then iterate by picking the next  $n$ .

1. Fill in the notebook functions `calculate_accuracy(y, X, w)` which computes the accuracy on the training/test dataset for any  $\mathbf{w}$  and `calculate_primal_objective(y, X, w, lambda_)` which computes the total primal objective (1).
2. Derive the SGD updates for the original SVM formulation and fill in the notebook function `calculate_stochastic_gradient()` which should return the stochastic gradient of the total cost function (loss plus regularizer) with respect to  $\mathbf{w}$ . Finally, use `sgd_for_svm_demo()` provided in the template for training.

## 2 Support Vector Machines using Coordinate Descent

As seen in class, another approach to train SVMs is by considering the dual optimization problem given by

$$\max_{\alpha \in \mathbb{R}^N} \alpha^\top \mathbf{1} - \frac{1}{2\lambda} \alpha^\top \mathbf{Y} \mathbf{X} \mathbf{X}^\top \mathbf{Y} \alpha \quad \text{such that} \quad 0 \leq \alpha_n \leq 1 \quad \forall n \in [N] \quad (2)$$

where  $\mathbf{1}$  is the vector of size  $N$  filled with ones,  $\mathbf{Y} := \text{diag}(\mathbf{y})$ , and  $\mathbf{X} \in \mathbb{R}^{N \times D}$  again collects all  $N$  data examples as its rows, as usual. In this approach we optimize over the dual variables  $\alpha$  and map the solutions back to the primal vector  $\mathbf{w}$  via the mapping we have seen in class:  $\mathbf{w}(\alpha) = \frac{1}{\lambda} \mathbf{X}^\top \mathbf{Y} \alpha$ .

### Problem 2 (Coordinate Descent for SVM):

In this part we will derive the coordinate descent (or rather ascent in our case) algorithm seen in class to solve the dual (2) of the SVM formulation. That is, at every iteration, we pick a coordinate  $n \in [N]$  uniformly at random, and fully optimize the objective (2) **with respect to that coordinate alone**. Hence one step of coordinate ascent corresponds to solving, for a given coordinate  $n \in [N]$  and our current vector  $\alpha \in [0, 1]^N$ , the one dimensional problem:

$$\max_{\gamma \in \mathbb{R}} f(\alpha + \gamma e_n) \quad \text{such that} \quad 0 \leq \alpha_n + \gamma \leq 1 \quad (3)$$

where  $f(\alpha) = \alpha^\top \mathbf{1} - \frac{1}{2\lambda} \alpha^\top \mathbf{Y} \mathbf{X} \mathbf{X}^\top \mathbf{Y} \alpha$  and  $e_n = [0, \dots, 1, \dots, 0]^\top$  (all zero vector except at the  $n^{\text{th}}$  position). We denote by  $\gamma^*$  the value of  $\gamma$  which maximises problem 3. The coordinate update is then  $\alpha^{\text{new}} = \alpha + \gamma^* e_n$ .

1. Solve problem 3 and give a closed form solution for the update  $\alpha^{\text{new}} = \alpha + \gamma^* e_n$ , this update should only involve  $\alpha$ ,  $\lambda$ ,  $x_n$ ,  $y_n$  and  $\mathbf{w}(\alpha)$  (*Hint*: Notice that  $\gamma \mapsto f(\alpha + \gamma e_n)$  is polynomial and **don't forget the constraints**!). Convince yourself that this update can be computed in  $O(D)$  time.
2. Find an efficient way of updating the corresponding  $\mathbf{w}(\alpha^{\text{new}})$ . This should be computable in  $O(D)$  time.
3. Fill in the notebook functions `calculate_coordinate_update()` which should compute the coordinate update for a single desired coordinate and `calculate_dual_objective()` which should return the objective (loss) for the dual problem (2).
4. Finally train your model using coordinate descent (here ascent) using the given function `sgd_for_svm_demo()` in the template. Compare to your SGD implementation. Which one is faster? (Compare the training objective values (1) for the  $\mathbf{w}$  iterates you obtain from each method). Is the gap going to 0?

## Theory Exercises

### Problem 3 (Kernels):

In class we have seen that many kernel functions  $k(\mathbf{x}, \mathbf{x}')$  can be written as inner products  $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$  for a suitably chosen feature map  $\phi(\cdot)$ . Let us say that such a kernel function is *valid*. We further discussed many operations on valid kernel functions that result again in valid kernel functions. Here are two more.

1. Let  $k_1(\mathbf{x}, \mathbf{x}')$  be a valid kernel function. Let  $f$  be a polynomial with positive coefficients. Show that  $k(\mathbf{x}, \mathbf{x}') = f(k_1(\mathbf{x}, \mathbf{x}'))$  is a valid kernel.
2. Show that  $k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}'))$  is a valid kernel assuming that  $k_1(\mathbf{x}, \mathbf{x}')$  is a valid kernel. *Hint*: The limit of valid kernels is still a valid kernel.