

## Homework 3

Lecturer: Inderjit Dhillon

Date Due: Sep 29, 2014

**Keywords:** *SVM, Coordinate Descent, OpenMP***Note:**

1. In this assignment we use the following two datasets: RCV1 and Covtype. These two datasets can be downloaded from the Assignments tab in Canvas (called `hw3-data.zip`).
2. You will run your code and collect timing results on TACC machines. See <https://www.tacc.utexas.edu/user-services/user-guides/stampede-user-guide> to know about using the Stampede cluster.
3. Submit your code through Canvas. Submit a hard copy of your answers in class.

**Problem 1: Dual Coordinate Descent for Linear SVM**

Given training data  $\{\mathbf{x}_i, y_i\}_{i=1}^n$ , each  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \{+1, -1\}$ , consider the following variant of SVM problem (without the bias term):

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i^2 \equiv g(\mathbf{w}) \\ \text{subject to} \quad & y_i \mathbf{w}^T \mathbf{x}_i \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad \forall i = 1, \dots, n. \end{aligned} \quad (1)$$

Note that in the above formulation we consider the squared hinge  $\xi_i^2$  instead of  $\xi_i$ . The dual problem has the following form:

$$\begin{aligned} \min_{\boldsymbol{\alpha} \in \mathbb{R}^n} \quad & \frac{1}{2} \boldsymbol{\alpha}^T (Q + \frac{1}{2C} I) \boldsymbol{\alpha} - \mathbf{e}^T \boldsymbol{\alpha} \equiv f(\boldsymbol{\alpha}), \\ \text{subject to} \quad & 0 \leq \alpha_i, i = 1, \dots, n, \end{aligned} \quad (2)$$

where  $Q_{ij} = y_i y_j \mathbf{x}_i^T \mathbf{x}_j$ . The optimal primal variable  $\mathbf{w}^*$  and optimal dual variable  $\boldsymbol{\alpha}^*$  satisfy

$$\mathbf{w}^* = \sum_{i=1}^n y_i \alpha_i^* \mathbf{x}_i.$$

Write a C/C++ program which implements the coordinate descent algorithm for solving the SVM dual problem (Algorithm 1), and parallelize it using OpenMP.

- (5 pt) Show that (2) is the dual problem of (1).
- (2 pt) Applying coordinate descent to (2) requires solving the following single variable sub-problem.

$$\delta^* = \arg \min_{\delta} f(\boldsymbol{\alpha} + \delta \mathbf{e}_i). \quad (3)$$

Derive the closed form solution to (3). What is the time complexity for directly computing the closed form solution?

---

**Algorithm 1** Co-ordinate Descent for Dual SVM
 

---

- Input: Data  $\{y_i, \mathbf{x}_i : i = 1, \dots, n\}$  and  $C$
- Initialize  $\boldsymbol{\alpha}$  and  $\mathbf{w}$  to zero
- For iter = 1, ..., (Outer iterations)
  - Choose a random permutation  $\pi$
  - For  $s = 1, \dots, n$  (Inner iterations)
    - \*  $i = \pi(s)$
    - \* Solve single variable problem for  $\alpha_i$ :

$$\delta^* = \arg \min_{\delta} f(\boldsymbol{\alpha} + \delta \mathbf{e}_i).$$

- \* Perform updates:
      - $\alpha_i = \alpha_i + \delta^*$
      - maintain  $\mathbf{w}$
    - Output timing, objective function value, and test accuracy
- 

- (c) (3 pt) Function value evaluation  $f(\boldsymbol{\alpha})$  involves a term  $\boldsymbol{\alpha}^T Q \boldsymbol{\alpha}$ , which requires  $O(n^2)$  time as  $Q$  is a dense matrix. Now assume that a  $d$  dimensional variable  $\mathbf{w} = \sum_{i=1}^n y_i \alpha_i \mathbf{x}_i$  is available. Give an alternative formulation of  $f(\boldsymbol{\alpha})$  using both  $\boldsymbol{\alpha}$  and  $\mathbf{w}$  such that the function value evaluation  $f(\boldsymbol{\alpha})$  costs only  $O(n + d)$ .
- (d) (2 pt) How to use the assumption in part (c) (i.e.,  $\mathbf{w} = \sum_{i=1}^n y_i \alpha_i \mathbf{x}_i$ ) to reduce the time complexity for computing the closed form solution for  $\delta^*$  to  $O(d_i)$ , where  $d_i$  is the number of non-zeros of  $\mathbf{x}_i$ ?
- (e) (1 pt) How would you maintain  $\mathbf{w}$  after the  $\alpha_i$  update to ensure the validity of the assumption  $\mathbf{w} = \sum_{i=1}^n y_i \alpha_i \mathbf{x}_i$ ? The time complexity for the maintenance should also be  $O(d_i)$ .
- (f) (18 pt) Implement the dual co-ordinate descent for solving (2) in C/C++ and parallelize the inner iterations using OpenMP. The pseudo code is given in Algorithm 1. At the end of each outer iteration, print the following quantities:
- (a) The wall time for co-ordinate descent updates (excluding preprocessing time and time for computing the measurements below).
  - (b) Compute and output the dual objective function value (using the current  $\boldsymbol{\alpha}, \mathbf{w}$ ).
  - (c) Compute and output the primal objective function value using  $\mathbf{w}$ .
  - (d) Compute and output the prediction accuracy using current  $\mathbf{w}$ .
  - (e) Compute and output  $\|\sum_{i=1}^n y_i \alpha_i \mathbf{x}_i - \mathbf{w}\|_2^2$ . This quantity should be zero if  $\mathbf{w}$  is correctly maintained. To ensure the value is zero in the multi-threaded environment, you might need to use `# pragma omp atomic`.

Report the 20-iteration running time, primal objective function, and prediction accuracy for both RCV1 and Covtype datasets with 1, 8, 16 threads using  $C = 0.1$ . What is the speedup you observe?

- (g) (4 pt) Remove the atomic flag (or locks, critical sections) in your code for correctly maintaining  $\mathbf{w} = \sum_{i=1}^n y_i \alpha_i \mathbf{x}_i$ .

For each dataset, and for the locking and non-locking cases, plot the wall time versus relative error  $((g(\mathbf{w}) - g(\mathbf{w}^*)) / g(\mathbf{w}^*))$ , where  $g(\cdot)$  is the primal SVM objective function (1). In each of the four plots, include curves for 1, 8, 16 threads using  $C = 0.1$ . Plot the relative error in log-scale. Comment on your observations.

Note that you can use the sequential version optimal solution as  $\mathbf{w}^*$ .

Below are some guidelines for your implementation.

- Fix the number of outer iterations to 20 and initialize  $\alpha$  and  $\mathbf{w}$  with the zero vector.
- Write a Makefile to compile your code using `icc`. The name of the output binary file should be `cd-svm`.
- Usage of the binary file: `./cd-svm C nr_threads traindata testdata`
- Datasets:
  - `hw3-data.zip` contains the text-format version of the three datasets used in HW3
  - For each dataset, there is a testing file (ended with `.t`) and a training file (ended with `.tr`). Each line of the files contain a data point, started with the label (+1 or -1) and followed by a sorted list of `< feature >:< value >` pairs. See SVMlight format described in <http://svmlight.joachims.org/> for detail.