

Dual Coordinate Descent-Group 31

AKYAM DHATRI NANDA - AI20BTECH11002
VOJESWITHA GOPIREDDY -AI20BTECH11024
VAISHNAVI W - AI20BTECH11025
MUCHARLA SUPRIYA - EE19BTECH11022
AKUNURI VISHWA VEENA - ES19BTECH11026

April 30, 2022

- Abstract
- Introduction
- Linear SVM Optimization Problem
 - ▶ Primal Optimization Problem
 - ▶ Dual Formulation
 - ▶ Dual Optimization Problem
- Dual Coordinate Descent
 - ▶ Outer Iteration
 - ▶ Inner Iteration
- Procedure
- Comparisons
- References

Abstract

Support Vector Machine (SVM) is a supervised machine learning algorithm used for both classification and regression tasks. The optimization problem is to maximize the margin with the constraint that no points lie in between the margin which simplifies to a quadratic programming problem.

Lagrange problem is typically solved in the dual form. One of the main challenges is the large computational cost of QPP. The long training time also prevents it from locating the optimal parameter set from a very fine grid of parameters over a large span. The Dual Coordinate Descent (DCD) is one of the efficient algorithms that directly solves the QPP by orderly optimizing the series of single variable sub-problems. We wish to explore the efficiency of the algorithm.

Introduction

Support Vector Machines are supervised learning algorithms for classification of the data. SVM though is a linear model, it can be extended to solve non-linear classification problems. In this paper we wish to explore dual coordinate descent, an optimisation algorithm, to solve linear SVM classification problem. We are applying coordinate descent technique to the dual optimisation problem of the Linear SVM, where we update one variable at a time by minimising a single variable sub-problem. In general implementations, accessing values per-instance is much easier than accessing data values of corresponding features. Hence, choosing the dual problem while implementing coordinate descent method is better option compared to the primal problem.

Linear SVM Optimisation Problem

Formulation of linear SVM for classification is as follows:

Consider a training set of data instances and labels $(x_i, y_i), i = 1, 2, \dots, l$, where $y_i \in \{-1, +1\}$ and $x_i \in \mathbb{R}^n$. And the classification function is $\text{sgn}(w^T x + b)$, where (w, b) are the weights and bias respectively.

Linear SVM Optimisation Problem

Primal Optimisation Problem

$$\begin{aligned} \min_{\omega, b} \quad & \frac{1}{2} \omega^T \omega + C \sum_{i=1}^l \xi_i \\ \text{subject to} \quad & y_i(\omega^T x_i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0, i = 1, 2, \dots, l \end{aligned} \tag{1}$$

where n is the number of features, $C \geq 0$ is the penalty parameter and b is the bias term.

Dual Formulation

For the optimisation problem, the Lagrangian is given by:

$$L(\omega, b, \xi, \alpha, \beta) = \frac{1}{2}\omega^T \omega + c \sum_{i=1}^I \xi_i + \sum_{i=1}^I \beta_i(-\xi_i) + \sum_{i=1}^I \alpha_i[(1 - \xi_i) - y_i(\omega^T x_i + b)] \quad (2)$$

The dual form of the problem is obtained by minimising Lagrangian with respect to ω , b and ξ . Each lagrange multiplier i.e., α_i , β_i is constrained to be greater than or equal to zero.

$$\max_{\alpha, \beta} f_d(\alpha, \beta) = \max_{\alpha, \beta} \min_{\omega, b, \xi} L(\omega, b, \xi, \alpha, \beta) \quad (3)$$

Dual Formulation

Setting derivative of Lagrangian wrt ω to zero we get,

$$\frac{\partial L}{\partial \omega} = \omega - \sum_{i=1}^l \alpha_i y_i x_i \quad (4)$$

$$\implies \omega = \sum_{i=1}^l \alpha_i y_i x_i \quad (5)$$

Setting derivative of Lagrangian wrt b to zero we get,

$$\frac{\partial L}{\partial b} = \sum_{i=1}^l \alpha_i y_i \quad (6)$$

$$\implies \sum_{i=1}^l \alpha_i y_i = 0 \quad (7)$$

Dual Formulation

Setting derivative of Lagrangian wrt ξ_i to zero we get,

$$\frac{\partial L}{\partial \xi_i} = C - \alpha_i - \beta_i \quad (8)$$

$$\implies C = \alpha_i + \beta_i \quad (9)$$

Since, we have $\beta_i = C - \alpha_i$, the constraint $\beta_i \geq 0$ becomes $\alpha_i \leq C$.
Therefore the constraints of the dual problem are

$$0 \leq \alpha_i \leq C \quad (10)$$

Dual Formulation

On substituting (6), (8), (10) in (2), the lagrangian simplifies to

$$f_d(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j \quad (11)$$

$$= -\left(\frac{1}{2} \alpha^T Q \alpha - e^T \alpha\right) \quad (12)$$

Linear SVM Optimisation Problem

Dual Optimisation Problem

$$\begin{aligned} \min_{\alpha} f(\alpha) &= \frac{1}{2} \alpha^T Q \alpha - \mathbf{e}^T \alpha \\ \text{subject to } &0 \leq \alpha_i \leq C, i = 1, 2, \dots, l \\ &\mathbf{y}^T \alpha = 0 \text{ (if bias is considered)} \end{aligned} \quad (13)$$

where $\mathbf{e} = [1, \dots, 1]^T$, $Q_{ij} = y_i y_j \mathbf{x}_i^T \mathbf{x}_j$ and α is the Lagrangian multiplier.

Dual Coordinate Descent

In coordinate descent method, each coordinate of the variable vector is updated successively, one at a time, unlike Gradient Descent method where all the coordinates are updated at once. In the above stated linear SVM dual optimization problem, we are going to update each coordinate of α (α_i). Once all the coordinates are updated, the process of updating repeats.

Outer Iteration

Let us consider $\alpha^0 \in \mathbb{R}'$ as the initial vector of Lagrangians. In the optimisation process, sequence of vectors $\{\alpha^k\}_{k=0}^{\infty}$ are generated. The process of generating α^{k+1} from α^k is called an outer iteration.

Inner Iteration

Each outer iteration consists of l inner iterations, where we update each coordinate of $\alpha^k = [\alpha_1^k, \alpha_2^k, \dots, \alpha_l^k]$ sequentially. After l inner iterations, we get α^{k+1}

$\alpha^{k,i}$ represents the vector generated at the end of k outer iterations followed by i inner iterations.

$$\alpha^{k,i} = [\alpha_1^{k+1}, \dots, \alpha_{i-1}^{k+1}, \alpha_i^k, \dots, \alpha_l^k]^T \quad \forall i = 2, \dots, l \quad (14)$$

Procedure

One variable sub-problem

Since we update one coordinate at a time in an iteration, the above dual optimization problem now reduces to a one-variable sub-problem in a variable d .

On updating $\alpha^{k,i}$ to $\alpha^{k,i+1}$, the i^{th} coordinate is updated while keeping others the same. The optimisation problem we now have is,

$$\begin{aligned} \min_d \quad & f(\alpha^{k,i} + d\mathbf{e}_i) \\ \text{subject to} \quad & 0 \leq \alpha_i^k + d \leq C \end{aligned} \tag{15}$$

Procedure

Objective Function

The objective function for the above optimisation problem (15) is a quadratic function in variable d . It is given as,

$$f(\alpha^{k,i} + d\mathbf{e}_i) = f(\alpha^{k,i}) + \nabla_i f(\alpha^{k,i})d + \frac{1}{2}Q_{ii}d^2 \quad (16)$$

where $\nabla_i f$ is the i^{th} component of ∇f

Procedure

Projection Gradient

$\nabla^P f(\alpha)$ is the projected gradient defined as below,

$$\nabla_i^P f(\alpha) = \begin{cases} \nabla_i f(\alpha) & 0 < \alpha_i < C \\ \min(0, \nabla_i f(\alpha)) & \alpha_i = 0 \\ \max(0, \nabla_i f(\alpha)) & \alpha_i = C \end{cases} \quad (17)$$

If projected gradient at $\alpha^{k,i}$ becomes 0, the updating of α_i is not necessary as the optimal α_i is reached, and move to the next index $i + 1$.

Procedure

Solving the sub-problem

Without any constraints, optimal d for the above objective function

$$d = -\frac{\nabla_i f(\alpha^{k,i})}{Q_{ii}} \quad (18)$$

From (15) subjecting d to the constraints, the obtained solution is

$$\alpha_i^{k,i+1} = \begin{cases} \alpha_i^{k,i} - \frac{\nabla_i f(\alpha^{k,i})}{Q_{ii}} & 0 \leq \alpha_i^k + d \leq C \\ C & \alpha_i^k + d \geq C \\ 0 & \alpha_i^k + d \leq 0 \end{cases} \quad (19)$$

which can also be written as,

$$\alpha_i^{k,i+1} = \min(\max(\alpha_i^{k,i} - \frac{\nabla_i f(\alpha^{k,i})}{Q_{ii}}, 0), C) \quad (20)$$

Procedure

Finding the Gradients (∇f)

The i^{th} component of gradient of objective function can be obtained as,

$$\begin{aligned}\nabla_i f(\boldsymbol{\alpha}) &= \sum_{j=1}^I Q_{ij} \alpha_j - 1 \\ &= \sum_{j=1}^I y_i y_j \mathbf{x}_i^T \mathbf{x}_j \alpha_j - 1\end{aligned}\quad (21)$$

The operations involved in the above expression to calculate gradient directly are computationally expensive, costs $O(I^n)$.

Procedure

Finding the Gradients (∇f)

However, since for linear SVM with soft margin we have

$$\omega = \sum_{j=1}^I y_j \alpha_j \mathbf{x}_j \quad (22)$$

Gradient now can be calculated in $O(n)$ time as follows

$$\nabla_i f(\alpha) = y_i \omega^T \mathbf{x}_i - 1 \quad (23)$$

Procedure

Updating the weights

From the primal-dual relationship, we have $\omega = \sum_{i=1}^l y_i \alpha_i \mathbf{x}_i$. Calculating ω using it is computationally expensive, takes $O(ln)$ time. An alternate way for updating ω after each inner iteration, taking $O(n)$ time is as follows,

$$\omega \leftarrow \omega + (\alpha_i^{k+1} - \alpha_i^k) y_i \mathbf{x}_i \quad (24)$$

At the end of all the outer iterations, we obtain the optimal value of ω of the primal problem (1)

Algorithm

- Given $\alpha = \alpha^0$ and initial weights $\omega = \sum_{i=1}^l y_i \alpha_i^0 \mathbf{x}_i$
- While α is not optimal (Outer Iteration)

For $i = 1, 2, \dots, l$ (Inner Iteration)

(a) $\alpha_{temp} = \alpha_i$

(b) $\nabla = y_i \omega^T \mathbf{x}_i - 1$

(c)
$$\nabla^P = \begin{cases} \min(\nabla, 0) & \alpha_i = 0 \\ \max(\nabla, 0) & \alpha_i = C \\ \nabla & 0 < \alpha_i < C \end{cases}$$

(d) If $\nabla^P \neq 0$

$$\alpha_i \leftarrow \min(\max(\alpha_i - \frac{\nabla}{Q_{ii}}, 0), C)$$

$$\omega \leftarrow \omega + (\alpha_i - \alpha_{temp}) y_i \mathbf{x}_i$$

Discussions

Bias Term

In our report, in the formulation of linear SVM optimisation problem, the bias term was not considered. With the inclusion of bias term, a linear constraint is introduced in the dual optimisation problem i.e.,

$$\mathbf{y}^T \boldsymbol{\alpha} = 0 \quad (25)$$

Dual coordinate descent for linear SVM without the bias term converges faster when compared with that with a bias term.

The stated reason, a linear constraint may make the sub-problem to be easily optimal already; α s are not moved/updated.

Results

We used LinearSVM from sklearn library in python over the dataset 'The Cleveland Heart Disease'. The class LinearSVM uses DCD when 'dual' is set to True. When fitting the dataset over a smaller number of samples, we observed that the time taken for the dual problem is lesser when compared to the primal problem. Further, we experimented with exclusion of bias term and random shuffling of sub-problems and they were computationally efficient. However, this wasn't the same when a larger sample size was used.

Results contd.

Sample Size	Method	Time(s)
n = 30	Dual - DCD	0.00961
	Dual - Rand	0.00832
	Dual - Bias	0.00848
	Primal	0.01073
n = 300	DCD	0.1136
	Primal	0.0251

Conclusions

For linear SVM, DCD is efficient when used without a bias term, thus avoiding linear constraint. For one-class SVMs and SVDD, which have an inherent linear constraint, novel coordinate descent methods have to be developed. Further studies shows that, the training speed increases by randomising the order of sub-problems in each iteration. Also, as the instances and number of features in each instance increase, the DCD method trains well over other solvers.