
Assignment 1

Ankit Kumar
170121

Chaitanya Pathare
19111061

Hemant Kumar
170297

Rahul Ninaji Pakhare
19111068

Raunak Kumar
19111069

Abstract

Given The Squared SVM Solver we need to add dual constraint to it and define the Lagrangian function to it. And then eliminate the primal variables using first order optimality. Solve the problem with any 3 methods and plot it in graph

1 Question

Given equation

$$\Rightarrow \operatorname{argmin}_{w \in R^d} 1/2 \mathbf{w}^2 + C \sum_{i=1}^n ([1 - y^i < \mathbf{w}, \mathbf{x}^i >]_+)^2 \quad (P1)$$

rewritten as new optimization problem

$$\Rightarrow \operatorname{argmin}_{w \in R^d} 1/2 \mathbf{w}^2 + C \sum_{i=1}^n \xi_i^2 \text{ for all } i \in [1, n]$$

$$s.t. \quad y^i < \mathbf{w}, \mathbf{x}^i > \geq 1 - \xi_i, \xi_i \geq 0 \quad \text{for all } i \in [1, n] \quad (P2)$$

Introducing dual variable α_i for each of the n constraints in (P2)

$$L(w, \xi, \alpha) = 1/2 \mathbf{w}^2 + C \sum_{i=1}^n \xi_i^2 + \sum_{i=1}^n \alpha_i (1 - \xi_i - y^i < \mathbf{w}, \mathbf{x}^i >) \quad , \alpha_i \geq 0 \quad (P3)$$

2 Question

To get the dual problem from the above equation (P3). We have appended one extra dimension to our data i.e if $X =$

$$X = \begin{bmatrix} x_{11} & x_{12} & x_{13} & \cdot & x_{1n} \\ x_{21} & \vdots & \vdots & \vdots & x_{2n} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ x_{n1} & \vdots & \vdots & \vdots & x_{nn} \end{bmatrix} \text{ and then transformed } X \rightarrow X^0 =$$

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & \cdot & x_{1n} & 1 \\ x_{21} & \cdot & \cdot & \cdot & x_{2n} & 1 \\ \cdot & \cdot & \ddots & \ddots & \cdot & 1 \\ \cdot & \cdot & \ddots & \ddots & \cdot & 1 \\ \cdot & \cdot & \ddots & \ddots & \cdot & 1 \\ x_{n1} & \cdot & \cdot & \cdot & x_{nn} & 1 \end{bmatrix} \text{ so that } w[d-1] = b, \{d \text{ is dimension of each data point in } X^0\}$$

put $\delta L / \delta \mathbf{w} = \mathbf{0}$ and $\delta L / \delta \xi_i = \mathbf{0}$ for all $i \in [1, n]$

$$\delta L / \delta \mathbf{w} = \mathbf{w} - \sum_{i=1}^n \alpha_i \mathbf{x}^i = \mathbf{0}$$

$$\text{or } \mathbf{w} = \sum_{i=1}^n \alpha_i \mathbf{x}^i$$

and

$$\delta L / \delta \xi_i = 2C\xi_i - \sum_{i=1}^n \alpha_i = \mathbf{0}$$

$$\text{or } \xi_i = 1/2C \sum_{i=1}^n \alpha_i$$

Optimization problem comes down to

$$\Rightarrow 1/2 \sum_{i=1}^n \sum_{j=1}^n (\alpha_i \alpha_j y^i y^j < \mathbf{x}^i \mathbf{x}^j >) + 1/4C \sum_{i=1}^n \alpha_i^2 + \sum_{i=1}^n \alpha_i - 1/2C \sum_{i=1}^n \alpha_i^2$$

$$- \sum_{i=1}^n \sum_{j=1}^n (\alpha_i \alpha_j y^i y^j < \mathbf{x}^i \mathbf{x}^j >)$$

$$\Rightarrow \text{argmin}(\sum_{i=1}^n \alpha_i - 1/4C(\sum_{i=1}^n \alpha_i^2) - 1/2 \sum_{i=1}^n \sum_{j=1}^n (\alpha_i \alpha_j y^i y^j < \mathbf{x}^i \mathbf{x}^j >)) \quad (D2)$$

$$\text{All above } \mathbf{w} \in R^{d+1}, \mathbf{x} \in R^{d+1}$$

3 Question

a. Mini Batch Stochastic Gradient descent on P1

Given,

$$P1 \quad f = \text{argmin}_{\mathbf{w} \in R^d} 1/2 \mathbf{w}^2 + C \sum_{i=1}^n ([1 - y^i < \mathbf{w}, \mathbf{x}^i >]_+)^2$$

now get gradient for P1

$$\nabla f = \mathbf{w} + 2C \sum_{i=1}^n ([1 - y^i < \mathbf{w}, \mathbf{x}^i >]_+) (-y^i \mathbf{x}^i)$$

Note: $[1 - y^i < \mathbf{w}, \mathbf{x}^i >]_+ = \{(1 - y^i < \mathbf{w}, \mathbf{x}^i >) \text{ if } y^i < \mathbf{w}, \mathbf{x}^i > < 1$

$$0 \text{ if } y^i < \mathbf{w}, \mathbf{x}^i > \geq 1\}$$

iterate till the time end for the X

Do

- 1.find gradient for each data point in the given Batch set B and step length eta
- 2.update w for these

i.e, $\mathbf{w} = \mathbf{w} - \eta 2C \sum_{i=x}^{x+B} ([1 - y^i < \mathbf{w}, \mathbf{x}^i >]_+) (-y^i \mathbf{x}^i)$

Differengt cobination of (batch B) and step length η we tried are as follows,all works fine

B	η
10	0.01
100	0.05
1000	0.001

b. Coordinate descent on P1

Given,

$$P1 \quad f = \operatorname{argmin}_{\mathbf{w} \in R^d} 1/2 \mathbf{w}^2 + C \sum_{i=1}^n ([1 - y^i < \mathbf{w}, \mathbf{x}^i >]_+)^2$$

now get gradient for P1

$$\nabla f = \mathbf{w} + 2C \sum_{i=1}^n ([1 - y^i < \mathbf{w}, \mathbf{x}^i >]_+) (-y^i \mathbf{x}^i)$$

Note: $[1 - y^i < \mathbf{w}, \mathbf{x}^i >]_+ = \begin{cases} (1 - y^i < \mathbf{w}, \mathbf{x}^i >) & \text{if } y^i < \mathbf{w}, \mathbf{x}^i > < 1 \\ 0 & \text{if } y^i < \mathbf{w}, \mathbf{x}^i > \geq 1 \end{cases}$

iterate till the time end for the X

Do

- 1.find gradient for each data point choosing the coordinate cyclically i,e j th= 1,2...d,1,2..
- 2.update w for these

i.e, $\mathbf{w} = \mathbf{w} - \eta 2C \sum_{i=1}^n ([1 - y^i < \mathbf{w}, \mathbf{x}_{j^{th}}^i >]_+) (-y^i \mathbf{x}_{j^{th}}^i)$

Note : η value was choosen very small or else the curve diverges.

c. Method used for D2 maximization

$$D2 \quad \Rightarrow \operatorname{argmax}(\sum_{i=1}^n \alpha_i - 1/4C(\sum_{i=1}^n \alpha_i^2) - 1/2 \sum_{i=1}^n \sum_{j=1}^n (\alpha_i \alpha_j y^i y^j < \mathbf{x}^i \mathbf{x}^j >))$$

concentrate on α_i

$$\operatorname{argmax}(\alpha_i - 1/4C(\alpha_i^2) - (\sum_{i \neq j} \alpha_j^2) - 1/2(\alpha_i^2 (\mathbf{x}^i)^2) - \alpha_i y^i \sum_{i \neq j} (\alpha_j y^j < \mathbf{x}^i \mathbf{x}^j >))$$

$$\text{Let } x = \alpha_i, \quad q = (\mathbf{x}^i)^2, \quad p = y^i \sum_{i \neq j} (\alpha_j y^j < \mathbf{x}^i \mathbf{x}^j >), \quad r = (\sum_{i \neq j} \alpha_j^2)$$

$$\operatorname{argmin}(x^2(1/2q + 1/4C) - x(1 - p) + r)$$

$$\text{Let } q^0/2 = (q/2 + \frac{1}{4C}) \quad \text{argmin}(q^0/2x^2 - x(1-p) + r)$$

minimum at $x' = (1-p)/q^0$

If $x' \in [0, \infty]$ then x' is solution

elseif $x' < 0$ solution is 0

else solution is ∞ .

Note:

$$p = y^i \sum_{i \neq j} (\alpha_j y^j < \mathbf{x}^i \mathbf{x}^j >)$$

$$= \mathbf{w}^T \mathbf{x}^i - \alpha_i y^i q^0$$

Algorithm minimize:

$$\text{initialize } \alpha^T = \{0, 0, \dots, 0\}_{1 \times n} \quad \mathbf{w} = \{0, 0, \dots, 0\}_{1 \times n}$$

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}^i \quad \text{for } i \in [1, n]$$

Note: We have appended one extra dimension to our data i.e if $\mathbf{X} =$

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & . & x_{1n} \\ x_{21} & \vdots & \vdots & \vdots & x_{2n} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ x_{n1} & \vdots & \vdots & \vdots & x_{nn} \end{bmatrix}$$

and then transformed $\mathbf{X} \rightarrow \mathbf{X}^0 =$

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & . & x_{1n} & 1 \\ x_{21} & . & . & . & x_{2n} & 1 \\ . & . & \ddots & \ddots & . & 1 \\ . & . & \ddots & \ddots & . & 1 \\ . & . & \ddots & \ddots & . & 1 \\ x_{n1} & . & . & . & x_{nn} & 1 \end{bmatrix} \quad \text{so that } \mathbf{w}[d-1] = b, \{d \text{ is}$$

dimension of each data point in $\mathbf{X}^0\}$. Iterate till time does not end, for each data point in \mathbf{X}^0

Do-

1. Calculate α_i if $\alpha_i \geq 0$ then $\alpha[i] = \alpha_i$

else $\alpha[i] = 0$

2. Update \mathbf{w} , $\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}^i$

\mathbf{w} will converge after sufficient number of iteration.

4 Question

Hyperparameters used

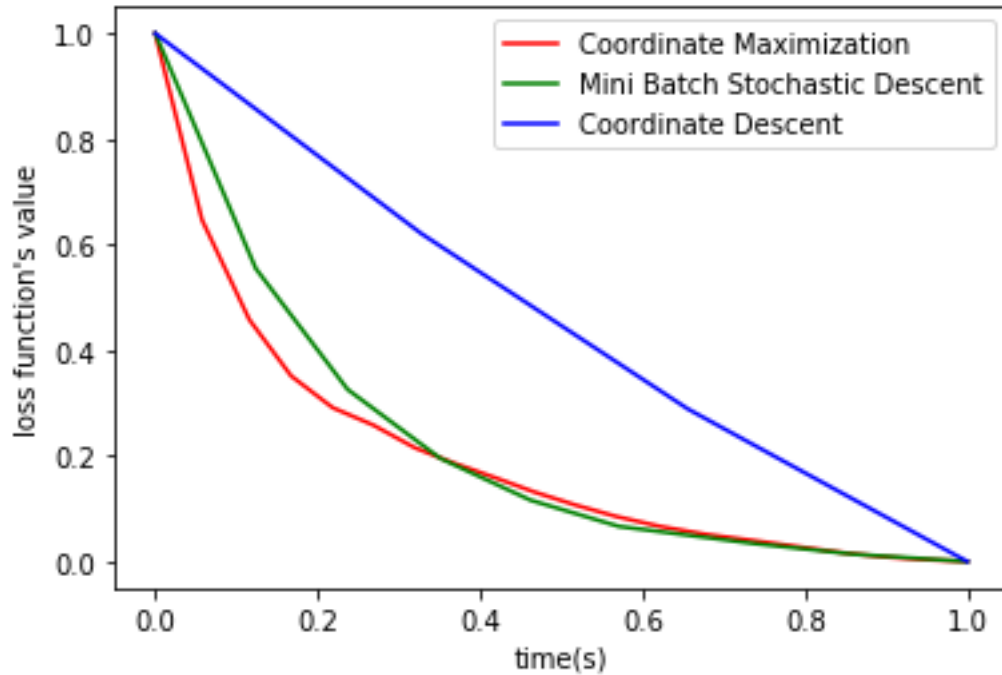
a. Mini Batch SGD on P1: tried different values for Batch size from 10 to 1000 all gives similar results with $\eta = 0.001$ or 0.01 .

b. Coordinate descent on P1: here the step length need to be very small or else the curve diverges. We tried with step length 0.1, 0.01, ... and it kept on diverging. Thus fixed step length to a very small number 0.000001 that lead to convergence of the curve.

c.Coordinate Maximization on D2: No hyperparameter used.

5 Question

Figure: Plot for convergence curves offered by those 3 method



Following graph was normalized using the minimum and maximum value for both the axis,as given below

a. Mini Batch SGD on P1: X axis (time): Min = 0.3228445053100586 Max = 3.0747430324554443

Y axis (loss function value): Min = 5607.35 Max = 6252.87

b. Coordinate Descent on P1: X axis (time): Min = 7.35390830039978 Max = 29.89153790473938

Y axis (loss function value): Min = 13471.7 Max = 18036.3

c. Coordinate Maximization on D2: X axis (time): Min = 0.4362635612487793 Max = 9.31814169883728

Y axis (loss function value): Min = 5236.12 Max = 9059.85

6 Question

Among the above three curves it is clear that Mini Batch Stochastic gradient descent gives the best convergence for the graph as seen in the plot.