

Homework report 1

Name: Ma xiaoqi
grasper: 308932

(the scan of pdf can be enlarged)

Problem 2

Problem 1: (\vec{x} means vector)

For n points of form: $(\vec{a}_1, b_1) \dots (\vec{a}_n, b_n)$

$$\therefore L_{SH}(\vec{x}, \mu) = \frac{1}{n} \sum_{i=1}^n g_i(\vec{x}, \mu)$$

$$g_i(\vec{x}, \mu) = \begin{cases} \frac{1}{2} - b_i(\vec{a}_i^T \vec{x} + \mu) & b_i(\vec{a}_i^T \vec{x} + \mu) \leq 0 \\ \frac{1}{2}(1 - b_i(\vec{a}_i^T \vec{x} + \mu))^2 & 0 < b_i(\vec{a}_i^T \vec{x} + \mu) \leq 1 \\ 0 & 1 \leq b_i(\vec{a}_i^T \vec{x} + \mu) \end{cases} \Rightarrow \begin{cases} I_L(i, i) = 1 \\ I_Q(i, i) = 1 \\ I_O(i, i) = 1 \end{cases}$$

(a) : In problem a, $\mu = 0$ $\tilde{A}\vec{x} = [\underline{b_1\vec{a}_1, \dots, b_n\vec{a}_n}]^T \vec{x}$ $[\vec{x}\vec{x}]_{n \times 1}$

$$\therefore \cancel{L_{SH}(\vec{x})} = \frac{1}{n} \text{ we use } \vec{T}^T = [1, 1, \dots, 1]_{n \times P} \vec{P} \vec{x}$$

$$\therefore L_{SH}(\vec{x}) = \frac{1}{n} \left\{ \frac{1}{2} - \vec{T}^T I_L \tilde{A} \vec{x} + \frac{1}{2} \| I_a \vec{T} - I_a \tilde{A} \vec{x} \|^2 \right\}$$

$$\therefore \nabla L_{SH}(\vec{x}) = \frac{1}{n} \left\{ \nabla \left(-\vec{T}^T I_L \tilde{A} \vec{x} \right) + \frac{1}{2} \nabla \| I_a \vec{T} - I_a \tilde{A} \vec{x} \|^2 \right\}$$

$$= \frac{1}{n} \left\{ -(\tilde{A}^T I_L^T \vec{T}) + \left[-(\tilde{A}^T I_a^T)(I_a \vec{T} - I_a \tilde{A} \vec{x}) \right] \right\}$$

$$= \frac{1}{n} \left\{ \tilde{A}^T I_a^T I_a \tilde{A} \vec{x} - \tilde{A}^T I_a^T I_a \vec{T} - \tilde{A}^T I_L^T \vec{T} \right\} \quad (I_a^T \cdot I_a = I_a)$$

$$= \frac{1}{n} \left\{ \tilde{A}^T I_a^T (\tilde{A} \vec{x} - \vec{T}) - \tilde{A}^T I_L^T \vec{T} \right\}$$

$$\therefore f(\vec{x}) = L_{SH}(\vec{x}) + \frac{\lambda}{2} \|\vec{x}\|^2 \quad \nabla f(\vec{x}) = \nabla L_{SH}(\vec{x}) + \lambda \vec{x}$$

$$\therefore \nabla f(\vec{x}) = \lambda \vec{x} + \frac{1}{n} \tilde{A}^T I_a (\tilde{A} \vec{x} - \vec{T}) - \frac{1}{n} \tilde{A}^T I_L \vec{T} \quad \star$$

If there exist L ($L > 0$), $\|\nabla f(\vec{x}) - \nabla f(\vec{y})\| \leq L \|\vec{x} - \vec{y}\|$

$$\therefore \|\nabla f(\vec{x}) - \nabla f(\vec{y})\| = \|\lambda \vec{x} - \lambda \vec{y} + \frac{1}{n} \tilde{A}^T I_a (\tilde{A} \vec{x} - \vec{T}) + \frac{1}{n} \tilde{A}^T I_L \vec{T}\|$$

$$= \|\lambda(\vec{x} - \vec{y}) + \frac{1}{n} \tilde{A}^T \tilde{A} (\vec{x} - \vec{y})\| \leq \lambda \|\vec{x} - \vec{y}\| + \frac{1}{n} \|\tilde{A}^T \tilde{A} (\vec{x} - \vec{y})\| \\ \leq [\lambda + \frac{1}{n} \|\tilde{A}^T\| \|A\|] \cdot \|\vec{x} - \vec{y}\|$$

$$\therefore \|A\| = \|\tilde{A}\| \quad \|A^T\| = \|\tilde{A}^T\|$$

$\therefore L = \lambda + \frac{1}{n} \|A^T\| \|A\|$ and f is L -Lipschitz

$$(b) \quad \nabla f(\vec{x}) = \lambda \vec{x} + \frac{1}{n} \tilde{A}^T I_a (\tilde{A} \vec{x} - \vec{T}) - \frac{1}{n} \tilde{A}^T I_L \vec{T}$$

$\therefore I_a = \mathbb{I}$, cut the constant term.

~~$\nabla^2 f(\vec{x}) = \nabla(\lambda \vec{x} + \frac{1}{n} \tilde{A}^T \mathbb{I} \tilde{A} \vec{x})$~~

$$\text{for the part } \nabla(\lambda \vec{x}) \quad \lambda \vec{x} = \begin{bmatrix} \lambda & & & \\ & \ddots & & \\ & & \lambda & \\ 0 & & & \ddots \end{bmatrix} \text{ diagonal}(\lambda) \quad \vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}_{P \times P} \quad \therefore \nabla(\lambda \vec{x}) = \lambda \mathbb{I} \vec{x} = \lambda \mathbb{I} \vec{x}$$

$$\text{for the part } \nabla(\frac{1}{n} \tilde{A}^T \mathbb{I} \tilde{A} \vec{x}) = \nabla(\frac{1}{n} \tilde{A}^T \tilde{A} \vec{x}) = \frac{1}{n} (\tilde{A} \tilde{A}^T)^T = \frac{1}{n} \tilde{A}^T (\tilde{A}^T)^T = \frac{1}{n} \tilde{A}^T \tilde{A}$$

~~$\therefore \tilde{A}^T \tilde{A} = A^T A$~~ $\because b \in \{1, -1\} \quad \therefore [b_1 \vec{a}_1, b_2 \vec{a}_2, \dots, b_n \vec{a}_n]^T$

~~$\therefore b_1 \cdot b_2 \cdots b_n = 1 \quad \therefore \tilde{A}^T \cdot \tilde{A} = A^T A$~~

$$\therefore \nabla^2 f(\vec{x}) = \lambda \mathbb{I} + \frac{1}{n} A^T A$$

$\nabla^2 f(\vec{x})$ does not depend on \vec{x} and f is twice-differentiable at \vec{x}

~~or fix~~ because the partial derivations are continuous at x
 (∇f)

(c) $f(\vec{x}) = LSH(\vec{x}, \mu) + \frac{\lambda}{2} \|\vec{x}\|^2$

So we need to improve the $LH(\vec{x})$ is convex. ($\mu=0$)

$$\nabla LH(\vec{x}) = \frac{1}{n} \tilde{A}^T I_Q (\tilde{A} \vec{x} - \vec{1}) - \frac{1}{n} \tilde{A}^T I_L \vec{1}$$

$$\nabla^2 LH(\vec{x}) = \frac{1}{n} (\tilde{A}^T I_Q \tilde{A})^T = \frac{1}{n} \tilde{A}^T \cdot (I_Q)^T (\tilde{A}^T)^T = \frac{1}{n} \tilde{A}^T I_Q \tilde{A}$$

This is the Hessian Matrix of ~~LH~~ $LH(\vec{x})$

So we need to prove Matrix $\tilde{A}^T I_Q \tilde{A}$ is positive semi-definite.

We set the Matrix $H = \tilde{A}^T I_Q \tilde{A}$

Then, ~~set the~~ For any vector \vec{P}, \vec{P}^T

$$\cancel{\vec{P} \cdot H \cdot \vec{P}^T} \quad \vec{P}^T H \vec{P} = \vec{P}^T \tilde{A}^T I_Q \tilde{A} \vec{P} = \|\vec{P}^T I_Q \tilde{A} \vec{P}\|^2 \geq 0$$

$\therefore H$ is the positive ~~definite~~ semi-definite

$LH(\vec{x})$ is convex

$f(\vec{x})$ is λ -strongly convex

(a) GD and GDstr

This problem asks us to implement the GD and GDstr and compare some difference between them. GDstr change the parameter alpha (step-size), and we can get the outcomes as following:

GD:

Set the max iteration is 4000, after 3995 iters,

$$f(x) = 0.039934508$$

Error w.r.t 0-1 loss: 0.029197080291970802

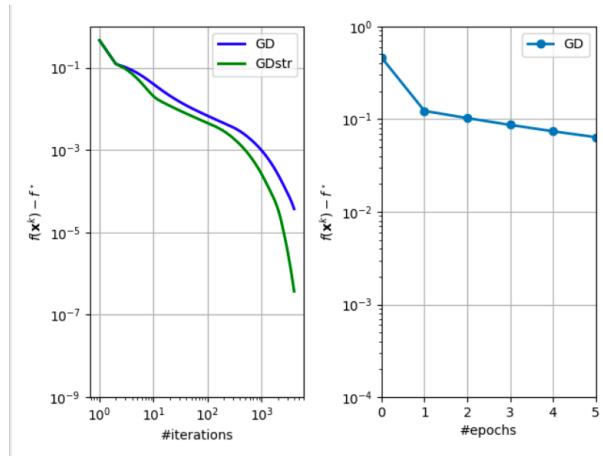
GDstr

Set the max iteration is 4000, after 3995 iters,

$$f(x) = 0.039897574$$

Error w.r.t 0-1 loss: 0.029197080291970802

Although the error is the same, we can find the GDstr converges faster than GD because the existence of strong convexity (from figure)



(b) AGD and AGDstr

Implementing the accelerate gradient descent method is to make a faster convergence than gradient decent, I get outcomes as following:

AGD:

Set the max iteration is 2300, after 2295 iters,

$$f(x) = 0.039897205$$

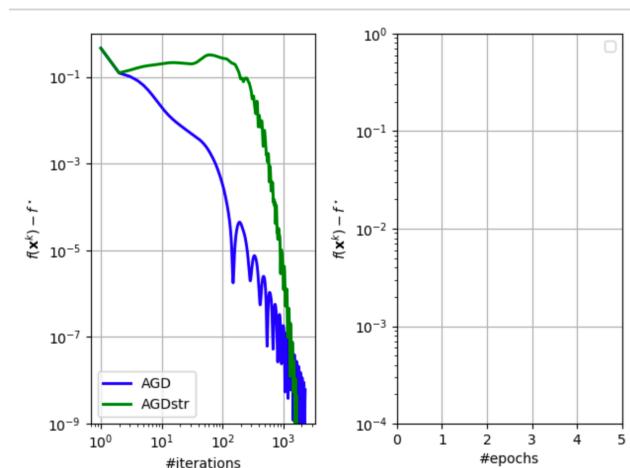
Error w.r.t 0-1 loss: 0.029197080291970802

AGDstr:

Set the max iteration is 1700, after 1696,

$$f(x) = 0.039897203$$

Error w.r.t 0-1 loss: 0.029197080291970802



However, after

about 10^{12} iterations, I

find the AGD algorithm performs oscillation by a wide margin, AGDstr also has oscillation but the frequency is low. At the same time, AGDstr just use less iterations than AGD to achieve the same performance.

(c) LSGD and LSAGD

Then, implement a line search procedure which aims at deterring the stepsize of alpha, I get the outcomes as following:

LSGD:

Set the max iteration is 400, after 395 iters,

$$f(x) = 0.039897202$$

Error w.r.t 0-1 loss: 0.029197080291970802

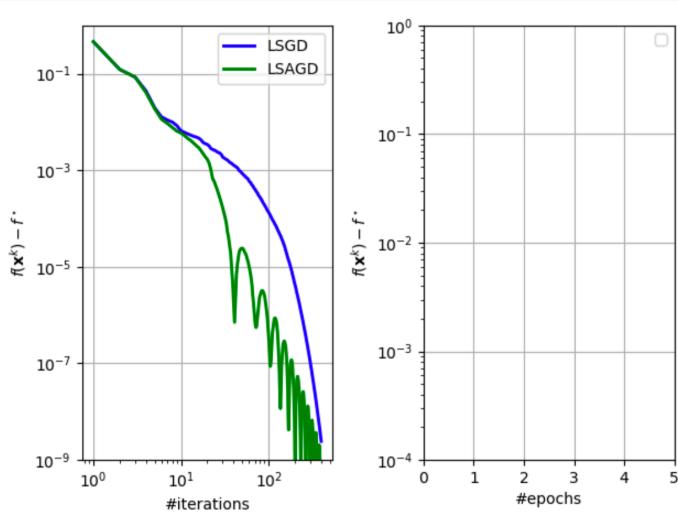
LSAGD:

Set the max iteration is 400, after 395 iters,

$$f(x) = 0.039897199$$

Error w.r.t 0-1 loss: 0.029197080291970802

I find the LSAGD has the faster convergence than LSGD, but LSAGD performs oscillation after 50 iterations, so we need improvement.



(d) AGDR and LSAGDR

In order to solve the problem that the the accelerated gradient descent can be oscillatory, implement the adaptive restart strategy. I get the outcomes as following:

AGDR: Set the max iteration is 300, after 295 iters,

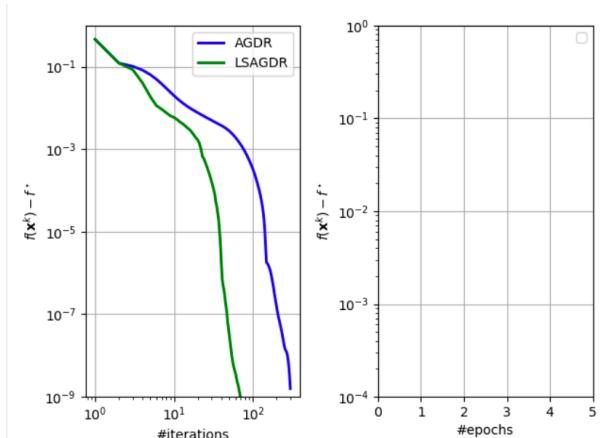
$$f(x) = 0.039897201$$

Error w.r.t 0-1 loss: 0.029197080291970802

LSAGDR: Set the max iteration is 100, after 295 iterations

$$f(x) = 0.039897199$$

Error w.r.t 0-1 loss: 0.029197080291970802



Compared with the AGD, the phenomenon of oscillation has disappeared and LSAGDR converges faster than AGDR

(e)(f)AdaGrad and ADAM

I implement the adaptive optimization method as following:

AdaGrad: Set the max iteration is 4000, after 3995 iterations

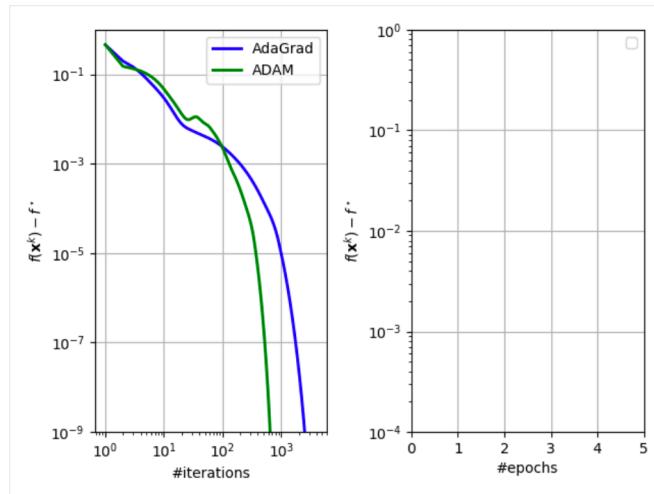
$$f(x) = 0.039897199$$

Error w.r.t 0-1 loss: 0.029197080291970802

ADAM: Set the max iteration is 4000, after 3995 iterations

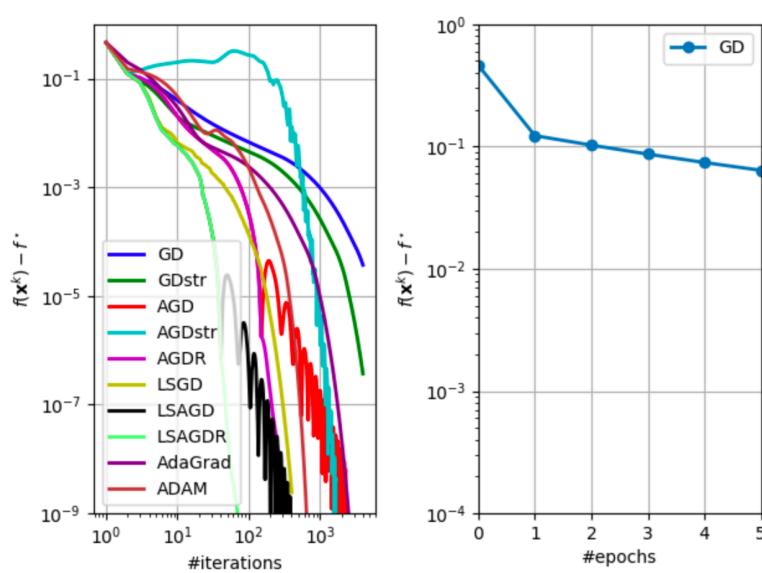
$$f(x) = 0.039897199$$

Error w.r.t 0-1 loss: 0.029197080291970802



ADAM is better than AdaGrad

Plot about program 2



Problem 3

$$(a) \text{① } \nabla f(\vec{x}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\vec{x})$$

$$E_{\nu_k} [\nabla f_{i_k}(\vec{x})] = \sum_{i=1}^n P(i_k=i) \nabla f_i(\vec{x}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\vec{x}) = \nabla f(\vec{x})$$

$\therefore \nabla f_{i_k}(\vec{x})$ is an unbiased estimation of $\nabla f(\vec{x})$

$$\text{② } \nabla f_{i_k}(\vec{x}) = \lambda \vec{x} + \vec{1}_{\{0 < b_i \vec{a}_i^\top \vec{x} \leq 0\}} \vec{a}_i (\vec{a}_i^\top \vec{x} - b_i) - \vec{1}_{\{b_i \vec{a}_i^\top \vec{x} \leq 0\}} b_i \vec{a}_i$$

$$\therefore \| \nabla f_{i_k}(\vec{x}) - \nabla f_{i_k}(\vec{y}) \| = \| \lambda(\vec{x} - \vec{y}) + \vec{a}_{i_k} (\vec{a}_{i_k}^\top \vec{x} - b_{i_k}) - \vec{a}_{i_k} (\vec{a}_{i_k}^\top \vec{y} - b_{i_k}) \|$$

$$= \| \lambda(\vec{x} - \vec{y}) + \vec{a}_{i_k} \vec{a}_{i_k}^\top (\vec{x} - \vec{y}) \|$$

$$\leq \lambda \| \vec{x} - \vec{y} \| + \| \vec{a}_{i_k} \vec{a}_{i_k}^\top \| \cdot \| \vec{x} - \vec{y} \|$$

$$\leq [\lambda + \| \vec{a}_{i_k} \vec{a}_{i_k}^\top \|^2] \cdot \| \vec{x} - \vec{y} \|$$

$\therefore L = \lambda + \| \vec{a}_{i_k} \vec{a}_{i_k}^\top \|^2$ and ∇f_{i_k} is Lipschitz continuous

Program 3(4)

(b)(c)(d) stochastic gradient method

In the stochastic gradient method, I use a random number uniformly, so for SGD, I always get the different outcomes about SGD

SGD (one experiment):

$$f(x) = 0.109049377$$

Error w.r.t 0-1 loss: 0.08029197080291971

SAG:

$$f(x) = 0.047288051$$

Error w.r.t 0-1 loss: 0.02919708029197082

SVR:

$$f(x) = 0.039897199$$

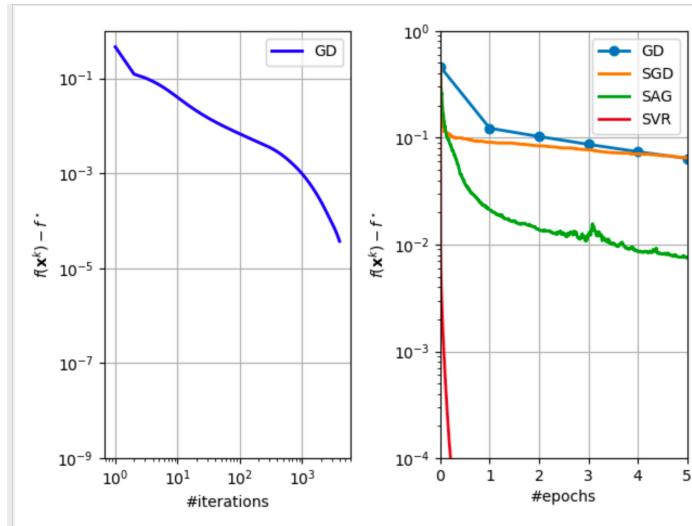
Error w.r.t 0-1 loss: 0.02919708029197082

SGD converges faster than GD but the value of function almost had no change with the change of epochs.

SAG use the average of previous gradients and its convergence rate is bigger than SAG and GD

SVR converges very fast, almost reach as a linear converge rate, so I find the effect of reducing the variance on convergence rate is obvious.

To sum up, SGD choose a single sample randomly as the sample population, so, the error sometimes become larger but when we are dealing with large amounts of data, SGD can help us get better performance.



Final results:

