

Problem 1. μ -strong convexity: for any $x, y \in X$ and $t \in [0, 1]$, $\exists \mu > 0$ s.t.

$$f(tx + (1-t)y) \leq t \cdot f(x) + (1-t)f(y) - \frac{\mu}{2} t(1-t) \|x-y\|^2 \dots \textcircled{1}$$

C1 $f(y) \geq f(x) + \nabla f(x)^T(y-x) + \frac{\mu}{2} \|y-x\|^2$

C2. $(\nabla f(x) - \nabla f(y))^T(x-y) \geq \mu \|x-y\|^2$

C3. $\nabla^2 f(x) - \mu I > 0$

$\textcircled{1} \Leftrightarrow \text{C1}$. $f(tx + (1-t)y) \leq t f(x) + f(y) - t f(y) - \frac{\mu}{2} t(1-t) \|x-y\|^2$

(calculate $\frac{d}{dt}|_{t=0}$) LHS: $\frac{d}{dt} f(tx + (1-t)y)|_{t=0} = \nabla f(y)^T(x-y)$

RHS: $\frac{d}{dt} [t f(x) + (1-t)f(y) - \frac{\mu}{2} t(1-t) \|x-y\|^2]|_{t=0} = [f(x) - f(y) - (1-t) \frac{\mu}{2} \|x-y\|^2]|_{t=0}$
 $= f(x) - f(y) - \frac{\mu}{2} \|x-y\|^2$

$\Leftrightarrow \nabla f(y)^T(x-y) \leq f(x) - f(y) - \frac{\mu}{2} \|x-y\|^2$

Rearrange Info. $\Leftrightarrow f(x) \geq f(y) + \nabla f(y)^T(x-y) + \frac{\mu}{2} \|x-y\|^2$

Interchange $x, y \Leftrightarrow f(y) \geq f(x) + \nabla f(x)^T(y-x) + \frac{\mu}{2} \|y-x\|^2$

$\text{C1} \Leftrightarrow \text{C2}$.
 $(\text{C1}) \begin{cases} f(x) \geq f(y) + \nabla f(y)^T(x-y) + \frac{\mu}{2} \|x-y\|^2 \dots a. \\ f(y) \geq f(x) + \nabla f(x)^T(y-x) + \frac{\mu}{2} \|y-x\|^2 \dots b. \end{cases}$

$(a. + b.) \Leftrightarrow \cancel{f(x)} + \cancel{f(x)} \geq \cancel{f(x)} + \cancel{f(x)} + f(y)^T(x-y) + \nabla f(x)^T(y-x) + \mu \|x-y\|^2$

(Rearrange) $\Leftrightarrow (\nabla f(x) - \nabla f(y))^T(x-y) \geq \mu \|x-y\|^2$

$\text{C1} \Leftrightarrow \text{C3}$. (C1) $f(y) \geq f(x) + \nabla f(x)^T(y-x) + \frac{\mu}{2} \|y-x\|^2 \dots c.$

(Taylor) $f(y) = f(x) + \nabla f(x)^T(y-x) + \frac{1}{2}(y-x)^T \nabla^2 f(x)(y-x) + o(\|y-x\|^2) \dots d.$

(by c, d) $\Leftrightarrow \frac{1}{2}(y-x)^T \nabla^2 f(x)(y-x) + o(\|y-x\|^2) \geq \frac{\mu}{2} \|y-x\|^2$

$(\times \frac{2}{\|y-x\|^2} \text{ on both side}) \Leftrightarrow \frac{(y-x)^T}{\|y-x\|} \nabla^2 f(x) \frac{(y-x)}{\|y-x\|} + o(1) \geq \mu$

$\Leftrightarrow \nabla^2 f(x) \succeq \mu I$

Problem 2.

$$f(x) = \begin{cases} \geq 5x^2 & x < 1 \dots f_1(x) & f_1(1) = 25 = f_1(1) & f_1'(x) = 50x & f_1'(1) = 50 = f_1'(1) & f_1'' = 50 \\ x^2 + 48x - 24 & 1 \leq x \leq 2 \dots f_2(x) & f_2(2) = 76 = f_2(2) & f_2'(x) = 2x + 48 & f_2'(2) = 50 = f_2'(2) & f_2'' = 2 \\ \geq 5x^2 - 48x + 72 & 2 < x \dots f_3(x) & f_3(2) = 76 = f_3(2) & f_3'(x) = 50x - 48 & f_3'(2) = 52 = f_3'(2) & f_3'' = 50 \end{cases}$$

Since the function pieces are continuous and differentiable at break points, so $f(x)$ is continuous and differentiable.

2. a. 2-strongly convex:

for $x < 1$, $\nabla^2 f(x) = 50 \geq 2$, for $1 \leq x \leq 2$, $\nabla^2 f(x) = 2 \geq 2$, for $x > 2$, $\nabla^2 f(x) = 50 \geq 2$
 else, check $(\nabla f(x) - \nabla f(y))^T(x-y) \geq \mu \|x-y\|^2$

Since $x < y$
↓

$$\text{for } x < 1, 1 \leq y \leq 2: (50x - (2y + 48))(x - y) = (50x - 2y - 48)(x - y) \stackrel{?}{\geq} \|x - y\|$$

$$\Rightarrow 50x - 2y - 48 \leq 2x - 2y \text{ holds for } x < 1 \neq$$

$$\text{for } x < 1, 2 < y: (50x - 50y + 48)(x - y) \stackrel{?}{\geq} \|x - y\|$$

$$\Rightarrow 50x - 50y + 48 \leq 2x - 2y$$

$$\Rightarrow 48x + 48 \leq 48y \text{ holds for } x < 1, y > 2 \neq$$

$$\text{for } 1 \leq x \leq 2, y > 2: (2x + 48 - 50y + 48)(x - y) \stackrel{?}{\geq} \|x - y\|$$

$$\Rightarrow 2x - 50y + 96 \leq 2x - 2y$$

$$\Rightarrow 48y \geq 96 \text{ holds for } y > 2 \neq$$

By all the cases above, $f(x)$ is 2 -strongly convex.

$$L\text{-smooth: } \nabla^2 f(x) \preceq LI \quad \cdot \quad \|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$$

$$\text{for } x < 1, \nabla^2 f(x) = 50 \leq 50, \text{ for } 1 \leq x \leq 2, \nabla^2 f(x) = 2 \leq 50, \text{ for } x > 2, \nabla^2 f(x) = 50 \leq 50$$

$$\text{for } x < 1, 1 \leq y \leq 2: \|50x - 2y - 48\| \leq 50\|x - y\|$$

$$\Rightarrow 50x - 50y \leq 50x - 2y - 48 \leq 50y - 50x$$

$$\Rightarrow \begin{cases} 50x - 50y \leq 50x - 2y - 48 \text{ holds for } y \geq 1 \\ 50y - 50x \geq 50x - 2y - 48 \end{cases}$$

$$\Rightarrow 100x \leq 52y + 48 \text{ holds for } x < 1, 1 \leq y \leq 2 \neq$$

$$\text{for } x < 1, y > 2: \|50x - 50y + 48\| \leq 50\|x - y\|$$

$$\Rightarrow \begin{cases} 50x - 50y + 48 \geq 50x - 50y \text{ always holds} \\ 50x - 50y + 48 \leq 50y - 50x \end{cases}$$

$$\Rightarrow 100x + 48 \leq 100y \text{ holds for } x < 1, y > 2 \neq$$

$$\text{for } 1 \leq x \leq 2, y > 2: \|2x + 48 - 50y + 48\| \leq 50\|x - y\|$$

$$\Rightarrow \begin{cases} 2x - 50y + 96 \geq 50x - 50y \text{ holds for } 1 \leq x \leq 2 \\ 2x - 50y + 96 \leq 50y - 50x \end{cases}$$

$$\Rightarrow 52x + 96 \leq 100y \text{ holds for } 1 \leq x \leq 2, y > 2 \neq$$

By all the six cases above, we know that $f(x)$ is L -smooth with $L = 50 \neq$

2.b. Since $f(x)$ is convex, so x with $f'(x) = 0$ implies it is a global optimizer.

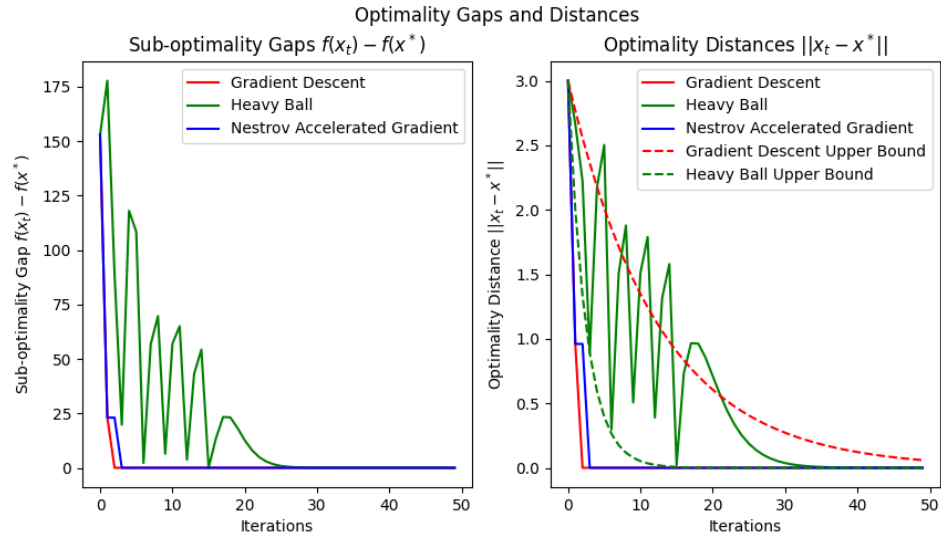
$$\text{for } x < 1, f'(x) = 50x = 0 \Leftrightarrow x = 0 \quad (\checkmark)$$

$$\text{for } 1 \leq x \leq 2, f'(x) = 2x + 48 = 0 \Leftrightarrow x = -24 \quad (x \notin [1, 2])$$

$$\text{for } x > 2, f'(x) = 50x - 48 = 0 \Leftrightarrow x = 0.96 \quad (x \notin (2, \infty))$$

So the global optimizer $x^* = 0$ with $f(x^*) = 0 \neq$

2. (c) Gradient Descent method performs a really small optimality distance compares to the upper bound, while heavy ball method seems always swinging above the upper bound. I think it is because this function is not really a quadratic function.



Problem 3.a. $f(x) = f(y) + \nabla f(y)^T(x-y)$ by convexity ... ①

$$f(y - \frac{1}{2L} \nabla f(y)) - f(y) \leq \nabla f(y)^T(-\frac{1}{2L} \nabla f(y)) + \frac{1}{2} \|-\frac{1}{2L} \nabla f(y)\|^2 \text{ by } L\text{-smoothness}$$

$$= -\frac{1}{2L} \|\nabla f(y)\|^2 \dots ②$$

By ① + ②: $f(y - \frac{1}{2L} \nabla f(y)) - f(x) \leq -\frac{1}{2L} \|\nabla f(y)\|^2 - \nabla f(y)^T(x-y) \dots ③$

3.b. $x_{t+1} = y_t - \frac{1}{2L} \nabla f(y_t) \Leftrightarrow f(y_t - \frac{1}{2L} \nabla f(y_t)) = f(x_{t+1}), \nabla f(y_t) = -L(x_{t+1} - y_t)$

Let $y = y_t, x = x_t$ in ③: $f(y_t - \frac{1}{2L} \nabla f(y_t)) - f(x_t) = f(x_{t+1}) - f(x_t)$

$$\leq -\frac{1}{2L} \|L(x_{t+1} - y_t)\|^2 + L(x_{t+1} - y_t)^T(x_t - y_t)$$

$$= -\frac{1}{2} \|x_{t+1} - y_t\|^2 + L(x_{t+1} - y_t)^T(x_t - y_t) \dots ④$$

Similarly, Let $y = y_t, x = x^*$ in ③, we can get

$$f(x_{t+1}) - f(x^*) \leq -\frac{1}{2} \|x_{t+1} - y_t\|^2 + L(x_{t+1} - y_t)^T(x^* - y_t) \dots ⑤$$

3.c. $\theta_t(\theta_t - 1) \times ③ + \theta_t \times ⑤$ LHS

$$= \theta_t^2(f(x_{t+1}) - f(x_t)) + \theta_t(f(x_{t+1}) - f(x^*))$$

$$= \theta_t^2 f(x_{t+1}) - \theta_t^2 f(x_t) + \theta_t f(x_{t+1}) - \theta_t f(x^*)$$

(by $\theta_t^2 - \theta_t - \theta_{t+1} = 0$)

$$= \theta_t^2 f(x_{t+1}) - \theta_t^2 f(x_t) + (\theta_t^2 - \theta_{t+1}) f(x_{t+1}) - (\theta_t^2 - \theta_{t+1}) f(x^*)$$

$$= \theta_t^2 (f(x_{t+1}) - f(x^*)) - \theta_{t+1} (f(x_t) - f(x^*))$$

(Let $\Delta_t := f(x_t) - f(x^*)$)

$$= \theta_t^2 \Delta_{t+1} - \theta_{t+1} \Delta_t$$
RHS

$$\leq -\frac{1}{2} \theta_t^2 \|x_{t+1} - y_t\|^2 + L(x_{t+1} - y_t)^T(\theta_t(\theta_t - 1)(x_t - y_t) + \theta_t(x^* - y_t))$$

$$= -\frac{1}{2} \theta_t^2 \|x_{t+1} - y_t\|^2 + \frac{1}{2} \theta_t^2 (x_{t+1} - y_t)^T((\theta_t - 1)x_t - \theta_t y_t + x^*)$$

$$= -\frac{1}{2} (\|\theta_t(x_{t+1} - y_t)\|^2 + 2\theta_t(x_{t+1} - y_t)^T(\theta_t y_t - (\theta_t - 1)x_t - x^*))$$

$$\Rightarrow \theta_t^2 \Delta_{t+1} - \theta_{t+1} \Delta_t \leq -\frac{1}{2} (\|\theta_t(x_{t+1} - y_t)\|^2 + 2\theta_t(x_{t+1} - y_t)^T(\theta_t y_t - (\theta_t - 1)x_t - x^*)) \dots ⑥$$

3.d. Complete the square in ⑥: $\|\theta_t(x_{t+1} - y_t)\|^2 + 2\theta_t(x_{t+1} - y_t)^T(\theta_t y_t - (\theta_t - 1)x_t - x^*)$

(Let $\varnothing_t := \theta_t y_t - (\theta_t - 1)x_t - x^*$)

$$= \|\theta_t(x_{t+1} - y_t)\|^2 + 2\theta_t(x_{t+1} - y_t)^T \varnothing_t + \|\varnothing_t\|^2 - \|\varnothing_t\|^2$$

$$= \|\theta_t(x_{t+1} - y_t) + \varnothing_t\|^2 - \|\varnothing_t\|^2$$

$$= \|\theta_t x_{t+1} - \theta_t y_t + \theta_t y_t - (\theta_t - 1)x_t - x^*\|^2 - \|\varnothing_t\|^2 \dots ⑦$$

Note that by the update rule (b.):

$$y_{t+1} = x_{t+1} - \frac{(1 - \theta_t)}{\theta_{t+1}} (x_{t+1} - x_t)$$

both $\times \theta_{t+1} \Rightarrow \theta_{t+1} y_{t+1} = \theta_{t+1} x_{t+1} + (\theta_t - 1)x_{t+1} - (\theta_t - 1)x_t$

$$\Rightarrow -(\theta_t - 1)x_t + \theta_t x_{t+1} = \theta_{t+1} y_{t+1} - (\theta_{t+1} - 1)x_{t+1} \dots ⑧$$

Plug ⑧ into ⑦:

$$⑦ = \|\theta_t x_{t+1} - (\theta_t - 1)x_t - x^*\|^2 - \|\varnothing_t\|^2$$

$$= \|\theta_{t+1} y_{t+1} - (\theta_{t+1} - 1)x_{t+1} - x^*\|^2 - \|\varnothing_t\|^2$$

$$= \|\varnothing_{t+1}\|^2 - \|\varnothing_t\|^2 \dots ⑨$$

Therefore, by ①, ②, ③, ④, we can get:

$$\theta_t^2 \Delta_{t+1} - \theta_{t+1}^2 \Delta_t \leq -\frac{1}{2} (\|\phi_{t+1}\|^2 - \|\phi_t\|^2) = \frac{1}{2} (\|\phi_t\|^2 - \|\phi_{t+1}\|^2) \quad \# \dots ⑤$$

3.e. Telescope summing ⑤ from $t=1$ to $T-1$

$$t=1: \theta_1^2 \Delta_2 - \theta_2^2 \Delta_1 \leq \frac{1}{2} (\|\phi_1\|^2 - \|\phi_2\|^2)$$

$$t=2: \theta_2^2 \Delta_3 - \theta_3^2 \Delta_2 \leq \frac{1}{2} (\|\phi_2\|^2 - \|\phi_3\|^2)$$

\vdots

$$+) t=T-1: \theta_{T-1}^2 \Delta_T - \theta_T^2 \Delta_{T-1} \leq \frac{1}{2} (\|\phi_{T-1}\|^2 - \|\phi_T\|^2)$$

$$\Rightarrow \theta_{T-1}^2 \Delta_T - \theta_0^2 \Delta_1 \leq \frac{1}{2} (\|\phi_1\|^2 - \|\phi_T\|^2) \leq \frac{1}{2} \|\phi_1\|^2$$

$$\Rightarrow \theta_{T-1}^2 \Delta_T \leq \frac{1}{2} \|\phi_1\|^2 + \theta_0^2 \Delta_1 \quad \theta_0=0$$

$$= \frac{1}{2} \|\theta_1 y_1 - (\theta_1 - 1) x_1 - x^*\|^2$$

$$= \frac{1}{2} \|x_1 - x^* - \theta_1 x_1 + \theta_1 x_1\|^2 \quad x_1 = y_1$$

$$\Rightarrow \Delta_T \leq \frac{1}{2\theta_{T-1}^2} \|x_1 - x^*\|^2 \dots ⑥$$

To Prove: $\theta_{t-1} \geq \frac{t}{2}$

Case $t=1$, $\theta_0=0 \geq 0$

Suppose it holds up to $t=n-1$, $\theta_{n-2} \geq \frac{n-1}{2}$

Case $t=n$, $\theta_{n-1} = \frac{1}{2} (1 + \sqrt{1 + 4\theta_{n-2}^2})$

$$\geq \frac{1}{2} (1 + \sqrt{1 + (n-1)^2}) \quad \text{by hypothesis}$$

$$> \frac{1}{2} (1 + \sqrt{n^2 - 2n + 1})$$

$$= \frac{1}{2} (1 + n - 1) = \frac{n}{2} \quad \#$$

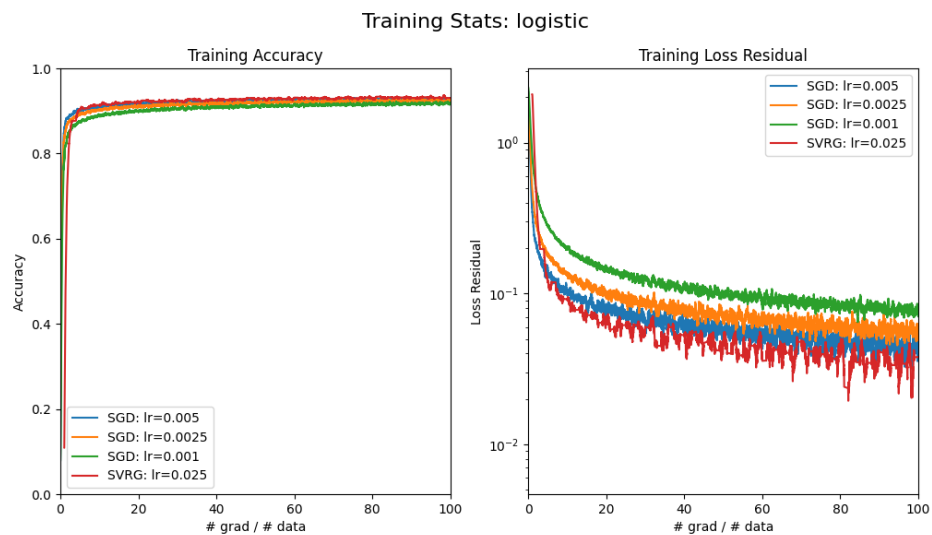
by induction, $\theta_{t-1} \geq \frac{t}{2}$ holds for $t \geq 1 \dots ⑦$

$$\text{by ③ and ④, } \Delta_T = f(x_t) - f(x^*) \leq \frac{1}{2\theta_{T-1}^2} \|x_1 - x^*\|^2 \leq \frac{1}{2} \left(\frac{2}{T}\right)^2 \|x_1 - x^*\|^2 = \frac{2}{T^2} \|x_1 - x^*\|^2 \quad \#$$

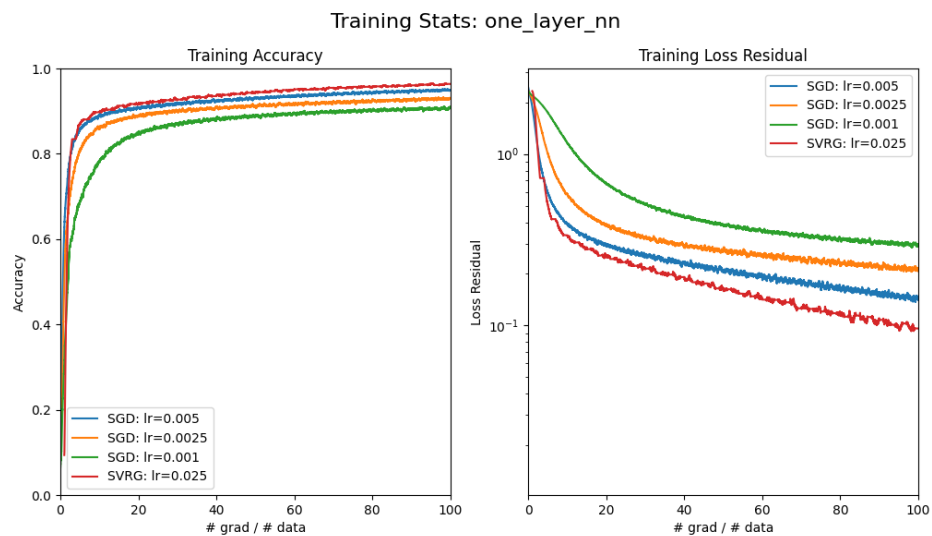
Problem 4.

μ	$f(x^*)$	x^*			
0	-0.12	1	0	0	0
0.1	-0.104694	$6.38 * 10^{-1}$	$3.61 * 10^{-1}$	$2.76 * 10^{-9}$	$2.02 * 10^{-10}$
1.0	-0.0715251	$2.36 * 10^{-1}$	$3.97 * 10^{-1}$	$3.66 * 10^{-1}$	$1.06 * 10^{-9}$
2.0	-0.0524063	0.16754123	0.31296852	0.35757121	0.16191904
5.0	-0.0232699	0.09715142	0.19490255	0.24962519	0.45832084
10.0	0.0126792	0.07368816	0.15554723	0.21364318	0.55712144

Problem 5. (a) logistic: SVRG loss outperforms SGD a sometimes.



Problem 5. (b) one-layer NN: SVRG outperforms SGD (in loss residual) more than pervious experiment.



Hyperparameters:

Batch size: 64

Learning rate: SGD = [0.005, 0.0025, 0.001], SVRG = 0.025, GD = 0.05

Iterations: 100

Logistic primal optimal (loss=0.21524): GD for 1000 iterations

NN primal optimal (loss=0.034793): SGD with lr = 0.005 for 1000 iterations

NN hidden layer: 128 nodes