Optimization in Machine Learning HWS 2024

For the exercise class on the 17.10.2023 at 12:00. Hand in your solutions by 10:15 in the lecture on Tuesday 15.10.2024.

Exercise 1 (Convergence Speed).

Proof that

(i) if we have

then $e(x_k)$ converges super-linearly.

Solution. We define $c_n := \sup_{k \ge n} \frac{e(x_{k+1})}{e(x_k)}$. Then

Thus we have super-linear convergence.

(ii) If for $c \in (0, 1)$ we have

then $e(x_k)$ converges linearly with rate c.

Solution. We again define $c_n := \sup_{k \ge n} \frac{e(x_{k+1})}{e(x_k)}$

$$\lim_{n \to \infty} c_n = \limsup_{k \to \infty} \frac{e(x_{k+1})}{e(x_k)} < c$$

thus there exists $N \ge 0$ such that for all $n \ge N$ we have $c_n \le c$ and therefore for all $n \ge N$

$$e(x_{n+1}) \le c_n e(x_n) \le c e(x_n).$$

(iii) If for $c \in (0, 1)$ we have

 $\limsup_{k \to \infty} \frac{e(x_{k+1})}{e(x_k)^2} < c,$

then $e(x_k)$ converges super-linearly with rate c.

Solution Sheet 3

 $\limsup_{k \to \infty} \frac{e(x_{k+1})}{e(x_k)} = 0,$

 $\lim_{n \to \infty} c_n = \limsup_{k \to \infty} \frac{e(x_{k+1})}{e(x_k)} = 0$

 $e(x_{k+1}) \le c_k e(x_k).$

 $\limsup_{k \to \infty} \frac{e(x_{k+1})}{e(x_k)} < c,$

(1 pt)

(3 Points)

(1 pt)

(1 pt)

Solution. We similarly define $c_n := \sup_{k \ge n} \frac{e(x_{k+1})}{e(x_k)^2}$ and again get $\lim_{n \to \infty} c_n < c$. Thus there exists $N \ge 0$ such that for all $n \ge K$ we have $c_n \le c$ and therefore for all $n \ge N$

$$e(x_{n+1}) \le c_n e(x_n)^2 \le c e(x_n)^2.$$

Exercise 2 (Sub-gradients).

Let $f, q : \mathbb{R}^d \to \mathbb{R}$ be convex functions.

(i) Prove that $\partial f(x)$ is a convex set for any $x \in \mathbb{R}^d$. (1 pt)

Solution. Let $g_1, g_2 \in \partial f(x)$. Then for any $\lambda \in [0, 1]$ and any $y \in \mathbb{R}^d$

$$f(y) = \lambda f(y) + (1 - \lambda) f(y)$$

$$\stackrel{g_1, g_2 \in \partial f(x)}{\geq} \lambda \Big(f(x) + \langle g_1, y - x \rangle \Big) + (1 - \lambda) \Big(f(x) + \langle g_2, y - x \rangle \Big)$$

$$= f(x) + \Big\langle \lambda g_1 + (1 - \lambda) g_2, y - x \Big\rangle$$

thus $\lambda g_1 + (1 - \lambda)g_2 \in \partial f(x)$ by definition.

(ii) For h(x) = f(Ax + b) prove $\partial h(x) \supseteq A^T \partial f(Ax + b)$. Prove equality for invertible A. (1 pt)

Solution. Let
$$g_x \in \partial f(x)$$
 i.e. $g_{Ax+b} \in \partial f(Ax+b)$. Then $A^T g_{Ax+b} \in \partial h(x)$ because

$$h(y) = f(Ay+b) \ge f(Ax+b) + \langle g_{Ax+b}, (Ay+b) - (Ax+b) \rangle$$

$$= h(x) + \langle A^T g_{Ax+b}, y - x \rangle.$$

If A is invertible, we have $f(x) = h(A^{-1}x - A^{-1}b)$ so by the previous statement with $\tilde{A} = A^{-1}$ and $\tilde{b} = -A^{-1}b$, we get the other direction.

Exercise 3 (Lasso).

Let

$$f(x) = \frac{1}{2} \|x - y\|^2 + \lambda \|x\|_1$$

for $x \in \mathbb{R}^d$ be the Lagrangian form of the least squares LASSO method.

(i) Compute a sub-gradient of f.

Solution. Using $\partial(g + \lambda h)(x) \supseteq \partial g(x) + \lambda \partial h(x)$, we only need to determine the subgradient of $g(x) := \frac{1}{2} ||x - y||^2$ and

$$h(x) := \|x\|_1 = \sum_{i=1}^{a} |x_i|$$

But $\nabla g(x) = x - y$ as g is differentiable. And since it is also convex, we have

$$\partial g(x) = \{\nabla g(x)\}\$$

(6 Points)

(2 pts)

(2 Points)

Now the subgradient of $h_i(x) = |x_i|$ is given by $sgn(x_i)e_i$, where $sgn(0) \in [-1, 1]$ can be selected arbitrarily, because

$$\begin{aligned} h_i(x) + \langle \operatorname{sgn}(x_i)e_i, y - x \rangle &= |x_i| + \operatorname{sgn}(x_i)y_i - \underbrace{\operatorname{sgn}(x_i)x_i}_{|x_i|} = \operatorname{sgn}(x_i)y_i \\ & \underbrace{\operatorname{sgn}(x_i)\in[-1,1]}_{\leq} |y_i| = h_i(y). \end{aligned}$$

So again

$$\partial h(x) \supseteq \sum_{i=1}^{d} \partial h_i(x) \ni (\operatorname{sgn}(x_1), \dots, \operatorname{sgn}(x_n))^T =: s(x).$$

So putting everything together we have

$$\partial f(x) \ni x - y + \lambda s(x).$$

(1 pt)

(1 pt)

(ii) Prove that f is convex.

Solution. As its sets of sub-gradients is nowhere empty, it is convex. \Box

(iii) Find a global minimum of f.

Solution. By the lecture it is sufficient to find a point x such that $0 \in \partial f(x)$. By the previous exercise we therefore want to solve

$$0 \stackrel{!}{=} x - y + \lambda s(x)$$

entry-wise this implies

$$\begin{aligned} x_i \stackrel{!}{=} y_i - \lambda \operatorname{sgn}(x_i) &= \begin{cases} y_i + \lambda & x_i < 0\\ y_i - \lambda [-1, 1] & x_i = 0\\ y_i - \lambda & x_i > 0 \end{cases} \\ &= \begin{cases} y_i + \lambda & y_i + \lambda < 0\\ 0 & y_i \in [-\lambda, \lambda]\\ y_i - \lambda & y_i - \lambda > 0 \end{cases} \\ &= \begin{cases} y_i + \lambda & y_i < -\lambda\\ 0 & y_i \in [-\lambda, \lambda]\\ y_i - \lambda & y_i > \lambda. \end{cases} \end{aligned}$$

(iv) Implement f as a sub-type of "DifferentiableFunction" (even though it is not) by returning a single sub-gradient and apply gradient descent to verify the global minimum https://classroom.github.com/a/Bm7FMb12 (2 pts).

Exercise 4 (Momentum Matrix).

let $D = \operatorname{diag}(\lambda_1, \ldots, \lambda_d), \alpha, \beta > 0$ and define

$$T = \begin{pmatrix} (1+\beta)\mathbb{I} - \alpha D & -\beta\mathbb{I} \\ \mathbb{I} & \mathbf{0} \end{pmatrix} \in \mathbb{R}^{2d \times 2d}$$

Prove there exists a regular $S \in \mathbb{R}^{2d \times 2d}$ such that

$$S^{-1}TS = \hat{T} = \begin{pmatrix} T_1 & & \\ & \ddots & \\ & & T_d \end{pmatrix}$$

with

$$T_i = \begin{pmatrix} 1 + \beta - \alpha \lambda_i & -\beta \\ 1 & 0 \end{pmatrix} \in \mathbb{R}^{2 \times 2}$$

Solution. We simply define for the standard basis $e_i \in \mathbb{R}^d$

$$S = \begin{pmatrix} e_1 & 0 & \dots & e_d & 0\\ 0 & e_1 & \dots & 0 & e_d \end{pmatrix} \in \mathbb{R}^{2d \times 2d}$$

in particular $S^T = S^{-1}$.

Exercise 5 (PL-Inequality).

Assume $f : \mathbb{R}^d \to \mathbb{R}$ is *L*-smooth and satisfies the Polyak-Łojasiewicz inequality

$$\|\nabla f(x)\|^2 \ge 2c(f(x) - f_*)$$
 (PL)

for some c > 0 and all $x \in \mathbb{R}^d$ with $f_* = \min_x f(x) > -\infty$.

(i) Prove that gradient descent with fixed step size $\alpha_k = \frac{1}{L}$ converges linearly in the sense

$$f(x_k) - f_* \le (1 - \frac{c}{L})^k (f(x_0) - f_*).$$
 (2 pts)

Solution. By L-smoothness and the descent lemma, we have

$$f(x_{k+1}) \le f(x_k) - \frac{1}{2L} \|\nabla f(x_k)\|^2 \stackrel{\text{(PL)}}{\le} f(x_k) - \frac{c}{L} (f(x_k) - f_*)$$

Subtracting f_* from both sides, we get

$$f(x_{k+1}) - f_* \le (1 - \frac{c}{L})(f(x_k) - f_*)$$

(ii) Prove that μ -strong-convexity and L-smoothness imply the PL-inequality. (2 pts)

(2 Points)

(6 Points)

Solution. Recall by the solution of sheet 1, exercise 6 (iii), and strong convexity we have

$$\begin{split} \mu \|x - y\|^2 &\leq D_f^{(B)}(x, y) + D_f^{(B)}(y, x) \\ &= \langle \nabla f(x) - \nabla f(y), x - y \rangle \\ &\stackrel{\text{C.S.}}{\leq} \|\nabla f(x) - \nabla f(y)\| \|x - y\| \end{split}$$

and therefore

$$\mu \|x - y\| \le \|\nabla f(x) - \nabla f(y)\|.$$
(1)

Finally we know by L-smoothness and $\nabla f(x_*) = 0$ where x_* is the minimum

$$f(x) - f(x_*) \stackrel{\nabla f(x_*)=0}{=} D_f^{(B)}(x, x_*) \stackrel{L-\text{smooth}}{\leq} \frac{L}{2} \|x - x_*\|^2 \stackrel{(1)}{\leq} \frac{L}{2\mu} \|\nabla f(x) - \underbrace{\nabla f(x_*)}_{=0} \|^2. \quad \Box$$

(iii) Use a graphing calculator to find c such that $f(x) = x^2 + 3\sin^2(x)$ satisfies the PL-condition (argue why $x \to \infty$ is not a problem) and prove it is not convex. (2 pts)

Solution. For $c = \frac{1}{6}$ we have the PL-condition



As $f'(x) = 2(x + 3\sin(x)\cos(x))$ and therefore

$$f'(x)^{2} = 4(x+3\underbrace{\sin(x)\cos(x)}_{\in [-1,1]})^{2} \stackrel{|x|\geq 3}{\geq} 4(|x|-3)^{2}$$

the x^2 dominates for large x, so if we make c small enough we can ensure the inequality for large x.

f is not convex because

$$f(\frac{1}{2}\pi + \frac{1}{2}0) = \frac{\pi^2}{4} + 3 > \frac{1}{2}\pi^2 = \frac{1}{2}f(\pi) + \frac{1}{2}f(0).$$

Exercise 6 (Weak PL-Inequality).

Assume $f : \mathbb{R}^d \to \mathbb{R}$ is *L*-smooth and satisfies the "weak PL inequality"

$$\|\nabla f(x)\| \ge 2c(f(x) - f_*)$$

for some c > 0 and all $x \in \mathbb{R}^d$ with $f_* = \min_x f(x) > -\infty$.

(i) Let $a_0 \in [0, \frac{1}{q}]$ for some q > 0 and assume for the sequence $(a_n)_{n \in \mathbb{N}}$ that it is positive and satisfies a diminishing contraction

$$0 \le a_{n+1} \le (1 - qa_n)a_n \qquad \forall n \ge 0.$$

Prove the convergence rate

$$a_n \le \frac{1}{nq+1/a_0} \le \frac{1}{(n+1)q}.$$
 (2 pts)

Solution. Divide the reordered contraction

$$a_n \ge a_{n+1} + qa_n^2$$

by $a_n a_{n+1}$ to obtain

$$\frac{1}{a_{n+1}} \ge \frac{1}{a_n} + q \underbrace{\frac{a_n}{a_{n+1}}}_{\ge 1} \ge \frac{1}{a_n} + q$$

which leads to

$$\frac{1}{a_n} - \frac{1}{a_0} = \sum_{k=0}^{n-1} \frac{1}{a_{k+1}} - \frac{1}{a_k} \ge nq.$$

Reordering we obtain our claim

$$a_n \le \frac{1}{nq + \frac{1}{a_0}} \stackrel{a_0 \le \frac{1}{q}}{\le} \frac{1}{(n+1)q}.$$

(ii) Prove that f is bounded. More specifically $e(x) := f(x) - f_* \le \frac{L}{2c^2}$ for all x. (1 pt)

Solution. Using Sheet 1 Exercise 1 (i), we get

$$f_* \le f(x) - \frac{1}{2L} \|\nabla f(x)\|^2$$

(5 Points)

and therefore

$$e(x) \ge \frac{1}{2L} \|\nabla f(x)\|^2 \stackrel{\text{weak PL}}{\ge} \frac{4c^2}{2L} e(x)^2$$

Dividing both sides by e(x) we obtain

$$1 \ge \frac{2c^2}{L}e(x)$$

and thus

$$e(x) \le \frac{L}{2c^2}.$$

(iii) For gradient descent $x_{n+1} - x_n = -\alpha_n \nabla f(x_n)$ with constant step size $\alpha_k = \frac{1}{L}$ prove the convergence rate

$$f(x_n) - f_* \le \frac{L}{2c^2(n+1)}.$$
 (2 pts)

Solution. Using L-smoothness, we have

$$f(x_{k+1}) \leq f(x_k) + \langle \nabla f(x_k), x_{k+1} - x_k \rangle + \frac{L}{2} ||x_{k+1} - x_k||^2$$

$$\leq f(x_k) - \underbrace{\alpha_k (1 - \frac{L}{2} \alpha_k)}_{=\frac{1}{2L}} \underbrace{||\nabla f(x_k)||^2}_{\geq 4c^2 e(x_k)^2}$$

If we subtract f_* from both sides and apply our weak PL inequality we get

$$e(x_{k+1}) \le e(x_k) - \frac{4c^2}{2L}e(x_k)^2 = (1 - \frac{2c^2}{L}e(x_k))e(x_k)$$

with $q = \frac{2c^2}{L}$ and $e(x_0) \le \frac{L}{2c^2} = \frac{1}{q}$ by (ii), we can apply (i) to obtain our claim.