

Constrained Optimization and SGD

Optimization for Machine Learning — Homework #2

Monday 12th June, 2023

The theory part can be handed-in physically during the exercise session, or digitally on Moodle. The programming part has to be sent on Moodle. *Group work is allowed (2 – 3 people), but submissions are personal..*

Part I: Theory

12 points

I.1. Constrained optimization

Exercise I.1 (Constrained optimization, 5 points). You encounter the following optimization problem on $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \in \mathbb{R}^2$:

$$\begin{array}{ll} \text{minimize} & f(x) := -(x_1 + x_2) \\ \text{subject to} & c(x) := x_1^2 + x_2^2 - 1 \leq 0 \end{array}$$

1. What is the dimension of the Lagrangian multiplier α ? Write down the Lagrangian $L(x, \alpha)$.
2. Show that the dual function $g(\alpha) := \min_x L(x, \alpha)$ is $g(\alpha) = -(\alpha + \frac{1}{2\alpha})$.
3. For feasible α and x (meaning $\alpha \geq 0$ and $c(x) \leq 0$), show that $g(\alpha) \leq f(x)$.
4. Solve the dual problem

$$\begin{array}{ll} \text{maximize} & g(\alpha) \\ \text{subject to} & \alpha \geq 0 \end{array}$$

5. What is the solution to the original problem?

Answer I.1. 1. The dimension of α is given by the number of constraints of the primal problem, here $\dim \alpha = 1$.

The constraint is in canonical form for minimization problem, i.e. it is non-negative when satisfied. Therefore, the Lagrangian is given by

$$\forall x \in \mathbb{R}^2, \forall \alpha \geq 0, \quad L(x, \alpha) := f(x) + \alpha c(x) = -(x_1 + x_2) + \alpha(x_1^2 + x_2^2 - 1)$$

Indeed, as the primal problem is to find a point x that minimizes $f(x)$, the product $\alpha c(x)$ will be positive when the constraint is *not* satisfied.

2. The dual function $g: [0, +\infty[\rightarrow \mathbb{R}$ is found by solving, for $\alpha \leq 0$, $\min_x L(x, \alpha)$. For $\alpha \geq 0$, the function $x \mapsto L(x, \alpha)$ is differentiable and convex (since f and c are, and $\alpha \geq 0$), therefore g is found by solving, for $\alpha \geq 0$,

$$0 \stackrel{!}{=} \nabla_x L(x^*, \alpha) = \begin{pmatrix} -1 + 2\alpha x_1^* \\ -1 + 2\alpha x_2^* \end{pmatrix} \implies x_1^* = x_2^* = \frac{1}{2\alpha}.$$

Then,

$$\begin{aligned} g(\alpha) &= L(x^*, \alpha) = -2 \left(\frac{1}{2\alpha} \right) + \alpha \left(\frac{1}{4\alpha^2} + \frac{1}{4\alpha^2} - 1 \right) \\ &= -\frac{1}{\alpha} + \frac{1}{2\alpha} - \alpha \\ &= -\left(\alpha + \frac{1}{2\alpha} \right). \end{aligned}$$

3. Let $\alpha \geq 0$ and $\tilde{x} \in \mathbb{R}^2$ such that $c(\tilde{x}) \leq 0$. Then, $L(\tilde{x}, \alpha) = f(\tilde{x}) + \alpha c(\tilde{x}) \leq f(\tilde{x})$, and $g(\alpha) = \min_x L(x, \alpha) \leq L(\tilde{x}, \alpha) \leq f(\tilde{x})$.
4. The function g differentiable and concave in α , since it is the point-wise minimum of a family of affine functions in α ¹. Hence, the dual maximisation problem can be solved using the first-order optimality condition for differentiable concave functions:

$$\begin{aligned} g(\alpha^*) &= \max_{\alpha} g(\alpha) \iff \nabla g(\alpha^*) = 0 \\ &\iff -1 - \frac{1}{2\alpha^{*2}} = 0 \\ &\iff \alpha^* = \pm \frac{\sqrt{2}}{2}. \end{aligned}$$

The dual feasibility condition requires that $\alpha^* \leq 0$, hence $\alpha^* = -\frac{\sqrt{2}}{2}$.

5. The functions f and f_1 are convex in x . Moreover, there exists a point $x \in \mathbb{R}^2$ such that the constraint is strictly satisfied, e.g. the point $\begin{pmatrix} 0 \\ 0 \end{pmatrix}$. Therefore, from the Slater's (strong) constraint qualification, the duality gap is 0 (strong duality holds), and $d^* = p^*$ where d^*

¹ $\min_x t_x(\alpha) = -\max_x (-t_x(\alpha))$, for affine t_x , $-t_x$ is also affine hence concave and the point-wise maximum of a collection of concave functions is concave.

is the dual optimal value and p^* the primal dual value. Therefore, the minimum value for f is

$$p^* = g(\alpha^*) = -\left(\frac{\sqrt{2}}{2} + \frac{\sqrt{2}}{2}\right) = -\sqrt{2}$$

Since $x \mapsto L(x, \alpha)$ is strictly convex, we have that

$$x^* \text{ solves the primal problem} \iff x^* = \operatorname{argmin}_x L(x, \alpha^*).$$

To find the solution to the primal problem, we therefore simply need to find the solution to $\operatorname{argmin}_x L(x, \alpha^*)$, which is given by (2) with $\alpha = \alpha^* = \frac{\sqrt{2}}{2}$. We find

$$x^* = \begin{pmatrix} \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{pmatrix}.$$

I.2. General analysis

Exercise I.2 (Inequality of c -strongly convex functions, 2 points). Recall that, given $c > 0$, a function $E: \Omega \subset \mathbb{R}^d \rightarrow \mathbb{R}$ is called c -strongly convex if Ω is convex and

$$\forall (x, y) \in \Omega^2, \quad E(y) \geq E(x) + \langle \nabla E(x), y - x \rangle + \frac{c}{2} \|y - x\|_2^2$$

A strongly convex function has a unique minimizer $x_* = \operatorname{argmin}_{x \in \Omega} E(x)$.

Show that, if E is c -strongly convex, it satisfies the following inequality:

$$\forall x \in \Omega, \quad 2c(E(x) - E_*) \leq \|\nabla E(x)\|_2^2$$

Hint: Study the function $q_x(y) = E(x) + \langle \nabla E(x), y - x \rangle + \frac{c}{2} \|y - x\|_2^2$.

Answer I.2. 1. Assume that $E: \Omega \rightarrow \mathbb{R}$ is c -strongly convex. Let $x \in \Omega$. Then,

$$\forall y \in \Omega, \quad E(y) \geq E(x) + \langle \nabla E(x), y - x \rangle + \frac{c}{2} \|y - x\|_2^2 =: q_x(y).$$

The quadratic function q_x has a unique minimizer y^* , given by (first-order optimality condition for convex function)

$$\nabla q_x(y^*) = 0 \implies \nabla E(x) + c(y^* - x) = 0$$

$$\implies y^* = x - \frac{1}{c} \nabla E(x)$$

$$\implies \min_y q_x(y) = q_x(y^*) = E(x) + \langle \nabla E(x), -\frac{1}{c} \nabla E(x) \rangle + \frac{c}{2} \left\| -\frac{1}{c} \nabla E(x) \right\|^2$$

$$= E(x) - \frac{1}{2c} \|\nabla E(x)\|^2$$

Then,

$$\begin{aligned}
\forall x \in \Omega, \forall y \in \Omega, E(y) \geq q_x(y) &\implies \min_y E(y) \geq \min_y q_x(y) \\
&\implies E_* \geq E(x) - \frac{1}{2c} \|\nabla E(x)\|^2 \\
&\implies 2c(E(x) - E_*) \leq \|\nabla E(x)\|^2.
\end{aligned}$$

This is true for all $x \in \Omega$, and thus the statement is proven.

I.3. SGD Analysis

Exercise I.3 (5 points). Recall that if we assume that F is strongly convex, we can show that the gradient descent converges with rate $\mathcal{O}(\rho^k)$ where $0 < \rho < 1$, and k is the number of iterations. This rate is called “linear convergence”. Assume we have a L -smooth and c -strongly convex function. Recall the expression for the convergence of stochastic gradient descent:

$$\begin{aligned}
\mathbb{E}[\|\nabla f(x_T)\|^2] &\leq 2 \left(\frac{2\sqrt{T+1}-1}{L} \right)^{-1} \cdot \left(\mathbb{E}[f(x_0)] - f^* + \frac{\log(T)+1}{L^2} \right) \\
&= \mathcal{O}\left(\frac{\log(T)}{L\sqrt{T}}\right),
\end{aligned}$$

where $f^* = \min_x f(x)$, α is the step size, L is the Lipschitz constant, and T is the total number of iterations.

1. How does this compare to the expression that we get for the gradient descent?
2. Derive the rate of convergence for the stochastic gradient descent for strongly convex functions.

Hints:

1. Start from the expression, valid when f is L -smooth:

$$\mathbb{E}[f(x_{k+1})] \leq \mathbb{E}[f(x_k)] - \frac{\alpha}{2} \mathbb{E}[\|\nabla f(x_k)\|^2] + \frac{\alpha^2 \sigma^2 L}{2}.$$

2. Apply the inequality for c -strongly convex functions (Polyak-Łojasewicz inequality):

$$\|\nabla f(x)\|^2 \geq 2c(f(x) - f^*), \quad c > 0$$

3. Subtract from both sides f^* .
4. Subtract from both sides the fixed point $\frac{\alpha^2 \sigma^2 L}{2c\alpha}$.
5. Apply the previous step recursively for T steps.
6. Use the inequality $(1 - c\alpha) \leq \exp(-c\alpha)$.

Answer I.3. 1. The rate

2.

Part II: Programming

8 points

Exercise II.1 (SVM with SGD, 8 points). In this exercise, (linear) SVM will be solved with SGD.

We consider the SVM problem with prediction $h(x, w) := \langle x, w \rangle + b$ and (differentiable) loss $\ell(\hat{y}, y) := \frac{1}{10} \ln(1 + \exp(10(1 - y \cdot \hat{y})))$. The composed loss is denoted by $f(w, (x, y)) := \ell(h(x, w), y)$. Given a training dataset $\{(x_i, y_i)\}_{i \in [n]}$, the empirical risk is $R_n(w) := \frac{1}{n} \sum_{i=1}^n f(w, (x_i, y_i))$.

The stochastic gradient descent (SGD) algorithm is given in Algorithm 1. The random variable ξ_k selects samples and their targets (i.e. elements from $\mathcal{X} \times \{-1, 1\}$).

Algorithm 1 Stochastic Gradient Descent [1, Algorithm 4.1].

- 1: Choose an initial iterate w_1
 - 2: **for** $k = 1, 2, \dots$ **do**
 - 3: Generate a realization of the random variable ξ_k with values in $\mathcal{X} \times \{-1, 1\}$ (e.g. batch of samples)
 - 4: Compute a stochastic vector $g(w_k, \xi_k)$
 - 5: Choose a step size $\alpha_k > 0$
 - 6: Set the new iterate as $w_{k+1} \leftarrow w_k - \alpha_k g(w_k, \xi_k)$.
 - 7: **end for**
-

1. What is the role of the factor 10 in ℓ ?
2. Is the function $w \mapsto R_n(w)$ (strongly) convex? L -smooth? What is the gradient $\nabla_w \ell(h(x, w), y)$?
3. Implement the SGD algorithm Algorithm 1, with data given by `utils.gen_linsep_data`. You can choose how to sample from the dataset, either one sample at a time or using a batch of samples.
4. Test different step sizes and show the convergence.

Answer II.1. 1. The factor $p = 10$ in ℓ controls how “close” the function ℓ is to $f: z \mapsto \max(0, 1 - z)$: with higher p , the function is closer to f . See it here.

2. We see that R_n is convex as addition of convex functions.

To see if the function R_n is strongly convex or L -smooth, one can first study the Hessian of $f: z \mapsto \frac{1}{10} \ln(1 + \exp(10(1 - z)))$. One computes $f''(z) = \frac{10 \exp(10(1-z))}{(1 + \exp(10(1-x)))^2}$.

The function f is *not* c -strongly convex: since $\lim_{z \rightarrow -\infty} f''(z) = \lim_{z \rightarrow +\infty} f''(z) = 0$, there is no $c > 0$ such that $\forall z \in \mathbb{R}, f''(z) \geq c$.

The function f is L -smooth: since $\lim_{z \rightarrow -\infty} f''(z) = \lim_{z \rightarrow +\infty} f''(z) = 0$ and f'' is continuous, there exists $L \in \mathbb{R}$ such that $\forall z \in \mathbb{R}, f''(z) \leq L$.

Since R_n is the composition of affine functions of w with f , we can draw the same conclusions for R_n (in 1D, $(f \circ g)'' = (f'' \circ g)(g')^2 + (f' \circ g)g'' = (f'' \circ g)(g')^2$ for affine

g , with a similar rule in dimension d). It is not c -strongly convex, and it is L -smooth. (Warning: the constant L is not necessarily the same one as the one for f).

To compute it more precisely, one can compute the Hessian of R_n and find the bounds (c, L) such that $cI \preceq \nabla^2 R_n(w) \preceq LI$.

The gradient of the loss ℓ with respect to w is (chain rule):

$$\nabla_w \ell(h(x; w), y) = \frac{\partial h(x; w)}{\partial w}^\top \ell'(h(x; w), y) = -x \cdot y \cdot \frac{\exp(10(1 - y \cdot h(x; w)))}{1 + \exp(10(1 - y \cdot h(x; w)))}.$$

Likewise,

$$\nabla_b \ell(h(x; w), y) = \frac{\partial h(x; w)}{\partial b}^\top \ell'(h(x; w), y) = -y \cdot \frac{\exp(10(1 - y \cdot h(x; w)))}{1 + \exp(10(1 - y \cdot h(x; w)))}.$$

References

- [1] Léon Bottou, Frank E. Curtis and Jorge Nocedal. ‘Optimization Methods for Large-Scale Machine Learning’. 2016. DOI: 10.1137/16M1080173. arXiv: 1606.04838.