

Exercise 9: Adaptive Methods

*Lecturer: Aurelien Lucchi***Problem 1 (Polyak's step size):**

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex function we aim to minimize. Assume that we have access to a subdifferential $\partial f(x)$ of f at each point $x \in \mathbb{R}^d$. We denote by x^* the minimizer of f . We consider the following iterative process to solve the problem

$$x^{t+1} = \arg \min \frac{1}{2} \|x - x^t\|^2 \quad \text{such that} \quad x \in \mathcal{S}_t := \{x \mid f(x^t) + \langle g^t, x - x^t \rangle \leq f(x^*)\}, \quad (1)$$

where $g^t \in \partial f(x^t)$ is any subgradient of f at point x^t .

- Using the definition of a convex function, show that \mathcal{S}_t is non-empty for all iterations t .
- Write the Lagrangian function $L(x, \lambda)$ for the problem (1).
- Write KKT conditions for the Lagrangian.
- Solving KKT conditions find the closed-form solution, i.e., derive the expression for x^{t+1} .
- Does the iteration process from d) remind you of an algorithm that has already been covered during the course? Write the name of the algorithm.

Problem 2 (Lipschitz Continuity Inequality):

Prove that for a L -smooth function $\frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \leq (f(\mathbf{x}_t) - f^*)$

Problem 3 (Gauss-Newton Decomposition):

Let $\ell : \mathbb{R}^m \times \mathcal{Y} \rightarrow \mathbb{R}$ be a loss function (e.g., a squared loss $\ell(\mathbf{z}, \mathbf{y}) = \frac{1}{2} \|\mathbf{z} - \mathbf{y}\|^2$). Here \mathbf{z} are predictions and \mathbf{y} is the ground truth. Let us consider some ML model with parameters \mathbf{w} : $\mathbf{f} : \mathbb{R}^d \times \mathcal{X} \rightarrow \mathbb{R}^m$ which takes features \mathbf{x} as input and outputs a prediction $\mathbf{f}(\mathbf{w}, \mathbf{x})$. The quality of the model given an input \mathbf{x} is $\ell(\mathbf{f}(\mathbf{w}, \mathbf{x}), \mathbf{y})$. Show that the Hessian of the function $g(\mathbf{w}) := \ell(\mathbf{f}(\mathbf{w}, \mathbf{x}), \mathbf{y})$ can be written as

$$\nabla^2 g(\mathbf{w}) = J_f(\mathbf{w}, \mathbf{x})^\top \nabla_{\mathbf{z}}^2 \ell(\mathbf{z}, \mathbf{y}) J_f(\mathbf{w}, \mathbf{x}) + \sum_{l=1}^m \frac{\partial \ell(\mathbf{z}, \mathbf{y})}{\partial \mathbf{z}_l} \nabla_{\mathbf{w}}^2 \mathbf{f}_l(\mathbf{w}, \mathbf{x}), \quad (2)$$

where $J_f(\mathbf{w}, \mathbf{x})$ is the Jacobian of \mathbf{f} , i.e. $[J_f(\mathbf{w}, \mathbf{x})]_{i,j} = \frac{\partial \mathbf{f}_i(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_j}$.