

Exercise sheet 1: Around gradient descent

M2 MIAGE ID Apprentissage

October 1st, 2024



Exercise 1.1: One-layer neural network (Exam 2021-2022)

In this exercise, we consider the special case of a dataset with scalar labels/outputs, i.e. of the form $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ with $\mathbf{x}_i \in \mathbb{R}^{d_x}$ and $y_i \in \mathbb{R}$ for every $i = 1, \dots, n$. We build a simple neural network with no activation function and one homogeneous linear layer to predict the value y_i from the vector \mathbf{x}_i , resulting in the model

$$\begin{aligned} h^{lin}(\cdot; \mathbf{w}) : \mathbb{R}^{d_x} &\longrightarrow \mathbb{R} \\ \mathbf{x} &\longmapsto \mathbf{W}_1 \mathbf{x}, \end{aligned} \quad (1)$$

with $\mathbf{W}_1 \in \mathbb{R}^{1 \times d_x}$. Letting $d = d_x$ and $\mathbf{w} = \mathbf{W}_1^T \in \mathbb{R}^d$, finding the best model amounts to solving

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f^{lin}(\mathbf{w}) := \frac{1}{2n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2. \quad (2)$$

- What class of problems does problem (2) belong to?
- The objective function f^{lin} is $\mathcal{C}_L^{1,1}$, i.e. its gradient is L -Lipschitz continuous. If L is known, how can its value be used in an algorithm such as gradient descent?
- Problem (2) is convex with a \mathcal{C}^1 objective function.
 - What can then be said about a point $\bar{\mathbf{w}}$ such that $\nabla f^{lin}(\bar{\mathbf{w}}) = \mathbf{0}_{\mathbb{R}^d}$?
 - What is the convergence rate of gradient descent on this problem?
 - What is the convergence rate of accelerated descent on a convex problem? Is it better or worse than that of the previous question?
- Suppose that the data is such that the objective f^{lin} is μ -strongly convex, in addition to the properties already mentioned above.
 - Let $\mathbf{w}, \mathbf{v} \in \mathbb{R}^d$ be two points such that $\nabla f^{lin}(\mathbf{w}) = \nabla f^{lin}(\mathbf{v}) = \mathbf{0}_{\mathbb{R}^d}$. What can we say about \mathbf{v} and \mathbf{w} ?
 - What is the convergence rate of accelerated gradient on this problem?

Exercise 1.2: Two-layer linear neural networks (exam 2021-2022)

We consider a dataset $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ where $\mathbf{x}_i \in \mathbb{R}^{d_x}$ and $\mathbf{y}_i \in \mathbb{R}^{d_y}$. We wish to learn a mapping from \mathbb{R}^{d_x} to \mathbb{R}^{d_y} that correctly outputs \mathbf{y}_i when given \mathbf{x}_i as an input. Our model will be that of a two-layer linear neural network :

$$\begin{aligned} \mathbf{h}(\cdot; \mathbf{w}) : \mathbb{R}^{d_x} &\longrightarrow \mathbb{R}^{d_y} \\ \mathbf{x} &\longmapsto \mathbf{W}_2(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2, \end{aligned} \quad (3)$$

where $\mathbf{W}_1 \in \mathbb{R}^{d_x \times m}$, $\mathbf{b}_1 \in \mathbb{R}^m$, $\mathbf{W}_2 \in \mathbb{R}^{m \times d_y}$ and $\mathbf{b}_2 \in \mathbb{R}^{d_y}$. We will consider \mathbf{h} as being parameterized by $\mathbf{w} \in \mathbb{R}^d$, with $d = d_x m + m + m d_y + d_y$ and \mathbf{w} concatenating all coefficients from $\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2$. Our goal is to determine a value of \mathbf{w} so that $\mathbf{h}(\mathbf{x}_i; \mathbf{w}) \approx \mathbf{y}_i$, which we formalize using the squared loss $(\mathbf{h}, \mathbf{y}) \mapsto \frac{1}{2} \|\mathbf{h} - \mathbf{y}\|^2$.

Overall, we obtain the following problem:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) := \frac{1}{2n} \sum_{i=1}^n \|\mathbf{h}(\mathbf{x}_i; \mathbf{w}) - \mathbf{y}_i\|^2. \quad (4)$$

It can be shown that the function f is \mathcal{C}^1 .

- Give a lower bound on the objective function of problem (4).
- In general, problem (4) is nonconvex. What does this imply about its local minima?
- Suppose that \mathbf{w}^* is a solution of (4). What can be said about the derivative of f at \mathbf{w}^* ?
- Write down the gradient descent iteration for problem (4) with an arbitrary stepsize.
- Given that the problem is nonconvex, what is the theoretical convergence rate of gradient descent applied to (4)?

Exercise 1.3: Matrix completion (exam 2022-2023)

Let $\mathbf{X} \in \mathbb{R}^{d \times d}$ be a data matrix such that only a subset of its entries $\mathcal{S} \subset \{1, \dots, d\}^2$ are known with $|\mathcal{S}| = n \leq d^2$. We consider the problem

$$\underset{\mathbf{W} \in \mathbb{R}^{d \times d}}{\text{minimize}} f(\mathbf{W}) := \frac{1}{2n} \sum_{(i,j) \in \mathcal{S}} ([\mathbf{W}]_{ij} - [\mathbf{X}]_{ij})^2. \quad (5)$$

- a) When $\mathcal{S} = \{1, \dots, d\}^2$, justify that $\mathbf{W}^* = \mathbf{X}$ is the unique solution of the problem.
 b) Problem (5) is convex in the coefficients of \mathbf{W} . Letting $\mathbf{w} \in \mathbb{R}^{d^2}$ denoting the column vector formed by stacking all columns of the matrix \mathbf{W} in order, we can reformulate the problem as

$$\underset{\mathbf{w} \in \mathbb{R}^{d^2}}{\text{minimize}} \hat{f}(\mathbf{w}) := \frac{1}{2n} \sum_{(i,j) \in \mathcal{S}} ([\mathbf{w}]_{i+(j-1)d} - [\mathbf{X}]_{ij})^2. \quad (6)$$

The function \hat{f} is convex and \mathcal{C}^1 .

- i) What convergence rate guarantee can we provide on gradient descent when applied to problem (6)? What quantity does this rate apply to?
 ii) What is the corresponding convergence rate for the accelerated gradient method due to Nesterov? Is it better than that of gradient descent?
 iii) When $n = d^2$, the function \hat{f} is a strongly convex quadratic function. Aside from Nesterov's method, what other approach can we use to obtain better convergence rates than gradient descent?
 c) We now suppose that the data matrix \mathbf{X} is symmetric, positive semidefinite and of rank $1 \ll d$. In this setting, rather than seeking an arbitrary matrix \mathbf{W} to approximate \mathbf{X} , we can force the matrix to be rank one by writing it $\mathbf{u}\mathbf{u}^T$ where $\mathbf{u} \in \mathbb{R}^d$. Problem (5) then becomes

$$\underset{\mathbf{u} \in \mathbb{R}^d}{\text{minimize}} \tilde{f}(\mathbf{u}) := \frac{1}{2n} \sum_{(i,j) \in \mathcal{S}} ([\mathbf{u}\mathbf{u}^T]_{ij} - [\mathbf{X}]_{ij})^2. \quad (7)$$

The objective function of problem (7) is \mathcal{C}^2 and nonconvex.

- i) State the first-order necessary optimality conditions for problem (7).
 ii) What is the convergence rate of gradient descent for this problem? What quantity does this rate apply to?
 iii) Under certain assumptions on \mathbf{X} and \mathcal{S} , one can show that all the local minima of this problem are global. In that case, what technique guarantees almost surely that gradient descent will converge to such a point?