

Tutorial 2: Around gradient descent

Optimization for data science, M2 MIAE ID/ID Apprentissage

November 10, 2025



Exercise 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

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

- a) What class of problems does problem (2) belong to?
- b) 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?
- c) Problem (2) is convex with a \mathcal{C}^1 objective function.
 - i) What can then be said about a point $\bar{\mathbf{w}}$ such that $\nabla f^{lin}(\bar{\mathbf{w}}) = \mathbf{0}_{\mathbb{R}^d}$?
 - ii) What is the convergence rate of gradient descent on this problem?
 - iii) What is the convergence rate of accelerated descent on a convex problem? Is it better or worse than that of the previous question?
- d) Suppose that the data is such that the objective f^{lin} is μ -strongly convex, in addition to the properties already mentioned above.
 - i) 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} ?
 - ii) What is the convergence rate of accelerated gradient on this problem?

Exercise 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 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

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

- When $\mathcal{S} = \{1, \dots, d\}^2$, justify that $\mathbf{W}^* = \mathbf{X}$ is the unique solution of the problem.
- 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 .

- What convergence rate guarantee can we provide on gradient descent when applied to problem (6)? What quantity does this rate apply to?
 - What is the corresponding convergence rate for the accelerated gradient method due to Nesterov? Is it better than that of gradient descent?
 - 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.

- State the first-order necessary optimality conditions for problem (7).
- What is the convergence rate of gradient descent for this problem? What quantity does this rate apply to?
- State the second-order necessary optimality conditions for problem (7).
- Under certain assumptions on \mathbf{X} and \mathcal{S} , one can show that all points satisfying second-order necessary optimality conditions are global minima of the problem.

Solutions

Solutions to Exercise 1

- a) The function $f(\mathbf{W})$ is always nonnegative (as a sum of squares, i.e. nonnegative numbers). When $n = d^2$, we have that

$$f(\mathbf{W}) = 0 \Leftrightarrow ([\mathbf{W}]_{ij} - [\mathbf{X}]_{ij})^2 = 0 \forall (i, j) \in \{1, \dots, d\}^2 \Leftrightarrow \mathbf{W} = \mathbf{X}.$$

As a result, the problem has a single global minimum given by $\mathbf{W}^* = \mathbf{X}$.

- b) Convex formulation

- i) Since the problem is convex, we know that after $K \geq 1$ iterations of gradient descent, the iterate \mathbf{w}_K satisfies

$$f^{lin}(\mathbf{w}_K) - \min_{\mathbf{w} \in \mathbb{R}^{d^2}} f^{lin}(\mathbf{w}) \leq \mathcal{O}\left(\frac{1}{K}\right).$$

Gradient descent thus converges at a rate $\frac{1}{K}$.

- ii) The rate for accelerated gradient on such a problem is $\frac{1}{K^2}$, which is a better rate as it converges more quickly to 0.
- iii) When f^{lin} is a strongly convex quadratic function, the heavy-ball method (aka Polyak's method) attains the optimal rate of convergence for strongly convex functions, which is better than gradient descent. *NB: The value of that rate is not required to answer the question.*

- c) (Nonconvex case)

- i) If $\bar{\mathbf{u}} \in \mathbb{R}^d$ is a local minima of problem (7), then $\nabla \tilde{f}(\bar{\mathbf{u}}) = \mathbf{0}$.
- ii) For this problem, after $K \geq 1$ iterations of gradient descent, we have

$$\min_{0 \leq k \leq K-1} \|\nabla f(\mathbf{w}_k)\| \leq \mathcal{O}\left(\frac{1}{\sqrt{K}}\right),$$

hence the convergence rate of gradient descent is in $\frac{1}{\sqrt{K}}$.

- iii) If $\bar{\mathbf{u}} \in \mathbb{R}^d$ is a local minima of problem (7), then $\nabla \tilde{f}(\bar{\mathbf{u}}) = \mathbf{0}$ and $\nabla^2 \tilde{f}(\bar{\mathbf{u}}) \succeq \mathbf{0}$.
- iv) Initializing gradient descent with a random point guarantees almost surely that it will converge to a second-order necessary point, hence a global minimum under the assumptions of this question.

Solutions to Exercise 2

- a) The value 0 is a lower bound on this objective function, since it is always nonnegative. Any value less than or equal to 0 also works.
- b) The local minima of a nonconvex problem are not necessarily global minima.

- c) By the first-order necessary conditions, if \mathbf{w}^* is a solution of (4), then its gradient is zero, that is $\nabla f(\mathbf{w}^*) = \mathbf{0}$.
- d) Using an arbitrary stepsize $\alpha_k > 0$, the k th iteration of gradient descent can be written as

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k).$$

- e) For a nonconvex problem such as (4), it can be guaranteed that, after $K \geq 1$ iterations of gradient descent, one has

$$\min_{0 \leq k \leq K-1} \|\nabla f(\mathbf{w}_k)\| \leq \mathcal{O}\left(\frac{1}{\sqrt{K}}\right).$$

Solutions to Exercise 3

- a) Problem (2) is a linear least-squares problem.
- b) If a Lipschitz constant L for the gradient is known, the stepsize can be chosen as the constant value $\alpha = \frac{1}{L}$. *NB: Other values less than $\frac{2}{L}$ would also guarantee decrease of the function value at every iteration.*
- c)

- i) Since the problem is convex, any point $\bar{\mathbf{w}}$ such that $\nabla f^{lin}(\bar{\mathbf{w}}) = \mathbf{0}_{\mathbb{R}^d}$ is a global minimum.
- ii) On such a convex problem, after $K \geq 1$ iterations of gradient descent, one obtains that

$$f(\mathbf{w}_k) - \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) \leq \mathcal{O}\left(\frac{1}{K}\right).$$

- iii) On a convex problem, after $K \geq 1$ iterations of accelerated gradient, one obtains that

$$f(\mathbf{w}_k) - \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) \leq \mathcal{O}\left(\frac{1}{K^2}\right),$$

which is better than the rate for gradient descent since it converges more rapidly towards 0.

- d)
- i) Since the function is strongly convex and continuously differentiable, it has a unique global minimum, which is the unique solution of the equation $\nabla f^{lin}(\mathbf{w}) = \mathbf{0}_{\mathbb{R}^d}$. Therefore, if \mathbf{w} and \mathbf{v} satisfy $\nabla f^{lin}(\mathbf{w}) = \nabla f^{lin}(\mathbf{v}) = \mathbf{0}_{\mathbb{R}^d}$, then we must have $\mathbf{v} = \mathbf{w}$.
- ii) On a strongly convex problem, after $K \geq 1$ iterations of accelerated gradient, one obtains that

$$f(\mathbf{w}_k) - \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) \leq \mathcal{O}\left((1 - \sqrt{\mu}L)^K\right).$$