

Assignment 1
CS60B: Algorithmic Machine Learning, Spring 2022
Due: April 13, 10PM

Guidelines for submitting the solutions:

- The assignments need to be submitted on Gradescope. Make sure you follow all the instructions - they are simple enough that exceptions will not be accepted.
 - Start each problem or sub-problem on a separate page even if it means having a lot of white-space and write/type in large font.
 - The solutions need to be submitted by 10 PM on the due date. No late submissions will be accepted.
 - Please adhere to the code of conduct outlined on the class page.
1. For a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, a point x_0 is a *local minimum* if there exists $\delta > 0$ such that for all x with $\|x - x_0\| < \delta$, $f(x_0) \leq f(x)$. Similarly, a point x_0 is a *local maximum* if the opposite holds: if there exists $\delta > 0$ such that for all x with $\|x - x_0\| < \delta$, $f(x_0) \geq f(x)$.

There exist smooth functions for which finding a *local minimum* is NP-hard. However, if f is differentiable, then one *necessary* condition for x_0 to be a local minimum is that $\nabla f(x_0) = 0$.

- (a) This is necessary but not sufficient! Give an example of a function in two dimensions and a point x_0 such that $\nabla f(x_0) = 0$ but x_0 is neither a local minimum nor a local maximum of f . [2 points]
- (b) In spite of the above, vanishing gradients is often desired. Suppose we have a β -smooth function $f : \mathbb{R}^d \rightarrow \mathbb{R}$. Show that gradient descent can be used to find a point w such that $\|\nabla f(w)\| \leq \varepsilon$. How many iterations of GD do you need for finding such a point? Your bound can depend on the starting point, $f(w_0)$, β , ε , and the value of the global optimum. [3 points]
- [Hint: Your point w could be any of the iterations of GD. Use our claim on monotonicity of GD.]
2. For $\alpha > 0$, a smooth function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is α -convex if for all $x, y \in \mathbb{R}^n$,

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\alpha}{2} \|y - x\|_2^2.$$

In this exercise we will show a better bound on convergence rate of gradient descent (GD) for such functions.

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a α -convex β -smooth function. Let x^* be an optimal minimizer for f . Consider the GD algorithm with starting point x_0 and step-size $t = 1/\beta$. Show that after k iterations we have

$$\|x_k - x^*\|_2^2 \leq \|x_0 - x^*\|_2^2 \left(1 - \frac{\alpha}{\beta}\right)^k.$$

In other words, the distance of the k 'th iterate to the optimal decreases exponentially in the number of iterations. [5 points]

[Hint: You may use the fact that $\nabla f(x^*) = 0$.]

Remark: Note that the number of iterations to get error ε for such functions is $O(\log(1/\varepsilon))$ for fixed x_0, α, β ; this is exponentially better than the bound for general convex functions which was $O(1/\varepsilon)$.

3. Show that for a differentiable convex function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, every local minimum, i.e., any point x with $\nabla f(x) = 0$, is also a global minimum. [3 points]
4. The goal of this exercise is to implement and compare GD, and Nesterov's accelerated gradient descent (NAGD) for *logistic regression*.

Suppose we have binary data, i.e., the labels are 0 or 1. When trying to *fit* binary labels with linear predictors, a commonly used loss function is the *logit loss* function: $\ell(a, b) = -a \log b - (1 - a) \log(1 - b)$ (also known as *cross-entropy* loss etc.).

One commonly used parameterized predictor family is $h_w(x) = \sigma(w \cdot x)$, where $\sigma(t) = 1/(1 + e^{-t})$ is the sigmoid function. Combining these two, given a dataset $(x_1, y_1), \dots, (x_n, y_n)$, the corresponding ERM optimization problem is to minimize

$$L(w) = (1/n) \sum_{i=1}^n \ell(y_i, \sigma(\langle w, x_i \rangle)).$$

The above formulation is called *logistic regression*. Your goal is to implement variants of GD for the above loss function and test it some random dataset.

Generating dataset. Set $n = 1000, d = 20$ (or higher too if you want). Pick a random 'hidden vector' $w_* \in \mathbb{R}^d$ of unit-norm. As in the workspace examples, generate random data as follows: For $i = 1, \dots, n$: Generate a random Gaussian vector $x_i \in \mathbb{R}^d$ and after that, independently set $y_i = 1$ with probability $\sigma(\langle w_*, x_i \rangle)$ and 0 with probability $1 - \sigma(\langle w_*, x_i \rangle)$. Sample code for this is provided at the bottom of GD workspace on edStem.

- (a) Write down a formula for the gradient of $L(w)$. [1 point]
- (b) Implement GD, NAGD with various step-sizes (say 3 step-sizes each) for several hundred iterations to get convergence and plot the loss values. Also plot the distance to the hidden vector w_* of the iterates under the three methods. [4 points]

(Don't ask why this is the **right** statistical model, but you can check the wiki page for details. We are only interested in optimization and not on the statistical aspects.)

Your submission on gradescope should be the following.

- i. Screenshots of the implementations of the three algorithms (the gradient subroutine and other important parts).
- ii. Plot of the values of f against the number of iteration for GD and NAGD on the same figure.
- iii. Plot of the distance of the current iterate w_k to w_* under GD and NAGD on the same figure.

You may use/copy any of the code from the workspaces as you need (especially for plotting and generating data). Make sure your plots are well-labeled (the axes, and with an appropriate legend).