

Optimization for Machine Learning HW 6

Shuyue Jia
BUID: U62343813

Due: 11/17/2023

All parts of each question are equally weighted. When solving one question/part, you may assume the results of all previous questions/parts. You may also assume all previous homework results and results from class or lecture notes, but please explain which result you are using when you use it.

1. We have seen that for H -smooth convex objectives, no first-order algorithm can faster than $O(H/T^2)$ in a *dimension-free* manner. That is, the “hard” function we studied is a very high dimensional function. However, if we restrict to considering $\mathcal{L} : \mathbb{R} \rightarrow \mathbb{R}$, the situation is quite different. Provide a first-order algorithm that, given a first-order oracle for a convex function $\mathcal{L} : \mathbb{R} \rightarrow \mathbb{R}$ such that \mathcal{L} is H -smooth and achieves its minimum at some $|w_*| \leq 1$, then after T iterations the algorithm outputs \hat{w} that satisfies $\mathcal{L}(\hat{w}) - \mathcal{L}(w_*) \leq O(H2^{-2T})$. Be careful: you need to have the 2 in the exponent since 2^{-T} is NOT $O(2^{-2T})$.

Solution:

The intuition is to design a **Bisection method** that meets the above requirements:

Algorithm 1 The Designed Bisection Algorithm

Require: First-order oracle for convex function, T, a, b

```
1: Initialize  $a \leftarrow -1, b \leftarrow 1$ 
2: for  $t = 1$  to  $T$  do
3:    $c \leftarrow \frac{a+b}{2}$    % Midpoint of the interval
4:    $\text{gradient}_c \leftarrow \text{oracle.gradient}(c)$    % Gradient at the midpoint
5:   if  $\text{gradient}_c > 0$  then
6:      $b \leftarrow c$ 
7:   else
8:      $a \leftarrow c$ 
9:   end if
10: end for
11: return  $\frac{a+b}{2}$    % Final approximation
```

Firstly, the bisection method halves the interval $[a, b]$ in each iteration, so after T iterations, the interval size is reduced to $(b - a)/2^T$.

Since the function \mathcal{L} is H -smooth, the absolute value of its gradient is bounded by H everywhere. Thus, the Lipschitz continuity of the gradient implies that $\mathcal{L}(c) - \mathcal{L}(w_*) \leq \frac{H}{2}(b - a)$ for any $c \in [a, b]$.

By combining steps 1 and 2, we have $\mathcal{L}(c) - \mathcal{L}(w_*) \leq \frac{H}{2^T}(b - a)$.

Then, let \hat{w} be the final approximation returned by the bisection method. Since \hat{w} lies in the final interval $[a, b]$, we can write $\mathcal{L}(\hat{w}) - \mathcal{L}(w_*) \leq \frac{H}{2^T}(b - a)$.

Next, by substituting $b - a$ with the initial interval size, we get $\mathcal{L}(\hat{w}) - \mathcal{L}(w_*) \leq \frac{H}{2^T}(b - a) = \frac{H}{2^T}(2) = \frac{H}{2^{T-1}}$.

Finally, since $2^{T-1} = 2^T/2$, we can rewrite the result as $\mathcal{L}(\hat{w}) - \mathcal{L}(w_*) \leq O(H2^{-2T})$, as desired. \square

2. Suppose you are trying to identify the bias of a coin. We model a coin flip as a “1” if it comes up heads, and “0” otherwise, and let p_* be the probability that it comes up heads. If $Z \in \{0, 1\}$ is the outcome of a coin flip, it holds that $p_* = \operatorname{argmin} \mathbb{E}[\ell(w, Z)] = \mathcal{L}(w)$ where $\ell(w, z) = (w - z)^2$. After observing T coin flips z_1, \dots, z_T , you make the natural prediction $\hat{p} = \frac{z_1 + \dots + z_T}{T}$. Show that $\mathbb{E}[\mathcal{L}(\hat{p}) - \mathcal{L}(p_*)] = \frac{p_*(1-p_*)}{T}$. Explain why this does *not* contradict our $\frac{1}{\sqrt{T}}$ lower bound for stochastic convex optimization?

Solution:

(1) **Proof.** The loss function is given by $\ell(w, z) = (w - z)^2$, and our goal is to find the probability p_* that minimizes the expected loss $\mathbb{E}[\ell(w, Z)]$, where Z , a Bernoulli random variable, is the outcome of a coin flip.

$$\mathbb{E}[(z_1 - p_*)^2] = \operatorname{Var}(Z) = p_*(1 - p_*). \quad (1)$$

The expected loss is given by

$$\mathbb{E}[\ell(w, Z)] = \sum_{z \in \{0, 1\}} \ell(w, z)P(Z = z). \quad (2)$$

For our coin flip, p_* minimizes this expectation. Now, let's calculate the expected loss for the natural prediction $\hat{p} = \frac{z_1 + \dots + z_T}{T}$ after observing T coin flips z_1, \dots, z_T .

Then, we will have:

$$\begin{aligned} \mathbb{E}[\mathcal{L}(\hat{p}) - \mathcal{L}(p_*)] &= \mathbb{E} \left[\left(\frac{z_1 + \dots + z_T}{T} - p_* \right)^2 - (p_* - p_*)^2 \right] \\ &= \frac{1}{T^2} \mathbb{E}[(z_1 + \dots + z_T - Tp_*)^2] \\ &= \frac{1}{T^2} \sum_{t=1}^T \mathbb{E}[(z_t - p_*)^2] \quad (\text{by linearity of expectation}) \\ &= \frac{1}{T^2} T \mathbb{E}[(z_1 - p_*)^2] \quad (\text{all terms are identical}) \\ &= \frac{1}{T} \mathbb{E}[(z_1 - p_*)^2] \\ &= \frac{1}{T} p_*(1 - p_*). \end{aligned} \quad (3)$$

□

(2) **Explanation** for why this does not contradict the $\frac{1}{\sqrt{T}}$ lower bound for stochastic convex optimization:

The lower bound of $\frac{1}{\sqrt{T}}$ is a convergence rate for optimization problems where the goal is to find the optimal parameter w_* in a convex function. In this case, we are not optimizing a convex function directly; instead, we are estimating a probability parameter p_* based on observed coin flips.

The rate $\frac{1}{\sqrt{T}}$ arises in optimization problems where the objective function is smooth and has a Lipschitz continuous gradient. In the coin flip scenario, the loss function is not differentiable at $p_* = 0$ and $p_* = 1$, so it doesn't satisfy the smoothness conditions typically assumed in optimization settings.