

Gradient descent

Daniel Hsu

September 2, 2023

1 Smooth functions

Smooth functions are functions whose derivatives (gradients) do not change too quickly. The change in the derivative is the second-derivative, so smoothness is a constraint on the second-derivatives of a function.

We say a twice-differentiable function $J: \mathbb{R}^d \rightarrow \mathbb{R}$ is β -smooth if the eigenvalues of its Hessian matrix at any point in \mathbb{R}^d are at most β .

A consequence of β -smoothness is the following. Recall that by Taylor's theorem, for any $w, \delta \in \mathbb{R}^d$, there exists $\tilde{w} \in \mathbb{R}^d$ on the line segment between w and $w + \delta$ such that

$$J(w + \delta) = J(w) + \nabla J(w)^\top \delta + \frac{1}{2} \delta^\top \nabla^2 J(\tilde{w}) \delta.$$

If J is β -smooth, then we can bound the third term from above as

$$\begin{aligned} \frac{1}{2} \delta^\top \nabla^2 J(\tilde{w}) \delta &\leq \frac{1}{2} \|\delta\|^2 \max_{u \in \mathbb{R}^d: \|u\|=1} u^\top \nabla^2 J(\tilde{w}) u \\ &\leq \frac{1}{2} \|\delta\|^2 \lambda_{\max}(\nabla^2 J(\tilde{w})) \\ &\leq \frac{1}{2} \|\delta\|^2 \beta. \end{aligned}$$

Therefore, if J is β -smooth, then for any $w, \delta \in \mathbb{R}^d$,

$$J(w + \delta) \leq J(w) + \nabla J(w)^\top \delta + \frac{\beta}{2} \|\delta\|^2.$$

2 Gradient descent on smooth objectives

Gradient descent starts with an initial point $w^{(0)} \in \mathbb{R}^d$, and for a given step size η , iteratively computes a sequence of points $w^{(1)}, w^{(2)}, \dots$ as follows. For $t = 1, 2, \dots$:

$$w^{(t)} = w^{(t-1)} - \eta \nabla J(w^{(t-1)}).$$

2.1 Motivation

The motivation for the gradient descent update is the following. Suppose we have a current point $w \in \mathbb{R}^d$, and we would like to locally change it from w to $w + \delta$ so as to decrease the objective value. How should we choose δ ?

In gradient descent, we consider the quadratic upper-bound that smoothness grants, i.e.,

$$J(w + \delta) \leq J(w) + \nabla J(w)^\top \delta + \frac{\beta}{2} \|\delta\|_2^2,$$

and then choose δ to minimize this upper-bound. The upper-bound is a convex quadratic function of δ , so its minimizer can be written in closed-form. The minimizer is the value of δ such that

$$\nabla J(w) + \beta \delta = 0.$$

In other words, it is $\delta^*(w)$, defined by

$$\delta^*(w) = -\frac{1}{\beta} \nabla J(w).$$

Plugging in $\delta^*(w)$ for δ in the quadratic upper-bound gives

$$\begin{aligned} J(w + \delta^*(w)) &\leq J(w) + \nabla J(w)^\top \delta^*(w) + \frac{\beta}{2} \|\delta^*(w)\|_2^2 \\ &= J(w) - \frac{1}{\beta} \nabla J(w)^\top \nabla J(w) + \frac{1}{2\beta} \|\nabla J(w)\|_2^2 \\ &= J(w) - \frac{1}{2\beta} \|\nabla J(w)\|_2^2. \end{aligned}$$

This inequality tells us that this local change to w will decrease the objective value as long as the gradient at w is non-zero. It turns out that if the function J is convex (in addition to β -smooth), then repeatedly making such local changes is sufficient to approximately minimize the function.

2.2 Analysis for smooth convex objectives

One of the simplest ways to mathematically analyze the behavior of gradient descent on smooth functions (with step size $\eta = 1/\beta$) is to monitor the change in a “potential function” during the execution of gradient descent. The potential function we will use is the squared Euclidean distance to a fixed vector $w^* \in \mathbb{R}^d$, which could be a minimizer of J (but need not be):

$$\Phi(w) = \frac{1}{2\eta} \|w - w^*\|_2^2.$$

The scaling by $\frac{1}{2\eta}$ is used just for notational convenience.

Let us examine the “drop” in the potential when we change a point w to $w + \delta^*(w)$ (as in gradient descent):

$$\begin{aligned} \Phi(w) - \Phi(w + \delta^*(w)) &= \frac{1}{2\eta} \|w - w^*\|_2^2 - \frac{1}{2\eta} \|w + \delta^*(w) - w^*\|_2^2 \\ &= \frac{\beta}{2} \|w - w^*\|_2^2 - \frac{\beta}{2} (\|w - w^*\|_2^2 + 2\delta^*(w)^\top(w - w^*) + \|\delta^*(w)\|_2^2) \\ &= -\beta\delta^*(w)^\top(w - w^*) - \frac{\beta}{2} \|\delta^*(w)\|_2^2 \\ &= \nabla J(w)^\top(w - w^*) - \frac{1}{2\beta} \|\nabla J(w)\|_2^2. \end{aligned}$$

In the last step, we have plugged in $\delta^*(w) = -\frac{1}{\beta}\nabla J(w)$. Now we use two key facts. The first is the inequality we derived above based on the smoothness of J :

$$J(w + \delta^*(w)) \leq J(w) - \frac{1}{2\beta} \|\nabla J(w)\|_2^2,$$

which rearranges to

$$-\frac{1}{2\beta} \|\nabla J(w)\|_2^2 \geq J(w + \delta^*(w)) - J(w).$$

The second comes from the first-order definition of convexity:

$$J(w^*) \geq J(w) + \nabla J(w)^\top(w^* - w),$$

which rearranges to

$$\nabla J(w)^\top(w - w^*) \geq J(w) - J(w^*).$$

So, we can bound the drop in potential as follows:

$$\begin{aligned}\Phi(w) - \Phi(w + \delta^*(w)) &= \nabla J(w)^\top (w - w^*) - \frac{1}{2\beta} \|\nabla J(w)\|_2^2 \\ &\geq (J(w) - J(w^*)) + (J(w + \delta^*(w)) - J(w)) \\ &= J(w + \delta^*(w)) - J(w^*).\end{aligned}$$

Let us write this inequality in terms of the iterates of gradient descent with $\eta = 1/\beta$:

$$\Phi(w^{(t-1)}) - \Phi(w^{(t)}) \geq J(w^{(t)}) - J(w^*).$$

Summing this inequality from $t = 1, 2, \dots, T$:

$$\sum_{t=1}^T \left(\Phi(w^{(t-1)}) - \Phi(w^{(t)}) \right) \geq \sum_{t=1}^T \left(J(w^{(t)}) - J(w^*) \right).$$

The left-hand side simplifies to $\Phi(w^{(0)}) - \Phi(w^{(T)})$. Furthermore, since $J(w^{(t)}) \geq J(w^{(T)})$ for all $t = 1, \dots, T$, the right-hand side can be bounded from below by

$$T \left(J(w^{(T)}) - J(w^*) \right).$$

So we are left with the inequality

$$J(w^{(T)}) - J(w^*) \leq \frac{1}{T} \left(\Phi(w^{(0)}) - \Phi(w^{(T)}) \right) = \frac{\beta}{2T} \left(\|w^{(0)} - w^*\|_2^2 - \|w^{(T)} - w^*\|_2^2 \right).$$

3 Gradient descent on non-smooth objectives

Gradient descent can also be used for non-smooth convex functions as long as the function itself does not change too quickly.

We say that a differentiable function $J: \mathbb{R}^d \rightarrow \mathbb{R}$ is *L-Lipschitz* if its gradient at any point in \mathbb{R}^d is bounded in Euclidean norm by L .

The motivation for gradient descent based on minimizing quadratic upper-bounds no longer applies. Indeed, the gradient at w could be very different from the gradient at a nearby w' , so the function value at $w - \eta \nabla J(w)$ could be worse than the function value at w . Therefore, we cannot expect to have the same convergence guarantee for non-smooth functions that we had for smooth functions.

Gradient descent, nevertheless, will produce a sequence $w^{(1)}, w^{(2)}, \dots$ such that the function value at these points is approximately minimal “on average”.

3.1 Motivation

A basic motivation for gradient descent for convex functions, that does not assume smoothness, comes from the first-order condition for convexity:

$$J(w^*) \geq J(w) + \nabla J(w)^\top(w^* - w),$$

which rearranges to

$$(-\nabla J(w))^\top(w^* - w) \geq J(w) - J(w^*).$$

Suppose $J(w) > J(w^*)$, so that moving from w to w^* would improve the function value. Then, the inequality implies that the negative gradient $-\nabla J(w)$ at w makes a positive inner product with the direction from w to w^* . This is the crucial property that makes gradient descent work.

3.2 Analysis

We again monitor the change in the potential function

$$\Phi(w) = \frac{1}{2\eta} \|w - w^*\|_2^2,$$

for a fixed vector $w^* \in \mathbb{R}^d$.

Again, let us examine the “drop” in the potential when we change a point w to $w - \eta \nabla J(w)$ (as in gradient descent):

$$\begin{aligned} \Phi(w) - \Phi(w - \eta \nabla J(w)) &= \frac{1}{2\eta} \|w - w^*\|_2^2 - \frac{1}{2\eta} \|w - \eta \nabla J(w) - w^*\|_2^2 \\ &= (-\nabla J(w))^\top(w - w^*) - \frac{\eta}{2} \|\nabla J(w)\|_2^2 \\ &\geq J(w) - J(w^*) - \frac{L^2\eta}{2}, \end{aligned}$$

where the inequality uses the convexity and Lipschitzness of J . In terms of the iterates of gradient descent, this reads

$$\Phi(w^{(t-1)}) - \Phi(w^{(t)}) \geq J(w^{(t-1)}) - J(w^*) - \frac{L^2\eta}{2}.$$

Summing this inequality from $t = 1, 2, \dots, T$:

$$\Phi(w^{(0)}) - \Phi(w^{(T)}) \geq \sum_{t=1}^T \left(J(w^{(t-1)}) - J(w^*) \right) - \frac{L^2\eta T}{2}.$$

Rearranging and dividing through by T (and dropping a term):

$$\frac{1}{T} \sum_{t=1}^T \left(J(w^{(t-1)}) - J(w^*) \right) \leq \frac{\|w^{(0)} - w^*\|_2^2}{2\eta T} + \frac{L^2\eta}{2}.$$

The left-hand side is the average sub-optimality relative to $J(w^*)$. Therefore, there exists some $t^* \in \{0, 1, \dots, T-1\}$ such that

$$J(w^{(t^*)}) - J(w^*) \leq \frac{1}{T} \sum_{t=1}^T \left(J(w^{(t-1)}) - J(w^*) \right) \leq \frac{\|w^{(0)} - w^*\|_2^2}{2\eta T} + \frac{L^2\eta}{2}.$$

The right-hand side is $O(1/\sqrt{T})$ when we choose $\eta = 1/\sqrt{T}$.