

# Gradient descent

Daniel Hsu (COMS 4771)

## 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.

For any  $\beta > 0$ , we say a twice-differentiable function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is  $\beta$ -smooth if the eigenvalues of its Hessian matrix at any point in  $\mathbb{R}^d$  are at most  $\beta$ .

### Example: logistic regression

Consider the empirical logistic loss risk on a training data set  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^d \times \{-1, +1\}$ :

$$\hat{\mathcal{R}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i \mathbf{x}_i^\top \mathbf{w}))$$

The Hessian of  $\hat{\mathcal{R}}$  at  $\mathbf{w}$

$$\nabla^2 \hat{\mathcal{R}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \sigma(y_i \mathbf{x}_i^\top \mathbf{w}) \mathbf{x}_i \mathbf{x}_i^\top$$

where  $\sigma(t) = 1/(1 + \exp(-t))$  is the sigmoid function. For any unit vector  $\mathbf{u}$

$$\mathbf{u}^\top \nabla^2 \hat{\mathcal{R}}(\mathbf{w}) \mathbf{u} = \frac{1}{n} \sum_{i=1}^n \sigma(y_i \mathbf{x}_i^\top \mathbf{w}) (\mathbf{x}_i^\top \mathbf{u})^2$$

$$\leq \frac{1}{4n} \sum_{i=1}^n (\mathbf{x}_i^\top \mathbf{u})^2$$

$$= \frac{1}{4} \mathbf{u}^\top \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{u}$$

$$\leq \frac{1}{4} \lambda_{\max} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right),$$

where  $\lambda_{\max}(\mathbf{M})$  is used to denote the largest eigenvalue of a symmetric matrix  $\mathbf{M}$ . So if  $\lambda_1$  is the largest eigenvalue of the empirical second moment matrix  $\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$ , then  $\hat{\mathcal{R}}$  is  $\beta$ -smooth for  $\beta = \lambda_1/4$ .

### Quadratic upper bound for smooth functions

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

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

If  $f$  is  $\beta$ -smooth, then we can bound the third term from above as

$$\begin{aligned}\frac{1}{2}\delta^\top \nabla^2 f(\tilde{\mathbf{w}})\delta &\leq \frac{1}{2}\|\delta\|_2^2 \max_{\mathbf{u} \in \mathbb{R}^d: \|\mathbf{u}\|_2=1} \mathbf{u}^\top \nabla^2 f(\tilde{\mathbf{w}})\mathbf{u} \\ &\leq \frac{1}{2}\|\delta\|_2^2 \lambda_{\max}(\nabla^2 f(\tilde{\mathbf{w}})) \\ &\leq \frac{1}{2}\|\delta\|_2^2 \beta.\end{aligned}$$

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

$$f(\mathbf{w} + \delta) \leq f(\mathbf{w}) + \nabla f(\mathbf{w})^\top \delta + \frac{\beta}{2}\|\delta\|_2^2.$$

## Gradient descent on smooth functions

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

$$\mathbf{w}^{(t)} := \mathbf{w}^{(t-1)} - \eta \nabla f(\mathbf{w}^{(t-1)}).$$

### Motivation for gradient descent on smooth functions

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

In gradient descent, we bound

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

and then choose  $\delta$  to minimize this upper-bound. The upper-bound is quadratic in  $\delta$ , so its minimizer can be written in closed-form. The minimizer is the value of

$$\nabla f(\mathbf{w}) + \beta \delta = \mathbf{0}.$$

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

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

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

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

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

### Analysis of gradient descent on smooth convex functions

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  $\mathbf{w}^* \in \mathbb{R}^d$ , which could be a minimizer of  $f$  (but need not be):

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

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  $\mathbf{w}$  to  $\mathbf{w} + \boldsymbol{\delta}^*(\mathbf{w})$  (as in gradient descent):

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

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

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

which rearranges to

$$\frac{1}{2\beta} \|\nabla f(\mathbf{w})\|_2^2 \leq f(\mathbf{w}) - f(\mathbf{w} + \boldsymbol{\delta}^*(\mathbf{w})).$$

The second comes from the first:

<https://eduassistpro.github.io/>

which rearranges to

Add WeChat <https://eduassistpro.github.io/>

(We’ll discuss this inequality more later.) So, we can bound the drop in potent

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

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

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

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

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

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

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

So we are left with the inequality

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

## Gradient descent on Lipschitz convex functions

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

For any  $L > 0$ , we say that a differentiable function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  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  $\mathbf{w}$  could be very different from the gradient at a nearby  $\mathbf{w}'$ , so the function value at  $\mathbf{w} - \eta \nabla f(\mathbf{w})$  could be worse than the function value at  $\mathbf{w}$ . 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  $\mathbf{w}^{(1)}, \mathbf{w}^{(2)}, \dots$  such that the function value at these points is approximately minimal *on average*.

## Motivation for gradient descent on Lipschitz convex functions

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

$$f(\mathbf{w}^*) \geq f(\mathbf{w}) + \nabla f(\mathbf{w})^\top (\mathbf{w}^* - \mathbf{w}),$$

which rearranges to

Suppose  $f(\mathbf{w}) > f(\mathbf{w}^*)$ , so that moving from  $\mathbf{w}$  to  $\mathbf{w}^*$  would improve the function value. Then, the inequality implies that the negative gradient  $-\nabla f(\mathbf{w})$  is a direction from  $\mathbf{w}$  to  $\mathbf{w}^*$ . This is the crucial property.

## Analysis of gradient descent

We again monitor the change in the potential function

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

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

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

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

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

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

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

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

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

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

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

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

The right-hand side is  $O(1/\sqrt{T})$  when we choose  $\eta = 1/\sqrt{T}$ .<sup>1</sup> Alternatively, if we compute the average point

$$\bar{\mathbf{w}}_T := \frac{1}{T} \sum_{t=1}^T \mathbf{w}^{(t-1)},$$

then by Jensen's inequality we have

$$f(\bar{\mathbf{w}}_T) = f\left(\frac{1}{T} \sum_{t=1}^T \mathbf{w}^{(t-1)}\right) \leq \frac{1}{T} \sum_{t=1}^T f(\mathbf{w}^{(t-1)}).$$

So the bound for  $\mathbf{w}^{(t^*)}$  also applies to  $\bar{\mathbf{w}}_T$ :

$$f(\bar{\mathbf{w}}_T) \leq \frac{\|\mathbf{w}^{(0)} - \mathbf{w}^*\|_2^2}{2\eta T} + \frac{L^2\eta}{2}.$$

**Assignment Project Exam Help**  
<https://eduassistpro.github.io/>  
**Add WeChat edu\_assist\_pro**

---

<sup>1</sup>A similar guarantee holds when the step size used for the  $t$ -th update is  $\eta_t = 1/\sqrt{t}$ .