

The method described in this section require a suitable starting point  $x^{(0)}$ . The starting point must lie in  $\text{dom}f$ , and in addition the sublevel set

$$S = \left\{x \in \text{dom}f : f(x) \leq f(x^{(0)})\right\}$$

must be closed. This condition is satisfied for all  $x^{(0)} \in \text{dom}f$  if the function  $f$  is closed. Continuous functions with  $\text{dom}(f) = \mathbb{R}^n$  are closed, so if  $\text{dom}(f) = \mathbb{R}^n$ , the initial sublevel set condition is satisfied by any  $x^{(0)}$ .

**Theorem 1.** Assume that  $f$  convex and differentiable, with  $\text{dom}(f) = \mathbb{R}^n$  and  $\nabla f$  is Lipschitz continuous with constant  $L > 0$ , i.e.

$$\|\nabla f(x) - \nabla f(y)\|_2 \leq L\|x - y\|_2 \quad \forall x, y$$

then the gradient descent with fixed step size  $t \leq 1/L$  satisfies

$$f(x^{(k)}) - f^* \leq \frac{\|x^{(0)} - x^*\|}{2tk}$$

We say that the gradient descent has convergence rate  $O(1/k)$ .

*Proof. Part I:* With  $\nabla f$  Lipschitz constant  $L$ , we have that

$$f(y) \leq f(x) + \nabla f(x)^T(y - x) + \frac{L}{2}\|x - y\|_2^2 \quad \forall x, y \quad (1)$$

Suppose we are at  $x$  in the gradient descent and the next iteration go to

$$x^+ = x - t\nabla f(x)$$

We can use the above inequality with  $y = x^+$  and

$$\begin{aligned} f(x^+) &\leq f(x) + \nabla f(x)^T(-t\nabla f(x)) + \frac{L}{2}\| -t\nabla f(x)\|_2^2 \\ &= f(x) - t\|\nabla f(x)\|_2^2 + \frac{Lt^2}{2}\|\nabla f(x)\|_2^2 \\ &= f(x) - \left(1 - \frac{Lt}{2}\right)t\|\nabla f(x)\|_2^2 \end{aligned}$$

If  $0 \leq t \leq 1/L$ , we get  $-t + \frac{Lt^2}{2} \leq \frac{-t}{2}$  which gives us that

$$f(x^+) \leq f(x) - \frac{t}{2}\|\nabla f(x)\|_2^2. \quad (2)$$

This result also implies the descent property of the gradient descent algorithm

$$f(x^+) \leq f(x).$$

**Part II:** Use convexity of  $f$ , we know that

$$\begin{aligned} f(x^*) &\geq f(x) + \nabla f(x)^T(x^* - x) \\ f(x) &\leq f(x^*) - \nabla f(x)^T(x^* - x) \end{aligned} \quad (3)$$

Plugin (3) into (2) and you get

$$\begin{aligned} f(x^+) &\leq f(x^*) + \nabla f(x)^T(x - x^*) - \frac{t}{2} \|\nabla f(x)\|_2^2 \\ f(x^+) - f(x^*) &\leq \nabla f(x)^T(x - x^*) - \frac{t}{2} \|\nabla f(x)\|_2^2 \\ &= \frac{1}{2t} (\|x - x^*\|_2^2 - \|x^+ - x^*\|_2^2) \end{aligned}$$

The last equality is true because

$$\begin{aligned} \frac{1}{2t} (\|x - x^*\|_2^2 - \|x - t\nabla f(x) - x^*\|_2^2) &= \frac{1}{2t} (\|x - x^*\|_2^2 - \|x - x^*\|_2^2 + 2t\nabla f(x)^T(x - x^*) - t^2\|\nabla f(x)\|_2^2) \\ &= \nabla f(x)^T(x - x^*) - \frac{t}{2} \|\nabla f(x)\|_2^2 \end{aligned}$$

Finally,

$$\begin{aligned} f(x^{(i)}) - f(x^*) &\leq \frac{1}{2t} (\|x^{(i-1)} - x^*\|_2^2 - \|x^{(i)} - x^*\|_2^2) \\ \sum_{i=1}^k (f(x^{(i)}) - f(x^*)) &\leq \frac{1}{2t} (\|x^{(0)} - x^*\|_2^2 - \|x^{(k)} - x^*\|_2^2) \leq \frac{1}{2t} \|x^{(0)} - x^*\|_2^2 \end{aligned}$$

because we've proved that  $f(x^{(0)}) \geq f(x^{(1)}) \geq \dots \geq f(x^{(k)})$ . Thus

$$f(x^{(k)}) - f(x^*) \leq \frac{1}{k} \sum_{i=1}^k (f(x^{(i)}) - f(x^*)) \leq \frac{\|x^{(0)} - x^*\|_2^2}{2tk}$$

□

*Remark 1.* We can show that in Theorem 1, the assumption that  $\nabla f$  is Lipschitz continuous with constant  $L > 0$  can be relaxed to that we only need Lipschitz gradient over the sublevel set

$$S = \{x \in \text{dom} f : f(x) \leq f(x^{(0)})\}.$$

**Theorem 2.** *If the sublevel sets contained in  $S$  are bounded, so in particular, if  $S$  is bounded. Then  $\nabla f$  is Lipschitz continuous with constant  $L > 0$  over  $S$ .*

*Proof.* If  $S$  is bounded, then the maximum eigenvalue of  $\nabla^2 f(x)$ , which is a continuous function of  $x$  on  $S$ , is also bounded above on  $S$ . i.e., there exist a constant  $L$  such that

$$\nabla^2 f(x) \preceq LI \quad \forall x \in S.$$

This upper bound on the Hessian implies for any  $x, y \in S$

$$f(y) \leq f(x) + \nabla f(x)^T(y - x) + \frac{L}{2} \|y - x\|_2^2$$

Therefore we get a similar condition to the original Lipschitz continuous assumption (1) except that it is on the sublevel set  $S$ , which is sufficient to prove Theorem 1 since this condition can also lead to the descent property on the sublevel set

$$f(x^{(1)}) \leq f(x^{(0)}) - \frac{t}{2} \|\nabla f(x^{(0)})\|_2^2 \quad \forall x \in S$$

□