

Mathematics of Image and Data Analysis

Math 5467

Gradient Descent

Instructor: Jeff Calder

Email: jcalder@umn.edu

<http://www-users.math.umn.edu/~jwcalder/5467>

Announcements

- HW2 due Feb 25

Last time

- PageRank

Today

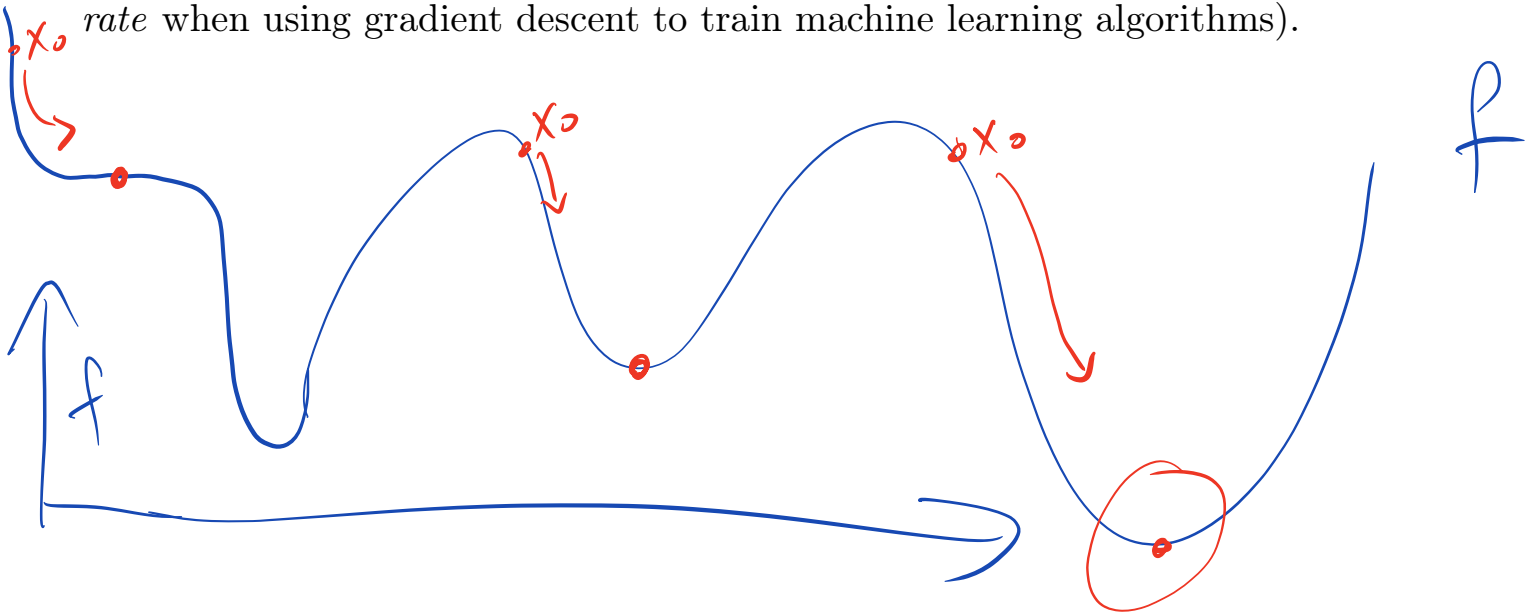
- Gradient Descent

Gradient Descent

Gradient descent is one of the most important algorithms in many areas of science and engineering. To minimize an objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, gradient descent iterates

$$(1) \quad x_{k+1} = x_k - \alpha \nabla f(x_k)$$

until convergence. The parameter $\alpha > 0$ is the time step (often called the *learning rate* when using gradient descent to train machine learning algorithms).



$$\nabla (v^T x) = v$$

Optimization interpretation

Exercise 1. Fix x_k and define

$$(2) \quad T(x) = \overbrace{f(x_k) + \nabla f(x_k)^T (x - x_k)}^{\text{hyper tangent plane at } x_k} + \underbrace{\frac{1}{2\alpha} \|x - x_k\|^2}_{\text{penalty}}.$$

If we define x_{k+1} as the minimizer of T , show that

$$x_{k+1} = x_k - \alpha \nabla f(x_k).$$

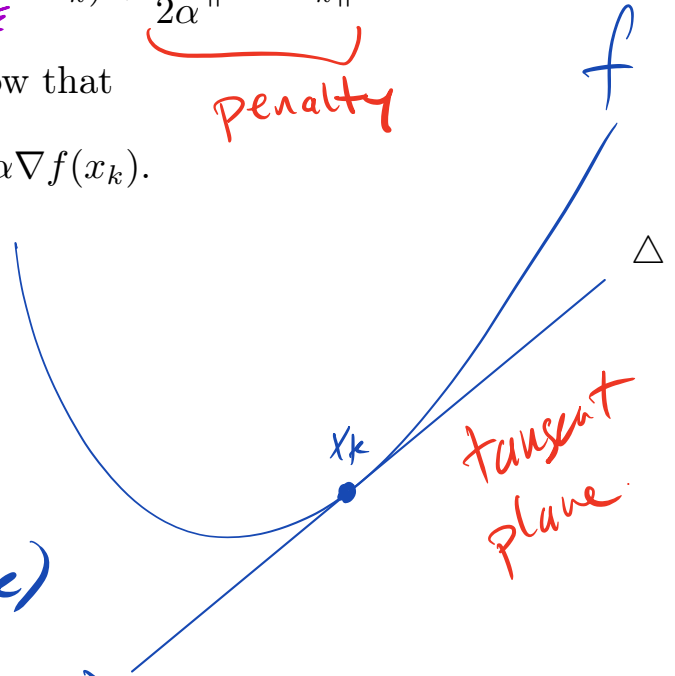
$$\nabla T(x) = 0$$

$$\nabla T = \nabla f(x_k) + \frac{1}{\alpha} (x - x_k)$$

$$\nabla T = 0$$

$$\hookrightarrow \frac{1}{\alpha} (x - x_k) = -\nabla f(x_k)$$

$$x_{k+1} = x = x_k - \alpha \nabla f(x_k).$$



Assumptions on f

We assume the objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a smooth function that admits a global minimizer $x_* \in \mathbb{R}^n$. That is

$$f(x_*) \leq f(x)$$

for all $x \in \mathbb{R}^n$. We denote the optimal value of f by $f_* := f(x_*)$.

Sublinear convergence rate

We say ∇f is *L-Lipschitz continuous* if

$$(3) \quad \|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\| \quad \text{for all } x, y \in \mathbb{R}^n.$$

Theorem 2. Assume ∇f is *L-Lipschitz* and that $\alpha \leq \frac{1}{L}$. Then for any integer $t \geq 1$ we have

$$(4) \quad \min_{0 \leq k \leq t} \|\nabla f(x_k)\|^2 \leq \frac{2(f(x_0) - f_*)}{\alpha t}.$$

Remark 3. The theorem says, with very few assumptions on f , that gradient descent converges at a rate of $O\left(\frac{1}{t}\right)$ to a critical point of f , in the sense that $\nabla f \sim \frac{1}{t} \rightarrow 0$. Since f is not assumed to be convex, critical points need not be minimizers and could also include saddle points.

$$f: \mathbb{R}^n \rightarrow \mathbb{R}, \quad x, y \in \mathbb{R}^n$$

Proof: Review Taylor expansion.

$$g(t) = f(x + t(y-x)), \quad t \in \mathbb{R}$$

$$g(1) = g(0) + \int_0^1 g'(t) dt \quad (\text{FTC})$$

$$= g(0) + \int_0^1 g'(0) dt + \int_0^1 g'(t) - g'(0) dt$$

$$g(1) = g(0) + g'(0) + R$$

$$R = \int_0^1 g'(t) - g'(0) dt$$

$$\begin{aligned}g'(t) &= \frac{d}{dt} f(x + t(y-x)) \\ &= \nabla f(x + t(y-x))^T (y-x)\end{aligned}$$

$$g'(0) = \nabla f(x)^T (y-x)$$

$$g(1) = g(0) + g'(0) + R$$

$$f(y) = f(x) + \nabla f(x)^T (y-x) + R$$

$$\begin{aligned}|R| &= \left| \int_0^1 g'(t) - g'(0) dt \right| \\ &\leq \int_0^1 |g'(t) - g'(0)| dt\end{aligned}$$

$$= \int_0^1 | \nabla f(x+t(y-x))^T (y-x) - \nabla f(x)^T (y-x) | dt$$

$$(\neq) = \int_0^1 | [\nabla f(x+t(y-x)) - \nabla f(x)]^T (y-x) | dt$$

Aside Cauchy-Schwarz: $v^T w \leq \underbrace{\|v\| \cdot \|w\|}_{=1}$

Proof: Assume $\|v\| = 1 = \|w\|$

$$\begin{aligned} 0 \leq \|v-w\|^2 &= \|v\|^2 - 2v^T w + \|w\|^2 \\ &= 2 - 2v^T w = 2(1 - v^T w) \end{aligned}$$

$$\Rightarrow 1 - v^T w \geq 0 \Rightarrow v^T w \leq 1$$

TUG

(*) Cauchy-Schwarz gives

$$|R| \leq \int_0^1 \underbrace{\|\nabla f(x + t(y-x)) - \nabla f(x)\|}_{\leq L \|x + t(y-x) - x\|} \cdot \|y-x\| dt$$

$$\leq L \|x + t(y-x) - x\|$$

since
 ∇f is
 L -Lipschitz

$$\leq L \|x-y\| \int_0^1 \|t(y-x)\| dt$$

$$= L \|x-y\|^2 \int_0^1 t dt = \frac{L}{2} \|x-y\|^2$$

$$f(y) = f(x) + \nabla f(x)^T (y-x) + R$$

where

$$|R| \leq \frac{L}{2} \|x - y\|^2$$

GD. $x_{k+1} - x_k = -\alpha \nabla f(x_k)$

Theorem 2. Assume ∇f is L -Lipschitz and that $\alpha \leq \frac{1}{L}$. Then for any integer $t \geq 1$ we have

$$(4) \quad \min_{0 \leq k \leq t} \|\nabla f(x_k)\|^2 \leq \frac{2(f(x_0) - f_*)}{\alpha t}.$$

Taylor expansion gives $x = x_k, y = x_{k+1}$,

$$f(x_{k+1}) \leq f(x_k) + \nabla f(x_k)^T (x_{k+1} - x_k) + \frac{L}{2} \|x_{k+1} - x_k\|^2$$

$$= f(x_k) - \alpha \nabla f(x_k)^T \nabla f(x_k) + \frac{L}{2} \|\alpha \nabla f(x_k)\|^2$$

$$= f(x_k) - \alpha \|\nabla f(x_k)\|^2 + \frac{L}{2} \alpha^2 \|\nabla f(x_k)\|^2$$

$$f(x_{k+1}) \leq f(x_k) - \left(\alpha - \frac{\alpha^2 L}{2}\right) \|\nabla f(x_k)\|^2$$

$$= f(x_k) - \alpha \left(1 - \frac{\alpha L}{2}\right) \|\nabla f(x_k)\|^2$$

Assume $\alpha \leq \frac{1}{L}$

$$\geq \frac{1}{2}$$

$$1 - \frac{\alpha L}{2} \geq \frac{1}{2}$$

$$\frac{\alpha L}{2} \leq \frac{1}{2}$$

$$\alpha \leq \frac{1}{L}$$

Then

$$f(x_{k+1}) - f(x_k) \leq -\frac{\alpha}{2} \|\nabla f(x_k)\|^2$$

$$\|\nabla f(x_k)\|^2 \leq \frac{-2}{\alpha} (f(x_{k+1}) - f(x_k))$$

$$= \frac{2}{\alpha} (f(x_k) - f(x_{k+1}))$$

$$\sum_{k=0}^t \|\nabla f(x_k)\|^2 \leq \frac{2}{\alpha} \sum_{k=0}^t (f(x_k) - f(x_{k+1}))$$

$$= \frac{2}{\alpha} (f(x_0) - f(x_{t+1}))$$

$$f(x_{t+1}) \geq f_* \quad \leq \frac{\alpha}{2} (f(x_0) - f_*)$$

$$\sum_{k=0}^t \|\nabla f(x_k)\|^2 \geq (t+1) \min_{0 \leq k \leq t} \|\nabla f(x_k)\|^2$$

$$\min_{0 \leq k \leq t} \|\nabla f(x_k)\|^2 \leq \frac{\alpha}{2} \frac{(f(x_0) - f_*)}{t+1}$$

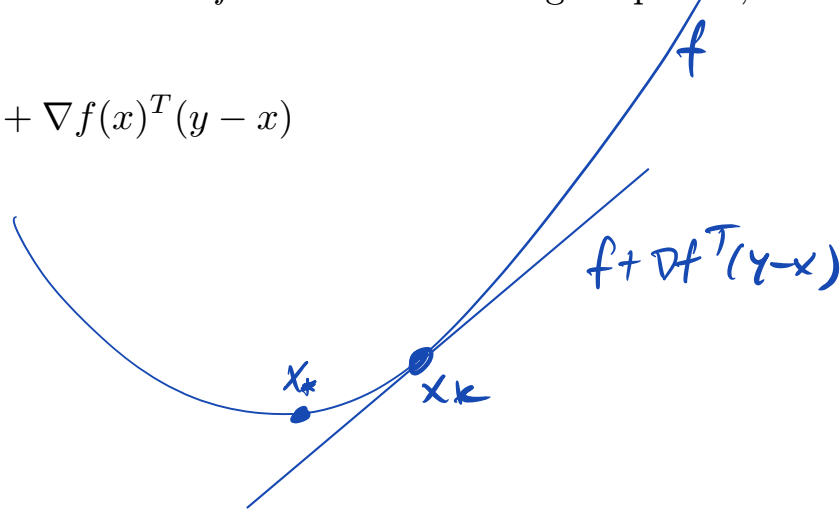


Convergence to a minimizer

To show that gradient descent converges to a global minimizer of f , we need to assume that f is *convex*, which for us means that f lies above its tangent planes, that is

$$(5) \quad f(y) \geq f(x) + \nabla f(x)^T (y - x)$$

for all $x, y \in \mathbb{R}^n$.



Other equivalent definitions of convexity include positive definiteness of the Hessian matrix $\nabla^2 f(x)$ for all x , and the convexity along lines definition

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$$

for all $x, y \in \mathbb{R}^n$ and $\lambda \in [0, 1]$.


Convergence to a minimizer

Theorem 4. Assume f is convex, ∇f is L -Lipschitz, and take $\alpha \leq \frac{1}{L}$. Then for any integer $t \geq 1$ we have

$$(6) \quad f(x_t) - f_* \leq \frac{\|x_0 - x_*\|^2}{2\alpha t},$$

where x_* is any minimizer of f .

Remark 5. Theorem 4 shows that the values $f(x_k)$ of gradient descent converge to the optimal value f_* at a rate of $O\left(\frac{1}{t}\right)$ when f is convex. This is an *extremely slow* convergence rate, known as *sublinear*. To get within $\varepsilon > 0$ of the optimal value requires $O(\varepsilon^{-1})$ iterations. So if you want 10^{-6} accuracy you need 10^6 iterations.



Proof: By convexity $(f(x_*) = f_*)$

$$f_* \geq f(x_k) + \nabla f(x_k)^T (x_* - x_k)$$

$$f(x_k) \leq f_* + \nabla f(x_k)^T (x_k - x_*)$$

From prev. proof $(\alpha \leq \frac{1}{L})$

$$f(x_{k+1}) \leq f(x_k) - \frac{\alpha}{2} \|\nabla f(x_k)\|^2$$

$$\leq f_* + \nabla f(x_k)^T (x_k - x_*) - \frac{\alpha}{2} \|\nabla f(x_k)\|^2$$

$$\|x - \gamma\|^2 = \|x\|^2 - 2x^T \gamma + \|\gamma\|^2$$

$$= f_* + \frac{1}{2\alpha} \left(2\alpha \nabla f(x_k)^T (x_k - x_*) - \alpha^2 \|\nabla f(x_k)\|^2 \right)$$

$$2x^T y - \|y\|^2 = \|x\|^2 - \|x - y\|^2$$

$$= f_* + \frac{1}{2\alpha} \left(\|x_k - x_*\|^2 - \underbrace{\|x_k - x_* - \alpha \nabla f(x_k)\|^2}_{\delta D = x_{k+1} - x_*} \right)$$

$$f(x_{k+1}) \leq f_* + \frac{1}{2\alpha} \left(\|x_k - x_*\|^2 - \|x_{k+1} - x_*\|^2 \right)$$

$$\sum_{k=0}^{t-1} (f(x_{k+1}) - f_*) \leq \frac{1}{2\alpha} \sum_{k=0}^{t-1} \left(\|x_k - x_*\|^2 - \|x_{k+1} - x_*\|^2 \right)$$

$$= \frac{1}{2\alpha} (\|x_0 - x_*\|^2 - \|x_t - x_*\|^2) \geq 0$$

$$\leq \frac{\|x_0 - x_*\|^2}{2\alpha}$$

Note $\sum_{k=0}^{t-1} (f(x_{k+1}) - f_*) \geq t(f(x_t) - f_*)$

Hence $f(x_t) - f_* \leq \frac{\|x_0 - x_*\|^2}{2\alpha t}$. \square

$$X_{k+1} = X_k - \alpha \nabla f(x_k)$$

If ∇f is L -Lipschitz then

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

$$\Rightarrow \mu \leq L$$

Linear convergence

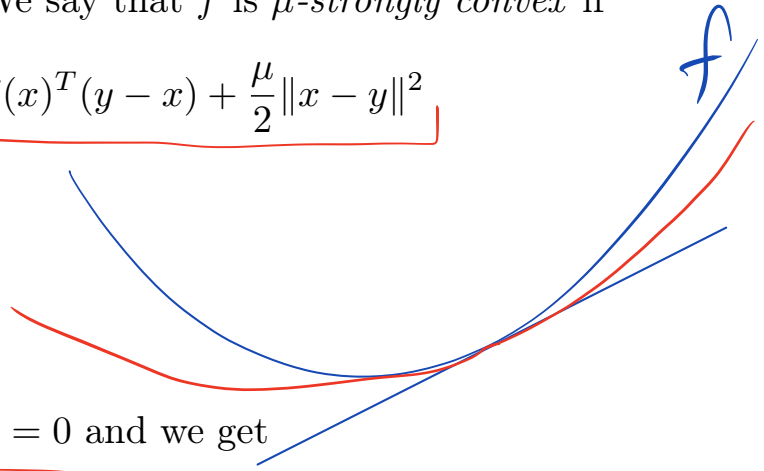
To obtain a better convergence rate, we need to make an additional assumption about how flat f can be at minima. We say that f is μ -strongly convex if

$$(7) \quad f(y) \geq \underbrace{f(x) + \nabla f(x)^T(y - x) + \frac{\mu}{2}\|x - y\|^2}_{\text{red bracket}} \quad \text{f}$$

for all $x, y \in \mathbb{R}^n$.

Note: If we take $x = x_*$ then $\nabla f(x_*) = 0$ and we get

$$(8) \quad f(y) \geq \underline{f_*} + \frac{\mu}{2}\|y - x_*\|^2.$$



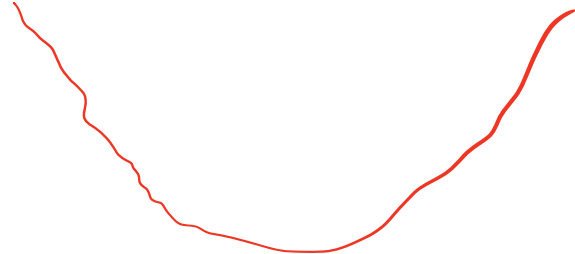
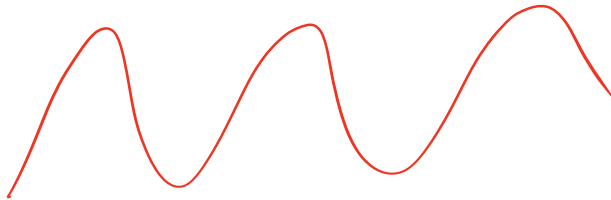
Polyak-Lojasiewicz (PL) inequality

If f is μ -strongly convex, then f satisfies the PL inequality

$$(9) \quad \frac{1}{2} \|\nabla f(x)\|^2 \geq \mu(f(x) - f_*)$$

for all $x \in \mathbb{R}^n$.

Remark 6. The PL inequality is weaker than strong convexity, and even nonconvex functions can satisfy it (as an exercise, show that $f(x) = x^2 + 3 \sin^2(x)$ satisfies the PL inequality (9) with $\mu = \frac{1}{32}$, but f is not convex).



Strong convexity

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|x-y\|^2$$

Minimize on both sides over $y \in \mathbb{R}^n$

$$f_* = \min f \geq f(x) + \min_{y \in \mathbb{R}^n} \left\{ \nabla f(x)^T (y-x) + \frac{\mu}{2} \|x-y\|^2 \right\}$$

$\nabla = 0$

$$\nabla f(x) + \mu(y-x) = 0$$

$$\Rightarrow y-x = -\frac{1}{\mu} \nabla f(x)$$

$$\begin{aligned} f_* &\geq f(x) + \nabla f(x)^T \left(-\frac{1}{\mu} \nabla f(x) \right) + \frac{\mu}{2} \left\| -\frac{1}{\mu} \nabla f(x) \right\|^2 \\ &= f(x) - \frac{1}{\mu} \|\nabla f(x)\|^2 + \frac{1}{2\mu} \|\nabla f(x)\|^2 \\ &= f(x) - \frac{1}{2\mu} \|\nabla f(x)\|^2 \end{aligned}$$

$$\frac{1}{2\mu} \|\nabla f(x)\|^2 \geq f(x) - f_*$$

$$\frac{1}{2} \|\nabla f(x)\|^2 \geq \mu (f(x) - f_*)$$

PL-inequality.

$$\mu \leq L \rightarrow \frac{1}{L} \leq \frac{1}{\mu}$$

Linear convergence

$$\alpha \leq \frac{1}{\mu} \Rightarrow \alpha\mu \leq 1$$

Theorem 7. Assume f satisfies the PL inequality (9), ∇f is L -Lipschitz, and take $\alpha \leq \frac{1}{L}$. Then for any integer $t \geq 0$ we have

$$(10) \quad \underline{f(x_t) - f_*} \leq (1 - \alpha\mu)^t (f(x_0) - f_*).$$

Proof: From previous proof, for $\alpha \leq \frac{1}{L}$

$$f(x_{k+1}) - f(x_k) \leq -\frac{\alpha}{2} \|\nabla f(x_k)\|^2$$

$$\leq -\alpha\mu (f(x_k) - f_*)$$

PL-inequality

$$\frac{1}{2} \|\nabla f(x)\|^2 \geq \mu (f(x) - f_*)$$

$$f(x_{k+1}) - f_* \leq f(x_k) - \alpha\mu(f(x_k) - f_*) - f_*$$

$$= f(x_k) - f_* - \alpha\mu(f(x_k) - f_*)$$

$$= (1 - \alpha\mu)(f(x_k) - f_*)$$

$$\leq (1 - \alpha\mu)^2(f(x_{k-1}) - f_*)$$

$$\leq (1 - \alpha\mu)^3(f(x_{k-2}) - f_*)$$

...

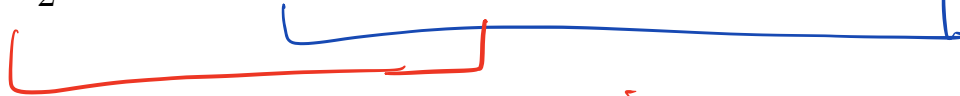
$$f(x_k) - f_* \leq (1 - \alpha\mu)^k(f(x_0) - f_*)$$



Convergence of minimizers

Remark 8. It is also natural to ask how quickly x_k is converging to x_* . For this, we require strong convexity. If f is μ -strongly convex then we have

$$\frac{\mu}{2} \|x_t - x_*\|^2 \leq f(x_t) - f_* \leq (1 - \alpha\mu)^t (f(x_0) - f_*).$$



Strong convexity

Gradient Descent Notebook ([.ipynb](#))