

[@bibtex](#)[Files](#)

Summary

A classical problem of function minimization is considered.

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

- The bottleneck (for almost all gradient methods) is choosing step-size, which can lead to the dramatic difference in method's behavior.
- One of the theoretical suggestions: choosing stepsize inversly proportional to the gradient Lipschitz constant $\eta_k = \frac{1}{L}$.
- In huge-scale applications the cost of iteration is usually defined by the cost of gradient calculation (at least $\mathcal{O}(p)$).
- If function has Lipschitz-continious gradient, then method could be rewritten as follows:

$$\begin{aligned} x_{k+1} &= x_k - \frac{1}{L} \nabla f(x_k) = \\ &= \arg \min_{x \in \mathbb{R}^n} \left\{ f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{L}{2} \|x - x_k\|_2^2 \right\} \end{aligned}$$

Intuition

Direction of local steepest descent

Let's consider a linear approximation of the differentiable function f along some direction h , $\|h\|_2 = 1$:

$$f(x + \eta h) = f(x) + \eta \langle f'(x), h \rangle + o(\eta)$$

We want h to be a decreasing direction:

$$f(x + \eta h) < f(x)$$

$$f(x) + \eta \langle f'(x), h \rangle + o(\eta) < f(x)$$

and going to the limit at $\eta \rightarrow 0$:

$$\langle f'(x), h \rangle \leq 0$$

Also from Cauchy–Bunyakovsky–Schwarz inequality:

$$|\langle f'(x), h \rangle| \leq \|f'(x)\|_2 \|h\|_2 \quad \rightarrow \quad \langle f'(x), h \rangle \geq -\|f'(x)\|_2 \|h\|_2 = -\|f'(x)\|_2$$

Thus, the direction of the antigradient

$$h = -\frac{f'(x)}{\|f'(x)\|_2}$$

gives the direction of the **steepest local** decreasing of the function f .

The result of this method is

$$x_{k+1} = x_k - \eta f'(x_k)$$

Gradient flow ODE

Let's consider the following ODE, which is referred as Gradient Flow equation.

$$\frac{dx}{dt} = -f'(x(t))$$

and discretize it on a uniform grid with η step:

$$\frac{x_{k+1} - x_k}{\eta} = -f'(x_k),$$

where $x_k \equiv x(t_k)$ and $\eta = t_{k+1} - t_k$ - is the grid step.

From here we get the expression for x_{k+1}

$$x_{k+1} = x_k - \eta f'(x_k),$$

which is exactly gradient descent.

Necessary local minimum condition

$$\begin{aligned}
f'(x) &= 0 \\
-\eta f'(x) &= 0 \\
x - \eta f'(x) &= x \\
x_k - \eta f'(x_k) &= x_{k+1}
\end{aligned}$$

This is, surely, not a proof at all, but some kind of intuitive explanation.

Minimizer of Lipschitz parabola

Some general highlights about Lipschitz properties are needed for explanation. If a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable and its gradient satisfies Lipschitz conditions with constant L , then $\forall x, y \in \mathbb{R}^n$:

$$|f(y) - f(x) - \langle \nabla f(x), y - x \rangle| \leq \frac{L}{2} \|y - x\|^2,$$

which geometrically means, that if we'll fix some point $x_0 \in \mathbb{R}^n$ and define two parabolas:

$$\phi_1(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle - \frac{L}{2} \|x - x_0\|^2,$$

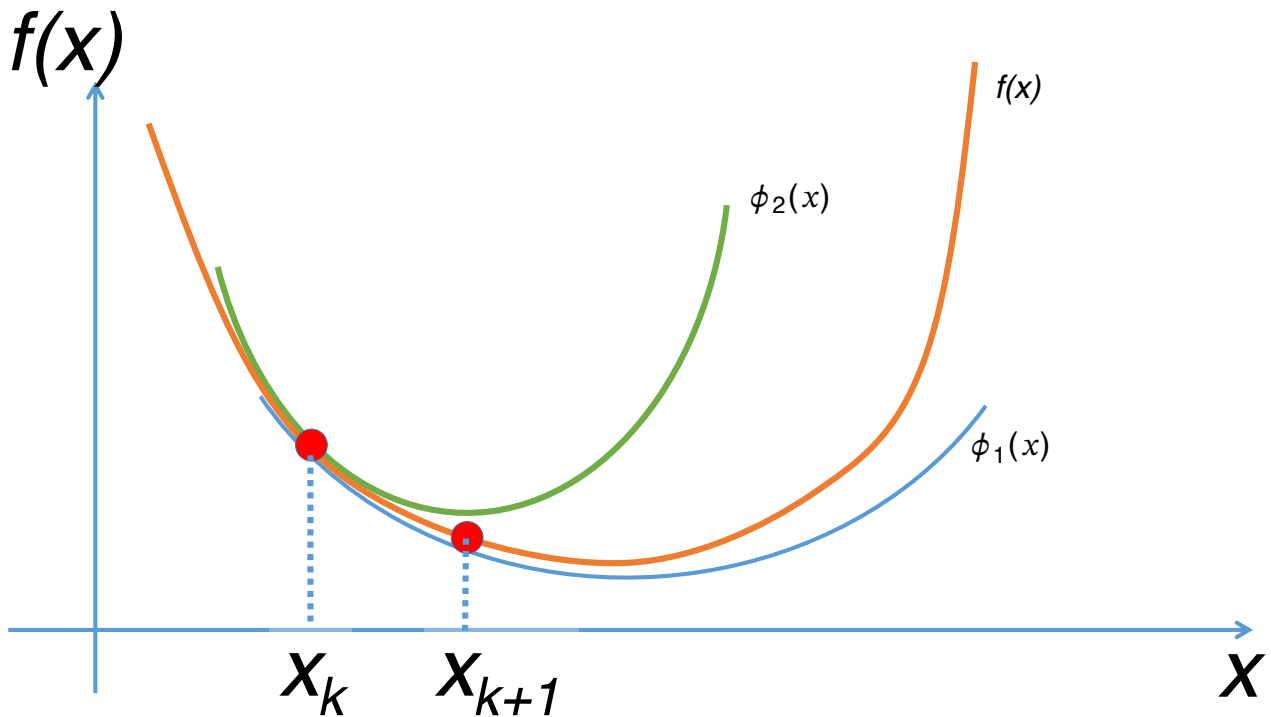
$$\phi_2(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle + \frac{L}{2} \|x - x_0\|^2.$$

Then

$$\phi_1(x) \leq f(x) \leq \phi_2(x) \quad \forall x \in \mathbb{R}^n.$$

Now, if we have global upper bound on the function, in a form of parabola, we can try to go directly to its minimum.

$$\begin{aligned}
\nabla \phi_2(x) &= 0 \\
\nabla f(x_0) + L(x^* - x_0) &= 0 \\
x^* &= x_0 - \frac{1}{L} \nabla f(x_0) \\
x_{k+1} &= x_k - \frac{1}{L} \nabla f(x_k)
\end{aligned}$$



This way leads to the $\frac{1}{L}$ stepsize choosing. However, often the L constant is not known.

But if the function is twice continuously differentiable and its gradient has Lipschitz constant L , we can derive a way to estimate this constant $\forall x \in \mathbb{R}^n$:

$$\|\nabla^2 f(x)\| \leq L$$

or

$$-LI_n \preceq \nabla^2 f(x) \preceq LI_n$$

Stepsize choosing strategies

Stepsize choosing strategy η_k significantly affects convergence. General [Line search](#) algorithms might help in choosing scalar parameter.

Constant stepsize

For $f \in C_L^{1,1}$:

$$\eta_k = \eta$$

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

With choosing $\eta = \frac{1}{L}$, we have:

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

Fixed sequence

$$\eta_k = \frac{1}{\sqrt{k+1}}$$

The latter 2 strategies are the simplest in terms of implementation and analytical analysis. It is clear that this approach does not often work very well in practice (the function geometry is not known in advance).

Exact line search aka steepest descent

$$\eta_k = \arg \min_{\eta \in \mathbb{R}^+} f(x_{k+1}) = \arg \min_{\eta \in \mathbb{R}^+} f(x_k - \eta \nabla f(x_k))$$

More theoretical than practical approach. It also allows you to analyze the convergence, but often exact line search can be difficult if the function calculation takes too long or costs a lot.

Interesting theoretical property of this method is that each following iteration is orthogonal to the previous one:

$$\eta_k = \arg \min_{\eta \in \mathbb{R}^+} f(x_k - \eta \nabla f(x_k))$$

Optimality conditions:

$$\nabla f(x_{k+1})^\top \nabla f(x_k) = 0$$

Goldstein-Armijo

Convergence analysis

Convex case

Lipschitz continuity of the gradient

Assume that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and differentiable, and additionally

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\| \quad \forall x, y \in \mathbb{R}^n$$

i.e., ∇f is Lipschitz continuous with constant $L > 0$.

Since ∇f Lipschitz with constant L , which means $\nabla^2 f \preceq LI$, we have $\forall x, y, z$:

$$(x - y)^\top (\nabla^2 f(z) - LI)(x - y) \leq 0$$

$$(x - y)^\top \nabla^2 f(z)(x - y) \leq L\|x - y\|^2$$

Now we'll consider second order Taylor approximation of $f(y)$ and Taylor's Remainder Theorem (we assume, that the function f is continuously differentiable), we have

$\forall x, y, \exists z \in [x, y]$:

$$\begin{aligned} f(y) &= f(x) + \nabla f(x)^\top (y - x) + \frac{1}{2}(x - y)^\top \nabla^2 f(z)(x - y) \\ &\leq f(x) + \nabla f(x)^\top (y - x) + \frac{L}{2}\|x - y\|^2 \end{aligned}$$

For the gradient descent we have $x = x_k, y = x_{k+1}, x_{k+1} = x_k - \eta_k \nabla f(x_k)$:

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) + \nabla f(x_k)^\top (-\eta_k \nabla f(x_k)) + \frac{L}{2}(\eta_k \nabla f(x_k))^2 \\ &\leq f(x_k) - \left(1 - \frac{L\eta}{2}\right)\eta \|\nabla f(x_k)\|^2 \end{aligned}$$

Optimal constant stepsize

Now, if we'll consider constant stepsize strategy and will maximize

$$\left(1 - \frac{L\eta}{2}\right)\eta \rightarrow \max_{\eta}, \text{ we'll get } \eta = \frac{1}{L}.$$

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

Convexity

$$f(x_k) \leq f(x^*) + \nabla f(x_k)^\top (x_k - x^*)$$

That's why we have:

$$\begin{aligned}
f(x_{k+1}) &\leq f(x^*) + \nabla f(x_k)^\top (x_k - x^*) - \frac{1}{2L} \|\nabla f(x_k)\|^2 \\
&= f(x^*) + \frac{L}{2} \left(\|x_k - x^*\|^2 - \|x_k - x^* - \frac{1}{L} \nabla f(x_k)\|^2 \right) \\
&= f(x^*) + \frac{L}{2} (\|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2)
\end{aligned}$$

Thus, summing over all iterations, we have:

$$\begin{aligned}
\sum_{i=1}^k (f(x_i) - f(x^*)) &\leq \frac{L}{2} (\|x_0 - x^*\|^2 - \|x_k - x^*\|^2) \\
&\leq \frac{L}{2} \|x_0 - x^*\|^2 = \frac{LR^2}{2},
\end{aligned}$$

where $R = \|x_0 - x^*\|$. And due to convexity:

$$f(x_k) - f(x^*) \leq \frac{1}{k} \sum_{i=1}^k (f(x_i) - f(x^*)) \leq \frac{LR^2}{2k} = \frac{R^2}{2\eta k}$$

Strongly convex case

If the function is strongly convex:

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x) + \frac{\mu}{2} \|y - x\|^2 \quad \forall x, y \in \mathbb{R}^n$$

...

$$\|x_{k+1} - x^*\|^2 \leq (1 - \eta\mu) \|x_k - x^*\|^2$$

Bounds

Conditions	$\ f(x_k) - f(x^*)\ \leq$	Type of convergence	$\ x_k - x^*\ \leq$
Convex Lipschitz-continuous function(G)	$\mathcal{O}\left(\frac{1}{k}\right) \frac{GR}{k}$	Sublinear	
Convex			

Lipschitz-continuous gradient (L)	$\mathcal{O}\left(\frac{1}{k}\right) \frac{LR^2}{k}$	Sublinear	
μ -Strongly convex Lipschitz-continuous gradient(L)		Linear	$(1 - \eta\mu)^k R^2$
μ -Strongly convex Lipschitz-continuous hessian(M)		Locally linear $R < \bar{R}$	$\frac{\bar{R}R}{\bar{R} - R} \left(1 - \frac{2\mu}{L + 3\mu}\right)$

- $R = \|x_0 - x^*\|$ - initial distance
- $\bar{R} = \frac{2\mu}{M}$

Materials

- [The zen of gradient descent. Moritz Hardt](#)
- [Great visualization](#)
- [Cheatsheet on the different convergence theorems proofs](#)