

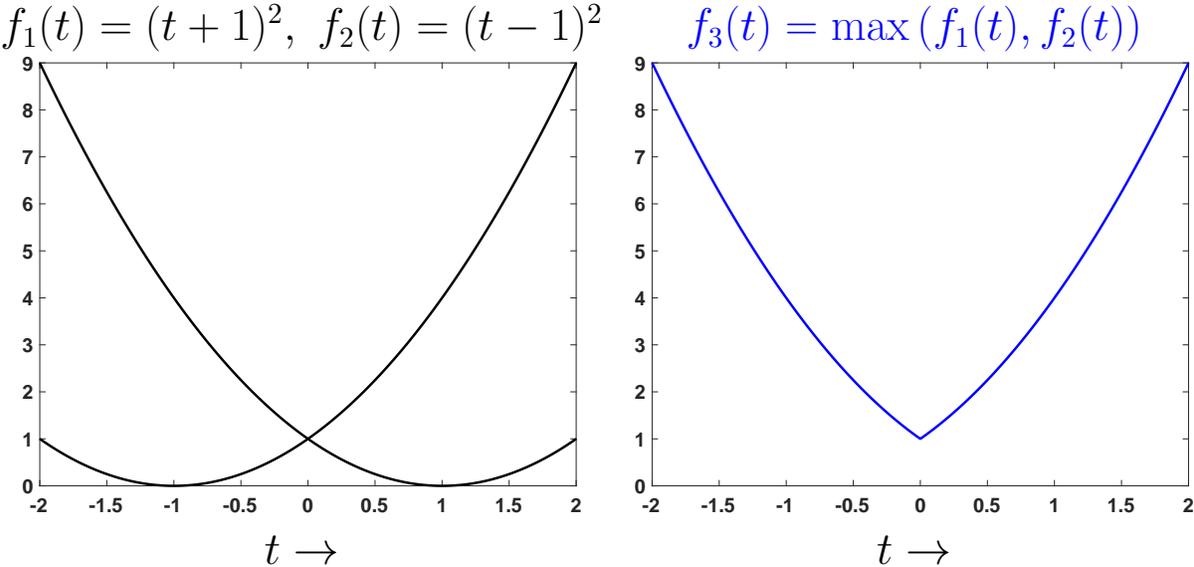
Nonsmooth optimization

Most of the theory and algorithms that we have explored for convex optimization have assumed that the functions involved are differentiable – that is, smooth.

This is not always the case in interesting scenarios. In fact, nonsmooth functions can arise quite naturally in applications. We already have looked at optimization programs involving the hinge loss $\max(\mathbf{a}^T \mathbf{x} + b, 0)$, the ℓ_1 norm, the ℓ_∞ norm, and the nuclear norm — none of these is differentiable. As another example, suppose f_1, \dots, f_Q are all perfectly smooth convex functions. Then the point-wise maximum

$$f(\mathbf{x}) = \max_{1 \leq q \leq Q} f_q(\mathbf{x})$$

is in general not smooth.

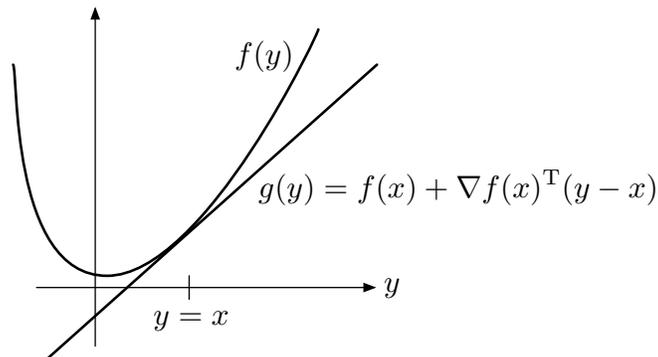


Fortunately, the theory for nonsmooth optimization is not too different than for smooth optimization. We really just need one new concept: that of a subgradient.

Subgradients

If you look back through the notes so far, you will see that the vast majority of the time we use the gradient of a convex function, it is in the context of the inequality

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle, \quad \text{for all } \mathbf{x}, \mathbf{y} \in \text{dom } f.$$



This is a very special property of convex functions, and it led to all kinds of beautiful results.

When a convex f is not differentiable at a point \mathbf{x} , we can more or less reproduce the entire theory using subgradients. A **subgradient** of f at \mathbf{x} is a vector \mathbf{g} such that

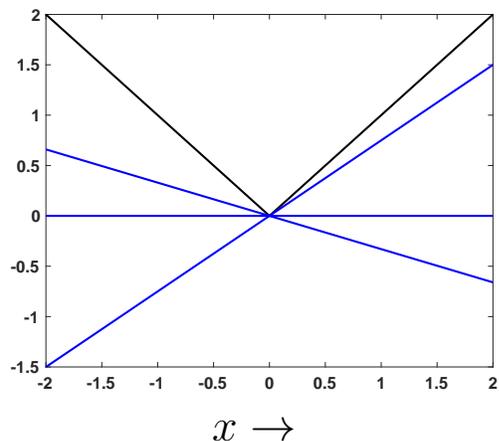
$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \mathbf{g} \rangle, \quad \text{for all } \mathbf{y} \in \text{dom } f.$$

Unlike gradients for smooth functions, there can be more than one subgradient of a nonsmooth function at a point. We call the collection of subgradients the **subdifferential** at \mathbf{x} :

$$\partial f(\mathbf{x}) = \{ \mathbf{g} : f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \mathbf{g} \rangle, \quad \text{for all } \mathbf{y} \in \text{dom } f \}.$$

Example:

$$f(x) = |x|, \quad \partial f(x) = \begin{cases} -1, & x < 0 \\ [-1, 1], & x = 0 \\ 1, & x > 0. \end{cases}$$



black: $f(x) = |x|$
blue: $f(0) + g(x-0)$ for a few $g \in \partial f(0)$

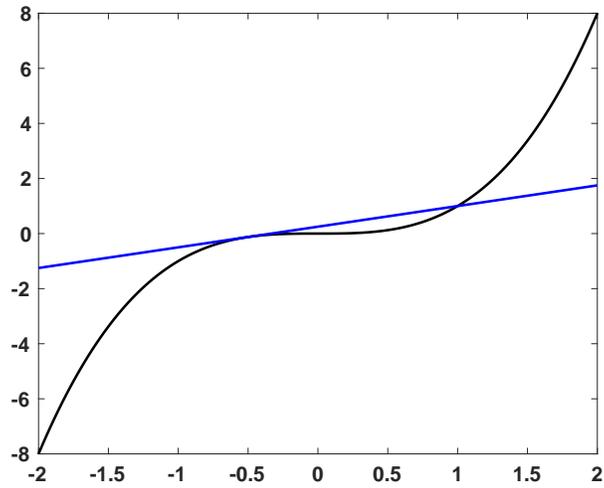
Facts for subdifferentials of convex functions:

1. If f is convex and differentiable at \mathbf{x} , then the subdifferential contains exactly one vector: the gradient,

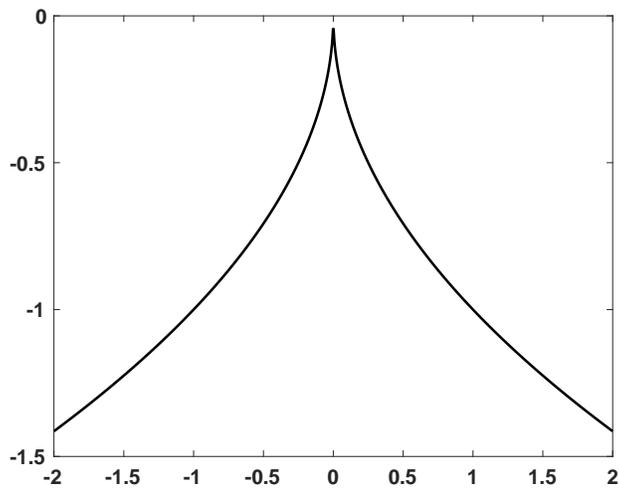
$$\partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}.$$

2. If f is convex on $\text{dom } f$, then the subdifferential is non-empty and bounded at all \mathbf{x} in the interior of $\text{dom } f$.

For non-convex f , these two points do not hold in general. The gradient at a point is not necessarily a subgradient:



and there can also be points where neither the gradient nor subgradient exist, e.g. $f(x) = -\sqrt{|x|}$ for $x \in \mathbb{R}$



Example: The ℓ_1 norm

Consider the function

$$f(\mathbf{x}) = \|\mathbf{x}\|_1.$$

The ℓ_1 norm is not differentiable at any \mathbf{x} that has at least one coordinate equal to zero. We will see that optimization problems involving the ℓ_1 norm very often have solutions that are sparse, meaning that they have many zeros. This is a big problem – the nonsmoothness is kicking in at exactly the points we are interested in.

What does the subdifferential $\partial\|\mathbf{x}\|_1$ look like in such a case? To see, recall that by definition, if a vector $\mathbf{u} \in \partial\|\mathbf{x}\|_1$, at the point \mathbf{x} , then we must have

$$\|\mathbf{y}\|_1 \geq \|\mathbf{x}\|_1 + \langle \mathbf{y} - \mathbf{x}, \mathbf{u} \rangle \quad (1)$$

for all $\mathbf{y} \in \mathbb{R}^N$. To understand what this means in terms of \mathbf{x} , it is useful to introduce the notation $\Gamma(\mathbf{x})$ to denote the set of indexes where \mathbf{x} is non-zero:

$$\Gamma(\mathbf{x}) = \{n : x_n \neq 0\}.$$

Using this, we can re-write the right-hand side of (1) as

$$\begin{aligned} \|\mathbf{x}\|_1 + \langle \mathbf{y} - \mathbf{x}, \mathbf{u} \rangle &= \sum_{n=1}^N |x_n| + \sum_{n=1}^N u_n(y_n - x_n) \\ &= \sum_{n \in \Gamma} (|x_n| - u_n x_n) + \sum_{n=1}^N u_n y_n. \end{aligned}$$

Note that if

$$u_n = \text{sign}(x_n) = \begin{cases} 1 & \text{if } x_n \geq 0, \\ -1 & \text{if } x_n < 0, \end{cases}$$

then $u_n x_n = |x_n|$. Thus, if $u_n = \text{sign}(x_n)$ for all $n \in \Gamma$, we have

$$\sum_{n \in \Gamma} |x_n| - u_n x_n = \sum_{n \in \Gamma} |x_n| - |x_n| = 0.$$

Thus, if we set $u_n = \text{sign}(x_n)$ for all $n \in \Gamma$, then (1) reduces to

$$\|\mathbf{y}\|_1 \geq \langle \mathbf{y}, \mathbf{u} \rangle.$$

As long as $|u_n| \leq 1$ for all n , then this will hold. Hence, if a vector \mathbf{u} satisfies

$$\begin{aligned} u_n &= \text{sign}(x_n) && \text{if } n \in \Gamma, \\ |u_n| &\leq 1 && \text{if } n \notin \Gamma, \end{aligned}$$

then $\mathbf{u} \in \partial\|\mathbf{x}\|_1$. It is not hard to show that for any \mathbf{u} that violates these conditions, we can construct a \mathbf{y} such that (1) is violated, and thus this is a complete description of all vectors in $\mathbf{u} \in \partial\|\mathbf{x}\|_1$.

Example: The ℓ_2 norm

While the function $\mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|_2^2$ is the prototypical differentiable ($\nabla f(\mathbf{x}) = \mathbf{x}$), smooth, and strongly convex function ($\nabla^2 f(\mathbf{x}) = \mathbf{I}$), the function $f(\mathbf{x}) = \|\mathbf{x}\|_2$ is not as nice; it is not strongly convex, and it is not differentiable at $\mathbf{x} = \mathbf{0}$ (to appreciate this latter point, consider that a 1D slice of the function $s(t) = \|t\mathbf{v}\|_2 = |t|\|\mathbf{v}\|_2$ looks like the absolute value function as function of t).

For $\mathbf{x} \neq \mathbf{0}$, an easy calculation¹ shows that

$$\nabla \|\mathbf{x}\|_2 = \frac{\mathbf{x}}{\|\mathbf{x}\|_2}.$$

At $\mathbf{x} = \mathbf{0}$, we know that $\mathbf{u} \in \partial \|\mathbf{x}\|_2$ if

$$\|\mathbf{y}\|_2 \geq \|\mathbf{0}\|_2 + \langle \mathbf{y} - \mathbf{0}, \mathbf{u} \rangle = \langle \mathbf{y}, \mathbf{u} \rangle \quad \text{for all } \mathbf{y} \in \mathbb{R}^N. \quad (2)$$

We can find \mathbf{u} that meet these conditions using the Cauchy-Schwarz inequality. Note that

$$\langle \mathbf{y}, \mathbf{u} \rangle \leq \|\mathbf{y}\|_2 \|\mathbf{u}\|_2,$$

so (2) will hold when $\|\mathbf{u}\|_2 \leq 1$. On the other hand, if $\|\mathbf{u}\|_2 > 1$, then for $\mathbf{y} = \mathbf{u}$, we have

$$\langle \mathbf{y}, \mathbf{u} \rangle = \|\mathbf{y}\|_2^2 > \|\mathbf{y}\|_2,$$

and (2) does not hold. Therefore

$$\partial \|\mathbf{x}\|_2 = \begin{cases} \{\mathbf{u} : \|\mathbf{u}\|_2 \leq 1\}, & \mathbf{x} = \mathbf{0} \\ \frac{\mathbf{x}}{\|\mathbf{x}\|_2}, & \mathbf{x} \neq \mathbf{0}. \end{cases}$$

¹Use the fact that $\frac{d}{dx} \sqrt{x^2 + a} = x / \sqrt{x^2 + a}$.

Example: Dual norms and $\partial\|\mathbf{x}\|$ at $\mathbf{x} = \mathbf{0}$

Norms in general are not differentiable at $\mathbf{x} = \mathbf{0}$, again because they look like an absolute value function along a line: $s(t) = \|t\mathbf{v}\| = |t| \cdot \|\mathbf{v}\|$ for any valid norm $\|\cdot\|$. We can generalize the result for the ℓ_2 norm at $\mathbf{x} = \mathbf{0}$ using the concept of a **dual norm**.

The dual norm $\|\cdot\|_*$ of a norm $\|\cdot\|$ is

$$\|\mathbf{y}\|_* = \max_{\|\mathbf{x}\| \leq 1} \langle \mathbf{x}, \mathbf{y} \rangle.$$

Since sublevel sets of norms in \mathbb{R}^N are compact, we know that this maximum is achieved, and it is an easy exercise to show that $\|\cdot\|_*$ is a valid norm. You can also verify the following easy facts at home

- the dual of $\|\cdot\|_2$ is again $\|\cdot\|_2$
- the dual of $\|\cdot\|_1$ is $\|\cdot\|_\infty$
- the dual of $\|\cdot\|_\infty$ is $\|\cdot\|_1$

It is also a fact (for norms on \mathbb{R}^N) that the dual of $\|\cdot\|_*$ is the original norm $\|\cdot\|$, i.e. $\|\mathbf{x}\|_{**} = \|\mathbf{x}\|$. We also have the generalized Cauchy-Schwarz inequality (see the Technical Details section)

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \cdot \|\mathbf{y}\|_*.$$

We can use these facts with an argument similar to the ℓ_2 case above to compute the subdifferential of any norm at $\mathbf{0}$ as

$$\partial\|\mathbf{0}\| = \{\mathbf{u} : \|\mathbf{u}\|_* \leq 1\}.$$

Properties of subdifferentials

Here are some properties of the subdifferential that we will state without proof (but are easy to prove). Below, we assume that all functions are well-defined on all of \mathbb{R}^N .

Summation: If $f(\mathbf{x}) = f_1(\mathbf{x}) + f_2(\mathbf{x})$, then

$$\partial f(\mathbf{x}) = \partial f_1(\mathbf{x}) + \partial f_2(\mathbf{x}).$$

That is, the set of all subgradients (at \mathbf{x}) of f is the set of vectors that can be written as a sum of a vector from $\partial f_1(\mathbf{x})$ plus a vector from $\partial f_2(\mathbf{x})$.

Chain rule for affine transformations: If $h(\mathbf{x}) = f(\mathbf{A}\mathbf{x} + \mathbf{b})$, then

$$\partial h(\mathbf{x}) = \mathbf{A}^T \partial f(\mathbf{A}\mathbf{x} + \mathbf{b}).$$

That is, we compute the subgradients of f at the point $\mathbf{A}\mathbf{x} + \mathbf{b}$, then map them through \mathbf{A}^T .

Max of functions: If $f(\mathbf{x}) = \max\{f_1(\mathbf{x}), \dots, f_M(\mathbf{x})\}$, then

$$\partial f(\mathbf{x}) = \text{conv} \left(\bigcup_{m \in \Gamma(\mathbf{x})} \partial f_m(\mathbf{x}) \right),$$

where $\Gamma(\mathbf{x}) = \{m : f_m(\mathbf{x}) = f(\mathbf{x})\}$, and conv takes the convex hull:

$$\text{conv}(\mathcal{X}) = \left\{ \sum_{p=1}^P \lambda_p \mathbf{x}_p, \mathbf{x}_p \in \mathcal{X}, \lambda_p \geq 0, \sum_{p=1}^P \lambda_p = 1, \forall P \right\}$$

Exercise: Compute $\partial f(\mathbf{x})$ for $f(\mathbf{x}) = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_1$.

Answer: Set $\Gamma(\mathbf{x}) = \{m : \mathbf{a}_m^\top \mathbf{x} \neq y_m\}$, where \mathbf{a}_m^\top is the m th row of \mathbf{A} . Then $\partial f(\mathbf{x})$ is the set of vectors that can be written as

$$\mathbf{u} = \sum_{m \in \Gamma(\mathbf{x})} \text{sgn}(\mathbf{a}_m^\top \mathbf{x} - y_m) \mathbf{a}_m + \sum_{m \notin \Gamma(\mathbf{x})} \beta_m \mathbf{a}_m$$

for any β_m with $|\beta_m| \leq 1$. Another way to write this is

$$\mathbf{y} = \mathbf{A}^\top \boldsymbol{\beta}, \quad \begin{cases} \beta_m = \text{sgn}(\mathbf{a}_m^\top \mathbf{x} - y_m), & m \in \Gamma(\mathbf{x}) \\ |\beta_m| \leq 1, & m \notin \Gamma(\mathbf{x}). \end{cases}$$

Exercise: Compute $\partial f(x)$ for $f(x) = \max(x, 0)$.

Answer:

$$\partial f(x) = \begin{cases} 0, & x < 0, \\ [0, 1], & x = 0, \\ 1, & x > 0. \end{cases}$$

Exercise: Compute $\partial f(x)$ for $f(x) = \max((x + 1)^2, (x - 1)^2)$.

Answer:

$$\partial f(x) = \begin{cases} 2(x - 1) & x < 0, \\ [-2, 2], & x = 0, \\ 2(x + 1), & x > 0. \end{cases}$$

Exercise: Compute $\partial f(\mathbf{x})$ for $f(\mathbf{x}) = \|\mathbf{x}\|_\infty$.

Optimality conditions for unconstrained optimization

With the right definition in place, it is very easy to re-derive the central mathematical results in this course for general² convex functions.

Let $f(\mathbf{x})$ be a general convex function. Then \mathbf{x}^* is a solution to the unconstrained problem

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad f(\mathbf{x})$$

if and only if

$$\mathbf{0} \in \partial f(\mathbf{x}^*).$$

The proof of this statement is so easy you could do it in your sleep. Suppose $\mathbf{0} \in \partial f(\mathbf{x}^*)$. Then

$$\begin{aligned} f(\mathbf{y}) &\geq f(\mathbf{x}^*) + \langle \mathbf{y} - \mathbf{x}^*, \mathbf{0} \rangle \\ &= f(\mathbf{x}^*) \end{aligned}$$

for all $\mathbf{y} \in \text{dom } f$. Thus \mathbf{x}^* is optimal. Likewise, if $f(\mathbf{y}) \geq f(\mathbf{x}^*)$ for all $\mathbf{y} \in \text{dom } f$, then of course it must also be true that $f(\mathbf{y}) \geq f(\mathbf{x}^*) + \langle \mathbf{y} - \mathbf{x}^*, \mathbf{0} \rangle$ for all \mathbf{y} , and so $\mathbf{0} \in \partial f(\mathbf{x}^*)$.

Example: The LASSO

Consider the ℓ_1 regularized least-squares problem

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \tau \|\mathbf{x}\|_1.$$

²Meaning not necessarily differentiable.

We can quickly translate the general result $\mathbf{0} \in \partial f(\mathbf{x}^*)$ into a useful set of optimality conditions. We need to compute the subdifferential of $f(\mathbf{x}) = \frac{1}{2}\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \tau\|\mathbf{x}\|_1$. The first term is smooth, so the subdifferential just contains the gradient:

$$\partial f(\mathbf{x}) = \mathbf{A}^\top(\mathbf{A}\mathbf{x} - \mathbf{y}) + \tau\partial\|\mathbf{x}\|_1.$$

As shown above $\partial\|\mathbf{x}\|_1$ is the set of all vectors \mathbf{u} such that

$$\begin{aligned} u_n &= \text{sign}(x_n) && \text{if } x_n \neq 0, \\ |u_n| &\leq 1 && \text{if } x_n = 0. \end{aligned}$$

Thus the optimality condition

$$\mathbf{0} \in \mathbf{A}^\top(\mathbf{A}\mathbf{x}^* - \mathbf{y}) + \tau\partial\|\mathbf{x}^*\|_1,$$

means that \mathbf{x}^* is optimal if and only if

$$\begin{aligned} \mathbf{a}_n^\top(\mathbf{y} - \mathbf{A}\mathbf{x}^*) &= \tau \text{sign } x_n^* && \text{if } x_n^* \neq 0, \\ |\mathbf{a}_n^\top(\mathbf{y} - \mathbf{A}\mathbf{x}^*)| &\leq \tau && \text{if } x_n^* = 0. \end{aligned}$$

where here \mathbf{a}_n is the n^{th} column of \mathbf{A} .

Note that this doesn't quite give us a closed form expression for \mathbf{x}^* (except when \mathbf{A} is an orthonormal matrix), but it is useful both algorithmically (for checking if a candidate \mathbf{x} is a solution) and theoretically (for understanding and analyzing the properties of the solution to this optimization problem.)

The subgradient method

The subgradient method is the non-smooth version of gradient descent. The basic algorithm is straightforward, consisting of the iteration

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{d}_k, \quad (3)$$

where \mathbf{d}_k is *any* subgradient at \mathbf{x}_k , i.e., $\mathbf{d}_k \in \partial f(\mathbf{x}_k)$. Of course, there could be many choices for \mathbf{d}_k at every step, and the progress you make at that iteration could vary dramatically with this choice. Making this determination, though, is often very difficult, and whether or not it can even be done it depends heavily on the particular problem. Thus the analytical results for the subgradient method just assume we have any subgradient at a particular step.

With the right choice of step sizes $\{\alpha_k\}$, some simple analysis (which we will get to in a minute) shows that the subgradient method converges. The convergence rate, though, is very slow. This is also evident in most practical applications of this method: it can take many iterations to arrive at a solution that is even close to optimal.

Here is what we know about this algorithm for solving the general unconstrained program

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} f(\mathbf{x}). \quad (4)$$

We will look at one particular case here; for more detailed results see [Nes04, Chapter 3]. Along with f being convex, we will assume that it has at least one minimizer. The results also assume that f is Lipschitz:

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq L \|\mathbf{x} - \mathbf{y}\|_2.$$

Note that here we are assuming that f is Lipschitz, not that f has a Lipschitz gradient map (since the gradient does not even necessarily

exist). A direct consequence of f being Lipschitz is that the norms of the subgradients are bounded:

$$\|\mathbf{d}\|_2 \leq L, \quad \text{for all } \mathbf{d} \in \partial f(\mathbf{x}), \quad \text{for all } \mathbf{x} \in \mathbb{R}^N. \quad (5)$$

The results below used pre-determined step sizes. Thus the iteration (3) does not necessarily decrease the functional $f(\mathbf{x})$ at every step. We will keep track of the best value we have up to the current iteration with

$$f_k^{\text{best}} = \min \{f(\mathbf{x}_i), \quad 0 \leq i < k\}.$$

We will let \mathbf{x}^* be any solution to (4) and set $f^* = f(\mathbf{x}^*)$.

Our analytical results stem from a careful look at what happens during a single iteration. Note that

$$\begin{aligned} \|\mathbf{x}_{i+1} - \mathbf{x}^*\|_2^2 &= \|\mathbf{x}_i - \alpha_i \mathbf{d}_i - \mathbf{x}^*\|_2^2 \\ &= \|\mathbf{x}_i - \mathbf{x}^*\|_2^2 - 2\alpha_i \langle \mathbf{x}_i - \mathbf{x}^*, \mathbf{d}_i \rangle + \alpha_i^2 \|\mathbf{d}_i\|_2^2 \\ &\leq \|\mathbf{x}_i - \mathbf{x}^*\|_2^2 - 2\alpha_i (f(\mathbf{x}_i) - f^*) + \alpha_i^2 \|\mathbf{d}_i\|_2^2, \end{aligned}$$

where the inequality follows from the definition of a subgradient:

$$f^* \geq f(\mathbf{x}_i) + \langle \mathbf{x}^* - \mathbf{x}_i, \mathbf{d}_i \rangle.$$

Rearranging the bound above we have

$$2\alpha_i (f(\mathbf{x}_i) - f^*) \leq \|\mathbf{x}_i - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{i+1} - \mathbf{x}^*\|_2^2 + \alpha_i^2 \|\mathbf{d}_i\|_2^2,$$

and so of course

$$2\alpha_i (f_i^{\text{best}} - f^*) \leq \|\mathbf{x}_i - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{i+1} - \mathbf{x}^*\|_2^2 + \alpha_i^2 \|\mathbf{d}_i\|_2^2.$$

Since f_i^{best} is monotonically decreasing, at iteration k we have

$$2\alpha_i (f_k^{\text{best}} - f^*) \leq \|\mathbf{x}_i - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{i+1} - \mathbf{x}^*\|_2^2 + \alpha_i^2 \|\mathbf{d}_i\|_2^2,$$

for all $i \leq k$. To understand what has happened after k iterations, we sum both sides of the expression above from $i = 0$ to $i = k - 1$. Notice that the two error terms on the right hand side give us the telescoping sum:

$$\begin{aligned} \sum_{i=0}^{k-1} (\|\mathbf{x}_i - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{i+1} - \mathbf{x}^*\|_2^2) &= \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_k - \mathbf{x}^*\|_2^2 \\ &\leq \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \end{aligned}$$

and so

$$f_k^{\text{best}} - f^* \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 + \sum_{i=0}^{k-1} \alpha_i^2 \|\mathbf{d}_i\|_2^2}{2 \sum_{i=0}^{k-1} \alpha_i} \quad (6)$$

We can now specialize this result to general step-size strategies.

Fixed step size. Suppose that $\alpha_k = \alpha > 0$ for all k . Then (6) becomes

$$f_k^{\text{best}} - f^* \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2k\alpha} + \frac{L^2\alpha}{2},$$

where we have also used the Lipschitz property (5). Note that in this case, no matter how small we choose α , **the subgradient algorithm is not guaranteed to converge**. This is, in fact, standard in practice as well. The problem is that, unlike gradients for smooth functions, the subgradients do not have to vanish as we approach the solution. Even at the solution, there can be subgradients that are large.

Fixed step length. A similar result holds if we always move the same amount, taking

$$\alpha_k = s / \|\mathbf{d}_k\|_2.$$

This means that

$$\|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2 = s.$$

Of course, with this strategy it is self-evident that it will never converge, since we move some fixed amount at every step. We can bound the suboptimality at step k as

$$f_k^{\text{best}} - f^* \leq \frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2ks} + \frac{Ls}{2},$$

which is not necessarily worse than the fixed step size result. In fact, notice that even though you are moving some fixed amount, you will never move too far from an optimal point.

Decreasing step size. The results above suggest that we might want to decrease the step size as k increases, so we can get rid of this constant offset term. To make the terms in (6) work out, we let $\alpha_k \rightarrow 0$, but not too fast. Specifically, we choose a sequence $\{\alpha_k\}$ such that

$$\sum_{k=1}^{\infty} \alpha_k = \infty, \quad \text{and} \quad \frac{\sum_{i=0}^{k-1} \alpha_i^2}{\sum_{i=0}^{k-1} \alpha_i} \rightarrow 0.$$

Looking at (6) above, we can see that under these conditions $f_k^{\text{best}} \rightarrow f^*$. It is an exercise (but a nontrivial one) to show that it is enough to choose $\{\alpha_k\}$ such that

$$\alpha_k \rightarrow 0 \text{ as } k \rightarrow \infty, \quad \text{and} \quad \sum_{k=1}^{\infty} \alpha_k = \infty. \quad (7)$$

To get an idea of the tradeoffs involved here, suppose that $\alpha_k = \alpha/(k+1)$. Then for large k , we have the approximations

$$\sum_{i=0}^{k-1} \alpha_i \sim \alpha \log k, \quad \text{and} \quad \sum_{i=0}^{k-1} \alpha_i^2 \sim \text{Const} = \alpha^2 \pi^2 / 6$$

that are good as upper and lower bounds to within constants. In this case, the convergence result (6) becomes

$$f_k^{\text{best}} - f^* \lesssim \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{\alpha \log k} + \text{Const} \cdot \frac{\alpha L^2}{\log k}.$$

So the convergence is extraordinarily slow – *logarithmic* in k .

You can get much better rates than this (but still not great) by decreasing the stepsize more slowly. Consider now $\alpha_k = \alpha/\sqrt{k+1}$. Then for large k

$$\sum_{i=0}^{k-1} \alpha_i \sim (\alpha + 1)\sqrt{k}, \quad \text{and} \quad \sum_{i=0}^{k-1} \alpha_i^2 \sim \alpha^2 \log k,$$

and so

$$f_k^{\text{best}} - f^* \lesssim \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{(\alpha + 1)\sqrt{k}} + \text{Const} \cdot \frac{\alpha L^2 \log k}{\sqrt{k}}.$$

This is something like $O(1/\sqrt{k})$ convergence. This means that if we want to guarantee $f_k^{\text{best}} - f^* \leq \epsilon$, we need $k = O(1/\epsilon^2)$ iterations.

In [Nes04, Chapter 3], it is shown that there is no better rate of convergence than $O(1/\sqrt{k})$ that holds uniformly across all problems.

Example. Consider the “ ℓ_1 approximation problem”

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_1.$$

We have already looked at the subdifferential of $\|\mathbf{x}\|_1$. Specifically, we showed that \mathbf{u} is a subgradient of $\|\mathbf{x}\|_1$ at \mathbf{x} if it satisfies

$$\begin{aligned} u_n &= \text{sign}(x_n) && \text{if } x_n \neq 0, \\ |u_n| &\leq 1 && \text{if } x_n = 0. \end{aligned}$$

In the exercise above, we also derived the subdifferential for $\|\mathbf{Ax} - \mathbf{b}\|_1$. We quickly re-derive it here using “guess and check”. First consider a vector \mathbf{z} that satisfies

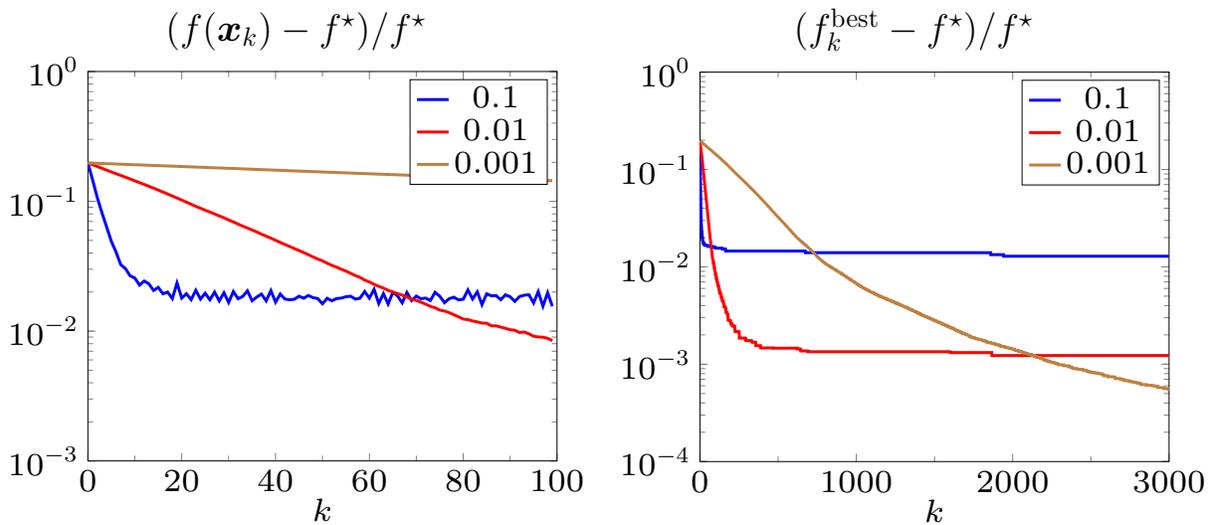
$$\begin{aligned} z_m &= \text{sign}(\mathbf{a}_m^T \mathbf{x} - b_m) && \text{if } \mathbf{a}_m^T \mathbf{x} - b_m \neq 0, \\ |z_m| &\leq 1 && \text{if } \mathbf{a}_m^T \mathbf{x} - b_m = 0. \end{aligned}$$

Now consider the vector $\mathbf{u} = \mathbf{A}^T \mathbf{z}$. Note that

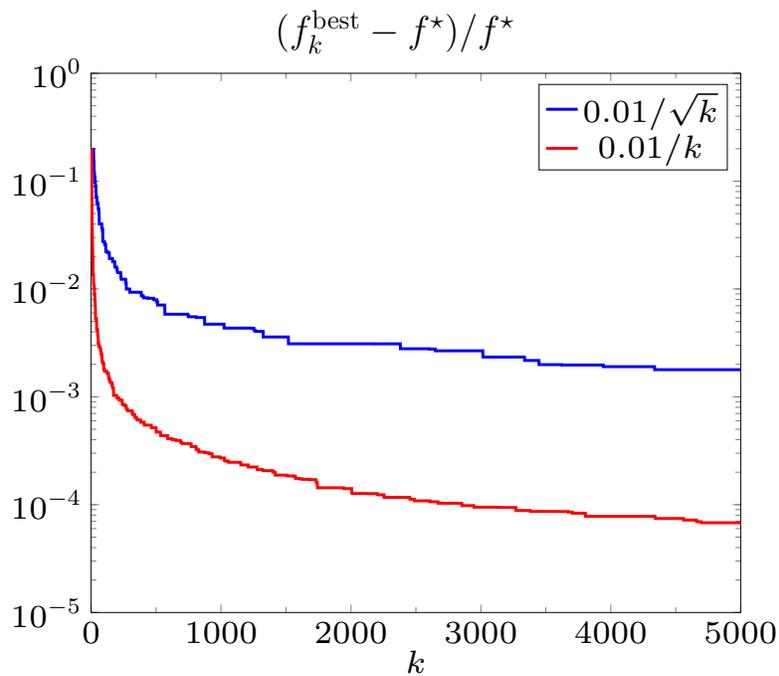
$$\begin{aligned} \mathbf{u}^T(\mathbf{y} - \mathbf{x}) &= \mathbf{z}^T \mathbf{A}(\mathbf{y} - \mathbf{x}) \\ &= \mathbf{z}^T(\mathbf{Ay} - \mathbf{b} + \mathbf{b} - \mathbf{Ax}) \\ &= \mathbf{z}^T(\mathbf{Ay} - \mathbf{b}) - \mathbf{z}^T(\mathbf{Ax} - \mathbf{b}) \\ &= \mathbf{z}^T(\mathbf{Ay} - \mathbf{b}) - \|\mathbf{Ax} - \mathbf{b}\|_1 \\ &\leq \|\mathbf{Ay} - \mathbf{b}\|_1 - \|\mathbf{Ax} - \mathbf{b}\|_1. \end{aligned}$$

Rearranging this shows that \mathbf{u} is a subgradient of $\|\mathbf{Ax} - \mathbf{b}\|_1$. Using this we can construct a subgradient at each step \mathbf{x}_k .

Below we illustrate the performance of this approach for a randomly generated example with $\mathbf{A} \in \mathbb{R}^{500 \times 100}$ and $\mathbf{b} \in \mathbb{R}^{1000}$. For three different sizes of fixed step length, $s = 0.1, 0.01, 0.001$, we make quick progress at the beginning, but then saturate, just as the theory predicts:



Here is a run using two different decreasing step size strategies: $\alpha_k = .01/\sqrt{k}$ and $\alpha_k = .01/k$.



As you can see, even though the theoretical worst case bound makes a stepsize of $\sim 1/\sqrt{k}$ look better, in this particular case, a stepsize $\sim 1/k$ actually performs better.

Qualitatively, the takeaways for the subgradient method are:

1. It is a natural extension of the gradient descent formulation
2. In general, it does not converge for fixed stepsizes.
3. If the stepsizes decrease, you can guarantee convergence.
4. Theoretical convergence rates are slow.
5. Convergence rates in practice are also very slow, but depend a lot on the particular example.

Technical Details: Dual norms and the generalized Cauchy-Schwarz inequality

Recall that the dual norm $\|\cdot\|_*$ to a norm $\|\cdot\|$ is

$$\|\mathbf{y}\|_* = \max_{\|\mathbf{x}\| \leq 1} \langle \mathbf{x}, \mathbf{y} \rangle.$$

First, note that dual norm is well-defined. For a fixed \mathbf{y} , we are maximizing a continuous function $\langle \mathbf{x}, \mathbf{y} \rangle$ over a compact set $\{\mathbf{x} : \|\mathbf{x}\| \leq 1\}$, and so by the Weierstrauss extreme value theorem, there exists a point in that set that achieves the max above. That is, there exists \mathbf{x}_y with $\|\mathbf{x}_y\| \leq 1$ and

$$\max_{\|\mathbf{x}\| \leq 1} \langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}_y, \mathbf{y} \rangle = \|\mathbf{y}\|_*.$$

Since $\langle \alpha \mathbf{x}, \mathbf{y} \rangle = \alpha \langle \mathbf{x}, \mathbf{y} \rangle$ and $\|\alpha \mathbf{x}\| = |\alpha| \|\mathbf{x}\|$, we know that $\|\mathbf{x}_y\| = 1$ (otherwise we could scale it to make the inner product bigger), and

$$\min_{\|\mathbf{x}\| \leq 1} \langle \mathbf{x}, \mathbf{y} \rangle = \langle -\mathbf{x}_y, \mathbf{y} \rangle = -\|\mathbf{y}\|_*.$$

This also gives the Generalized Cauchy Schwarz inequality, as

$$\langle \mathbf{x}, \mathbf{y} \rangle = \|\mathbf{x}\| \cdot \left\langle \frac{\mathbf{x}}{\|\mathbf{x}\|}, \mathbf{y} \right\rangle \leq \|\mathbf{x}\| \cdot \max_{\|\mathbf{x}'\| \leq 1} \langle \mathbf{x}', \mathbf{y} \rangle = \|\mathbf{x}\| \cdot \|\mathbf{y}\|_*.$$

We can replace the max with min above (while of course changing the direction of the inequality) to get

$$-\|\mathbf{x}\| \cdot \|\mathbf{y}\|_* \leq \langle \mathbf{x}, \mathbf{y} \rangle \leq \|\mathbf{x}\| \cdot \|\mathbf{y}\|_*,$$

with equality achieved on the right when $\mathbf{x}/\|\mathbf{x}\| = \mathbf{x}_y$ and equality achieved on the left when $\mathbf{x}/\|\mathbf{x}\| = -\mathbf{x}_y$.

References

- [Nes04] Y. Nesterov. *Introductory Lectures on Convex Optimization*. Springer Science+Business Media, 2004.