

Subgradients

S. Boyd, J. Duchi, and L. Vandenberghe
Notes for EE364b, Stanford University, Spring 2014-15

April 1, 2018

1 Definition

We say a vector $g \in \mathbf{R}^n$ is a *subgradient* of $f : \mathbf{R}^n \rightarrow \mathbf{R}$ at $x \in \mathbf{dom} f$ if for all $z \in \mathbf{dom} f$,

$$f(z) \geq f(x) + g^T(z - x). \quad (1)$$

If f is convex and differentiable, then its gradient at x is a subgradient. But a subgradient can exist even when f is not differentiable at x , as illustrated in figure 1. The same example shows that there can be more than one subgradient of a function f at a point x .

There are several ways to interpret a subgradient. A vector g is a subgradient of f at x if the affine function (of z) $f(x) + g^T(z - x)$ is a global underestimator of f . Geometrically, g is a subgradient of f at x if $(g, -1)$ supports **epi** f at $(x, f(x))$, as illustrated in figure 2.

A function f is called *subdifferentiable* at x if there exists at least one subgradient at x . The set of subgradients of f at the point x is called the *subdifferential* of f at x , and is denoted $\partial f(x)$. A function f is called subdifferentiable if it is subdifferentiable at all $x \in \mathbf{dom} f$.

Example. *Absolute value.* Consider $f(z) = |z|$. For $x < 0$ the subgradient is unique: $\partial f(x) = \{-1\}$. Similarly, for $x > 0$ we have $\partial f(x) = \{1\}$. At $x = 0$ the subdifferential is defined by the inequality $|z| \geq gz$ for all z , which is satisfied if and only if $g \in [-1, 1]$. Therefore we have $\partial f(0) = [-1, 1]$. This is illustrated in figure 3.

2 Basic properties

The subdifferential $\partial f(x)$ is always a closed convex set, even if f is not convex. This follows from the fact that it is the intersection of an infinite set of halfspaces:

$$\partial f(x) = \bigcap_{z \in \mathbf{dom} f} \{g \mid f(z) \geq f(x) + g^T(z - x)\}.$$

In addition, if f is continuous at x , then the subdifferential $\partial f(x)$ is bounded. Indeed, choose some $\epsilon > 0$ such that that $-\infty < \underline{f} \leq f(y) \leq \bar{f} < \infty$ for all $y \in \mathbf{R}^n$ such that $\|y - x\|_2 \leq \epsilon$.

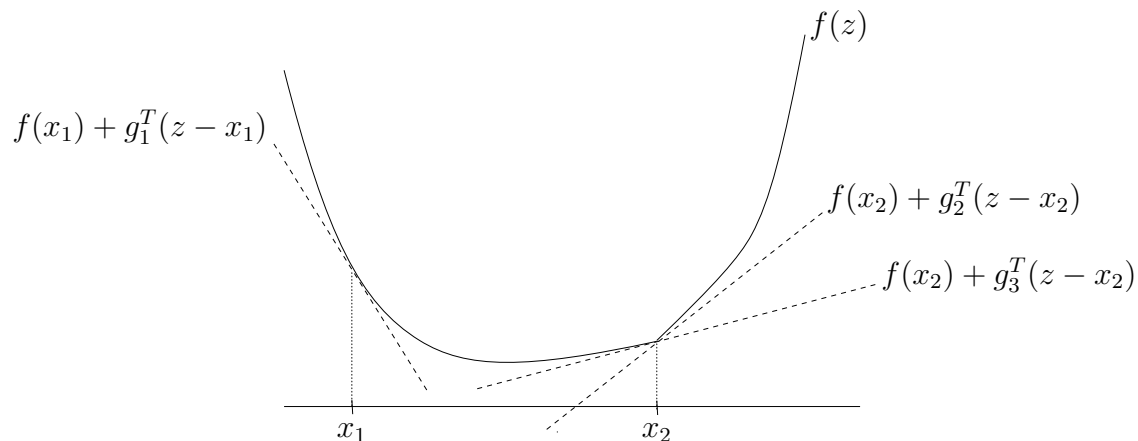


Figure 1: At x_1 , the convex function f is differentiable, and g_1 (which is the derivative of f at x_1) is the unique subgradient at x_1 . At the point x_2 , f is not differentiable. At this point, f has many subgradients: two subgradients, g_2 and g_3 , are shown.

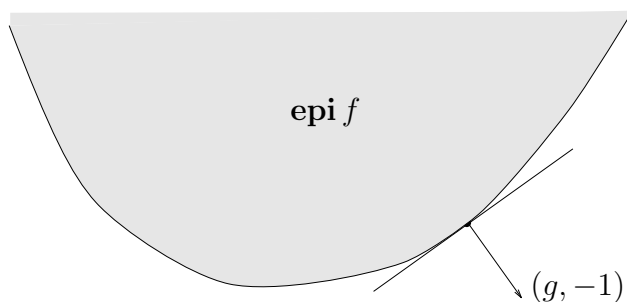


Figure 2: A vector $g \in \mathbf{R}^n$ is a subgradient of f at x if and only if $(g, -1)$ defines a supporting hyperplane to $\mathbf{epi} f$ at $(x, f(x))$.

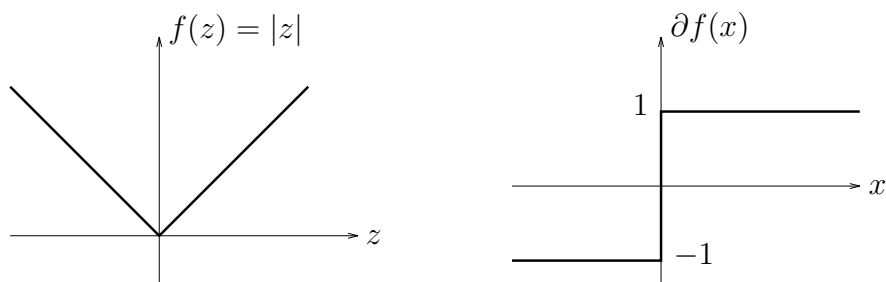


Figure 3: The absolute value function (left), and its subdifferential $\partial f(x)$ as a function of x (right).

If $\partial f(x)$ is unbounded, then there is a sequence $g_n \in \partial f(x)$ such that $\|g_n\|_2 \rightarrow \infty$. Taking the sequence $y_n = x + \epsilon g_n / \|g_n\|_2$, we find that $f(y_n) \geq f(x) + g_n^T(y_n - x) = f(x) + \epsilon \|g_n\|_2 \rightarrow \infty$, which is a contradiction to $f(y_n)$ being bounded.

2.1 Existence of subgradients

If f is convex and $x \in \mathbf{int\,dom}\,f$, then $\partial f(x)$ is nonempty and bounded. To establish that $\partial f(x) \neq \emptyset$, we apply the supporting hyperplane theorem to the convex set $\mathbf{epi}\,f$ at the boundary point $(x, f(x))$, to conclude the existence of $a \in \mathbf{R}^n$ and $b \in \mathbf{R}$, not both zero, such that

$$\begin{bmatrix} a \\ b \end{bmatrix}^T \left(\begin{bmatrix} z \\ t \end{bmatrix} - \begin{bmatrix} x \\ f(x) \end{bmatrix} \right) = a^T(z - x) + b(t - f(x)) \leq 0$$

for all $(z, t) \in \mathbf{epi}\,f$. This implies $b \leq 0$, and that

$$a^T(z - x) + b(f(z) - f(x)) \leq 0$$

for all z . If $b \neq 0$, we can divide by b to obtain

$$f(z) \geq f(x) - (a/b)^T(z - x),$$

which shows that $-a/b \in \partial f(x)$. Now we show that $b \neq 0$, *i.e.*, that the supporting hyperplane cannot be vertical. If $b = 0$ we conclude that $a^T(z - x) \leq 0$ for all $z \in \mathbf{dom}\,f$. This is impossible since $x \in \mathbf{int\,dom}\,f$.

This discussion shows that a convex function has a subgradient at x if there is at least one nonvertical supporting hyperplane to $\mathbf{epi}\,f$ at $(x, f(x))$. This is the case, for example, if f is continuous. There are pathological convex functions which do not have subgradients at some points, but we will assume in the sequel that all convex functions are subdifferentiable (at every point in $\mathbf{dom}\,f$).

2.2 Subgradients of differentiable functions

If f is convex and differentiable at x , then $\partial f(x) = \{\nabla f(x)\}$, *i.e.*, its gradient is its only subgradient. Conversely, if f is convex and $\partial f(x) = \{g\}$, then f is differentiable at x and $g = \nabla f(x)$.

2.3 The minimum of a nondifferentiable function

A point x^* is a minimizer of a function f (not necessarily convex) if and only if f is subdifferentiable at x^* and

$$0 \in \partial f(x^*),$$

i.e., $g = 0$ is a subgradient of f at x^* . This follows directly from the fact that $f(x) \geq f(x^*)$ for all $x \in \mathbf{dom}\,f$. And clearly if f is subdifferentiable at x^* with $0 \in \partial f(x^*)$, then $f(x) \geq f(x^*) + 0^T(x - x^*) = f(x^*)$ for all x .

While this simple characterization of optimality via the subdifferential holds for nonconvex functions, it is not particularly useful in that case, since we generally cannot find the subdifferential of a nonconvex function.

The condition $0 \in \partial f(x^*)$ reduces to $\nabla f(x^*) = 0$ when f is convex and differentiable at x^* .

2.4 Directional derivatives and subgradients

For convex functions f , the *directional derivative* of f at the point $x \in \mathbf{R}^n$ in the direction v is

$$f'(x; v) \triangleq \lim_{t \searrow 0} \frac{f(x + tv) - f(x)}{t}.$$

This quantity always exists for convex f , though it may be $+\infty$ or $-\infty$. To see the existence of the limit, we use that the ratio $(f(x + tv) - f(x))/t$ is non-decreasing in t . For $0 < t_1 \leq t_2$, we have $0 \leq t_1/t_2 \leq 1$, and

$$\begin{aligned} \frac{f(x + t_1 v) - f(x)}{t_1} &= \frac{f(\frac{t_1}{t_2}(x + t_2 v) + (1 - \frac{t_1}{t_2})x) - f(x)}{t_1} \\ &\leq \frac{\frac{t_1}{t_2} f(x + t_2 v)}{t_1} + \frac{(1 - \frac{t_1}{t_2}) f(x) - f(x)}{t_1} = \frac{f(x + t_2 v) - f(x)}{t_2}, \end{aligned}$$

so the limit in the definition of $f'(x; v)$ exists.

The directional derivative $f'(x; v)$ possesses several interesting properties as well. First, it is convex in v , and if f is finite in a neighborhood of x , then $f'(x; v)$ exists. Additionally, f is differentiable at x if and only if for some g (which is $\nabla f(x)$) and all $v \in \mathbf{R}^n$ we have $f'(x; v) = g^T v$, that is, if and only if $f'(x; v)$ is a linear function of v .¹ For general convex f , $f'(x; v)$ is *positively homogeneous* in v , meaning that for $\alpha \geq 0$, we have $f'(x; \alpha v) = \alpha f'(x; v)$ (replace t by t/α in the defining limit).

The directional derivative $f'(x; v)$ satisfies the following general formula for convex f :

$$f'(x; v) = \sup_{g \in \partial f(x)} g^T v.$$

To see this inequality, note that $f'(x; v) \geq \sup_{g \in \partial f(x)} g^T v$ by the definition of a subgradient: $f(x + tv) - f(x) \geq t g^T v$ for any $t \in \mathbf{R}$ and $g \in \partial f(x)$, so $f'(x; v) \geq \sup_{g \in \partial f(x)} g^T v$. For the other direction, we claim that all affine functions that are below the function $v \mapsto f'(x; v)$ may be taken to be linear. Specifically, suppose that $(g, r) \in \mathbf{R}^n \times \mathbf{R}$ and $g^T v - r \leq f'(x; v)$ for all v . Then $r \geq 0$, as taking $v = 0$ gives $-r \leq f'(x; 0) = 0$. By the positive homogeneity of $f'(x; v)$, we see that for any $t \geq 0$ we have $t g^T v - r \leq f'(x; tv) = t f'(x; v)$, and thus we have

$$g^T v - \frac{r}{t} \leq f'(x; v) \quad \text{for all } t > 0.$$

¹This is simply the standard definition of differentiability.

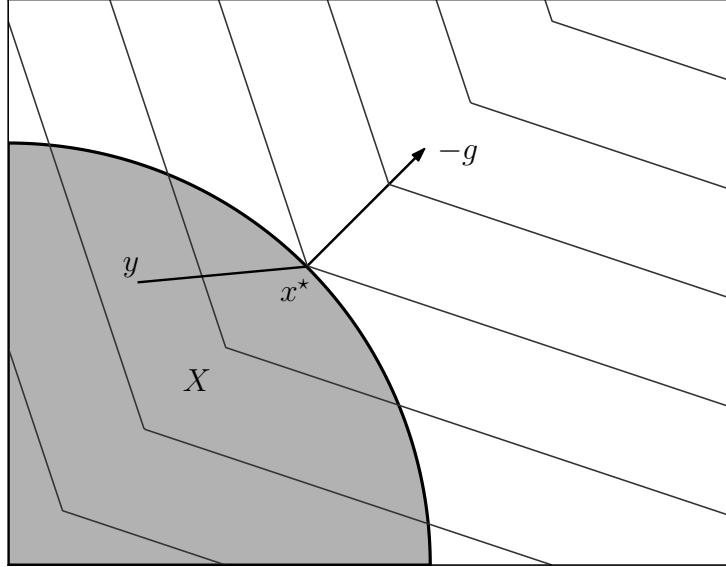


Figure 4: The point x^* minimizes f over X (the shown level curves) if and only if for some $g \in \partial f(x^*)$, $g^T(y - x^*) \geq 0$ for all $y \in X$. Note that not all subgradients satisfy this inequality.

Taking $t \rightarrow +\infty$ gives that any affine minorizer of $f'(x; v)$ may be taken to be linear. As any (closed) convex function can be written as the supremum of its affine minorants, we have

$$f'(x; v) = \sup \{g^T v \mid g^T \Delta \leq f'(x; \Delta) \text{ for all } \Delta \in \mathbf{R}^n\}.$$

On the other hand, if $g^T \Delta \leq f'(x; \Delta)$ for all $\Delta \in \mathbf{R}^n$, then we have $g^T \Delta \leq f(x + \Delta) - f(x)$, so that $g \in \partial f(x)$, and we may as well have taken the preceding supremum only over $g \in \partial f(x)$.

2.5 Constrained minimizers of nondifferentiable functions

There is a somewhat more complex version of the result that $0 \in \partial f(x)$ if and only if x minimizes f for constrained minimization. Consider finding the minimizer of a subdifferentiable function f over a (closed) convex set X . Then we have that x^* minimizes f if and only if there exists a subgradient $g \in \partial f(x^*)$ such that

$$g^T(y - x^*) \geq 0 \text{ for all } y \in X.$$

See Fig. 4 for an illustration of this condition.

To see this result, first suppose that $g \in \partial f(x^*)$ satisfies the preceding condition. Then by definition, $f(x) \geq f(x^*) + g^T(x - x^*) \geq f(x^*)$ for $x \in X$. The converse is more subtle, and we show it under the assumption that $x^* \in \mathbf{int dom} f$, though x^* may be on the boundary of

X . We suppose that $f(x) \geq f(x^*)$ for all $x \in X$. In this case, for any $x \in X$, the directional derivative

$$f'(x^*; x - x^*) = \lim_{t \searrow 0} \frac{f(x^* + t(x - x^*)) - f(x^*)}{t} \geq 0,$$

that is, for any x , the direction $\Delta = x - x^*$ pointing into X satisfies $f'(x^*; \Delta) \geq 0$. By our characterization of the directional derivative earlier, we know that $f'(x^*; \Delta) = \sup_{g \in \partial f(x^*)} g^T \Delta \geq 0$. Thus, defining the ball $\mathbf{B}_\epsilon = \{y + x^* \in \mathbf{R}^n \mid \|y\|_2 \leq \epsilon\}$, we have

$$\inf_{x \in X \cap \mathbf{B}_\epsilon} \sup_{g \in \partial f(x^*)} g^T (x - x^*) \geq 0.$$

As $\partial f(x^*)$ is bounded, we may swap the min and max (see, for example, Exercise 5.25 of [BV04]), finding that there must exist some $g \in \partial f(x^*)$ such that

$$\inf_{x \in X \cap \mathbf{B}_\epsilon} g^T (x - x^*) \geq 0.$$

But any $y \in X$ may be written as $t(x - x^*) + x^*$ for some $t \geq 0$ and $x \in X \cap \mathbf{B}_\epsilon$, which gives the result.

For fuller explanations of these inequalities and derivations, see also the books by Hiriart-Urruty and Lemaréchal [HUL93, HUL01].

3 Calculus of subgradients

In this section we describe rules for constructing subgradients of convex functions. We will distinguish two levels of detail. In the ‘weak’ calculus of subgradients the goal is to produce one subgradient, even if more subgradients exist. This is sufficient in practice, since subgradient, localization, and cutting-plane methods require only *a* subgradient at any point.

A second and much more difficult task is to describe the complete set of subgradients $\partial f(x)$ as a function of x . We will call this the ‘strong’ calculus of subgradients. It is useful in theoretical investigations, for example, when describing the precise optimality conditions.

3.1 Nonnegative scaling

For $\alpha \geq 0$, $\partial(\alpha f)(x) = \alpha \partial f(x)$.

3.2 Sum and integral

Suppose $f = f_1 + \dots + f_m$, where f_1, \dots, f_m are convex functions. Then we have

$$\partial f(x) = \partial f_1(x) + \dots + \partial f_m(x).$$

This property extends to infinite sums, integrals, and expectations (provided they exist).

3.3 Affine transformations of domain

Suppose f is convex, and let $h(x) = f(Ax + b)$. Then $\partial h(x) = A^T \partial f(Ax + b)$.

3.4 Pointwise maximum

Suppose f is the pointwise maximum of convex functions f_1, \dots, f_m , *i.e.*,

$$f(x) = \max_{i=1, \dots, m} f_i(x),$$

where the functions f_i are subdifferentiable. We first show how to construct a subgradient of f at x .

Let k be any index for which $f_k(x) = f(x)$, and let $g \in \partial f_k(x)$. Then $g \in \partial f(x)$. In other words, to find a subgradient of the maximum of functions, we can choose one of the functions that achieves the maximum at the point, and choose any subgradient of that function at the point. This follows from

$$f(z) \geq f_k(z) \geq f_k(x) + g^T(z - x) = f(x) + g^T(z - x).$$

More generally, we have

$$\partial f(x) = \mathbf{Co} \cup \{\partial f_i(x) \mid f_i(x) = f(x)\},$$

i.e., the subdifferential of the maximum of functions is the convex hull of the union of subdifferentials of the ‘active’ functions at x .

Example. *Maximum of differentiable functions.* Suppose $f(x) = \max_{i=1, \dots, m} f_i(x)$, where f_i are convex and differentiable. Then we have

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

At a point x where only one of the functions, say f_k , is active, f is differentiable and has gradient $\nabla f_k(x)$. At a point x where several of the functions are active, $\partial f(x)$ is a polyhedron.

Example. ℓ_1 -norm. The ℓ_1 -norm

$$f(x) = \|x\|_1 = |x_1| + \dots + |x_n|$$

is a nondifferentiable convex function of x . To find its subgradients, we note that f can be expressed as the maximum of 2^n linear functions:

$$\|x\|_1 = \max\{s^T x \mid s_i \in \{-1, 1\}\},$$

so we can apply the rules for the subgradient of the maximum. The first step is to identify an active function $s^T x$, *i.e.*, find an $s \in \{-1, +1\}^n$ such that $s^T x = \|x\|_1$. We can choose $s_i = +1$ if $x_i > 0$, and $s_i = -1$ if $x_i < 0$. If $x_i = 0$, more than one function

is active, and both $s_i = +1$, and $s_i = -1$ work. The function $s^T x$ is differentiable and has a unique subgradient s . We can therefore take

$$g_i = \begin{cases} +1 & x_i > 0 \\ -1 & x_i < 0 \\ -1 \text{ or } +1 & x_i = 0. \end{cases}$$

The subdifferential is the convex hull of all subgradients that can be generated this way:

$$\partial f(x) = \{g \mid \|g\|_\infty \leq 1, g^T x = \|x\|_1\}.$$

3.5 Supremum

Next we consider the extension to the supremum over an infinite number of functions, *i.e.*, we consider

$$f(x) = \sup_{\alpha \in \mathcal{A}} f_\alpha(x),$$

where the functions f_α are subdifferentiable. We only discuss the weak property.

Suppose the supremum in the definition of $f(x)$ is attained. Let $\beta \in \mathcal{A}$ be an index for which $f_\beta(x) = f(x)$, and let $g \in \partial f_\beta(x)$. Then $g \in \partial f(x)$. If the supremum in the definition is not attained, the function may or may not be subdifferentiable at x , depending on the index set \mathcal{A} .

Assume however that \mathcal{A} is compact (in some metric), and that the function $\alpha \mapsto f_\alpha(x)$ is upper semi-continuous for each x . Then

$$\partial f(x) = \mathbf{Co} \cup \{\partial f_\alpha(x) \mid f_\alpha(x) = f(x)\}.$$

Example. *Maximum eigenvalue of a symmetric matrix.* Let $f(x) = \lambda_{\max}(A(x))$, where $A(x) = A_0 + x_1 A_1 + \cdots + x_n A_n$, and $A_i \in \mathbf{S}^m$. We can express f as the pointwise supremum of convex functions,

$$f(x) = \lambda_{\max}(A(x)) = \sup_{\|y\|_2=1} y^T A(x) y.$$

Here the index set \mathcal{A} is $\mathcal{A} = \{y \in \mathbf{R}^n \mid \|y\|_2 = 1\}$.

Each of the functions $f_y(x) = y^T A(x) y$ is affine in x for fixed y , as can be easily seen from

$$y^T A(x) y = y^T A_0 y + x_1 y^T A_1 y + \cdots + x_n y^T A_n y,$$

so it is differentiable with gradient $\nabla f_y(x) = (y^T A_1 y, \dots, y^T A_n y)$.

The active functions $y^T A(x) y$ are those associated with the eigenvectors corresponding to the maximum eigenvalue. Hence to find a subgradient, we compute an eigenvector y with eigenvalue λ_{\max} , normalized to have unit norm, and take

$$g = (y^T A_1 y, y^T A_2 y, \dots, y^T A_n y).$$

The ‘index set’ in this example is $\{y \mid \|y\| = 1\}$ is a compact set. Therefore

$$\partial f(x) = \mathbf{Co} \{\nabla f_y \mid A(x) y = \lambda_{\max}(A(x)) y, \|y\| = 1\}.$$

Example. *Maximum eigenvalue of a symmetric matrix, revisited.* Let $f(A) = \lambda_{\max}(A)$, where $A \in \mathbf{S}^n$, the symmetric n -by- n matrices. Then as above, $f(A) = \lambda_{\max}(A) = \sup_{\|y\|_2=1} y^T A y$, but we note that $y^T A y = \mathbf{Tr}(A y y^T)$, so that each of the functions $f_y(A) = y^T A y$ is linear in A with gradient $\nabla f_y(A) = y y^T$. Then using an identical argument to that above, we find that

$$\partial f(A) = \mathbf{Co} \{ y y^T \mid \|y\|_2 = 1, y^T A y = \lambda_{\max}(A) \} = \mathbf{Co} \{ y y^T \mid \|y\|_2 = 1, A y = \lambda_{\max}(A) y \},$$

the convex hull of the outer products of maximum eigenvectors of the matrix A .

3.6 Minimization over some variables

The next subgradient calculus rule concerns functions of the form

$$f(x) = \inf_y F(x, y)$$

where $F(x, y)$ is subdifferentiable and jointly convex in $x \in \mathbf{R}^n$ and $y \in \mathbf{R}^m$.

Suppose that the infimum over y in the definition of $f(x)$ is attained on the set $Y_x \subset \mathbf{R}^m$ (where $Y_x \neq \emptyset$), so that $F(x, y) = f(x)$ for $y \in Y_x$. By definition, a vector $g \in \mathbf{R}^n$ is a subgradient of f if and only if

$$f(x') \geq f(x) + g^T(x' - x) = F(x, y) + g^T(x' - x)$$

for all $x' \in \mathbf{R}^n$ and any $y \in Y_x$. This is equivalent to

$$F(x', y') \geq F(x, y) + g^T(x' - x) = F(x, y) + \begin{bmatrix} g \\ 0 \end{bmatrix}^T \left(\begin{bmatrix} x' \\ y' \end{bmatrix} - \begin{bmatrix} x \\ y \end{bmatrix} \right)$$

for all $(x', y') \in \mathbf{R}^n \times \mathbf{R}^m$ and $x, y \in Y_x$. In particular, we have the result that

$$\partial f(x) = \{ g \in \mathbf{R}^n \mid (g, 0) \in \partial F(x, y) \text{ for some } y \in Y_x \}.$$

That is, there exist $g \in \mathbf{R}^n$ such that $(g, 0) \in \partial F(x, y)$ for some $y \in Y_x$, and any such g is a subgradient of f at x (as long as the infimum is attained and $x \in \mathbf{int dom } f$).

3.7 Optimal value function of a convex optimization problem

Suppose $f : \mathbf{R}^m \times \mathbf{R}^p \rightarrow \mathbf{R}$ is defined as the optimal value of a convex optimization problem in standard form, with $z \in \mathbf{R}^n$ as optimization variable,

$$\begin{aligned} & \text{minimize} && f_0(z) \\ & \text{subject to} && f_i(z) \leq x_i, \quad i = 1, \dots, m \\ & && A z = y. \end{aligned} \tag{2}$$

In other words, $f(x, y) = \inf_z F(x, y, z)$ where

$$F(x, y, z) = \begin{cases} f_0(z) & f_i(z) \leq x_i, \quad i = 1, \dots, m, \quad A z = y \\ +\infty & \text{otherwise,} \end{cases}$$

which is jointly convex in x, y, z . Subgradients of f can be related to the dual problem of (2) as follows.

Suppose we are interested in subdifferentiating f at (\hat{x}, \hat{y}) . We can express the dual problem of (2) as

$$\begin{aligned} & \text{maximize} && g(\lambda) - x^T \lambda - y^T \nu \\ & \text{subject to} && \lambda \succeq 0. \end{aligned} \tag{3}$$

where

$$g(\lambda) = \inf_z \left(f_0(z) + \sum_{i=1}^m \lambda_i f_i(z) + \nu^T A z \right).$$

Suppose strong duality holds for problems (2) and (3) at $x = \hat{x}$ and $y = \hat{y}$, and that the dual optimum is attained at λ^*, ν^* (for example, because Slater's condition holds). From the global perturbation inequalities we know that

$$f(x, y) \geq f(\hat{x}, \hat{y}) - \lambda^{*T}(x - \hat{x}) - \nu^{*T}(y - \hat{y})$$

In other words, the dual optimal solution provides a subgradient:

$$-(\lambda^*, \nu^*) \in \partial f(\hat{x}, \hat{y}).$$

4 Quasigradients

If $f(x)$ is quasiconvex, then g is a *quasigradient* at x_0 if

$$g^T(x - x_0) \geq 0 \Rightarrow f(x) \geq f(x_0),$$

Geometrically, g defines a supporting hyperplane to the sublevel set $\{x \mid f(x) \leq f(x_0)\}$. Note that the set of quasigradients at x_0 form a cone.

Example. *Linear fractional function.* $f(x) = \frac{a^T x + b}{c^T x + d}$. Let $c^T x_0 + d > 0$. Then $g = a - f(x_0)c$ is a quasigradient at x_0 . If $c^T x + d > 0$, we have

$$a^T(x - x_0) \geq f(x_0)c^T(x - x_0) \implies f(x) \geq f(x_0).$$

Example. *Degree of a polynomial.* Define $f : \mathbf{R}^n \rightarrow \mathbf{R}$ by

$$f(a) = \min\{i \mid a_{i+2} = \dots = a_n = 0\},$$

i.e., the degree of the polynomial $a_1 + a_2 t + \dots + a_n t^{n-1}$. Let $a \neq 0$, and $k = f(a)$, then $g = \text{sign}(a_{k+1})e_{k+1}$ is a quasigradient at a

To see this, we note that

$$g^T(b - a) = \text{sign}(a_{k+1})b_{k+1} - |a_{k+1}| \geq 0$$

implies $b_{k+1} \neq 0$.

References

- [BV04] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [HUL93] J. Hiriart-Urruty and C. Lemaréchal. *Convex Analysis and Minimization Algorithms I & II*. Springer, New York, 1993.
- [HUL01] Jean-Baptiste Hiriart-Urruty and Claude Lemaréchal. *Fundamentals of Convex Analysis*. Springer, 2001.