

Subgradients

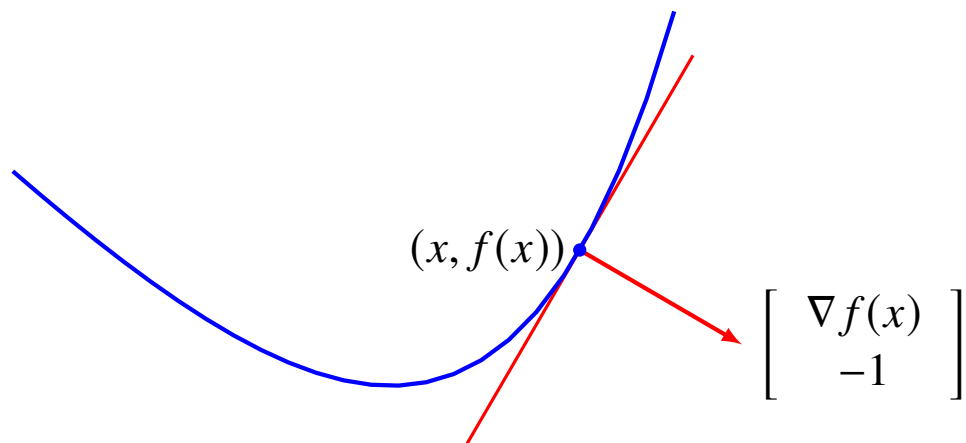
Acknowledgement: slides are based on Prof. Lieven Vandenberghes.

- definition
- subgradient calculus
- duality and optimality conditions
- directional derivative

Basic inequality

recall the basic inequality for differentiable convex functions:

$$f(y) \geq f(x) + \nabla f(x)^T (y - x) \quad \text{for all } y \in \text{dom } f$$



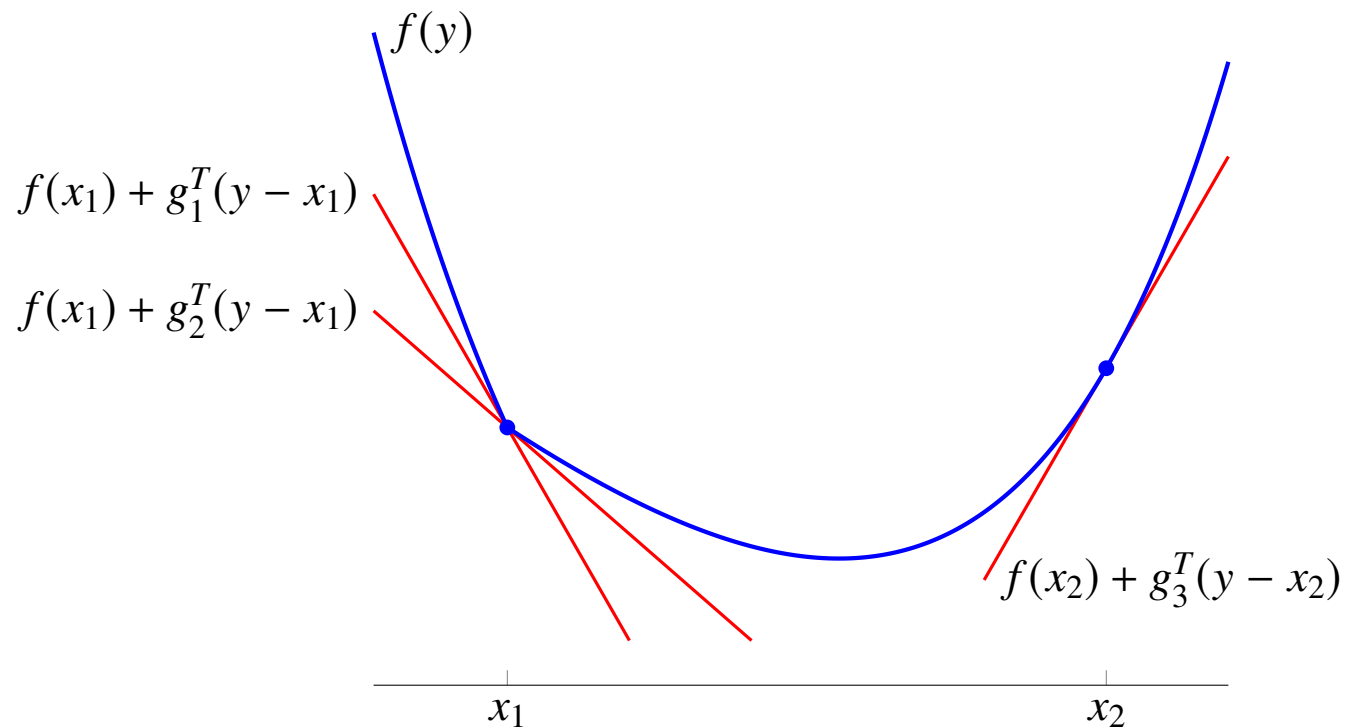
- the first-order approximation of f at x is a global lower bound
- $\nabla f(x)$ defines a non-vertical supporting hyperplane to $\text{epi } f$ at $(x, f(x))$:

$$\begin{bmatrix} \nabla f(x) \\ -1 \end{bmatrix}^T \left(\begin{bmatrix} y \\ t \end{bmatrix} - \begin{bmatrix} x \\ f(x) \end{bmatrix} \right) \leq 0 \quad \text{for all } (y, t) \in \text{epi } f$$

Subgradient

g is a *subgradient* of a convex function f at $x \in \text{dom } f$ if

$$f(y) \geq f(x) + g^T(y - x) \quad \text{for all } y \in \text{dom } f$$



g_1, g_2 are subgradients at x_1 ; g_3 is a subgradient at x_2

Properties

- $f(x) + g^\top(y - x)$ is a global lower bound on $f(y)$
- g defines non-vertical supporting hyperplane to $\text{epi} f$ at $(x, f(x))$:

$$\begin{bmatrix} g \\ -1 \end{bmatrix}^\top \left(\begin{bmatrix} y \\ t \end{bmatrix} - \begin{bmatrix} x \\ f(x) \end{bmatrix} \right) \leq 0, \forall (y, t) \in \text{epi} f$$

- If f is convex and differentiable, then $\nabla f(x)$ is a subgradient of g at x
- algorithms for non-differentiable convex optimization
- unconstrained optimality: x minimizes $f(x)$ iff $0 \in \partial f(x)$
- KKT conditions with non-differentiable functions.

Subdifferential

the *subdifferential* $\partial f(x)$ of f at x is the set of all subgradients:

$$\partial f(x) = \{g \mid g^T(y - x) \leq f(y) - f(x), \forall y \in \text{dom } f\}$$

Properties

- $\partial f(x)$ is a closed convex set (possibly empty)

this follows from the definition: $\partial f(x)$ is an intersection of halfspaces

- if $x \in \text{int dom } f$ then $\partial f(x)$ is nonempty and bounded

proof on next two pages

Proof: we show that $\partial f(x)$ is nonempty when $x \in \text{int dom } f$

- $(x, f(x))$ is in the boundary of the convex set $\text{epi } f$
- therefore there exists a supporting hyperplane to $\text{epi } f$ at $(x, f(x))$:

$$\exists(a, b) \neq 0, \quad \begin{bmatrix} a \\ b \end{bmatrix}^T \left(\begin{bmatrix} y \\ t \end{bmatrix} - \begin{bmatrix} x \\ f(x) \end{bmatrix} \right) \leq 0 \quad \forall (y, t) \in \text{epi } f$$

- $b > 0$ gives a contradiction as $t \rightarrow \infty$
- $b = 0$ gives a contradiction for $y = x + \epsilon a$ with small $\epsilon > 0$
- therefore $b < 0$ and $g = \frac{1}{|b|}a$ is a subgradient of f at x

Proof: $\partial f(x)$ is bounded when $x \in \text{int dom } f$

- for small $r > 0$, define a set of $2n$ points

$$B = \{x \pm r e_k \mid k = 1, \dots, n\} \subset \text{dom } f$$

and define $M = \max_{y \in B} f(y) < \infty$

- for every $g \in \partial f(x)$, there is a point $y \in B$ with

$$r \|g\|_\infty = g^T(y - x)$$

(choose an index k with $|g_k| = \|g\|_\infty$, and take $y = x + r \text{sign}(g_k)e_k$)

- since g is a subgradient, this implies that

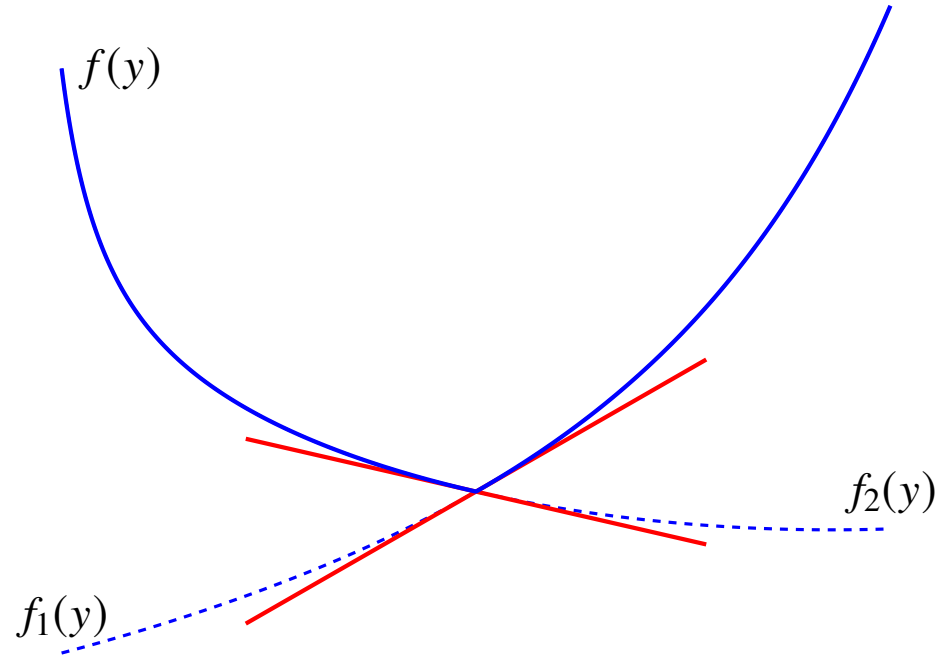
$$f(x) + r \|g\|_\infty = f(x) + g^T(y - x) \leq f(y) \leq M$$

- we conclude that $\partial f(x)$ is bounded:

$$\|g\|_\infty \leq \frac{M - f(x)}{r} \quad \text{for all } g \in \partial f(x)$$

Example

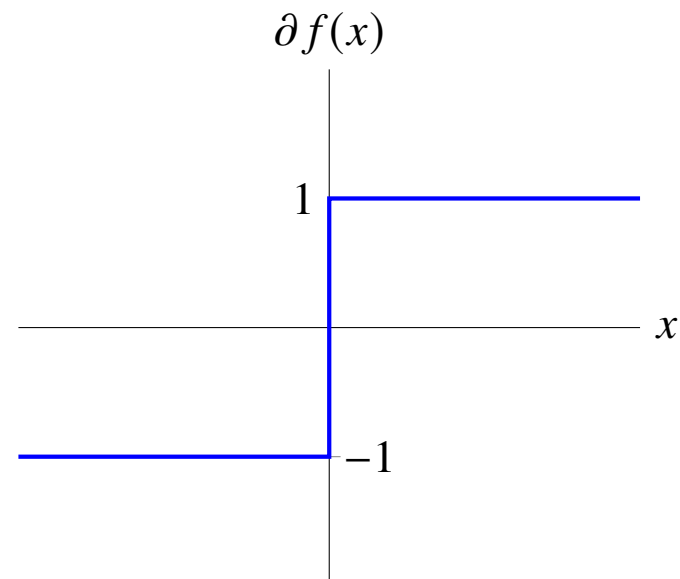
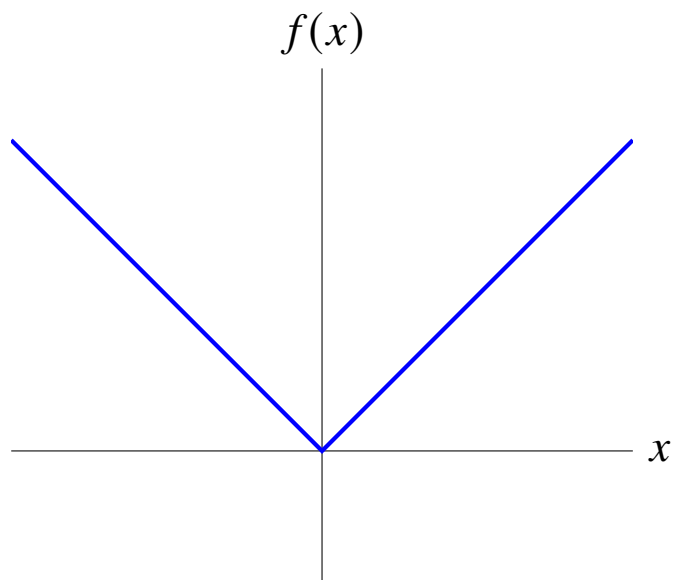
$$f(x) = \max \{f_1(x), f_2(x)\} \quad \text{with } f_1, f_2 \text{ convex and differentiable}$$



- if $f_1(\hat{x}) = f_2(\hat{x})$, subdifferential at \hat{x} is line segment $[\nabla f_1(\hat{x}), \nabla f_2(\hat{x})]$
- if $f_1(\hat{x}) > f_2(\hat{x})$, subdifferential at \hat{x} is $\{\nabla f_1(\hat{x})\}$
- if $f_1(\hat{x}) < f_2(\hat{x})$, subdifferential at \hat{x} is $\{\nabla f_2(\hat{x})\}$

Examples

Absolute value $f(x) = |x|$



Euclidean norm $f(x) = \|x\|_2$

$$\partial f(x) = \left\{ \frac{1}{\|x\|_2} x \right\} \quad \text{if } x \neq 0, \quad \partial f(x) = \{g \mid \|g\|_2 \leq 1\} \quad \text{if } x = 0$$

Monotonicity

the subdifferential of a convex function is a *monotone operator*:

$$(u - v)^T(x - y) \geq 0 \quad \text{for all } x, y, u \in \partial f(x), v \in \partial f(y)$$

Proof: by definition

$$f(y) \geq f(x) + u^T(y - x), \quad f(x) \geq f(y) + v^T(x - y)$$

combining the two inequalities shows monotonicity

Examples of non-subdifferentiable functions

the following functions are not subdifferentiable at $x = 0$

- $f : \mathbf{R} \rightarrow \mathbf{R}, \text{dom } f = \mathbf{R}_+$

$$f(x) = 1 \quad \text{if } x = 0, \quad f(x) = 0 \quad \text{if } x > 0$$

- $f : \mathbf{R} \rightarrow \mathbf{R}, \text{dom } f = \mathbf{R}_+$

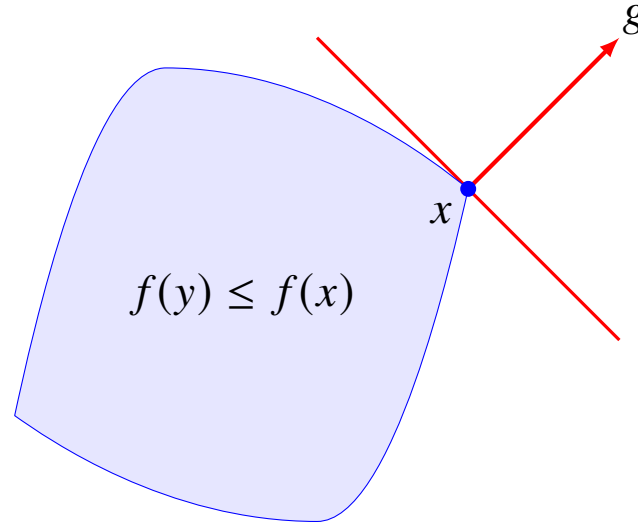
$$f(x) = -\sqrt{x}$$

the only supporting hyperplane to $\text{epi } f$ at $(0, f(0))$ is vertical

Subgradients and sublevel sets

if g is a subgradient of f at x , then

$$f(y) \leq f(x) \implies g^T(y - x) \leq 0$$



the nonzero subgradients at x define supporting hyperplanes to the sublevel set

$$\{y \mid f(y) \leq f(x)\}$$

Outline

- definition
- **subgradient calculus**
- duality and optimality conditions
- directional derivative

Subgradient calculus

Weak subgradient calculus: rules for finding *one* subgradient

- sufficient for most nondifferentiable convex optimization algorithms
- if you can evaluate $f(x)$, you can usually compute a subgradient

Strong subgradient calculus: rules for finding $\partial f(x)$ (*all* subgradients)

- some algorithms, optimality conditions, etc., need entire subdifferential
- can be quite complicated

we will assume that $x \in \text{int dom } f$

Basic rules

Differentiable functions: $\partial f(x) = \{\nabla f(x)\}$ if f is differentiable at x

Nonnegative linear combination

if $f(x) = \alpha_1 f_1(x) + \alpha_2 f_2(x)$ with $\alpha_1, \alpha_2 \geq 0$, then

$$\partial f(x) = \alpha_1 \partial f_1(x) + \alpha_2 \partial f_2(x)$$

(right-hand side is addition of sets)

Affine transformation of variables: if $f(x) = h(Ax + b)$, then

$$\partial f(x) = A^T \partial h(Ax + b)$$

Pointwise maximum

$$f(x) = \max \{f_1(x), \dots, f_m(x)\}$$

define $I(x) = \{i \mid f_i(x) = f(x)\}$, the ‘active’ functions at x

Weak result:

to compute a subgradient at x , choose any $k \in I(x)$, any subgradient of f_k at x

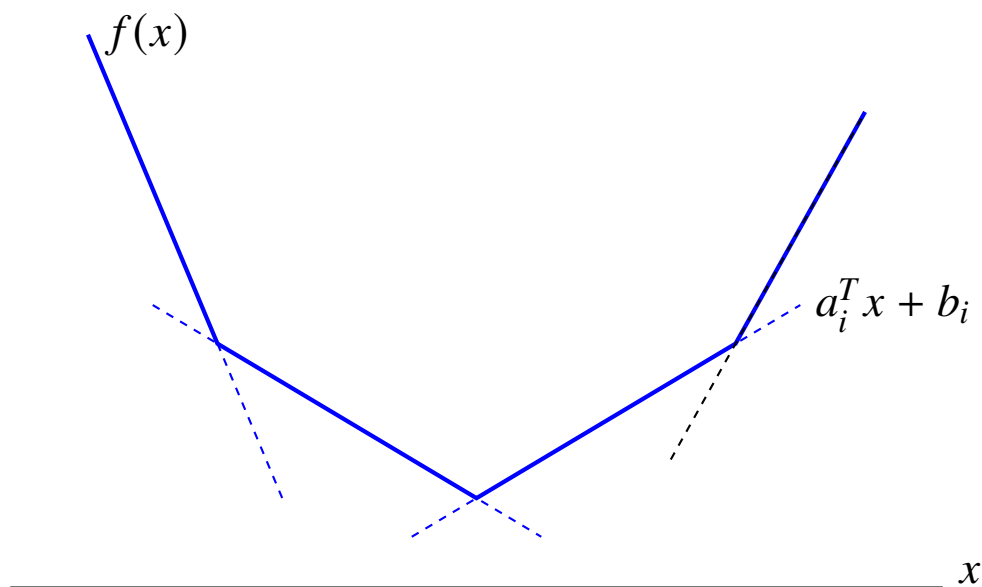
Strong result

$$\partial f(x) = \text{conv} \bigcup_{i \in I(x)} \partial f_i(x)$$

- the convex hull of the union of subdifferentials of ‘active’ functions at x
- if f_i ’s are differentiable, $\partial f(x) = \text{conv} \{\nabla f_i(x) \mid i \in I(x)\}$

Example: piecewise-linear function

$$f(x) = \max_{i=1,\dots,m} (a_i^T x + b_i)$$



the subdifferential at x is a polyhedron

$$\partial f(x) = \text{conv} \{a_i \mid i \in I(x)\}$$

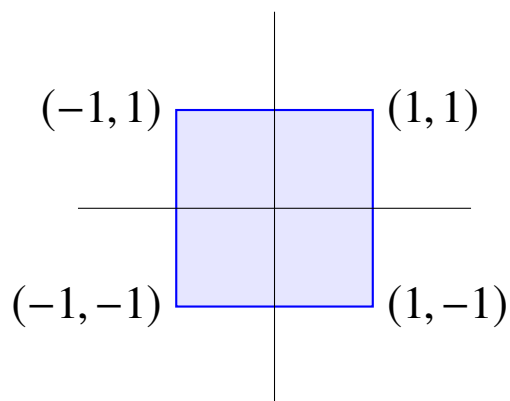
with $I(x) = \{i \mid a_i^T x + b_i = f(x)\}$

Example: ℓ_1 -norm

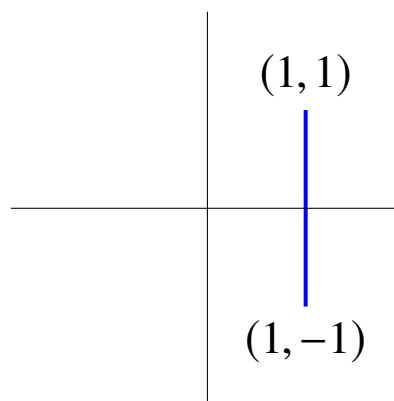
$$f(x) = \|x\|_1 = \max_{s \in \{-1,1\}^n} s^T x$$

the subdifferential is a product of intervals

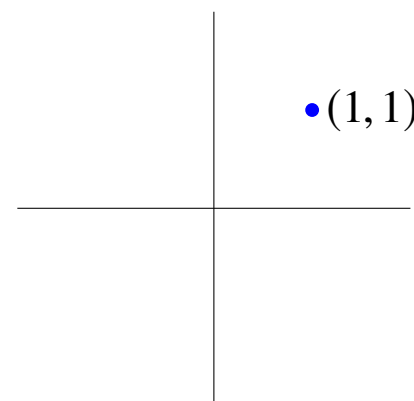
$$\partial f(x) = J_1 \times \cdots \times J_n, \quad J_k = \begin{cases} [-1, 1] & x_k = 0 \\ \{1\} & x_k > 0 \\ \{-1\} & x_k < 0 \end{cases}$$



$$\partial f(0,0) = [-1, 1] \times [-1, 1]$$



$$\partial f(1,0) = \{1\} \times [-1, 1]$$



$$\partial f(1,1) = \{(1,1)\}$$

Pointwise supremum

$$f(x) = \sup_{\alpha \in \mathcal{A}} f_{\alpha}(x), \quad f_{\alpha}(x) \text{ convex in } x \text{ for every } \alpha$$

Weak result: to find a subgradient at \hat{x} ,

- find *any* β for which $f(\hat{x}) = f_{\beta}(\hat{x})$ (assuming maximum is attained)
- choose *any* $g \in \partial f_{\beta}(\hat{x})$

(Partial) strong result: define $I(x) = \{\alpha \in \mathcal{A} \mid f_{\alpha}(x) = f(x)\}$

$$\text{conv} \bigcup_{\alpha \in I(x)} \partial f_{\alpha}(x) \subseteq \partial f(x)$$

equality requires extra conditions (for example, \mathcal{A} compact, f_{α} continuous in α)

Exercise: maximum eigenvalue

Problem: explain how to find a subgradient of

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

where $A(x) = A_0 + x_1 A_1 + \cdots + x_n A_n$ with symmetric coefficients A_i

Solution: to find a subgradient at \hat{x} ,

- choose *any* unit eigenvector y with eigenvalue $\lambda_{\max}(A(\hat{x}))$
- the gradient of $y^T A(x) y$ at \hat{x} is a subgradient of f :

$$(y^T A_1 y, \dots, y^T A_n y) \in \partial f(\hat{x})$$

Minimization

$$f(x) = \inf_y h(x, y), \quad h \text{ jointly convex in } (x, y)$$

Weak result: to find a subgradient at \hat{x} ,

- find \hat{y} that minimizes $h(\hat{x}, y)$ (assuming minimum is attained)
- find subgradient $(g, 0) \in \partial h(\hat{x}, \hat{y})$

Proof: for all x, y ,

$$\begin{aligned} h(x, y) &\geq h(\hat{x}, \hat{y}) + g^T(x - \hat{x}) + 0^T(y - \hat{y}) \\ &= f(\hat{x}) + g^T(x - \hat{x}) \end{aligned}$$

therefore

$$f(x) = \inf_y h(x, y) \geq f(\hat{x}) + g^T(x - \hat{x})$$

Exercise: Euclidean distance to convex set

Problem: explain how to find a subgradient of

$$f(x) = \inf_{y \in C} \|x - y\|_2$$

where C is a closed convex set

Solution: to find a subgradient at \hat{x} ,

- if $f(\hat{x}) = 0$ (that is, $\hat{x} \in C$), take $g = 0$
- if $f(\hat{x}) > 0$, find projection $\hat{y} = P(\hat{x})$ on C and take

$$g = \frac{1}{\|\hat{y} - \hat{x}\|_2}(\hat{x} - \hat{y}) = \frac{1}{\|\hat{x} - P(\hat{x})\|_2}(\hat{x} - P(\hat{x}))$$

Composition

$$f(x) = h(f_1(x), \dots, f_k(x)), \quad h \text{ convex and nondecreasing, } f_i \text{ convex}$$

Weak result: to find a subgradient at \hat{x} ,

- find $z \in \partial h(f_1(\hat{x}), \dots, f_k(\hat{x}))$ and $g_i \in \partial f_i(\hat{x})$
- then $g = z_1 g_1 + \dots + z_k g_k \in \partial f(\hat{x})$

reduces to standard formula for differentiable h, f_i

Proof:

$$\begin{aligned} f(x) &\geq h\left(f_1(\hat{x}) + g_1^T(x - \hat{x}), \dots, f_k(\hat{x}) + g_k^T(x - \hat{x})\right) \\ &\geq h(f_1(\hat{x}), \dots, f_k(\hat{x})) + z^T \left(g_1^T(x - \hat{x}), \dots, g_k^T(x - \hat{x})\right) \\ &= h(f_1(\hat{x}), \dots, f_k(\hat{x})) + (z_1 g_1 + \dots + z_k g_k)^T (x - \hat{x}) \\ &= f(\hat{x}) + g^T (x - \hat{x}) \end{aligned}$$

Optimal value function

define $f(u, v)$ as the optimal value of convex problem

$$\begin{array}{ll}\text{minimize} & f_0(x) \\ \text{subject to} & f_i(x) \leq u_i, \quad i = 1, \dots, m \\ & Ax = b + v\end{array}$$

(functions f_i are convex; optimization variable is x)

Weak result: suppose $f(\hat{u}, \hat{v})$ is finite and strong duality holds with the dual

$$\begin{array}{ll}\text{maximize} & \inf_x \left(f_0(x) + \sum_i \lambda_i (f_i(x) - \hat{u}_i) + v^T (Ax - b - \hat{v}) \right) \\ \text{subject to} & \lambda \geq 0\end{array}$$

if $\hat{\lambda}, \hat{v}$ are optimal dual variables (for right-hand sides \hat{u}, \hat{v}) then $(-\hat{\lambda}, -\hat{v}) \in \partial f(\hat{u}, \hat{v})$

Proof: by weak duality for problem with right-hand sides u, v

$$\begin{aligned} f(u, v) &\geq \inf_x \left(f_0(x) + \sum_i \hat{\lambda}_i (f_i(x) - u_i) + \hat{v}^T (Ax - b - v) \right) \\ &= \inf_x \left(f_0(x) + \sum_i \hat{\lambda}_i (f_i(x) - \hat{u}_i) + \hat{v}^T (Ax - b - \hat{v}) \right) \\ &\quad - \hat{\lambda}^T (u - \hat{u}) - \hat{v}^T (v - \hat{v}) \\ &= f(\hat{u}, \hat{v}) - \hat{\lambda}^T (u - \hat{u}) - \hat{v}^T (v - \hat{v}) \end{aligned}$$

Expectation

$$f(x) = \mathbf{E} h(x, u) \quad u \text{ random, } h \text{ convex in } x \text{ for every } u$$

Weak result: to find a subgradient at \hat{x} ,

- choose a function $u \mapsto g(u)$ with $g(u) \in \partial_x h(\hat{x}, u)$
- then, $g = \mathbf{E}_u g(u) \in \partial f(\hat{x})$

Proof: by convexity of h and definition of $g(u)$,

$$\begin{aligned} f(x) &= \mathbf{E} h(x, u) \\ &\geq \mathbf{E} \left(h(\hat{x}, u) + g(u)^T (x - \hat{x}) \right) \\ &= f(\hat{x}) + g^T (x - \hat{x}) \end{aligned}$$

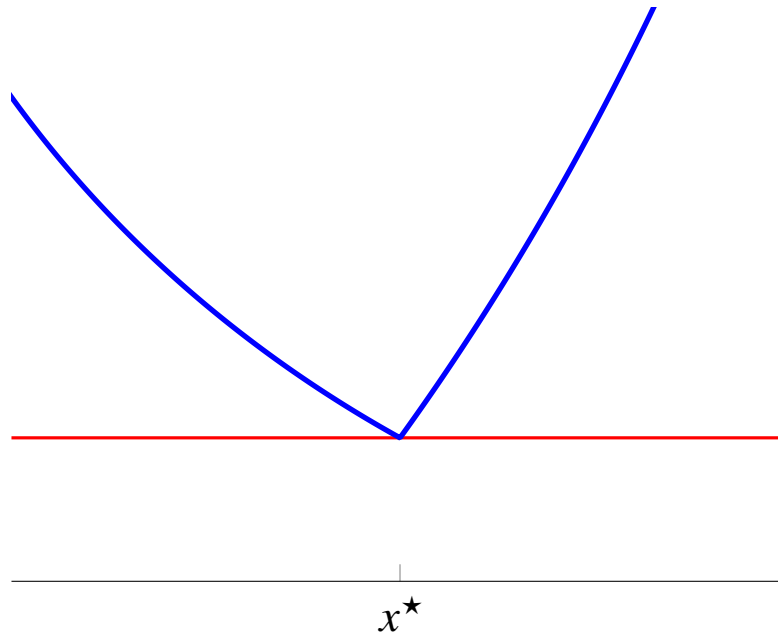
Outline

- definition
- subgradient calculus
- **duality and optimality conditions**
- directional derivative

Optimality conditions — unconstrained

x^\star minimizes $f(x)$ if and only

$$0 \in \partial f(x^\star)$$



this follows directly from the definition of subgradient:

$$f(y) \geq f(x^\star) + 0^T(y - x^\star) \quad \text{for all } y \quad \Longleftrightarrow \quad 0 \in \partial f(x^\star)$$

Example: piecewise-linear minimization

$$f(x) = \max_{i=1,\dots,m} (a_i^T x + b_i)$$

Optimality condition

$$0 \in \text{conv} \{a_i \mid i \in I(x^\star)\} \quad \text{where } I(x) = \{i \mid a_i^T x + b_i = f(x)\}$$

- in other words, x^\star is optimal if and only if there is a λ with

$$\lambda \geq 0, \quad \mathbf{1}^T \lambda = 1, \quad \sum_{i=1}^m \lambda_i a_i = 0, \quad \lambda_i = 0 \text{ for } i \notin I(x^\star)$$

- these are the optimality conditions for the equivalent linear program

$$\begin{array}{ll} \text{minimize} & t \\ \text{subject to} & Ax + b \leq t\mathbf{1} \end{array}$$

$$\begin{array}{ll} \text{maximize} & b^T \lambda \\ \text{subject to} & A^T \lambda = 0 \\ & \lambda \geq 0, \quad \mathbf{1}^T \lambda = 1 \end{array}$$

Optimality conditions — constrained

$$\begin{array}{ll}\text{minimize} & f_0(x) \\ \text{subject to} & f_i(x) \leq 0, \quad i = 1, \dots, m\end{array}$$

assume $\text{dom } f_i = \mathbf{R}^n$, so functions f_i are subdifferentiable everywhere

Karush–Kuhn–Tucker conditions

if strong duality holds, then x^\star , λ^\star are primal, dual optimal if and only if

1. x^\star is primal feasible
2. $\lambda^\star \geq 0$
3. $\lambda_i^\star f_i(x^\star) = 0$ for $i = 1, \dots, m$
4. x^\star is a minimizer of $L(x, \lambda^\star) = f_0(x) + \sum_{i=1}^m \lambda_i^\star f_i(x)$:

$$0 \in \partial f_0(x^\star) + \sum_{i=1}^m \lambda_i^\star \partial f_i(x^\star)$$

Outline

- definition
- subgradient calculus
- duality and optimality conditions
- **directional derivative**

Directional derivative

Definition (for general f): the *directional derivative* of f at x in the direction y is

$$\begin{aligned} f'(x; y) &= \lim_{\alpha \searrow 0} \frac{f(x + \alpha y) - f(x)}{\alpha} \\ &= \lim_{t \rightarrow \infty} \left(t \left(f\left(x + \frac{1}{t}y\right) - f(x) \right) \right) \end{aligned}$$

(if the limit exists)

- $f'(x; y)$ is the right derivative of $g(\alpha) = f(x + \alpha y)$ at $\alpha = 0$
- $f'(x; y)$ is homogeneous in y :

$$f'(x; \lambda y) = \lambda f'(x; y) \quad \text{for } \lambda \geq 0$$

Directional derivative of a convex function

Equivalent definition (for convex f): replace \lim with \inf

$$\begin{aligned} f'(x; y) &= \inf_{\alpha > 0} \frac{f(x + \alpha y) - f(x)}{\alpha} \\ &= \inf_{t > 0} \left(t f\left(x + \frac{1}{t}y\right) - t f(x) \right) \end{aligned}$$

Proof

- the function $h(y) = f(x + y) - f(x)$ is convex in y , with $h(0) = 0$
- its perspective $th(y/t)$ is nonincreasing in t (ECE236B ex. A2.5); hence

$$f'(x; y) = \lim_{t \rightarrow \infty} th(y/t) = \inf_{t > 0} th(y/t)$$

Properties

consequences of the expressions (for convex f)

$$\begin{aligned} f'(x; y) &= \inf_{\alpha > 0} \frac{f(x + \alpha y) - f(x)}{\alpha} \\ &= \inf_{t > 0} \left(t f\left(x + \frac{1}{t}y\right) - t f(x) \right) \end{aligned}$$

- $f'(x; y)$ is convex in y (partial minimization of a convex function in y, t)
- $f'(x; y)$ defines a lower bound on f in the direction y :

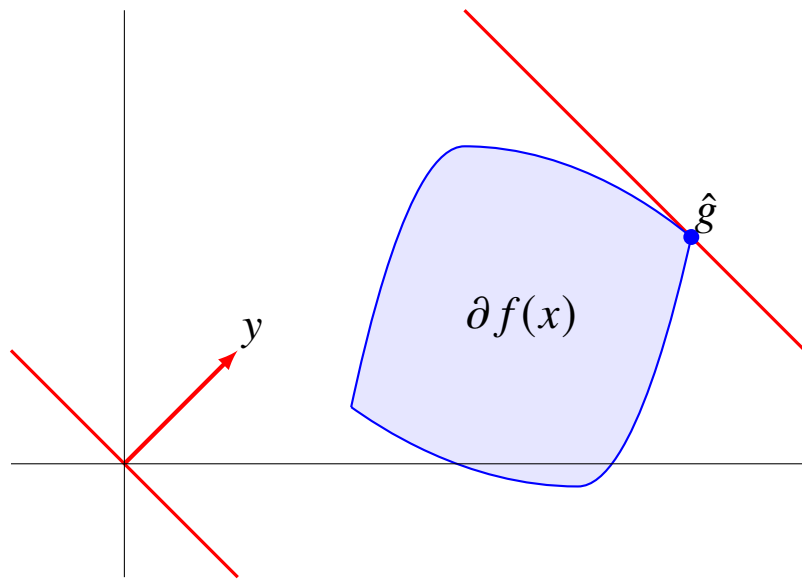
$$f(x + \alpha y) \geq f(x) + \alpha f'(x; y) \quad \text{for all } \alpha \geq 0$$

Directional derivative and subgradients

for convex f and $x \in \text{int dom } f$

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

$$\hat{f}'(x, y) = g^T y$$



$f'(x; y)$ is *support function* of $\partial f(x)$

- generalizes $f'(x; y) = \nabla f(x)^T y$ for differentiable functions
- implies that $f'(x; y)$ exists for all $x \in \text{int dom } f$, all y (see page 2.4)

Proof: if $g \in \partial f(x)$ then from page 2.29

$$f'(x; y) \geq \inf_{\alpha > 0} \frac{f(x) + \alpha g^T y - f(x)}{\alpha} = g^T y$$

it remains to show that $f'(x; y) = \hat{g}^T y$ for at least one $\hat{g} \in \partial f(x)$

- $f'(x; y)$ is convex in y with domain \mathbf{R}^n , hence subdifferentiable at all y
- let \hat{g} be a subgradient of $f'(x; y)$ at y : then for all $v, \lambda \geq 0$,

$$\lambda f'(x; v) = f'(x; \lambda v) \geq f'(x; y) + \hat{g}^T (\lambda v - y)$$

- taking $\lambda \rightarrow \infty$ shows that $f'(x; v) \geq \hat{g}^T v$; from the lower bound on page 2.30,

$$f(x + v) \geq f(x) + f'(x; v) \geq f(x) + \hat{g}^T v \quad \text{for all } v$$

hence $\hat{g} \in \partial f(x)$

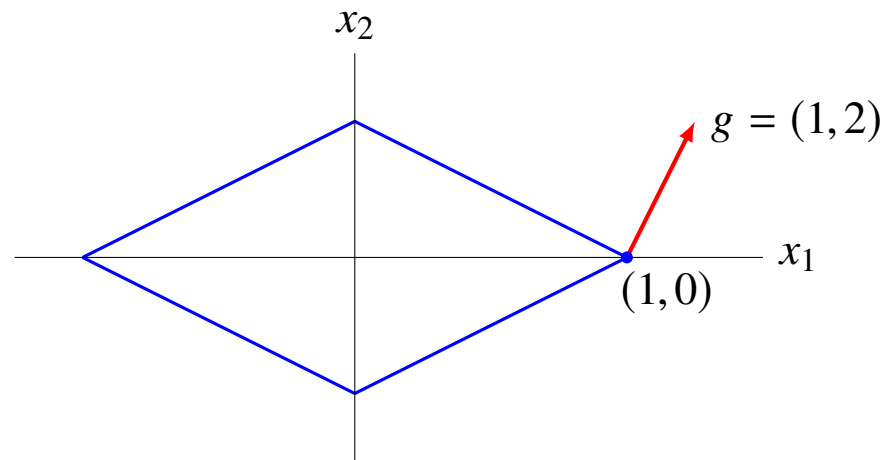
- taking $\lambda = 0$ we see that $f'(x; y) \leq \hat{g}^T y$

Descent directions and subgradients

y is a *descent direction* of f at x if $f'(x; y) < 0$

- the negative gradient of a differentiable f is a descent direction (if $\nabla f(x) \neq 0$)
- negative subgradient is **not** always a descent direction

Example: $f(x_1, x_2) = |x_1| + 2|x_2|$



$g = (1, 2) \in \partial f(1, 0)$, but $y = (-1, -2)$ is not a descent direction at $(1, 0)$

Steepest descent direction

Definition: (normalized) steepest descent direction at $x \in \text{int dom } f$ is

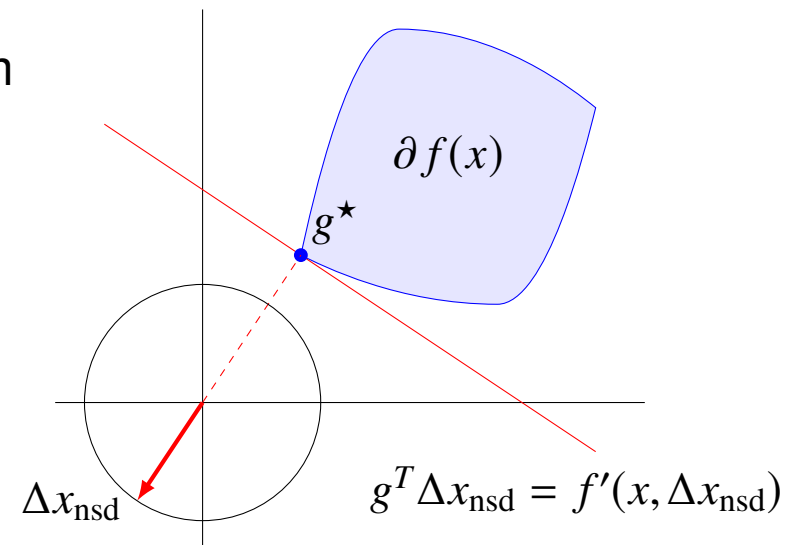
$$\Delta x_{\text{nsd}} = \underset{\|y\|_2 \leq 1}{\operatorname{argmin}} f'(x; y)$$

Δx_{nsd} is the primal solution y of the pair of dual problems (BV §8.1.3)

$$\begin{array}{ll} \text{minimize (over } y) & f'(x; y) \\ \text{subject to} & \|y\|_2 \leq 1 \end{array}$$

$$\begin{array}{ll} \text{maximize (over } g) & -\|g\|_2 \\ \text{subject to} & g \in \partial f(x) \end{array}$$

- dual optimal g^\star is subgradient with least norm
- $f'(x; \Delta x_{\text{nsd}}) = -\|g^\star\|_2$
- if $0 \notin \partial f(x)$, $\Delta x_{\text{nsd}} = -g^\star / \|g^\star\|_2$
- Δx_{nsd} can be expensive to compute



Subgradients and distance to sublevel sets

if f is convex, $f(y) < f(x)$, $g \in \partial f(x)$, then for small $t > 0$,

$$\begin{aligned}\|x - tg - y\|_2^2 &= \|x - y\|_2^2 - 2tg^T(x - y) + t^2\|g\|_2^2 \\ &\leq \|x - y\|_2^2 - 2t(f(x) - f(y)) + t^2\|g\|_2^2 \\ &< \|x - y\|_2^2\end{aligned}$$

- $-g$ is descent direction for $\|x - y\|_2$, for **any** y with $f(y) < f(x)$
- in particular, $-g$ is descent direction for distance to any minimizer of f

References

- A. Beck, *First-Order Methods in Optimization* (2017), chapter 3.
- D. P. Bertsekas, A. Nedić, A. E. Ozdaglar, *Convex Analysis and Optimization* (2003), chapter 4.
- J.-B. Hiriart-Urruty, C. Lemaréchal, *Convex Analysis and Minimization Algorithms* (1993), chapter VI.
- Yu. Nesterov, *Lectures on Convex Optimization* (2018), section 3.1.
- B. T. Polyak, *Introduction to Optimization* (1987), section 5.1.