

Gradient Characterization of Convexity

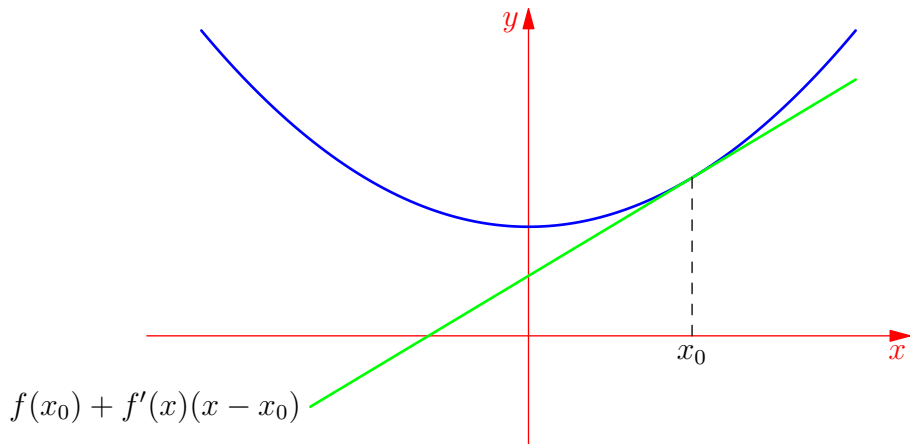
Theorem

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable. Then f is convex iff

$$f(x + v) \geq f(x) + \nabla f(x)^T v$$

hold for all $x, v \in \mathbb{R}^d$.

Gradient Approximation Gives Global Underestimator



Subgradients

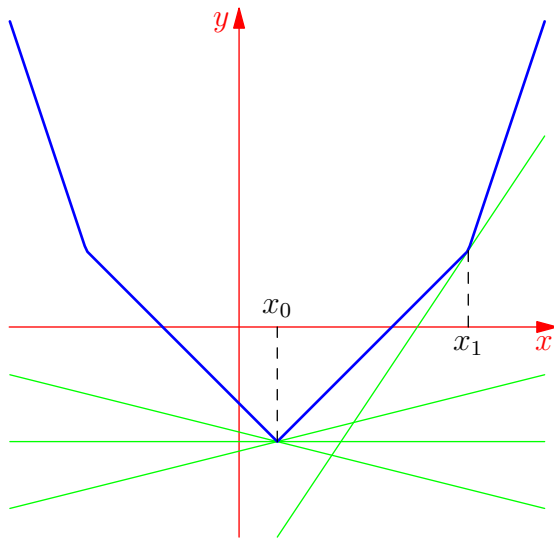
Definition (Subgradient, Subdifferential, Subdifferentiable)

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$. We say that $g \in \mathbb{R}^d$ is a *subgradient* of f at $x \in \mathbb{R}^d$ if

$$f(x + v) \geq f(x) + g^T v$$

for all $v \in \mathbb{R}^d$. The *subdifferential* $\partial f(x)$ is the set of all subgradients of f at x . We say that f is *subdifferentiable* at x if $\partial f(x) \neq \emptyset$ (i.e., if there is at least one subgradient).

Subgradients at x_0 and x_1



Facts About Subgradients

- 1 If f is convex and differentiable at x then $\partial f(x) = \{\nabla f(x)\}$.

Facts About Subgradients

- 1 If f is convex and differentiable at x then $\partial f(x) = \{\nabla f(x)\}$.
- 2 If f is convex then $\partial f(x) \neq \emptyset$ for all x .

Facts About Subgradients

- 1 If f is convex and differentiable at x then $\partial f(x) = \{\nabla f(x)\}$.
- 2 If f is convex then $\partial f(x) \neq \emptyset$ for all x .
- 3 The subdifferential $\partial f(x)$ is a convex set. Thus the subdifferential can contain 0, 1, or infinitely many elements.

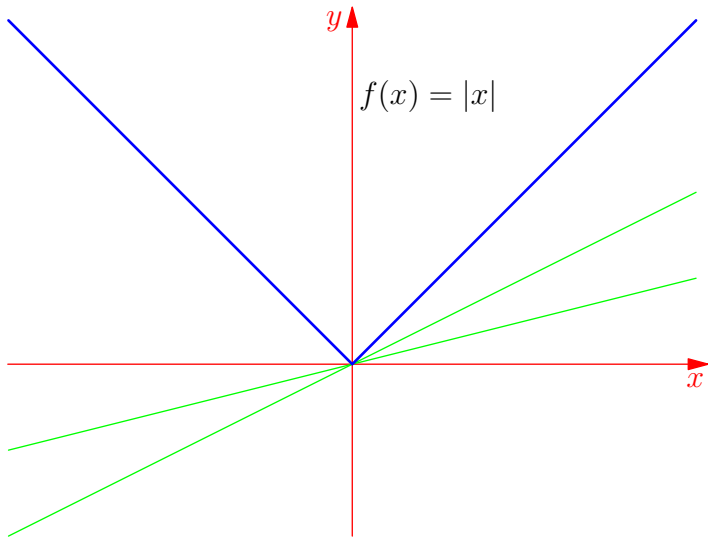
Facts About Subgradients

- 1 If f is convex and differentiable at x then $\partial f(x) = \{\nabla f(x)\}$.
- 2 If f is convex then $\partial f(x) \neq \emptyset$ for all x .
- 3 The subdifferential $\partial f(x)$ is a convex set. Thus the subdifferential can contain 0, 1, or infinitely many elements.
- 4 If the zero vector is a subgradient of f at x , then x is a global minimum.

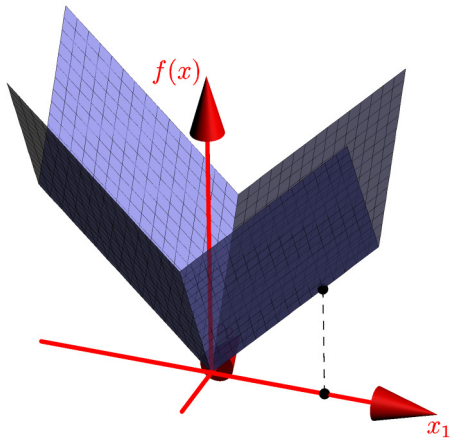
Facts About Subgradients

- ① If f is convex and differentiable at x then $\partial f(x) = \{\nabla f(x)\}$.
- ② If f is convex then $\partial f(x) \neq \emptyset$ for all x .
- ③ The subdifferential $\partial f(x)$ is a convex set. Thus the subdifferential can contain 0, 1, or infinitely many elements.
- ④ If the zero vector is a subgradient of f at x , then x is a global minimum.
- ⑤ If g is a subgradient of f at x , then $(g, -1)$ is orthogonal to the underestimating hyperplane $\{(x + v, f(x) + g^T v) \mid v \in \mathbb{R}^d\}$ at $(x, f(x))$.

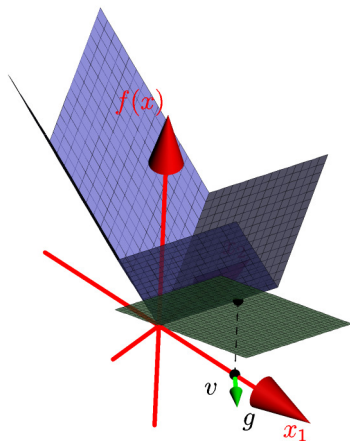
Compute the Subdifferentials of $f(x) = |x|$



Compute $\partial f(3, 0)$ For $f(x_1, x_2) = |x_1| + 2|x_2|$

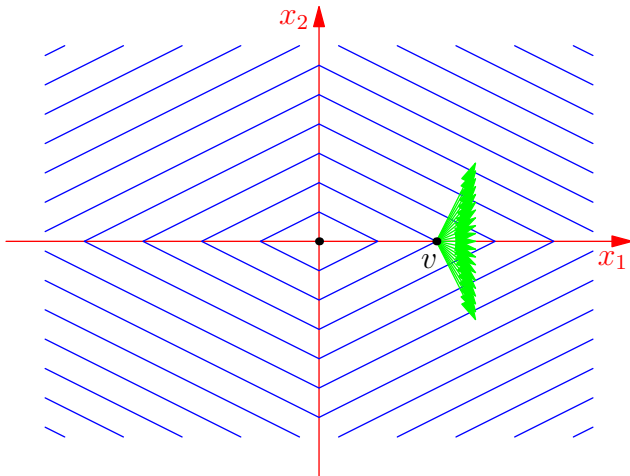


Compute $\partial f(3, 0)$ For $f(x_1, x_2) = |x_1| + 2|x_2|$



Compute $\partial f(3, 0)$ For $f(x_1, x_2) = |x_1| + 2|x_2|$

$$\partial f(3, 0) = \{(1, b)^T \mid b \in [-2, 2]\}$$

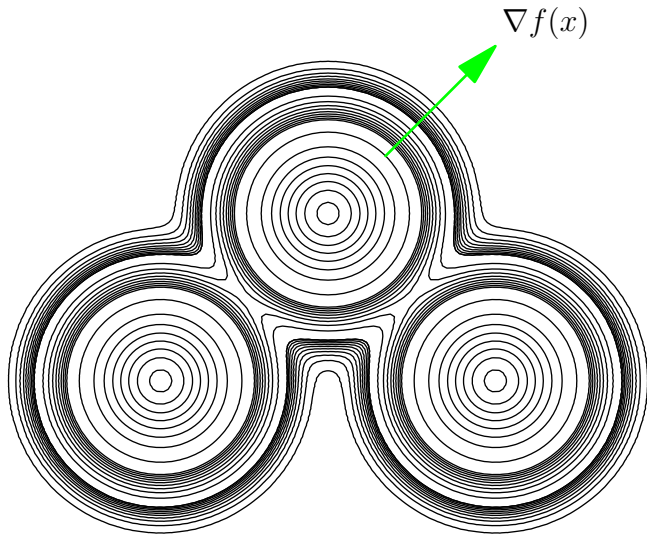


Gradient Lies Normal To Contours

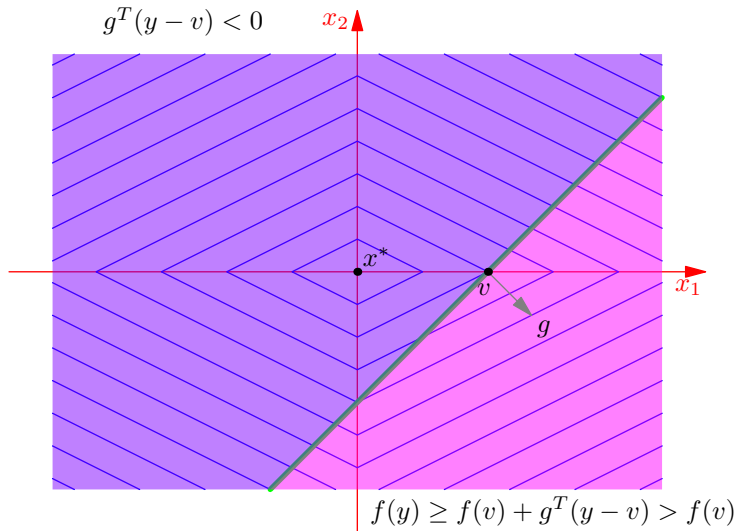
Theorem

If $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is continuously differentiable and $x_0 \in \mathbb{R}^d$ with $\nabla f(x_0) \neq 0$ then $\nabla f(x_0)$ is normal to the level set $S = \{x \in \mathbb{R}^d \mid f(x) = f(x_0)\}$.

Gradient Lies Normal To Contours



Normal Plane to Subgradient Splits Space



Subgradient Descent

- 1 Let $x^{(0)}$ denote the initial point.
- 2 For $k = 1, 2, \dots$
 - 1 Assign $x^{(k)} = x^{(k-1)} - \alpha_k g$, where $g \in \partial f(x^{(k-1)})$ and α_k is the step size.
 - 2 Set $f_{\text{best}}^{(k)} = \min_{i=1, \dots, k} f(x^{(i)})$. (Used since this isn't a descent method.)

Convergence of Subgradient Descent

Theorem

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be convex and Lipschitz with constant G , and let x^* be a minimizer. For a fixed step size t , the subgradient method satisfies:

$$\lim_{k \rightarrow \infty} f(x_{\text{best}}^{(k)}) \leq f(x^*) + G^2 t / 2.$$

For step sizes respecting the Robbins-Monro conditions,

$$\lim_{k \rightarrow \infty} f(x_{\text{best}}^{(k)}) = f(x^*).$$

Robbins-Monro conditions

- For convergence guarantee, use decreasing step sizes (dampens noise in step direction).
- Let η_t be the step size at the t 'th step.

Robbins-Monro Conditions

Many classical convergence results depend on the following two conditions:

$$\sum_{t=1}^{\infty} \eta_t^2 < \infty \quad \sum_{t=1}^{\infty} \eta_t = \infty$$

- As fast as $\eta_t = O\left(\frac{1}{t}\right)$ would satisfy this... but should be faster than $O\left(\frac{1}{\sqrt{t}}\right)$.
- A useful reference for practical techniques: Leon Bottou's "Tricks":
<http://research.microsoft.com/pubs/192769/tricks-2012.pdf>