

Optimization for Machine Learning

CSCI-599

Lecture 3: Constrained and Non-Smooth Gradient Descent

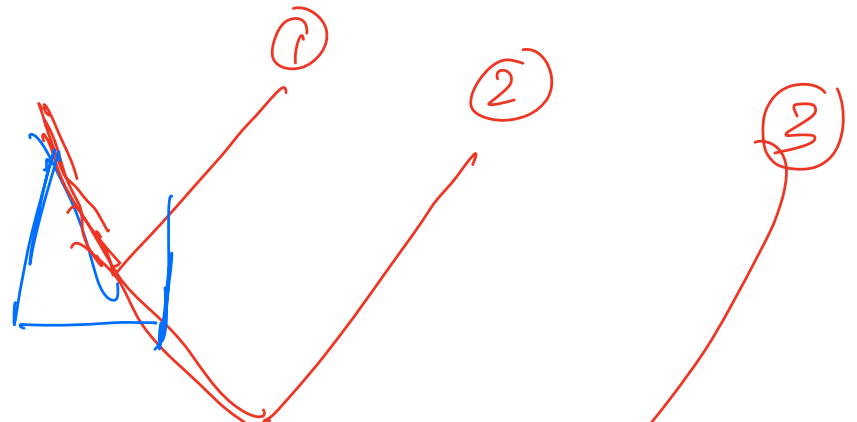
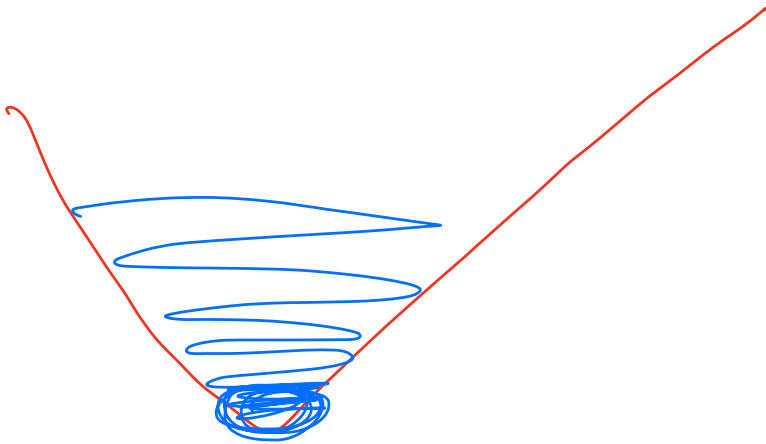
Sai Praneeth Karimireddy

USC – <https://spkreddy.org/optmlspring2025.html>

February 3, 2025

Convergence in $\mathcal{O}(1/\varepsilon)$ steps

Same as vanilla case for smooth functions, but now for any h for which we can compute the proximal mapping.



Recap

$$x_{t+1} = x_t - \gamma \nabla f(x_t) \quad \text{GD}$$

①

$$f \text{ is } B\text{-Lip} \Leftrightarrow |f(x) - f(y)| \leq B \|x - y\|_2 \quad \forall x, y$$
$$\|\nabla f(x)\|_2 \leq B \quad \forall x$$

$$\gamma = \frac{R}{B\sqrt{T}} \Rightarrow f(\bar{x}_T) - f^* \leq \frac{BR}{\sqrt{T}}$$
$$\bar{x}_T = \frac{1}{T} \sum_{t=1}^T x_t$$

②

$$f \text{ is } L\text{-smooth} \Rightarrow \|\nabla f(x_t)\|^2 \leq 2L (f(x_t) - f(x_{t+1}))$$

$$\gamma = \frac{1}{L} \Rightarrow f(x_T) - f^* \leq \frac{LR}{T}$$

③ f is L -smooth + μ -strong convex \Rightarrow

$$\gamma = \frac{1}{L} \Rightarrow \|x_T - x^*\|_2^2 \leq \left(1 - \frac{\mu}{L}\right)^T \|x_0 - x^*\|_2^2$$

$$f(x) = \underbrace{\frac{1}{2} x^T A x - b^T x}_{\text{CVX}} + \underbrace{\lambda \|x\|_1}_{\text{CVX}}$$

ReLU

Lasso Regression

$$f(x) = \max(\underbrace{Ax - b}_x, 0)$$

Subgradients

if f is diff \Rightarrow

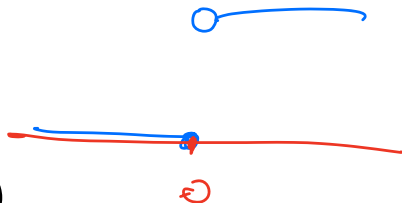
What if f is not differentiable?

$\partial f(x)$ is a subgrad

Definition

$g \in \mathbb{R}^d$ is a **subgradient** of f at x if

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

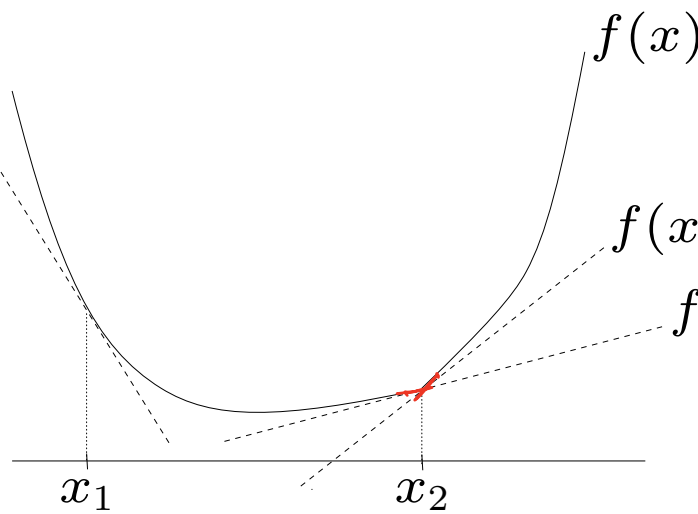


$f(x) = |x|$

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

$f(x_1) + g_1^T(x - x_1)$
 $= \partial f(x)$

$f(x_2) + g_2^T(x - x_2)$
 $f(x_2) + g_3^T(x - x_2)$

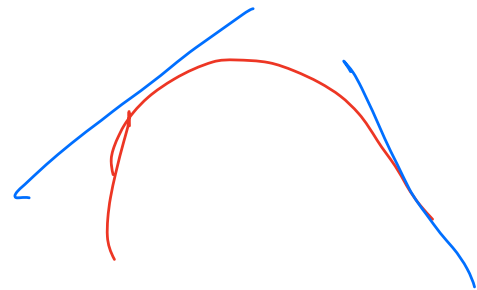
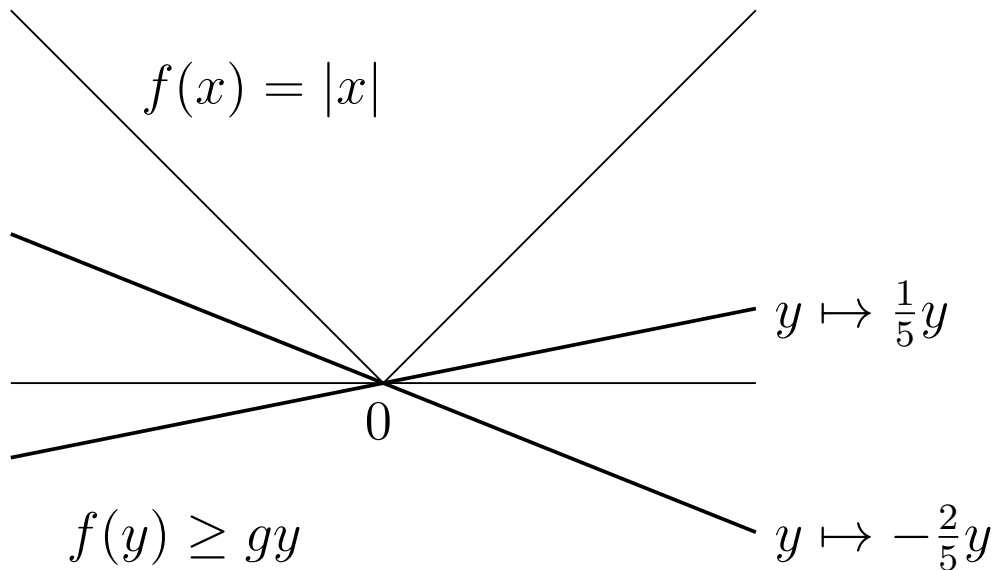


$\partial f(x) \subseteq \mathbb{R}^d$ is the **subdifferential**, the set of subgradients of f at x .

Subgradients II

Example:

$$f(x) = -x^2$$

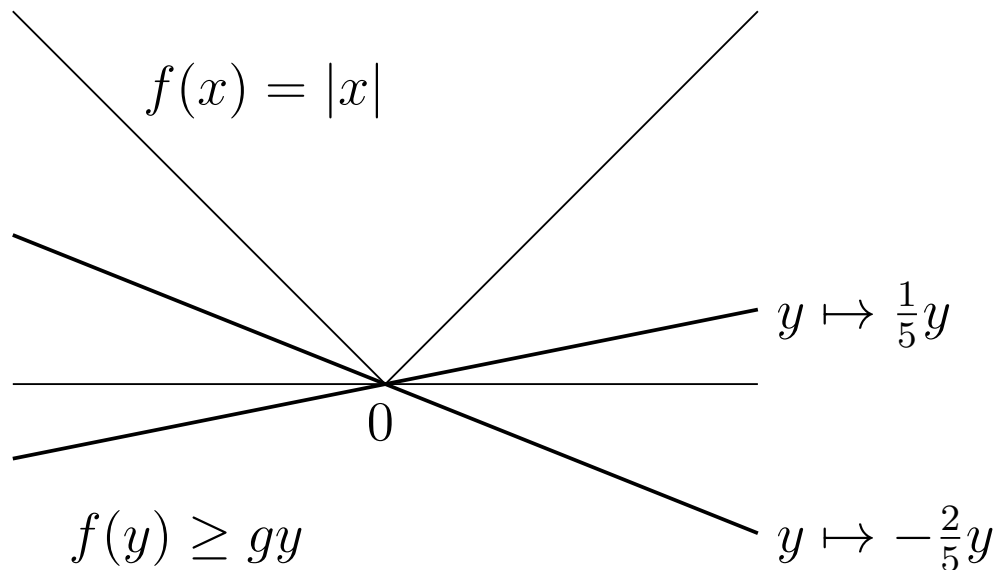


Subgradient condition at $x = 0$:

f is convex \Rightarrow subgrad exists

Subgradients II

Example:



Subgradient condition at $x = 0$: $f(y) \geq f(0) + g(y - 0) = gy$.

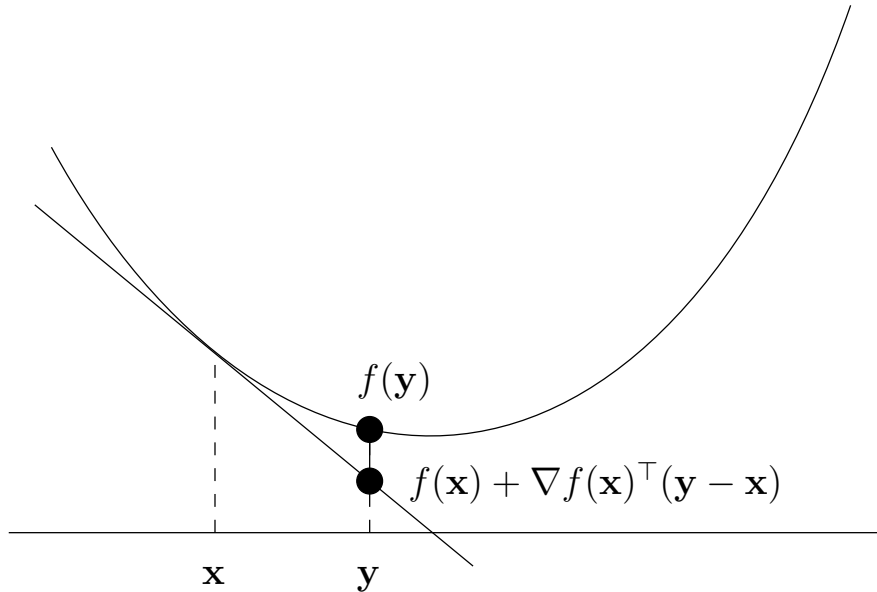
$$\partial f(0) = [-1, 1]$$

Subgradients III

Lemma (Exercise 28)

If $f : \text{dom}(f) \rightarrow \mathbb{R}$ is differentiable at $\mathbf{x} \in \text{dom}(f)$, then $\partial f(\mathbf{x}) \subseteq \{\nabla f(\mathbf{x})\}$.

Either exactly one subgradient $\nabla f(\mathbf{x}) \dots$



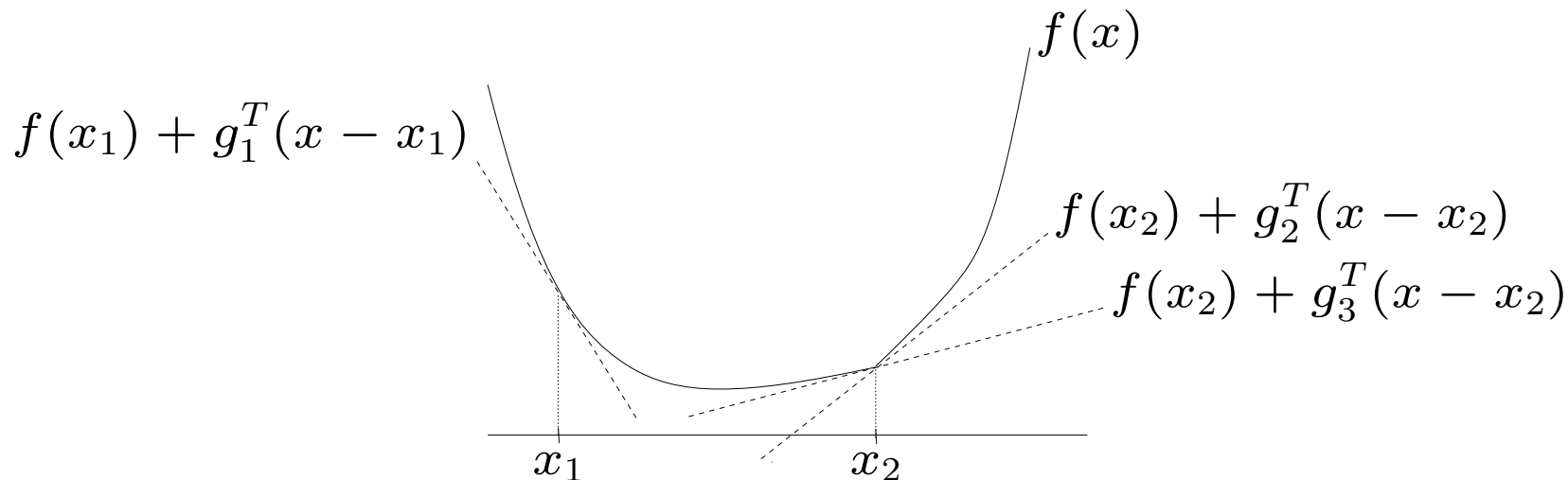
...or no subgradient at all.

Subgradient characterization of convexity

“convex = subgradients everywhere”

Lemma (Exercise 29)

A function $f : \text{dom}(f) \rightarrow \mathbb{R}$ is convex if and only if $\text{dom}(f)$ is convex and $\partial f(\mathbf{x}) \neq \emptyset$ for all $\mathbf{x} \in \text{dom}(f)$.



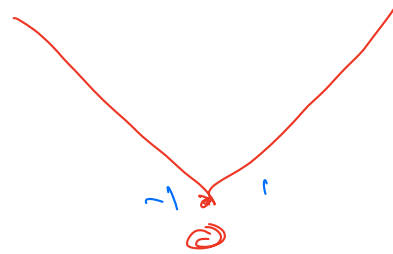
→ let $g(x_+)$ be "a" subgradient of f at x_+

— $x_{++} = x_+ - \gamma g(x_+)$ subgradient descent

x^* is opt of $\min_x f(x)$

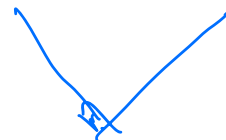
if $\nabla f(x^*) = 0$

$f(x) = |x|$



$0 \in g(x^*) \Rightarrow x^*$ is optimum

Convex and Lipschitz = bounded subgradients



Lemma (Exercise 30)

Let $f : \mathbf{dom}(f) \rightarrow \mathbb{R}$ be convex, $\mathbf{dom}(f)$ open, $B \in \mathbb{R}_+$. Then the following two statements are equivalent.

- (i) $\|\mathbf{g}\| \leq B$ for all $\mathbf{x} \in \mathbf{dom}(f)$ and all $\mathbf{g} \in \partial f(\mathbf{x})$.
- (ii) $|f(\mathbf{x}) - f(\mathbf{y})| \leq B\|\mathbf{x} - \mathbf{y}\|$ for all $\mathbf{x}, \mathbf{y} \in \mathbf{dom}(f)$.

Subgradient optimality condition

Lemma

Suppose that $f : \text{dom}(f) \rightarrow \mathbb{R}$ and $\mathbf{x} \in \text{dom}(f)$. If $\mathbf{0} \in \partial f(\mathbf{x})$, then \mathbf{x} is a global minimum.

Proof.

By definition of subgradients, $\mathbf{g} = \mathbf{0} \in \partial f(\mathbf{x})$ gives

Subgradient optimality condition

Lemma

Suppose that $f : \text{dom}(f) \rightarrow \mathbb{R}$ and $\mathbf{x} \in \text{dom}(f)$. If $\mathbf{0} \in \partial f(\mathbf{x})$, then \mathbf{x} is a global minimum.

Proof.

By definition of subgradients, $\mathbf{g} = \mathbf{0} \in \partial f(\mathbf{x})$ gives

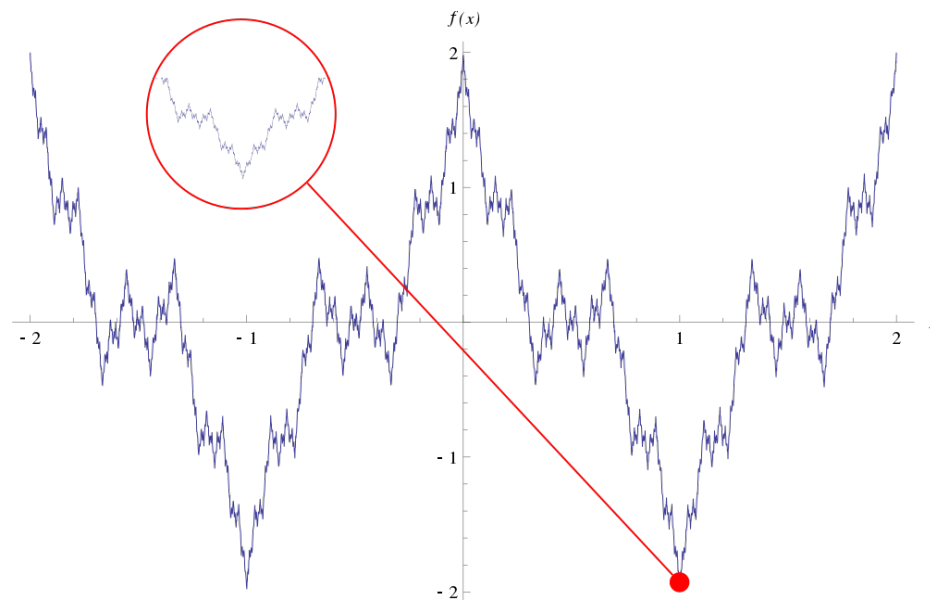
$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{y} - \mathbf{x}) = f(\mathbf{x})$$

for all $\mathbf{y} \in \text{dom}(f)$, so \mathbf{x} is a global minimum. □

Differentiability of convex functions

How “wild” can a non-differentiable convex function be?

Weierstrass function: a function that is continuous **everywhere** but differentiable **nowhere**



<https://commons.wikimedia.org/wiki/File:WeierstrassFunction.svg>

Differentiability of convex functions



Theorem ([?, Theorem 25.5])

A *convex* function $f : \mathbf{dom}(f) \rightarrow \mathbb{R}$ is differentiable *almost everywhere*.

In other words:

- ▶ Set of points where f is non-differentiable has measure 0 (no volume).
- ▶ For all $\mathbf{x} \in \mathbf{dom}(f)$ and all $\varepsilon > 0$, there is a point \mathbf{x}' such that $\|\mathbf{x} - \mathbf{x}'\| < \varepsilon$ and f is differentiable at \mathbf{x}' .

The subgradient descent algorithm

Subgradient descent: choose $\mathbf{x}_0 \in \mathbb{R}^d$.

Let $\mathbf{g}_t \in \partial f(\mathbf{x}_t)$

$\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma_t \mathbf{g}_t$

for **times** $t = 0, 1, \dots$, and **stepsizes** $\gamma_t \geq 0$.

Stepsize can vary with time!

This is possible in (projected) gradient descent as well, but so far, we didn't need it.

Lipschitz convex functions: $\mathcal{O}(1/\varepsilon^2)$ steps

Theorem

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and B -Lipschitz continuous with a global minimum \mathbf{x}^* ; furthermore, suppose that $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq R$. Choosing the constant stepsize

$$\gamma := \frac{R}{B\sqrt{T}},$$

subgradient descent yields

$$\frac{1}{T} \sum_{t=0}^{T-1} f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{RB}{\sqrt{T}}.$$

Proof is identical to the one of Theorem 2.1, except...

- ▶ In vanilla analysis, now use $\mathbf{g}_t \in \partial f(\mathbf{x}_t)$ instead of $\mathbf{g}_t = \nabla f(\mathbf{x}_t)$.
- ▶ Inequality $f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*)$ now follows from subgradient property instead of first-order characterization of convexity.

$$\begin{aligned} \sum g_t^T (x_t - x^*) &\leq \frac{1}{2\gamma} \left(\|x_0 - x^*\|^2 - \|x_T - x^*\|^2 \right) \\ &\leq \frac{1}{2\gamma} \left(\|x_0 - x^*\|^2 - \|x_T - x^*\|^2 \right) \\ &\quad + \frac{\gamma}{2} \sum \|g_t\|^2 \leq B \end{aligned}$$

Optimality of first-order methods

With all the convergence rates we have seen so far, a very natural question to ask is if these rates are best possible or not. Surprisingly, the rate can indeed not be improved in general.

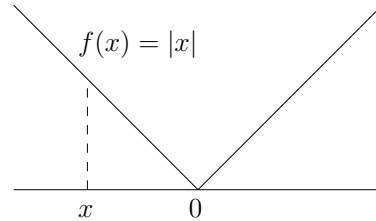
Theorem (Nesterov)

For any $T \leq d - 1$ and starting point \mathbf{x}_0 , there is a function f in the problem class of B -Lipschitz functions over \mathbb{R}^d , such that any (sub)gradient method has an objective error at least

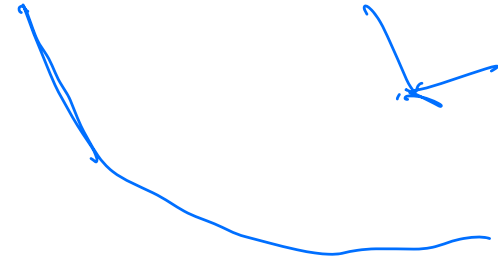
$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \geq \frac{RB}{2(1 + \sqrt{T + 1})} .$$

Smooth (non-differentiable) functions?

They don't exist (Exercise 31)!



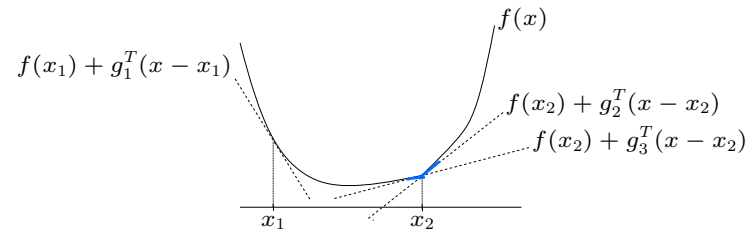
smoothness
differentiability



At 0, graph can't be below a tangent paraboloid.

Can we still improve over $O(1/\varepsilon^2)$ steps for Lipschitz functions?

Yes, if we also require strong convexity (graph is above not too flat tangent paraboloids).



Optimality

x^* is opt of

$$\min_x f(x)$$

f is convex

if $0 \in \text{subgrad}(x^*)$

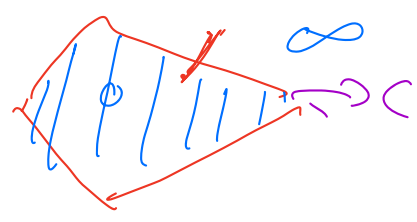
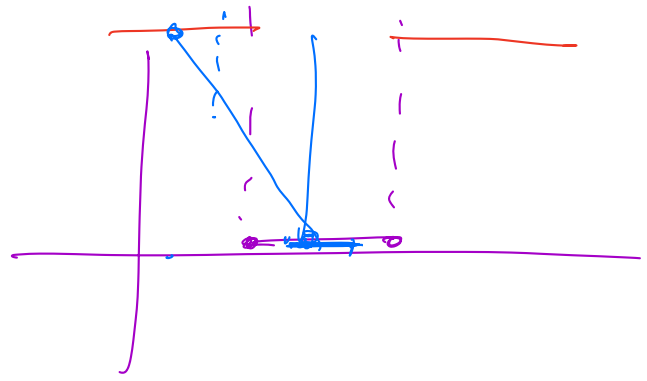
$$\min_{x \in C} f(x)$$

constrained opt

$$|f(x) - f(y)| \leq B \|x - y\|_2$$

$$\tilde{f}(x) = f(x) + \mathbb{1}(x \in C)$$

$$\min_x \tilde{f}(x) \Leftrightarrow \min_{x \in C} f(x)$$



we need Lip assumption
+ smooth

f is convex is not sufficient for convergence