

Lecture 17: Clarke & his differential.

Announcement.

- * hw2 tomorrow
- * notes

Why Clarke differential:

- * Need a model of a gradient
- * Doesn't match python; hard to predict what is used in 10 years.

Subgradients

We saw it in optimization:

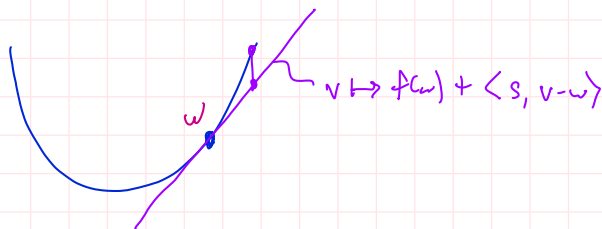
$$\partial_s f(w) = \text{"set of lower bounding slopes at } w\text{"}$$

$$= \{s \in \mathbb{R}^d : \forall v. f(v) \geq f(w) + \langle s, v-w \rangle\}.$$

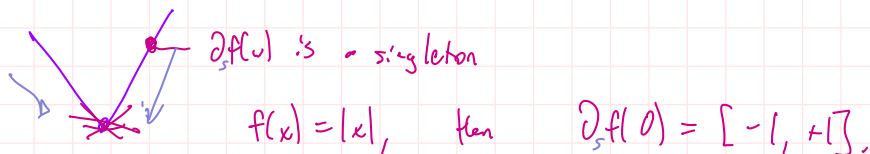


Examples

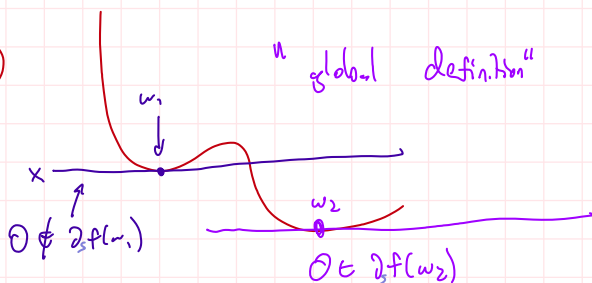
①



②



③



Properties (for convex $f: \mathbb{R}^d \rightarrow \mathbb{R}$).

- * $\partial_s f$ is nonempty everywhere, and is closed convex.
- * Related to directional derivatives

$$f'(w; v) = \lim_{h \downarrow 0} \frac{f(w+hu) - f(w)}{h} = \sup \{ \langle v, s \rangle : s \in \partial_s f(w) \}$$

- * Mean value theorems; see for instance

Hiriart-Urruty & Lemaréchal "fundamentals of convex analysis" (many pictures)

Borwein & Lewis "Convex analysis & non-linear optimization"

- * We used it in our first opt proof.

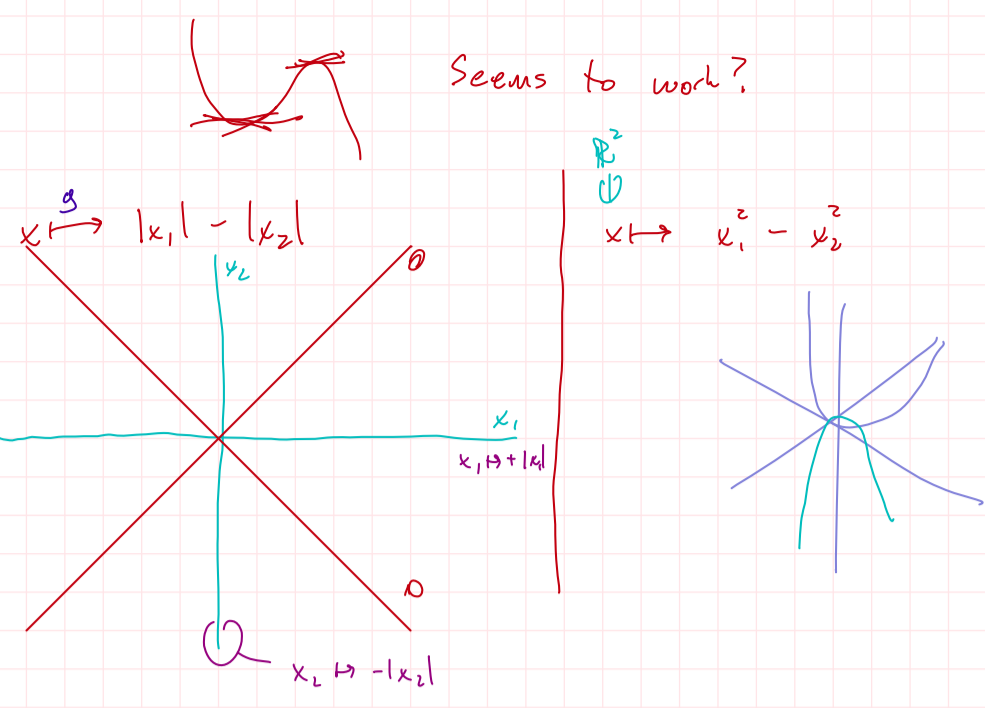
- * Proposition (Jensen's inequality). If $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex & X is a rv. then $\mathbb{E} f(X) \geq f(\mathbb{E} X)$.

Proof: $\forall s \in \partial f(\mathbb{E} X), \mathbb{E} f(X) \geq \mathbb{E} [f(\mathbb{E} X) + \langle s, X - \mathbb{E} X \rangle]$

$$= f(\mathbb{E} X) + \langle s, \mathbb{E} X - \mathbb{E} X \rangle. //$$

Clarke differential.

First try: "local" sub- supergradient



Want $0 \in \partial g(0)$, fails for "local sub- supergradient"

First definition of Clarke starts from directional derivative

Recall directional derivative: $f'(w;v) = \lim_{h \rightarrow 0} \frac{f(w+hv) - f(w)}{h}$

Then $\partial_0 f(w) = \{s \in \mathbb{R}^d : \forall v, f'(w;v) \geq \langle v, s \rangle\}$. ✓

We'll use a related form.

Defn: $\partial f(w) = \text{conv} \left(\left\{ \lim_{i \rightarrow \infty} \nabla f(w_i) : \begin{array}{l} w_i \rightarrow w \\ \nabla f(w_i) \text{ exists} \\ \lim \nabla f(w_i) \text{ exists} \end{array} \right\} \right)$

Example: also natural from directional derivative view, with then \rightarrow boundary pos-ns of subgradients.

Example 2: $g(x) = |x_1| - |x_2|$
 claim: $\partial g(0) = \text{conv}(\{(\pm 1, \pm 1)\})$

$\nabla g(x) = \begin{cases} (s_{g_1}(x_1), s_{g_2}(x_2)) & x_1 \neq 0 \text{ and } x_2 \neq 0 \\ \text{not differentiable} & x_1 = 0 \text{ or } x_2 = 0 \end{cases}$

Convex hull
 $\text{conv}(S) = \left\{ \sum_{i=1}^N \alpha_i u_i : N \geq 1, u_1, \dots, u_N \in S, \alpha_i \in \Delta, \sum \alpha_i = 1, \alpha_i \geq 0 \right\}$

Remark: Considering any sequence $w_i \rightarrow 0$, Clarke differential at 0 seems unstable.

What we want:

- 1 nonempty everywhere
- 2 chain rules & gradient flow

1 Nonempty everywhere

Defn: $f: U \rightarrow \mathbb{R}$ (over an open domain U) is locally Lipschitz if $\forall x \in U, \exists$ open $S \ni x$ s.t. f is Lipschitz over S ($\exists M, \forall y, z \in S, |f(y) - f(z)| \leq M \|y - z\|$).

picture: $\forall x$, can "zoom in" so that f appears Lipschitz.

- Ex: $|x|$ is Lip \Rightarrow locally Lip
- x^2 is locally Lipschitz but not Lipschitz
- All convex $f: \mathbb{R}^d \rightarrow \mathbb{R}$ are locally Lipschitz
- $1/x$ over $(0, \infty)$ is locally Lip
- $x \sin(1/x)$ over $(0, \infty)$ not locally Lip but uniformly cont.
- all standard networks are locally Lipschitz.

Theorem (Rademacher) $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is locally Lip \Rightarrow differentiable almost everywhere

Properties: $f: \mathbb{R}^d \rightarrow \mathbb{R}$

- 1 f locally Lipschitz $\Rightarrow \partial f$ nonempty everywhere
- 2 f convex $\Rightarrow \partial f = \partial_s f$
- 3 f differentiable $\Rightarrow \partial f(w) = \{\nabla f(w)\}$.

Remark: Not what pytorch computes.

$x \mapsto \sigma(x) - \sigma(-x) = x$ (maybe pytorch correct?)
 $x \mapsto \sigma(\sigma(x) - \sigma(-x) + 1) - 1$ (maybe not?)

