

Lec 15: nonsmooth, pos-hom, margins

Announcements:

- * Schedule until Thanksgiving mostly decided.
- * Thursday in Siebel 1105.
- * HW1 graded; check Zinei comments

Plan until Thanksgiving

- * Optimization: 4-5 more lectures; non-smooth, positive homogeneity, margins.
- * Generalization: cover all main topics before Thanksgiving; thus fair game for homework 4.

After Thanksgiving:

- * HW4 & project due.
- * 3 "extra" lectures.

Next 1.5 lectures:

goal is nonsmooth analog GF .

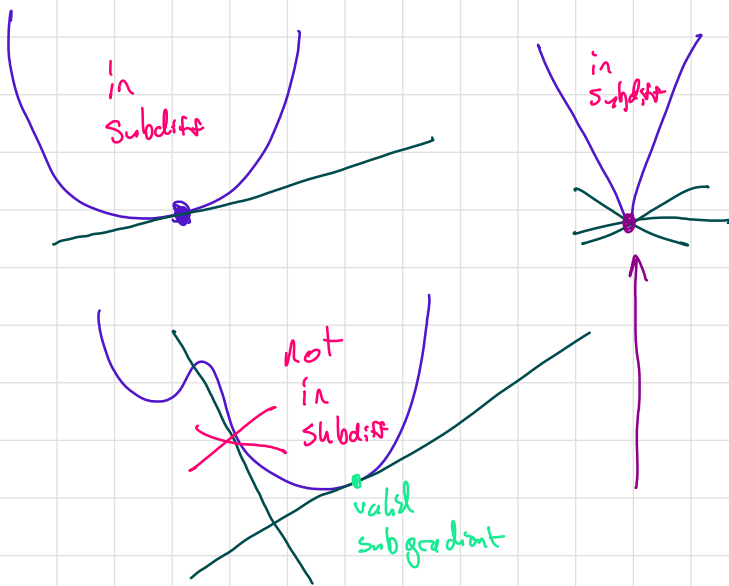
Progression: subgradients \rightarrow Clarke differential
 \rightarrow corresponding flow
(positive homogeneity, along way)

Subgradients

Define subdifferential / subgradient

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

("set of lower bounding tangents").



Properties: (for convex F .)

* $\partial_c F$ nonempty everywhere.
 (Convex analysis texts allow $F: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\pm\infty\}$,
 in which case in general $\partial_c F$ nonempty for $\text{ri}(\text{dom}(F))$.)

* $\partial_c F$ closed convex set.

* Related to directional derivatives

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

* Most other things you expect from gradients, hold with subgradients of convex functions.

Proposition (Jensen's inequality). Let convex $F: \mathbb{R}^d \rightarrow \mathbb{R}$
 & r.v. $X \in \mathbb{R}^d$ be given.

Then $F(\mathbb{E}X) \leq \mathbb{E}F(X)$.

Proof. $\exists s \in \partial_c F(\mathbb{E}X)$ (nonempty because convex $F: \mathbb{R}^d \rightarrow \mathbb{R}$).

$$\mathbb{E}[F(X)] \geq \mathbb{E}\left[F(\mathbb{E}X) + \langle s, X - \mathbb{E}X \rangle \right] \\
= F(\mathbb{E}X) + \langle s, \mathbb{E}X - \mathbb{E}X \rangle$$

defn subgrad applied to integrand.

Pretty books on convexity:

① Fundamentals of convex analysis } pretty; but long
 Hiriart-Urruty & Lemaréchal

② Convex analysis & ("?") } efficient
 Borwein & Lewis.

Clarke differential.

Defn. $F: U \rightarrow \mathbb{R}$ is locally Lipschitz if every open neighborhood is Lipschitz
 $(\forall x \in U, \exists \text{ open } S \ni x \exists M \text{ s.t. } \forall y \in S, |F(x) - F(y)| \leq M \|x - y\|)$

E.g. $x \mapsto |x|$ Lipschitz; γ_x is locally Lipschitz over $(0, \infty)$, γ_y is Lipschitz over $[1, \infty)$.

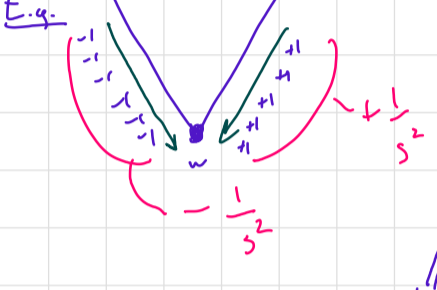
$x \mapsto x \sin(1/x)$ over $[0, \infty)$: continuous but not locally Lipschitz.

Example: All standard multilayer networks (except threshold) are locally Lipschitz.

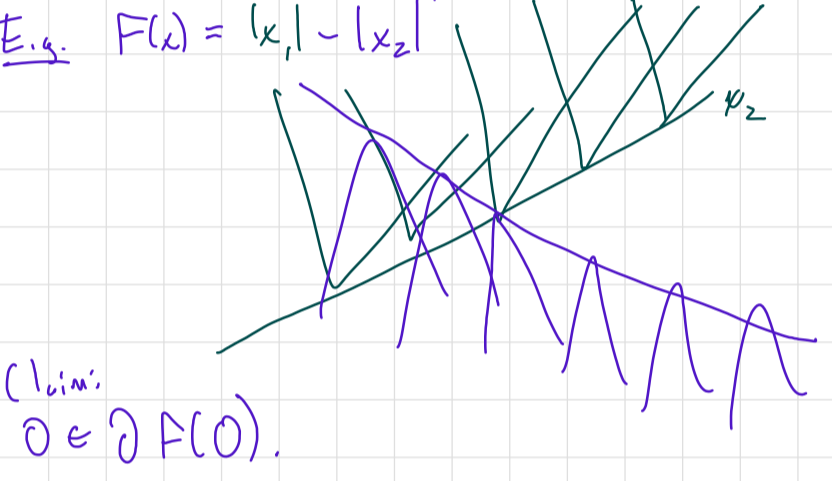
Theorem (Rademacher). $F: U \rightarrow \mathbb{R}$ locally Lipschitz \Rightarrow differentiable almost everywhere.

Definition. Clarke differential:

$$\partial F(w) = \text{conv} \left(\left\{ \lim_{s \rightarrow \infty} \nabla F(w_s) : \begin{array}{l} w_s \rightarrow w, \\ \nabla F(w_s) \text{ exists,} \\ \lim_{s \rightarrow \infty} \nabla F(w_s) \text{ exists.} \end{array} \right\} \right)$$



another see with gradients: $-1, -1, -1, -1, +1, -1, +1, -1, 4, \dots$



Claim: $0 \in \partial F(0)$.

Main properties:

- * F locally Lipschitz $\Rightarrow \partial F$ exists everywhere.
- * $F: \mathbb{R}^d \rightarrow \mathbb{R}$ convex $\Rightarrow \partial F = \partial_c F$ everywhere.
- * $F: \mathbb{R}^d \rightarrow \mathbb{R}$ differentiable everywhere $\Rightarrow \partial F(w) = \{ \nabla F(w) \}$ everywhere.

Replace GF: ~~$\dot{w}(t) = \frac{dw}{dt}(t) = -\nabla F(w(t)) \forall t \geq 0$~~

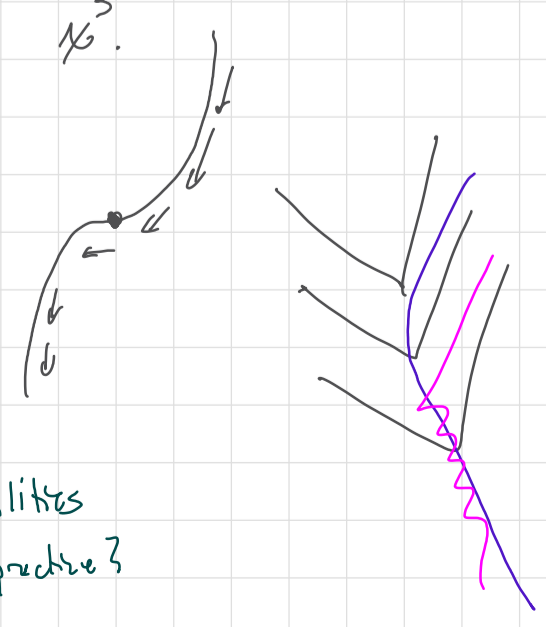
Differential inclusion: $\dot{w}(t) = \frac{dw}{dt}(t) \in -\partial F(w(t)) \forall \text{ a.e. } t \geq 0$

Good news:

- * We got a "chain rule" e.g. $\forall \text{ a.e. } t \geq 0 \quad \frac{d}{dt} F(w(t)) = \langle -v, \dot{w}(t) \rangle$ where $v \in \partial F(w(t))$
- under assumptions (e.g., holds for all standard types of networks)
- * Just like regular GF: function values decrease, pass near opt stationary points.

Bad news:

- * Does not correspond pytorch etc.
- not just crazy stuff like $\sigma(\alpha(x)) - \sigma(-x)$, but even x^3 .



* Non differentiability don't matter practice?