

# Lecture 18: gradient flow & differential inclusion

GF

Announcement:

\* hw2 ?? 50%

Gradient Flow. (differentiable)

$$GD: w' := w - \eta \nabla \hat{R}(w)$$

$$\frac{w' - w}{\eta} = -\nabla \hat{R}(w)$$

$$\xrightarrow{\eta \downarrow 0} \frac{d}{dt} w = \dot{w}_t = -\nabla \hat{R}(w_t)$$

"Gradient flow" ODE; for GF,

we assume solutions exist & are unique,

meaning  $(w_t)_{t \geq 0}$ ,  $t \in \mathbb{R}$  is the "solution".

Why GF:

\* "free" smoothness equality.  
(which implies  $\hat{R}(w_t)$  is non-increasing).

\* nice first thing to check

key byproduct

\* GD on  $\beta$ -smooth  $\hat{R}$

$$\hat{R}(w_{i+1}) \leq \hat{R}(w_i) - \eta \left(1 - \frac{\eta\beta}{2}\right) \|\nabla \hat{R}(w_i)\|^2$$

$$\Rightarrow \frac{\hat{R}(w_{i+1}) - \hat{R}(w_i)}{\eta} \leq -\left(1 - \frac{\eta\beta}{2}\right) \|\nabla \hat{R}(w_i)\|^2$$

\* GF gives us this "for free"

$$\begin{aligned} \frac{d}{dt} \hat{R}(w_t) &= \langle \nabla \hat{R}(w_t), \frac{d}{dt} w_t \rangle \\ &= -\|\nabla \hat{R}(w_t)\|^2 \end{aligned}$$

needed  
continuity/  
assumptions  
for existence  
& uniqueness

Why GF not nice:

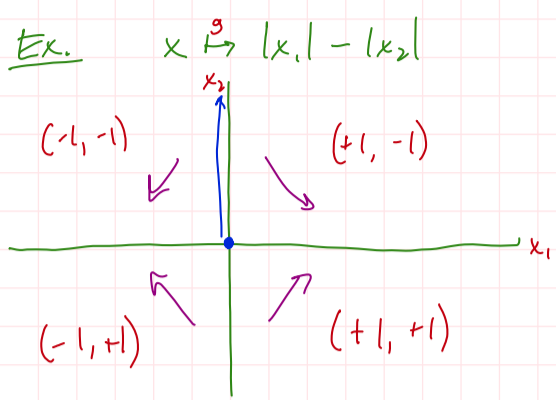
\* Not an algorithm  
(no real "running time").

\* can be hard to convert "rate"  
to GD (I have scary examples).

\* "circuitous"

# Clarke differential & GF

$\hat{R}$  diff & ( $\nabla \hat{R}$  uniformly continuous)?  $\Rightarrow$  differential equation  $\dot{w}_t = -\nabla \hat{R}(w_t)$   $\hookrightarrow$  unique solutions  
 $\hat{R}$  locally Lipschitz  $\Rightarrow$  differential inclusion  $\dot{w}_t \in -\partial \hat{R}(w_t)$  may fail to have unique solutions



Consider  $\dot{z}_t \in -\partial g(z_t)$   
 Claim  $\forall r \geq 0$   
 $t \in [0, r] \begin{cases} \dot{w}_t = (0, 0) \\ \dot{w}_t = (0, 0) \end{cases}$   
 $t \in (r, \infty) \begin{cases} \dot{w}_t = (0, t-r) \\ \dot{w}_t = (0, t) \end{cases}$

for any  $a, b$ ,  $w_b - w_a$

$$\begin{aligned}
 \textcircled{1} \quad r \in [a, b] &\Rightarrow w_b - w_a = \int_a^b \dot{w}_t dt \\
 &= \int_a^r \dot{w}_s ds + \int_r^b \dot{w}_s ds \\
 &= 0 + \int_r^b (0, +1) \\
 &= 0 + (b-r)(0, +1).
 \end{aligned}$$

- $\textcircled{2} \quad a \geq r$
  - $\textcircled{3} \quad b \leq r$
- } similar

- Remarks
- $\textcircled{1}$  (essentially univariate: consider  $z \mapsto -|z|$ .)
  - $\textcircled{2}$  Instability also exists with GD.
  - $\textcircled{3}$  Ruiteng Lyu & friends "simplicity obs: ..." last appendix

"Fix": "Clarke differential inclusion is well-behaved for almost all time".

Lemma (Bolte, Lewis, Ji, others). Consider  $\dot{z}_t \in -\partial f(z_t)$ .

Suppose  $f$  is locally Lipschitz and  $\phi$ -minimal definable. Then "always" satisfied for finite-width networks

- for almost every  $t \geq 0$
- $\textcircled{1}$  chain rule:  $\frac{d}{dt} f(z_t) = \langle s, z_t \rangle \quad \forall s \in \partial f(z_t)$
  - $\textcircled{2}$  minimum norm selection:  $\dot{z}_t = -\bar{\partial} f(z_t) = -\arg \min \{ \|s\| \mid s \in \partial f(z_t) \}$
  - $\textcircled{3}$   $\frac{d}{dt} f(z_t) = -\|\bar{\partial} f(z_t)\|^2$ .  
captures main benefit of GF.

Consider again  $g(x) = |x_1| - |x_2|$ . Then letting  $x_t$  be a solution to Clarke flow,

$$\begin{aligned}
 g(x_t) - g(x_0) &= \int_0^t \frac{d}{ds} g(x_s) ds \stackrel{\text{(using lemma)}}{=} - \int_0^t \|\bar{\partial} g(x_s)\|_2^2 ds \\
 &\stackrel{\text{assume } t \geq r}{=} - \int_{[0, r]} \|\bar{\partial} g(x_s)\|_2^2 ds - \int_{(r, \infty)} \|\bar{\partial} g(x_s)\|_2^2 ds \\
 &= -(t-r).
 \end{aligned}$$

## Descent proofs.

Proposition. Suppose  $\hat{R}$  is convex and  $w_t \in -\partial \hat{R}(w_t)$ ,  
&  $z \in \mathbb{R}^d$  arbitrary. Then  $\forall t$

$$\frac{1}{2} \|w_t - z\|^2 + t \hat{R}(w_t) \leq \frac{1}{2} \|w_0 - z\|^2 + t \hat{R}(z).$$

Proof.

$$\begin{aligned} \frac{1}{2} \|w_t - z\|^2 - \frac{1}{2} \|w_0 - z\|^2 &= \int_0^t \frac{d}{ds} \frac{1}{2} \|w_s - z\|^2 ds \\ &= \int_0^t \langle \dot{w}_s, w_s - z \rangle ds \\ &= \int_0^t -\langle \bar{\partial} \hat{R}(w_s), w_s - z \rangle ds \\ &\leq \int_0^t (\hat{R}(z) - \hat{R}(w_s)) ds. \end{aligned}$$

$$\Rightarrow \frac{1}{2} \|w_t - z\|^2 + \int_0^t \hat{R}(w_s) ds \leq \frac{1}{2} \|w_0 - z\|^2 + \int_0^t \hat{R}(z) ds, \quad \& \text{ use } \hat{R}(w_t) \leq \hat{R}(w_s) \quad \forall s \leq t. //$$