

Hw 1 Released

Simon OH → Thursday 10:30AM
This Week

Clarke Differential

W

Subdifferential and Subgradient

Definition: Given $f: \mathbb{R}^d \rightarrow \mathbb{R}$, for every x , the subdifferential set is defined as

$\partial_s f(x) \triangleq \{s \in \mathbb{R}^d : \forall x' \in \mathbb{R}^d, f(x') \geq f(x) + s^\top (x' - x)\}$. The elements in the subdifferential set are subgradients.

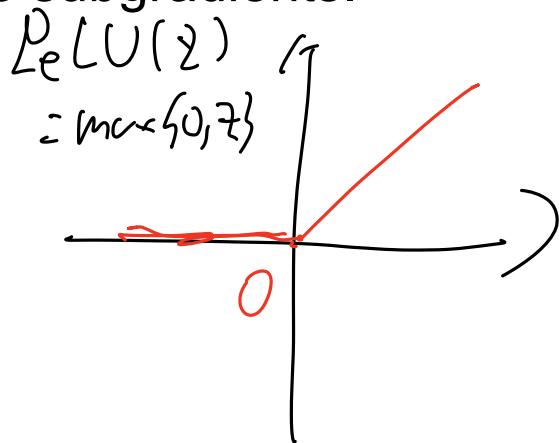
$$\min f(w)$$

$$G): \quad x_{t+1} = x_t - \eta_t \nabla f(x_t)$$

Subgradient descent

$$x_{t+1} = x_t - \eta_t g_t$$

$$g_t \in \partial_s f(x_t)$$



$$\partial_s f(0) = [0, 1]$$

Subdifferential and Subgradient

Definition: Given $f: \mathbb{R}^d \rightarrow \mathbb{R}$, for every x , the subdifferential set is defined as

$\partial_s f(x) \triangleq \{s \in \mathbb{R}^d : \forall x' \in \mathbb{R}^d, f(x') \geq f(x) + s^\top (x' - x)\}$. The elements in the subdifferential set are subgradients.

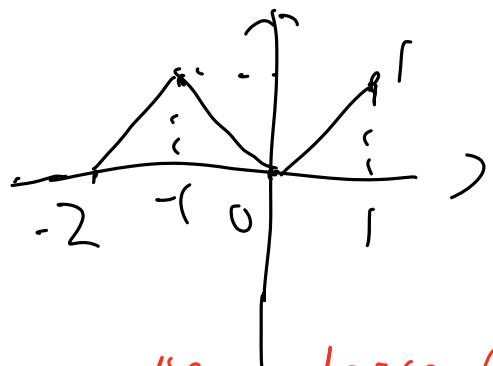
- If f is convex $\rightarrow \partial_s f$ exists everywhere
- If f is convex & differentiable
 $\partial_s f = \{ \nabla f \}$
- $\mathcal{O}\left(\frac{1}{\sqrt{f}}\right)$ in convex

Subdifferential is not enough

Definition: Given $f: \mathbb{R}^d \rightarrow \mathbb{R}$, for every x , the subdifferential set is defined as

$\partial_s f(x) \triangleq \{s \in \mathbb{R}^d : \forall x' \in \mathbb{R}^d, f(x') \geq f(x) + s^\top (x' - x)\}$. The elements in the subdifferential set are subgradients.

Problem: N/N is not convex



subgradient descent
is not well-defined

$$\begin{aligned} x &= -1 \\ \text{we need } s, \forall x' & \\ f(x') &\geq 1 + s^\top (x' - (-1)) \\ \text{let } x' = -2 & \\ 0 &\geq 1 + s(-1) \Rightarrow s \leq 1 \\ \text{let } x' = 1 & \\ 1 &\geq 1 + 2s \Rightarrow s \leq 0 \end{aligned}$$

Clarke Differential

$$x_{f+} = x - y f_f \\ y \in \partial f(x)$$

Definition: Given $f: \mathbb{R}^d \rightarrow \mathbb{R}$, for every x , the Clarke differential is defined as

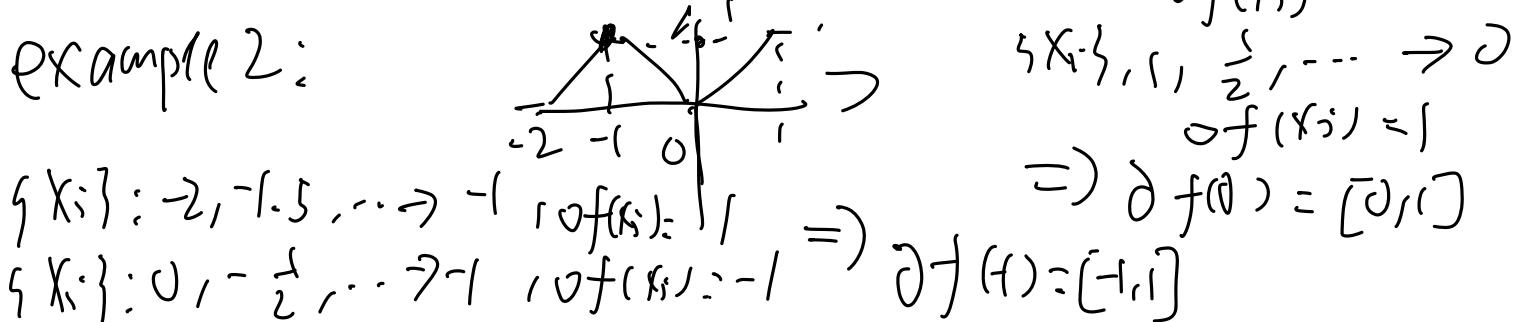
$$\partial f(x) \triangleq \text{conv} \left(\{s \in \mathbb{R}^d : \exists \underbrace{\{x_i\}_{i=1}^\infty}_{\rightarrow x}, \{\nabla f(x_i)\}_{i=1}^\infty \rightarrow s\} \right).$$

The elements in the subdifferential set are subgradients.

$$\text{Conv}(S) = \left\{ v : v = \sum_{i=1}^n \lambda_i u_i, u_i \in S \right. \\ \left. \lambda_i \geq 0, \sum_{i=1}^n \lambda_i = 1 \right\}$$

Example 1 $\text{ReLU}(z) = \max\{0, z\}$

Example 2:



When does Clarke differential exists

Definition (Locally Lipschitz): $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is locally Lipschitz if $\forall x \in \mathbb{R}^d$, there exists a neighborhood S of x , such that f is Lipschitz in S .

$$\text{local} \quad \forall x, x' \in S, |f(x) - f(x')| \leq L \|x - x'\|$$

- If f is locally Lipschitz \Rightarrow ∂f exists everywhere
- If f is convex $\Rightarrow \partial f = \partial_S f$
- If f is differentiable $\Rightarrow \partial f = \{\nabla f\}$
Satisfies chain rule

Positive Homogeneity

Definition: $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is positive homogeneous of degree L if $f(\alpha x) = \alpha^L f(x)$ for any $\alpha \geq 0$. \leftarrow homogeneous

(1) ReLU: $\tilde{f}(\alpha \cdot z) = \alpha \cdot f(z)$

$$f(z) = \max\{0, z\}$$

(2) monomials of degree L : $\prod_{i=1}^d x_i^{p_i}, \sum p_i = L$

$$\prod_{i=1}^d (\alpha x_i)^{p_i} = \alpha^{\sum p_i} \cdot \prod x_i^{p_i} = \alpha^L \prod x_i^{p_i}$$

(3) Norm: $\|(\alpha \cdot x)\| = \alpha \cdot \|x\|$

$$x \in \mathbb{R}^d, w_1 \in \mathbb{R}^{m \times d}$$

$$w_2, \dots, w_t \in \mathbb{R}^{m \times m}$$

$$W_{H+1} \in \mathbb{R}^m$$

(4) Multi-layer ReLU NN

$$f(x, w_1, \dots, w_{H+1}) = W_{H+1} \sigma(W_H \dots \sigma(W_1 x) \dots)$$

1) for each layer : (- homogeneous)

$$\begin{aligned} f(x, w_1, \dots, \alpha w_h, \dots, w_{H+1}) &= W_{H+1} \sigma(\dots \sigma(\alpha w_h \sigma(\dots \sigma(w_1 x) \dots))) \\ &= \alpha f(x, w_1, \dots, w_{H+1}) \end{aligned}$$

2) for all layers

$$f(x, \alpha w_1, \dots, \alpha w_{H+1}) = \alpha^{H+1} f(x, w_1, \dots, w_{H+1})$$

\Rightarrow ($H+1$) - homogeneous

Positive Homogeneity

Fact: , $\forall h$

$$\left\langle w_h, \frac{\partial f(x, w_1, \dots, w_{h-1})}{\partial w_h} \right\rangle = \underbrace{f(x, w_1, \dots, w_{h-1})}_{\text{independent of } h}$$

Pf: $A_n = \text{diag}(\underbrace{g'(w_1 g(\dots g(w_1 x) \dots))}_{m \times m} \in \mathbb{R}^{m \times m})$

$g': 0 \text{ or } 1$ pattern whether the activation is on or not

$$f(x, w_1, \dots, w_{h-1}) = w_{h-1} A_{h-1} w_{h-2} \dots A_1 w_1 x$$
$$\frac{\partial f}{\partial w_h} = \underbrace{(w_{h-1} A_{h-1} \dots w_1 A_1)^T}_{1 \times m} \underbrace{(A_{h-1} w_{h-1} \dots w_1 x)^T}_{m \times 1}$$
$$g(z) = g'(z) \cdot z$$

$$\begin{aligned}
\left\langle W_h, \frac{\partial f}{\partial w_n} \right\rangle &= \left\langle W_h, (W_{H+1} A_H \cdots W_{n+1} A_n)^T (A_{n+1} \cdots W_1 x)^T \right\rangle \\
&\stackrel{A^T B^T}{=} \text{tr} ((W_{H+1} A_H \cdots W_n)^T (A_{n+1} \cdots W_1 x)^T) \\
&\stackrel{=(BA)^T}{=} \text{tr} ((A_{n+1} \cdots W_1 x)^T (W_{H+1} \cdots A_n W_n)^T) \\
&\stackrel{\text{tr}(AB) = \text{tr}(BA)}{=} \text{tr} (W_{H+1} A_H \cdots A_1 W_1 x) \\
&= f(x, W_1, \dots, W_{H+1})
\end{aligned}$$

□

Positive Homogeneity and Clark Differential

clark differential plays

Lemma: Suppose $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is Locally Lipschitz and L -positively homogeneous. For any $x \in \mathbb{R}^d$ and $s \in \partial f(x)$, we have $\langle s, x \rangle = Lf(x)$.

View $\chi = \mathcal{U}_h : f \text{-hom}$

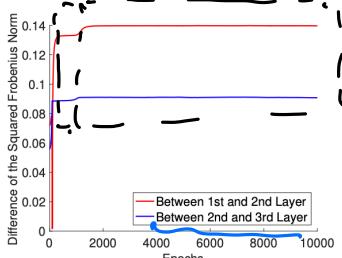
Norm Preservation

$$f(x, w_1, w_2, w_3) = w_3 \sigma(w_2 \sigma(w_1 x)),$$

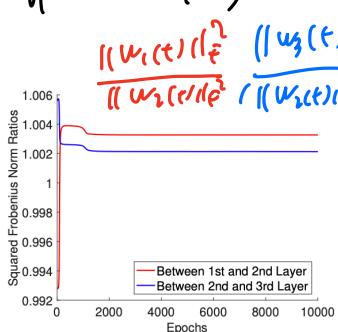
quadratic loss

$$\|w_1\|_F^2, \|w_2\|_F^2, \|w_3\|_F^2, \|A\|_F^2 = \sum_{i,j} A_{ij}^2$$

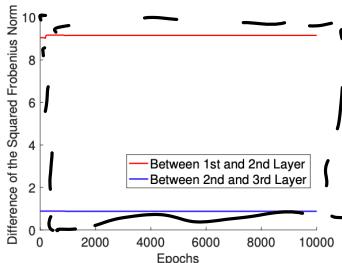
$$\begin{aligned} w_3 &\propto r_0 \\ w_2 &\propto r_0 \\ w_2(t+\epsilon) &= w_2(t) + \frac{\partial f}{\partial w_2} \end{aligned}$$



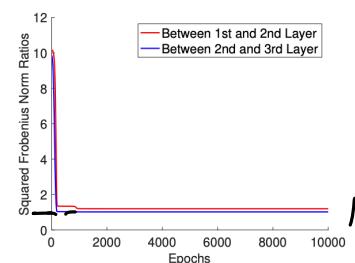
(a) Balanced initialization, squared norm differences.



(b) Balanced initialization, squared norm ratios.



(c) Unbalanced Initialization, squared norm differences.



(d) Unbalanced initialization, squared norm ratios.

- $\|w_1(0)\|_F^2 \approx \|w_2(0)\|_F^2 \approx \|w_3(0)\|_F^2$
- $\|w_1(t)\|_F^2 - \|w_2(t)\|_F^2$
- $\|w_2(t)\|_F^2 - \|w_3(t)\|_F^2$

$$\begin{aligned} \|w_1(0)\|_F^2 &\neq \|w_2(0)\|_F^2 \quad \text{small} \\ \|w_1(t)\|_F^2 &\neq \|w_2(t)\|_F^2 \end{aligned}$$

Gradient flow and gradient inclusion

Discrete-time dynamics can be complex. Let's use continuous-time dynamics to simplify:

$$\text{Gradient flow: } x_{t+1} = x_t - \eta \nabla f(x_t) \Rightarrow \frac{x(t)}{dt} = -\nabla f(x(t))$$

$$\text{Gradient inclusion: } \frac{dx(t)}{dt} \in -\partial f(x(t))$$

$$\frac{x_{t+1} - x_t}{\eta} = -\nabla f(x_t)$$

$$\text{Let } \eta \rightarrow 0$$

Norm preservation by gradient inclusion

Theorem (Du, Hu, Lee '18) Suppose $\alpha > 0$,
 $f(x; (W_{H+1}, \dots, \alpha W_i, \dots, W_1)) = \alpha f(x, (W_{H+1}, \dots, W_1))$, i.e.,
predictions are 1-homogeneous in each layer. Then for every pair
of layers $(j, h) \in [H+1] \times [H+1]$, the gradient inclusion
maintains: for all $t \geq 0$,

$$\frac{1}{2} \|W_h(t)\|_F^2 - \frac{1}{2} \|W_h(0)\|_F^2 = \frac{1}{2} \|W_h(t)\|_F^2 - \frac{1}{2} \|W_h(0)\|_F^2.$$
$$\Leftrightarrow \|W_h(t)\|_F^2 - \|W_h(0)\|_F^2 = \|W_h(t)\|_F^2 - \|W_h(0)\|_F^2$$

- If $\|W_h(0)\|_F^2 \text{ small } \rightarrow \|W_h(t)\|_F^2 \approx \|W_h(0)\|_F^2$

Norm preservation by gradient inclusion

Pf: loss = $L(f(x, w_1, \dots, w_{H+1}), y)$

$$\frac{d W_h(f)}{d f} = - L'(f(x, w_1, \dots, w_{H+1}), y) \cdot \frac{\partial f(x, w_1, \dots, w_H)}{\partial w_h}$$

$$= \frac{1}{2} \|W_h(t)\|_F^2 - \frac{1}{2} \|W_h(0)\|_F^2$$

$$= \int_0^t \frac{d}{dt} \frac{1}{2} \|W_h(\tau)\|_F^2 d\tau$$

$$= \int_0^t \langle W_h, \frac{d W_h(\tau)}{d \tau} \rangle d\tau \quad (\text{chain rule})$$

$$= \int_0^t \underbrace{\langle W_h, 1 \rangle}_{\text{scalar}} \underbrace{- L'(f(x, w_1, \dots, w_{H+1}), y) \cdot \frac{\partial f(x, -w_H)}{\partial w_h}}_{\text{scalar}} d\tau$$

of h = $= - \int_0^t L'(f(x, w_1, \dots, w_{H+1}), y) \cdot f(x, w_1, \dots, w_{H+1}) d\tau$

Optimization Methods for Deep Learning

W

Gradient descent for non-convex optimization

$\|A\|_2$: condition number, largest eigenvalue / absolute value

Descent Lemma: Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be twice differentiable, and $\|\nabla^2 f\|_2 \leq \beta$. Then setting the learning rate $\eta = 1/\beta$, and applying gradient descent, $x_{t+1} = x_t - \eta \nabla f(x_t)$, we have:

$$f(x_t) - f(x_{t+1}) \geq \frac{1}{2\beta} \|\nabla f(x_t)\|_2^2. \quad \text{if } f(x_t) \uparrow$$

Pf: by Taylor expansion & mean-value theorem
[Δ too large]

$$f(x+\Delta) = f(x) + \Delta^\top \nabla f(x) + \frac{1}{2} \Delta^\top \nabla^2 f(y) \Delta$$

$$\Delta^\top \nabla^2 f(y) \Delta \leq \|\nabla^2 f(y)\|_2 \cdot \|\Delta\|_2^2 \quad (\Delta \leq \beta \|\Delta\|_2^2) \quad (y = \lambda x + (1-\lambda)(x+\Delta))$$

$$\text{Let } \Delta = -\eta \cdot \nabla f(x)$$

$$\begin{aligned} f(x+\eta \nabla f(x)) &\leq f(x) - \eta \cdot \|\nabla f(x)\|_2^2 + \frac{1}{2} \beta \eta^2 \|\nabla f(x)\|_2^2 \\ &= f(x) - \frac{1}{2} \cdot \eta \cdot \|\nabla f(x)\|_2^2 \end{aligned}$$

Converging to stationary points

Theorem: In $T = O(\frac{\beta}{\epsilon^2})$ iterations, we have $\|\nabla f(x)\|_2 \leq \epsilon$.

Gradient Descent for Quadratic Functions

Problem: $\min_x \frac{1}{2} x^\top A x$ with $A \in \mathbb{R}^{d \times d}$ being positive-definite.

Theorem: Let λ_{\max} and λ_{\min} be the largest and the smallest eigenvalues of A . If we set $\eta \leq \frac{1}{\lambda_{\max}}$, we have

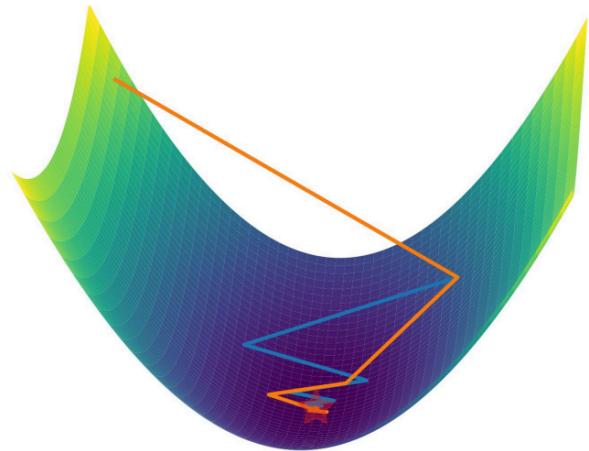
$$\|x_t\|_2 \leq (1 - \eta \lambda_{\min})^t \|x_0\|_2$$

Momentum: Heavy-Ball Method (Polyak '64)

Problem: $\min_x f(x)$

Method: $v_{t+1} = -\nabla f(x_t) + \beta v_t$

$$x_{t+1} = x_t + \eta v_{t+1}$$



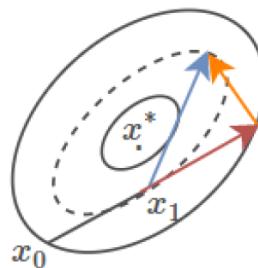
Momentum: Nesterov Acceleration (Nesterov '89)

Problem: $\min_x f(x)$

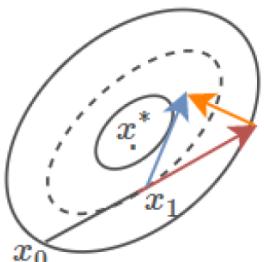
Method: $v_{t+1} = -\nabla f(x_t + \beta v_t) + \beta v_t$

$$x_{t+1} = x_t + \eta v_{t+1}$$

Polyak's Momentum

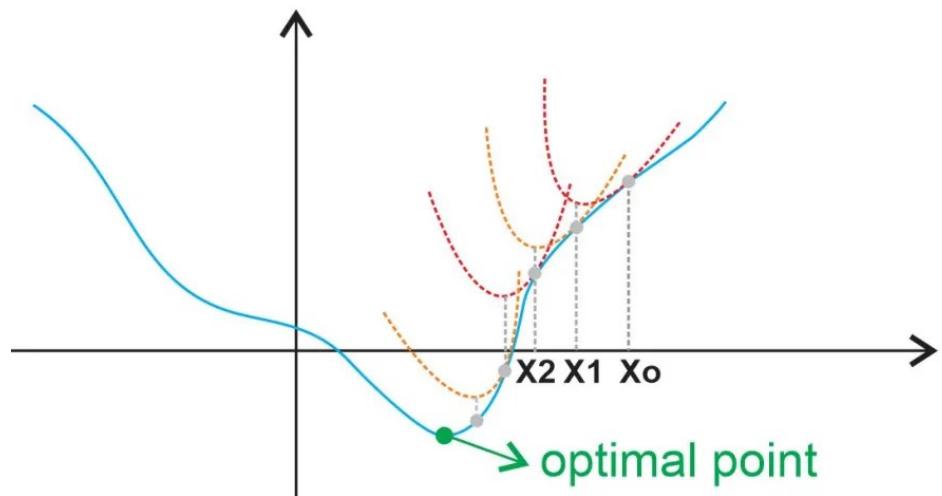


Nesterov Momentum



Newton's Method

Newton's Method: $x_{t+1} = x_t - \eta(\nabla^2 f(x_t))^{-1} \nabla f(x_t)$



AdaGrad (Duchi et al. '11)

Newton Method: $x_{t+1} = x_t - \eta(\nabla^2 f(x_t))^{-1} \nabla f(x_t)$

AdaGrad: separate learning rate for every parameter

$$x_{t+1} = x_t - \eta(G_{t+1} + \epsilon I)^{-1} \nabla f(x_t), (G_t)_{ii} = \sqrt{\sum_{j=1}^{t-1} (\nabla f(x_t)_i)^2}$$

RMSProp (Hinton et al. '12)

AdaGrad: separate learning rate for every parameter

$$x_{t+1} = x_t - \eta(G_{t+1} + \epsilon I)^{-1} \nabla f(x_t), \quad (G_t)_{ii} = \sqrt{\sum_{j=1}^{t-1} (\nabla f(x_t)_i)^2}$$

RMSProp: exponential weighting of gradient norms

$$\begin{aligned} x_{t+1} &= x_t - \eta(G_{t+1} + \epsilon I)^{-1/2} \nabla f(x_t), \\ (G_{t+1})_{ii} &= \beta(G_t)_{ii} + (1 - \beta)(\nabla f(x_t)_i)^2 \end{aligned}$$

AdaDelta (Zeiler '12)

RMSProp:

$$x_{t+1} = x_t - \eta(G_{t+1} + \epsilon I)^{-1/2} \nabla f(x_t),$$
$$(G_{t+1})_{ii} = \beta(G_t)_{ii} + (1 - \beta)(\nabla f(x_t)_i)^2$$

AdaDelta:

$$x_{t+1} = x_t - \eta \Delta x_t,$$
$$\Delta x_t = \sqrt{u_t + \epsilon} \cdot (G_{t+1} + \epsilon I)^{-1/2} \nabla f(x_t)$$
$$(G_{t+1})_{ii} = \rho(G_t)_{ii} + (1 - \rho)(\nabla f(x_t)_i)^2,$$
$$u_{t+1} = \rho u_t + (1 - \rho) \|\Delta x_t\|_2^2$$

Adam (Kingma & Ba '14)

Momentum:

$$v_{t+1} = -\nabla f(x_t) + \beta v_t, \quad x_{t+1} = x_t + \eta v_{t+1}$$

RMSProp: exponential weighting of gradient norms

$$x_{t+1} = x_t - \eta(G_{t+1} + \epsilon I)^{-1} \nabla f(x_t),$$

$$(G_t)_{ii} = \beta(G_t)_{ii} + (1 - \beta)(\nabla f(x_t)_i)^2$$

Adam

$$v_{t+1} = \beta_1 v_t + (1 - \beta_1) \nabla f(x_t)$$

$$(G_{t+1})_{ii} = \beta_2(G_t)_{ii} + (1 - \beta_2)(\nabla f(x_t)_i)^2$$

$$x_{t+1} = x_t - \eta(G_{t+1} + \epsilon I)^{-1/2} v_{t+1}$$

Default choice nowadays.

Other Optimizers

- AdamW
- NAdam
- RAdam
- GGT
- K-FAC
- ...