

HW 1 released, template
(Notes) on approximation released

Clarke Differential



Clarke Differential

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

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

The elements in the subdifferential set are subgradients.

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 .

\Rightarrow Clarke differential exists

Positive Homogeneity

→ motivate ReLU

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

(1) ReLU : $\sigma(\alpha z) = \alpha \cdot \sigma(z)$

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

$$\begin{aligned} \prod_{i=1}^d (\alpha x_i)^{p_i} &= \alpha^{\sum_{i=1}^d p_i} \prod_{i=1}^d x_i^{p_i} \\ &= \alpha^L \cdot \prod_{i=1}^d x_i^{p_i} \end{aligned}$$

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

Positive Homogeneity

(4) Multi-layer ReLU

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

for one-layer

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

for all-layer

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

\Rightarrow $(H+1)$ -homogeneous function

Say $W_h \in \mathbb{R}^{m \times m}$

Positive Homogeneity

- independent of (h)
- hold for ReLU

Fact: $\forall h = 1, \dots, H+1$

$$\left\langle W_h, \frac{\partial f(x, W_1, \dots, W_{H+1})}{\partial W_h} \right\rangle = \underline{f(x, W_1, \dots, W_{H+1})}$$

A, B matrix: $\langle A, B \rangle = \sum_{ij} A_{ij} B_{ij}$

Pf: $A_h = \text{diag} (\sigma'(W_h \sigma(\dots \sigma(W_1 x) \dots))) \in \mathbb{R}^{m \times m}$

($\sigma' = 0$ or 1) \Rightarrow matter whether activation is on or off

$\sigma(z) = z \cdot \sigma'(z)$: holds for ReLU

$$f(x, W_1, \dots, W_{H+1}) = W_{H+1} A_H W_H \dots A_1 W_1 x$$

$$\frac{\partial f}{\partial W_h} = (W_{H+1} A_H \dots W_{h+1} A_h)^T (A_{h-1} W_h \dots W_1 x)^T$$

Property $\left\langle W_h, \frac{\partial f}{\partial W_h} \right\rangle = f(x, W_1, \dots, W_{H+1})$

Positive Homogeneity and Clark Differential

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 $\langle s, x \rangle = Lf(x)$.

a?

⊙ ~~~~~

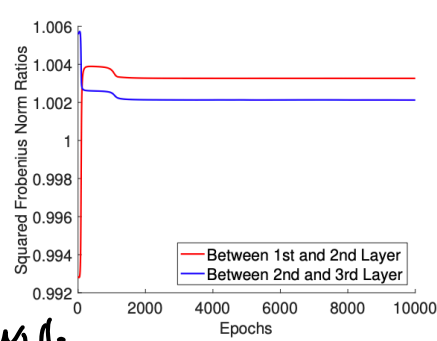
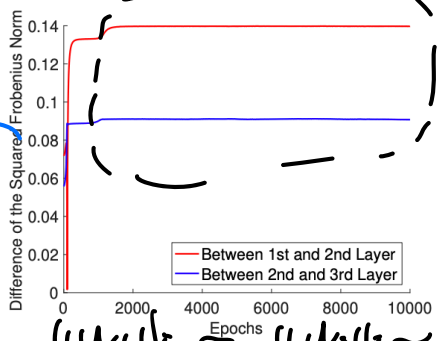
(HW 1)

Norm Preservation

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

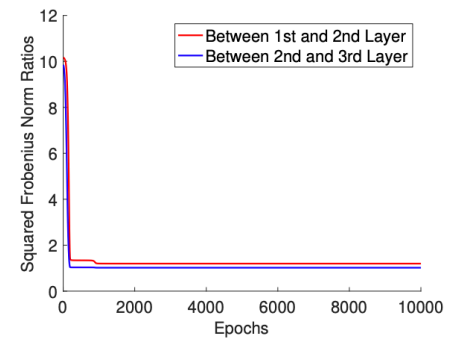
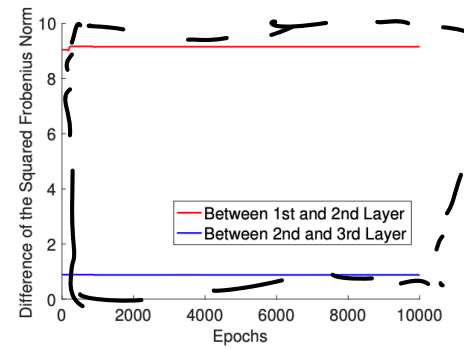
quadratic loss

$$\|w_1\|_F \ll \|w_2\|_F \ll \|w_3\|_F^2$$



(a) Balanced initialization, squared norm differences.

(b) Balanced initialization, squared norm ratios.



(c) Unbalanced Initialization, squared norm differences.

(d) Unbalanced initialization, squared norm ratios.

$$\text{---} : \|w_1\|_F^2 - \|w_2\|_F^2$$

$$\text{---} : \|w_2\|_F^2 - \|w_3\|_F^2$$

Given A matrix, $\|A\|_F^2 = \sum_{i,j} A_{ij}^2$

Gradient flow and gradient inclusion

widely used in CVX / opt theory

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

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

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

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

let $\eta \rightarrow 0 \Rightarrow \frac{dx(t)}{dt}$

$$f(x, W_1, W_2, W_3) = W_3 G(W_2 G(W_1 x)), \quad W_3 \times 10, W_2/10$$

$$W_2(t+1) = W_2(t) - \eta \frac{\partial f}{\partial W_2}$$

Norm preservation by gradient inclusion

No assumption on loss

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 $(i, j) \in [H+1] \times [H+1]$, the gradient inclusion maintains: for all $t \geq 0$,

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

• if $\|W_i(0)\|_F^2$ small for all i same for all layers
 $\Rightarrow \|W_i(t)\|_F^2 \approx \|W_j(t)\|_F^2 \Rightarrow$ balance

(pf sketch: 1) $\frac{dW_i(t)}{dt}$ formula
 2) $\frac{1}{2} \|W_i(t)\|_F^2 - \frac{1}{2} \|W_i(0)\|_F^2 = \int_0^t \frac{d}{dt} \frac{1}{2} \|W_i(t)\|_F^2 dt$
 \Rightarrow independent of i

Optimization Methods for Deep Learning



Gradient descent for non-convex optimization

$\|A\|_2$: operator norm, largest absolute eigenvalue

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:

$$f(x_t) \downarrow \quad f(x_t) - f(x_{t+1}) \geq \frac{1}{2\beta} \|\nabla f(x_t)\|_2^2 \quad \begin{array}{l} f(x_t) \uparrow \\ \Rightarrow \text{learning rate} \\ \text{too large} \end{array}$$

Pf: by Taylor expansion & mean-value theorem

$$f(x+\Delta) = f(x) + \Delta^T \nabla f(x) + \frac{1}{2} \Delta^T \nabla^2 f(y) \Delta \quad \text{for some } y$$

$$\Delta^T \nabla^2 f(y) \Delta \leq \|\nabla^2 f(y)\|_2 \cdot \|\Delta\|_2^2 \leq \beta \|\Delta\|_2^2$$

$$\text{let } \Delta = -\eta \nabla f(x_t)$$

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

Converging to stationary points

an approximate stationary point

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

$$P_f: f(x_{t+1}) \leq f(x_t) - \frac{\eta}{2} \|\nabla f(x_t)\|_2^2$$

Sum over $t = 0, \dots, T-1$

$$\sum_{t=1}^T f(x_t) \leq \sum_{t=0}^{T-1} f(x_t) - \frac{\eta}{2} \sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2$$

$$\Rightarrow f(x_T) \leq f(x_0) - \frac{\eta}{2} \sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2$$

$$\Rightarrow \frac{\eta}{2} \sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2 \leq f(x_0) - f(x_T)$$

$$\frac{\eta}{2} T \cdot \min_{0 \leq t \leq T-1} \|\nabla f(x_t)\|_2^2 \leq f(x_0) - \min_x f(x)$$

$$\Rightarrow \min_{0 \leq t \leq T-1} \|\nabla f(x_t)\|_2 \leq \sqrt{\frac{2\beta(f(x_0) - \min_x f(x))}{(T-1)}} = \epsilon$$

scale $\frac{1}{T}$

Gradient Descent for Quadratic Functions

optimal $x=0$

$\lambda_{\min}(A) > 0$

Problem: $\min_x \frac{1}{2} x^T A x$ 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

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

$$\begin{aligned} \|x_{t+1}\|_2 &= \|x_t - \eta A x_t\|_2 \\ &= \|(I - \eta A)x_t\|_2 \\ &\leq \|I - \eta A\|_2 \|x_t\|_2 \\ &\leq (1 - \eta \lambda_{\min}) \|x_t\|_2 \\ &\leq (1 - \eta \lambda_{\min})^{t+1} \|x_0\|_2 \end{aligned}$$

To make $\|x_t\|_2 \leq \epsilon$
when $\eta = \frac{1}{\lambda_{\max}}$

need $\left(\frac{\lambda_{\max}}{\lambda_{\min}} \log\left(\frac{1}{\epsilon}\right) \right)$ steps

$\kappa = \frac{\lambda_{\max}}{\lambda_{\min}}$ condition number

Momentum: Heavy-Ball Method (Polyak '64)

Problem: $\min_x f(x)$

$\beta < 1$

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

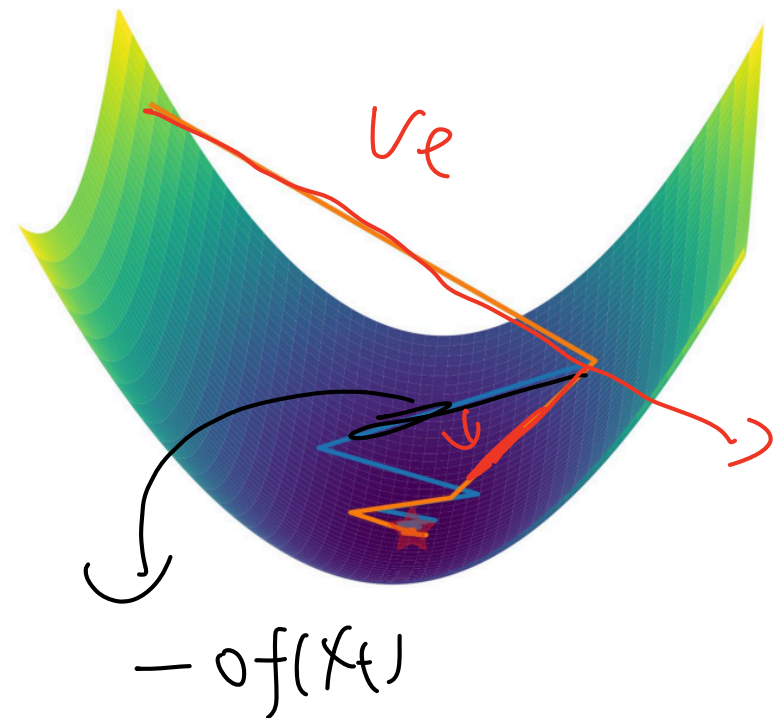
$$x_{t+1} = x_t + \eta v_{t+1}$$

For quadratic optimization
provable improvement

$$O\left(\sqrt{\kappa} \cdot \log\left(\frac{1}{\epsilon}\right)\right)$$

vs. $O\left(\kappa \cdot \log\left(\frac{1}{\epsilon}\right)\right)$

doesn't work for
general convex function



Momentum: Nesterov Acceleration (Nesterov '89)

Problem: $\min_x f(x)$

\mathcal{K} : $\frac{\text{Smoothness}}{\text{S.C. parameter}}$

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

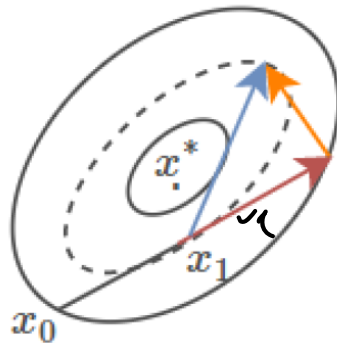
$$x_{t+1} = x_t + \eta v_{t+1}$$

For general strongly convex function

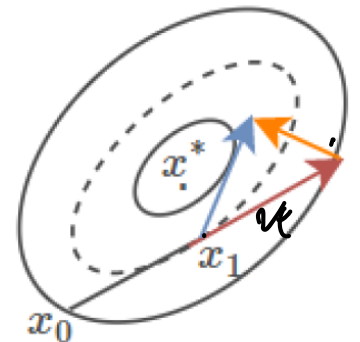
$$O\left(\sqrt{\kappa} \log\left(\frac{1}{\epsilon}\right)\right)$$

continuous approximation

Polyak's Momentum



Nesterov Momentum



Newton's Method

2nd order method

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

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

$\Leftrightarrow f(x+\Delta) \approx f(x) + \Delta^T \nabla f(x) + \frac{1}{2} \Delta^T \nabla^2 f(x) \Delta$

$\Rightarrow \Delta = -\nabla f(x)$

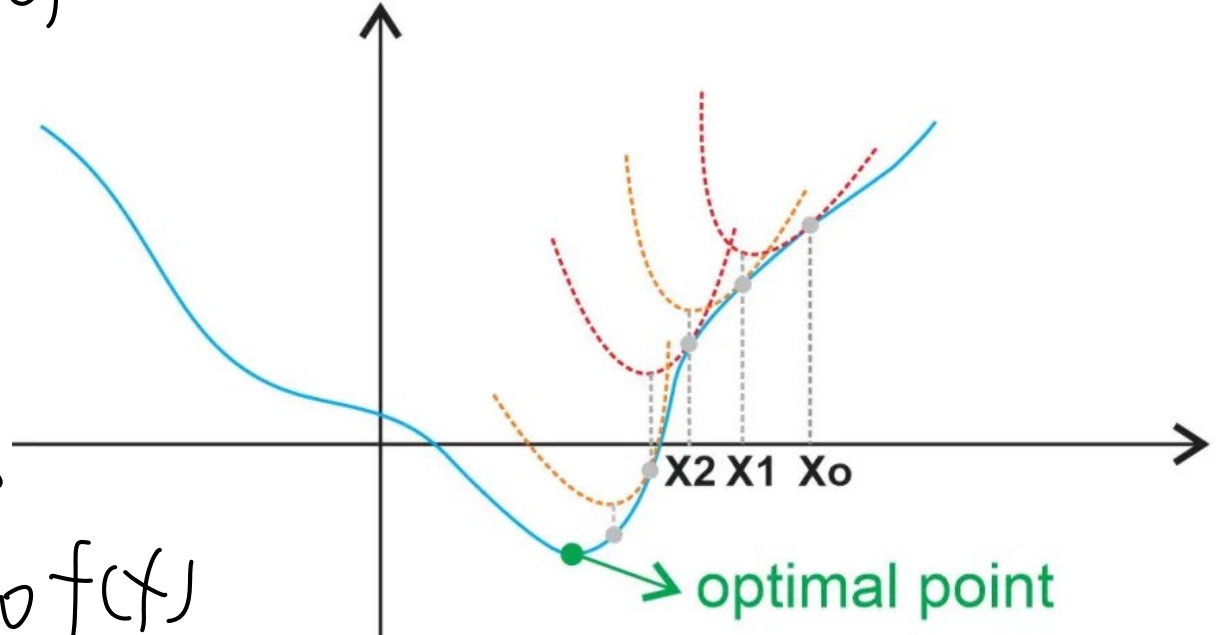
• Newton

$f(x+\Delta) \approx f(x) + \Delta^T \nabla f(x) + \frac{1}{2} \Delta^T \nabla^2 f(x) \Delta$

$\Rightarrow \Delta = -(\nabla^2 f(x))^{-1} \nabla f(x)$

$\mathcal{O}(\log \log(\frac{1}{\epsilon}))$, 1st dependence K

• Problem: invert $\nabla^2 f(x) \Rightarrow \mathcal{O}(d^3)$



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

positive-definite, G_t : diagonal

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

- dynamic learning rate for each coordinate

- default value works well

$$\eta = 0.01, \quad \epsilon = 10^{-8}$$

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

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

exponential average

$$0 < \beta < 1$$

AdaDelta (Zeiler '12)

RMSProp:

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

$\frac{\partial f}{\partial x} / \left(\frac{\partial^2 f}{\partial x^2} \right)^{1/2}$
unitless

AdaDelta:

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

unit observation

Newton:

$\frac{\partial}{\partial x} / \frac{\partial^2 f}{\partial x^2}$ unit of x

\Rightarrow approximate Newton

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

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

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

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

Default choice nowadays.

Are these actually useful

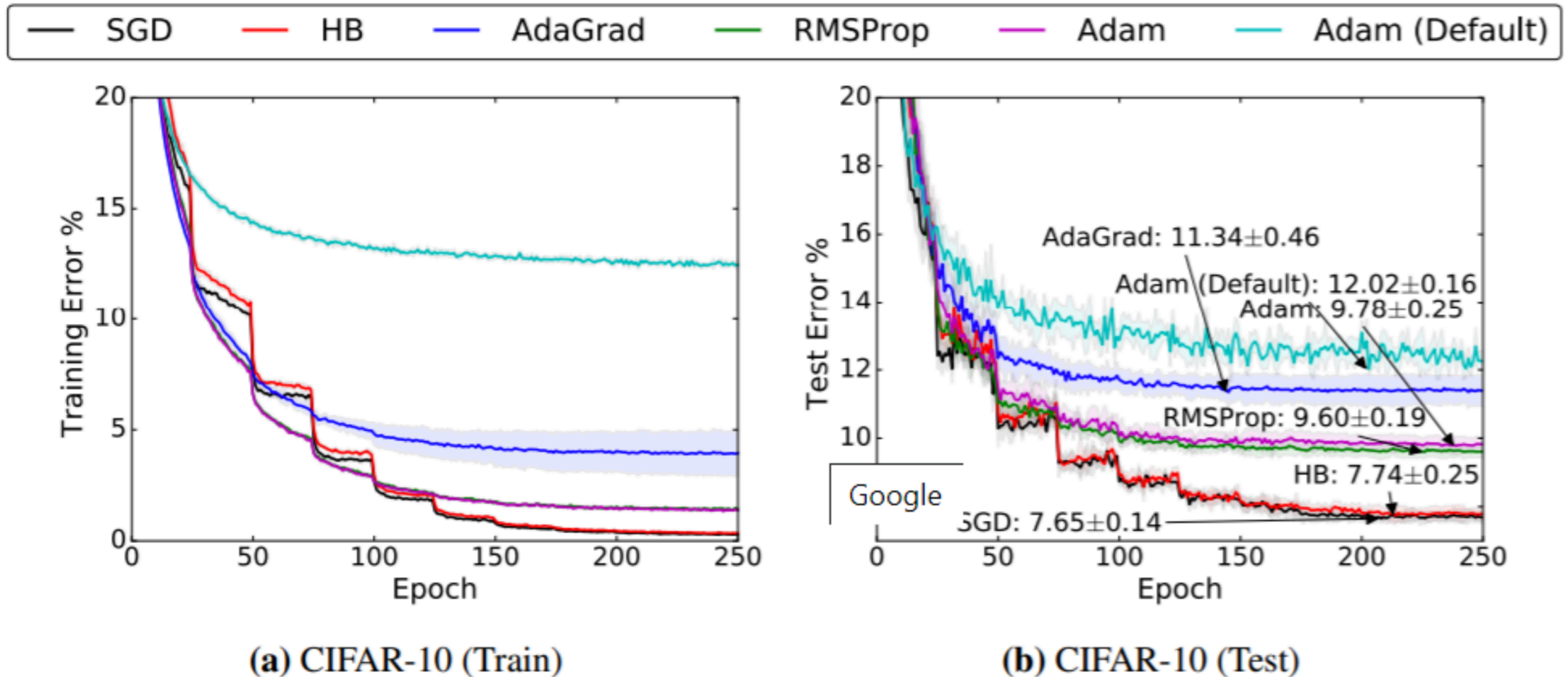


Figure 1: Training (left) and top-1 test error (right) on CIFAR-10. The annotations indicate where the best performance is attained for each method. The shading represents \pm one standard deviation computed across five runs from random initial starting points. In all cases, adaptive methods are performing worse on both train and test than non-adaptive methods.

Wilson, Roelofs, Stern, Srebro, Recht '18