

Large Stepsize GD for Logistic Loss

Non-Monotonicity of the Loss Improves Optimization Efficiency

Jingfeng Wu¹, Peter Bartlett^{1,3}, Matus Telgarsky², Bin Yu¹

¹UC Berkeley, ²New York University, ³Google DeepMind

Background

$$w_+ = w - \eta \nabla L(w)$$

How to choose **stepsize**?

Descent lemma

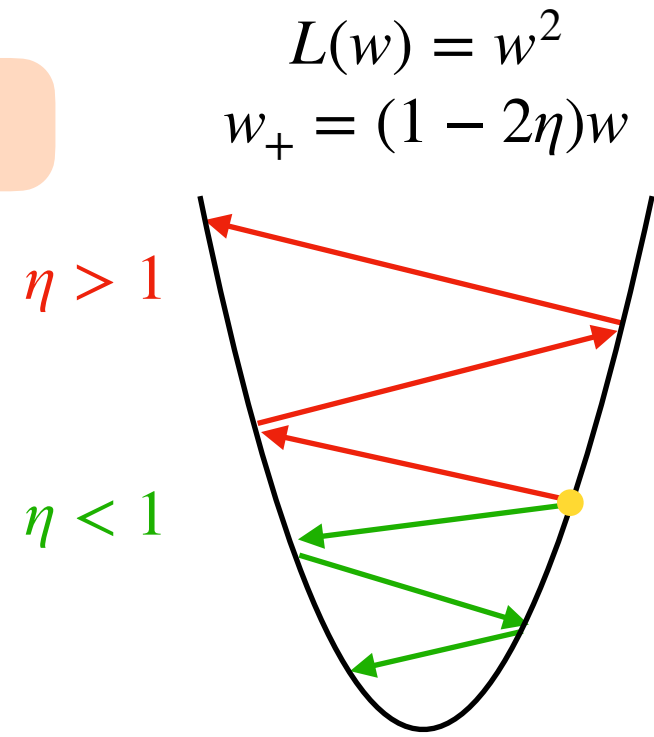
For **small** η , $L(w_t)$ decreases **monotonically**

For **large** η , $L(w_t)$ **diverges** for quadratics

$$L(w_+) = L(w - \eta \nabla L(w))$$

$$= L(w) - \eta \|\nabla L(w)\|^2 + \frac{\eta^2}{2} \nabla L(w)^\top \nabla^2 L(w) \nabla L(w) + O(\eta^3)$$

$$\leq L(w) - \eta \left(1 - \frac{\eta}{2} \|\nabla^2 L(w)\|_2\right) \|\nabla L(w)\|^2 + O(\eta^3)$$

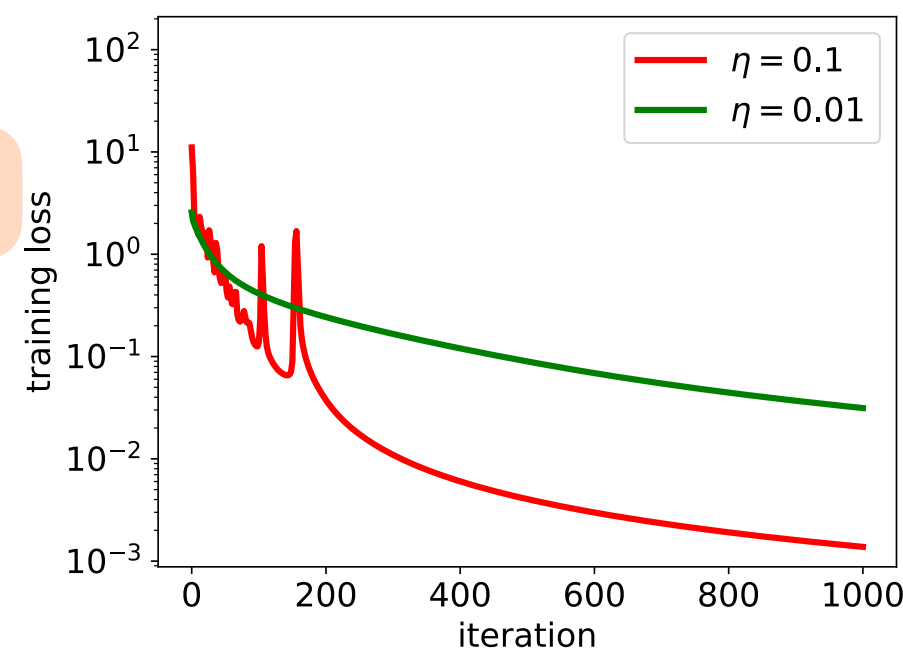


Edge of Stability (EoS)

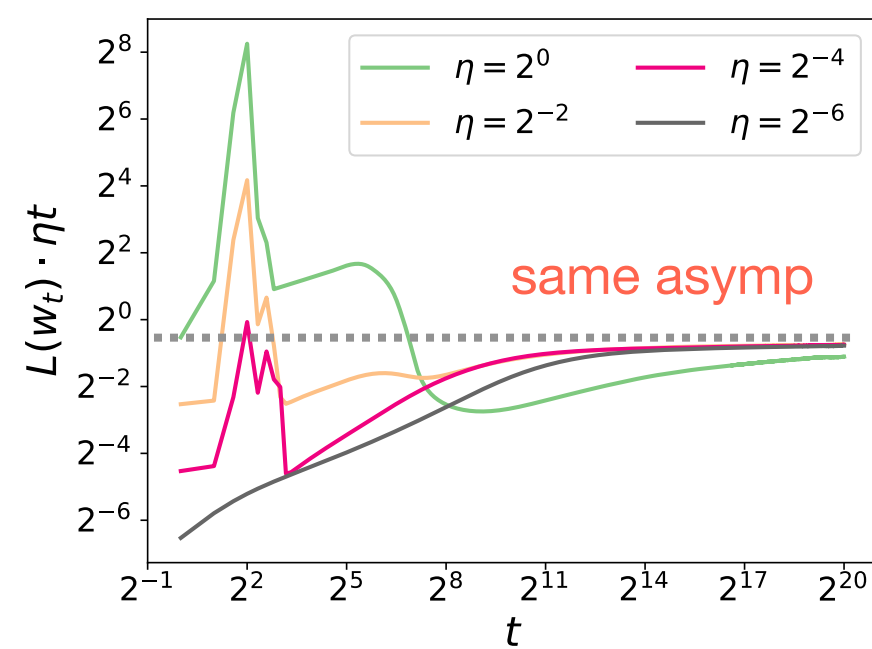
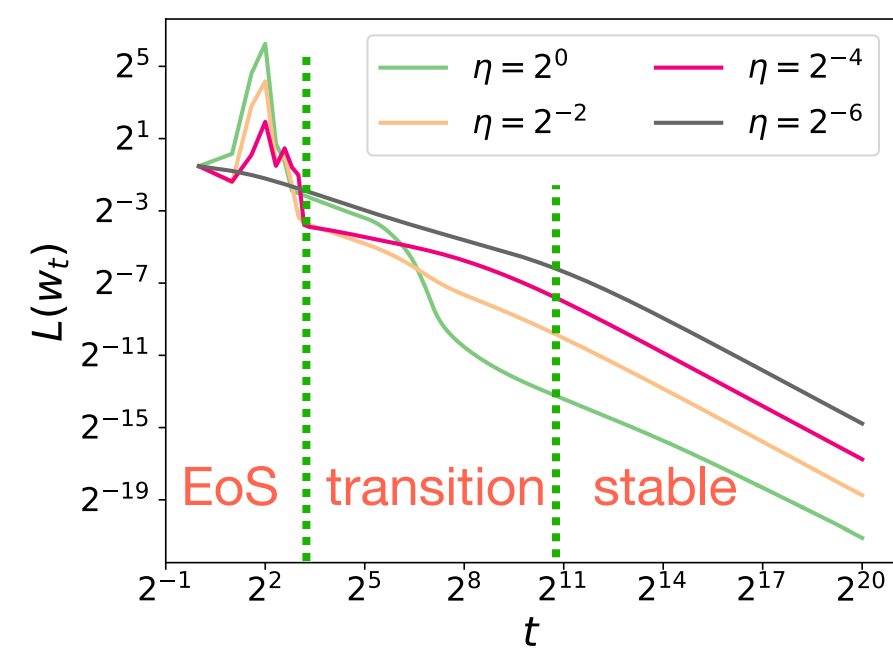
large stepsize works better;

“spikes” or “edge of stability”

unexplained by descent lemma



3-layer net + 1,000 samples from MNIST



logistic regression + 1,000 samples from MNIST “0” or “8”

A EoS Theory in Logistic Regression

classification data $(x_i, y_i)_{i=1}^n$, $\|x_i\| \leq 1$, $y_i \in \{\pm 1\}$

logistic loss + linear model

$$L(w) := \frac{1}{n} \sum_i \ln(1 + \exp(-y_i x_i^\top w))$$

Assume [linear separability]:

\exists vector w_* such that $y x^\top w_* > \gamma > 0$

Theorem

• **EoS phase.** For every t

$$\frac{1}{t} \sum_{k=0}^{t-1} L(w_k) \leq \tilde{O}\left(\frac{1 + \eta^2}{\eta t}\right)$$

• **Stable phase.** If $L(w_s) \leq 1/\eta$ for some s , then $L(w_{s+t}) \downarrow$ for $t \geq 0$ and

$$L(w_{s+t}) \leq \tilde{O}\left(\frac{F(w_s)}{\eta t}\right), \quad F(w_s) := \mathbb{E} \exp(-y x^\top w_s)$$

• **Phase transition.** We have $L(w_s) \leq 1/\eta$ and $F(w_s) \leq 1$ for

$$s \leq \tau := \Theta(\max\{\eta, n, n/\eta \ln(n/\eta)\})$$

Benefits of large stepsizes

- Asymptotic $\tilde{O}(1/\eta t)$ for **every** η (beyond $1/\text{smoothness}$)
- Larger $\eta \Rightarrow$ smaller const factor, but longer EoS
- Given #steps $T \geq \Omega(n)$, if choose $\eta = \Theta(T)$, then

$$\tau \leq T/2 \text{ and } L(w_T) \leq \tilde{O}(1/T^2)$$

“acceleration” by EoS
w/o momentum or varying stepsizes

- Theorem. In general, if not enter EoS, then $L(w_T) \geq \Omega(1/T)$

Extensions — SGD — General Loss Functions — Neural Tangent Kernel

Theorem. SGD for logistic regression

Let $(w_k)_{k=1}^n$ be iterates of **const stepsize online SGD** for **logistic regression** on iid data from a **separable distribution**. Then for **every stepsize** η , w.p. $\geq 1 - \delta$:

$$\frac{1}{n} \sum_{k=1}^n \mathbb{E} \ln(1 + \exp(-y x^\top w_k)) \lesssim \frac{\ln^2(\gamma^2 \eta n) + \eta^2}{\gamma^2 \eta n} + \frac{(\ln(\gamma^2 \eta n) + \eta) \ln(1/\delta)}{\gamma n}$$

$$\frac{1}{n} \sum_{k=1}^n \Pr(y x^\top w_k \leq 0) \lesssim \frac{\ln(\gamma^2 \eta n) + \eta}{\gamma^2 \eta n} + \frac{\ln(1/\delta)}{n}$$

large stepsize works but no acceleration (upto log)

A general loss function $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$

A. **Regularity.** Assume ℓ is \mathcal{C}^2 , convex, \downarrow , and $\ell(+\infty) = 0$,

define $\rho(\lambda) := \min_{z \in \mathbb{R}} \lambda \ell(z) + z^2$, $\lambda \geq 1$ minimizer away from init

B. **Lipschitzness.** Assume $g(\cdot) := |\ell'(\cdot)| \leq C_g$ prevents GD from diverging

C. **Self-boundedness.** Assume $g(\cdot) \leq C_\beta \ell(\cdot)$ and for entering stable phase

$$\ell(z) \leq \ell(x) + \ell'(z-x) + C_\beta g(x)(z-x)^2, \text{ for } |z-x| \leq 1$$

D. **Exp-tail.** Assume $\ell(\cdot) \leq C_e g(\cdot)$ unnecessary but improves transition time

$$L(w) := \mathbb{E} \ell(y f_x(w)), \quad f_x(w) := \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \max\{x^\top w^{(s)}, 0\}, \quad w \in \mathbb{R}^{md}$$

Assume NTK init: $w_0 \sim \mathcal{N}(0, I_{md})$, $(a_s)_{s=1}^m$ random from $\{\pm 1\}$ & fixed

Assume: “separable” in NTK RKHS holds generically

Theorem. GD for NN/general losses

Assume ℓ satisfies A-B. Fix T , assume $m \geq \Omega(R^2)$ for $R := \Theta(\sqrt{\rho(\eta T)} + \eta)$.

• **Lazy training.** For $t \leq T$, we have $\|w_t - w_0\| \leq R$

• **EoS phase.** For $t \leq T$, we have $\frac{1}{t} \sum_{k=0}^{t-1} L(w_k) \leq O\left(\frac{\rho(\eta t) + \eta^2}{\eta t}\right)$

• **Stable phase.** Assume ℓ also satisfies C. If $L(w_s) \leq \Theta(1/(\eta + n))$ for some s ,

$$\text{then } L(w_{s+t}) \downarrow \text{ and } L(w_{s+t}) \leq O\left(\frac{\rho(\eta t)}{\eta t}\right), \quad s+t \leq T$$

• **Phase transition.** We have $L(w_s) \leq \Theta(1/(\eta + n))$ for some $s \leq \tau$, where

$$\tau := \Theta(\max\{\psi^{-1}(\eta + n), \eta(\eta + n)\}), \quad \psi(\lambda) := \lambda/\rho(\lambda)$$

or $\tau := \Theta(\max\{\eta, n \ln(n)\})$ if ℓ also satisfies D

large stepsize accelerates NTK & general losses