#### Reimagining Gradient Descent Large Stepsize, Oscillation, Acceleration

Jingfeng Wu



#### Gradient descent

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

"GD  $\approx$  discrete time gradient flow"



Cauchy, 1847

$$dw = -\nabla L(w)dt \implies dL(w) = \nabla L(w)^{\mathsf{T}}dw$$
$$= -\|\nabla L(w)\|^{2}dt$$
$$\implies L(w) \downarrow$$

#### because of stepsize, $GD \neq discrete$ time flow

## Small stepsize for stability $L(w_{+}) = L(w - \eta \nabla L(w))$ $= L(w) - \eta \|\nabla L(w)\|^2 + \frac{\eta^2}{2} \nabla L(w)^{\mathsf{T}} \nabla^2 L(v) \nabla L(w)$ $\leq L(w) - \eta \|\nabla L(w)\|^2 \left( \left\|1 - \frac{\eta}{2} \|\nabla^2 L(v)\|_2 \right) \right)$ $L(w) = w^2$ $\eta < \frac{2}{\sup \|\nabla^2 L(\cdot)\|}$ $w_{+} = (1 - 2\eta)w$ **Descent lemma:** for small $\eta$ , $L(w_t)$ decreases monotonically for large $\eta$ , $L(w_t)$ diverges in "bad" cases

## Classical theory

Let *L* be 1-smooth with a finite minimizer  $w^*$ . For GD with  $\eta = 1$ ,

descent lemma  $L(w_t) \downarrow$ convexity  $L(w_t) - \min L \leq \frac{\|w_0 - w^*\|^2}{2t}$ 

 $\alpha$ -strong convexity  $L(w_t) - \min L \le e^{-\alpha t}(L(w_0) - \min L)$ 

Nesterov's momentum accelerates GD to

$$O(1/t^2)$$
 and  $O(e^{-\sqrt{\alpha}t})$ 

these are minimax optimal among first-order methods

## Experiment (3-layer net, MNIST)



large stepsize is

- unstable
- but faster

#### "edge of stability"

Cohen, Kaur, Li, Kolter, and Talwalkar. "Gradient descent on neural networks typically occurs at the edge of stability." ICLR 2021



# (1/3) Seeking "simplest" answer

linear regression

unstable

convergence

impossible

logistic regression

observable

& provable

• • • • • •

deep learning

unstable convergence observed



Peter Bartlett



Matus Telgarsky



Bin Yu

Wu, Bartlett\*, Telgarsky\*, and Yu\*. "Large stepsize gradient descent for logistic loss: nonmonotonicity of the loss improves optimization efficiency." COLT 2024

## Logistic regression

$$L(w) = \frac{1}{n} \sum_{i=1}^{n} \ln(1 + \exp(-y_i x_i^{\mathsf{T}} w))$$
$$w_{t+1} = w_t - \eta \nabla L(w_t)$$

**Assumption** (bounded + separable)

- $||x_i|| \le 1, y_i \in \{\pm 1\}, i = 1, ..., n$
- $\exists$  unit vector  $w^*$ ,  $\min_i y_i x_i^\top w^* \ge \gamma > 0$

**Classical theory** 

For  $\eta = \Theta(1)$ ,  $L(w_t) \downarrow$  and  $L(w_t) = \tilde{O}(1/t)$ 

improved to  $\tilde{O}(1/t^2)$  by Nesterov



#### MNIST "O" vs "8"



Stable phase:  $L(w_t) \downarrow$  from t and onwards EoS phase: otherwise

#### Theorem

#### **Phase transition.** GD exists EoS in au steps for

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

**Stable phase.** From au and onwards

$$L(w_{\tau+t}) = \tilde{O}\left(\frac{1}{\eta t}\right)$$

- 1. Convergence for **every**  $\eta$
- 2. Large  $\eta$ : faster in stable phase but stays longer in EoS

3. Given #steps  $T \ge \Theta(n)$ , if choose  $\eta = \Theta(T)$ , then

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

acceleration by large stepsize



#### Proof

$$||w_{t+1} - u||^{2} = ||w_{t} - u||^{2} + 2\eta \langle \nabla L(w_{t}), u - w_{t} \rangle + \eta^{2} ||\nabla L(w_{t})||^{2}$$
$$= ||w_{t} - u||^{2} + 2\eta \langle \nabla L(w_{t}), u_{1} - w_{t} \rangle$$
$$+ \eta^{2} \left( \langle \nabla L(w_{t}), 2u_{2}/\eta \rangle + ||\nabla L(w_{t})||^{2} \right)$$

$$\langle \nabla L(w), w^* \rangle < 0 \implies \leq 0 \text{ if } u_2 = w^* \cdot \Theta(\eta)$$

$$\|\nabla L(w)\| \le 1$$

$$\leq \|w_t - u\|^2 + 2\eta \langle \nabla L(w_t), u_1 - w_t \rangle$$
  
$$\leq \|w_t - u\|^2 + 2\eta (L(u_1) - L(w_t))$$

Telescoping the sum...

#### Two extensions



 $\lambda \rightarrow \infty$ 



unstable convergence under finite minimizer



## (2/3) Large stepsize for adaptive GD

self-bounded $\|\nabla^2 L\| \le L$ 



large stepsizes for GD variants



Ruiqi Zhang



Licong Lin



Peter Bartlett

Zhang, **Wu**, Lin, Bartlett. "Minimax optimal convergence of gradient descent in logistic regression via large and adaptive stepsizes." ICML 2025

Adaptive GD

$$L(w) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i x_i^{\mathsf{T}} w) \qquad \ell(t) = \ln(1 + \exp(-t))$$

$$w_{t+1} = w_t - \eta \left( (-\ell^{-1})' \circ L(w_t) \right) \nabla L(w_t)$$
$$\approx w_t - \frac{\eta}{L(w_t)} \nabla L(w_t)$$

$$w_{t+1} = w_t - \eta \nabla \phi(w_t) \qquad \phi(w) = -\ell^{-1}(L(w))$$
  
 
$$\approx \ln \sum \exp(-y_i x_i^{\mathsf{T}} w)$$
  
[Ji & Telgarsky, 2021]

For  $\eta = \Theta(1)$ ,  $L(w_t) \downarrow$  and  $L(w_t) \leq \exp(-\Theta(t))$ 

large stepsize makes adaptive GD even faster



#### Theorem

Assume separability with margin  $\gamma$ . For  $t \ge 1/\gamma^2$ , we have

$$L(\bar{w}_t) \le \exp\left(-\Theta(\gamma^2 \eta t)\right)$$
, where  $\bar{w}_t = \frac{1}{t} \sum_{k=1}^t w_k$ 

1. Arbitrarily small error in  $1/\gamma^2$  steps

$$\lim_{\eta \to \infty} L(\bar{w}_t) = 0 \quad \text{for} \quad t = 1/\gamma^2$$

2. Averaged iterate, no "stable phase"

3. small < large < small adaptive << large adaptive  $\tilde{O}(1/\epsilon) \quad \tilde{O}(1/\epsilon^{1/2}) \quad O(\ln(1/\epsilon)) \quad O(1)$ 

## Theorem (lower bound)

 $\forall w_0, \exists (x_i, y_i)_{i=1}^n$  with margin  $\gamma$  such that: for any first-order batch method

$$\min_{i} y_i x_i^{\mathsf{T}} w_t > 0 \implies t \ge \Omega(1/\gamma^2)$$

First-order batch method:

 $w_t \in w_0 + \text{span}\{ \nabla L(w_0), ..., \nabla L(w_{t-1}) \}$ 

where  $L(w) = \hat{\mathbb{E}}\ell(yx^{\mathsf{T}}w)$  for any  $\ell$ 

adaptive GD + large stepsize = minimax optimal

#### (3/3) Large stepsize under finite minimizer

minimizer at  $\infty$ 

 $\lim_{\lambda \to \infty} L(\lambda w^*) = 0$ 



unstable convergence under finite minimizer



Pierre Marion



Peter Bartlett

Wu\*, Marion\*, and Bartlett. "Large stepsizes accelerate gradient descent for regularized logistic regression." arXiv:2506.02336

#### Regularized logistic regression

$$\tilde{L}(w) = L(w) + \frac{\lambda}{2} \|w\|^2$$
$$w_{t+1} = w_t - \eta \nabla \tilde{L}(w_t)$$

 $\lambda$ -strongly convex,  $\Theta(1)$ -smooth,  $\kappa = \Theta(1/\lambda)$ finite minimizer  $w_{\lambda}$ ,  $||w_{\lambda}|| = O(\ln(1/\lambda))$ 

#### **Classical theory**

 $\tilde{O}(1/\lambda)$ 

 $L(w) = \frac{1}{n} \sum_{i} \ell(y_i x_i^{\mathsf{T}} w)$ 

For  $\eta = \Theta(1)$ ,  $\tilde{L}(w_t) \downarrow$  and  $\tilde{L}(w_t) - \min \tilde{L} \leq \epsilon$  for  $t = O(\kappa \ln(1/\epsilon))$ 

improved to  $\tilde{O}(1/\lambda^{1/2})$  by Nesterov

#### Theorem (small $\lambda$ )

Assume separability and

$$\lambda \leq \Theta\left(\frac{1}{n\ln n}\right) \quad \eta \leq \Theta\left(\min\left\{\frac{1}{\lambda^{1/2}}, \frac{1}{n\lambda}\right\}\right)$$

**Phase transition.** GD exists EoS in  $\tau$  steps for

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

**Stable phase.** From  $\tau$  and onward

$$\tilde{L}(w_{\tau+t}) - \min \tilde{L} \lesssim \exp(-\lambda \eta t)$$

for small  $\lambda$ , large stepsize GD matches Nesterov

#### Theorem (general $\lambda$ )

Assume separability and

$$\lambda \leq \Theta(1), \quad \eta \leq \Theta(1/\lambda^{1/3})$$

**Phase transition.** GD exists EoS in  $\tau$  steps for

$$\tau := \Theta(\eta^2)$$

**Stable phase.** From  $\tau$  and onward

$$\tilde{L}(w_{\tau+t}) - \min \tilde{L} \lesssim \exp(-\lambda \eta t)$$

for general  $\lambda$ , large stepsizes is faster than small stepsizes  $\tilde{O}(1/\lambda^{2/3})$   $\tilde{O}(1/\lambda)$ 

#### A new picture



**EoS.**  $\tilde{L} \approx L, R \leq \Theta(1)$ , "overshoot"  $||w_{\lambda}|| = O(\ln(1/\lambda))$ **Stable.** "move back"  $\sup ||w_t|| = \Theta(\eta) = \operatorname{poly}(1/\lambda)$ 

## Stepsize diagram



# (4/3) More results

generalization

• SGD

- networks in kernel regime
- two-layer networks with linear teacher
- other loss functions
  implicit bias

Wu, Bartlett\*, Telgarsky\*, and Yu\*. "Large stepsize gradient descent for logistic loss: nonmonotonicity of the loss improves optimization efficiency." COLT 2024

Zhang, **Wu**, Lin, Bartlett. "Minimax optimal convergence of gradient descent in logistic regression via large and adaptive stepsizes." ICML 2025

**Wu**\*, Marion\*, and Bartlett. "Large stepsizes accelerate gradient descent for regularized logistic regression." arXiv:2506.02336

Cai, **Wu**, Mei, Lindsey, and Bartlett. "Large stepsize GD for non-homogeneous two-layer networks: margin improvement and fast optimization." NeurIPS 2024

Cai\*, Zhou\*, **Wu**, Mei, Lindsey, and Bartlett. "Implicit bias of gradient descent for nonhomogeneous deep networks." ICML 2025

## Contribution

Provable unstable convergence in three cases

Next: a general theory?





theory >

divergent unstable

convergent

stable convergent

 $\eta = o(1)$ , gradient flow

 $\eta = \infty$