

Large Stepsizes Accelerate GD for Regularized Logistic Regression

Jingfeng Wu

with Pierre Marion and Peter Bartlett



Gradient descent

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

“GD \approx discrete time gradient flow”

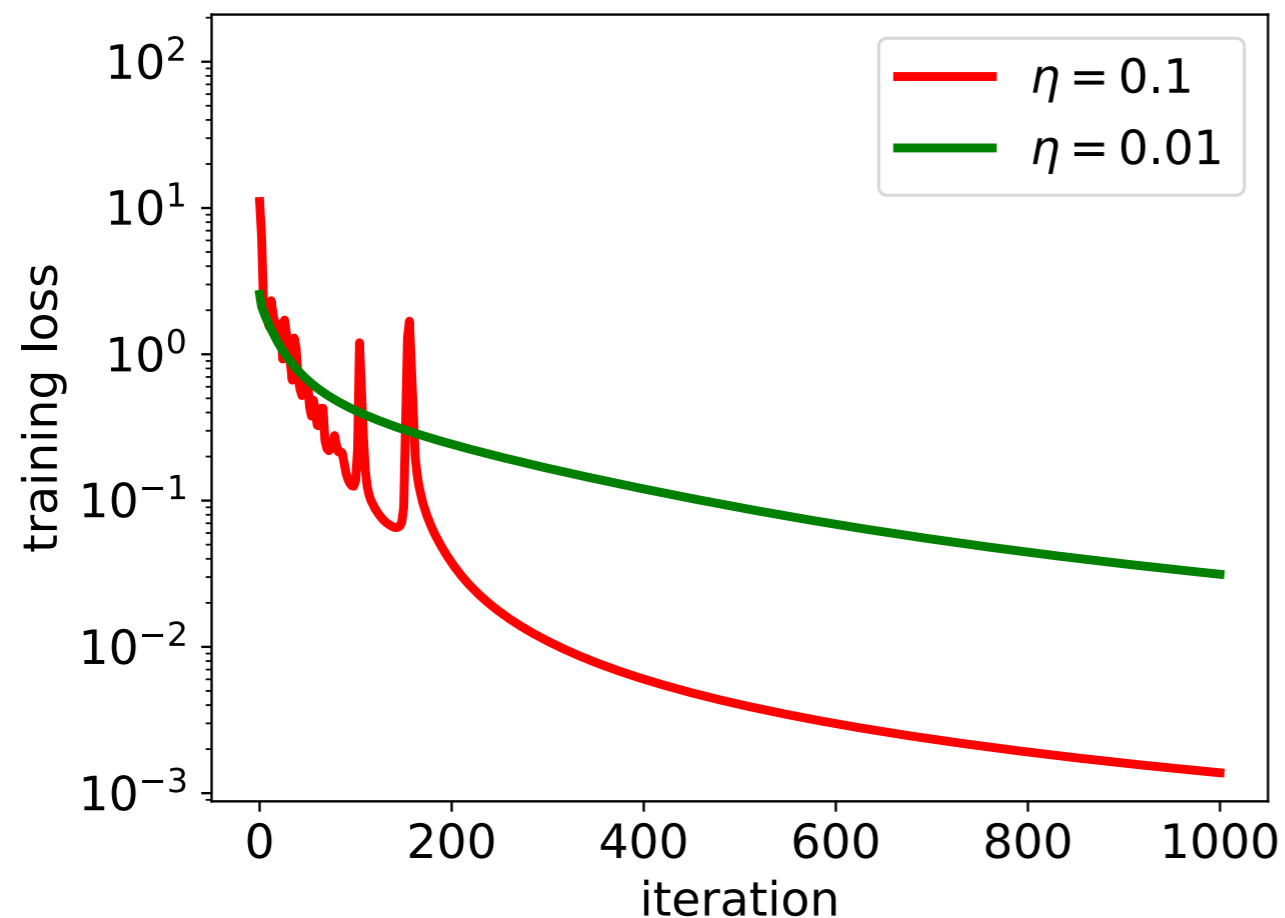
$$\begin{aligned} dw = -\nabla L(w)dt &\Rightarrow dL(w) = \nabla L(w)^\top dw \\ &= -\|\nabla L(w)\|^2 dt \\ &\Rightarrow L(w) \downarrow \end{aligned}$$

small $\eta \Rightarrow$ convex optimization theory



Cauchy, 1847

Experiment (3-layer net, MNIST)



large stepsize is

- unstable
- but faster

this work analyzes benefits of large stepsize

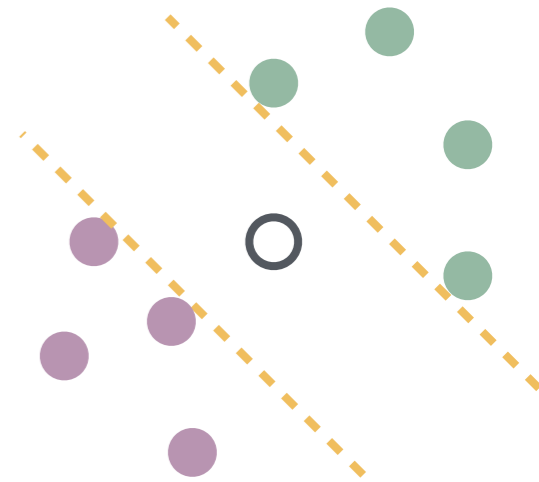
Regularized logistic regression

$$\tilde{L}(w) = L(w) + \frac{\lambda}{2} \|w\|^2 \quad L(w) = \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i x_i^\top w))$$

$$w_{t+1} = w_t - \eta \nabla \tilde{L}(w_t)$$

Assumption (bounded + separable)

- $\|x_i\| \leq 1, y_i \in \{\pm 1\}, i = 1, \dots, n$
- \exists unit vector $w^*, \min_i y_i x_i^\top w^* \geq \Theta(1)$



“almost surely” separable when overparameterized

Classical theory prediction

- $\Theta(1)$ -smooth
- λ -strongly convex
- condition number $\kappa = \Theta(1/\lambda)$
- finite minimizer w_λ , $\|w_\lambda\| = O(\ln(1/\lambda))$

Classical theory

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

$$\tilde{L}(w_t) - \min \tilde{L} \leq \epsilon \text{ for } t = \tilde{O}(1/\lambda)$$

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

Theorem (small λ)

$$\eta_{\max} = \Theta(1/\lambda^{1/2})$$

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 unstable phase in τ steps for

$$\tau := \max\{\eta, n, n/\eta \ln(n/\eta)\} \quad \tau = \Theta(1/\lambda^{1/2})$$

Stable phase. From τ and onward

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

$$t = \Theta(\ln(1/\epsilon)/\lambda^{1/2})$$

for small λ , large stepsize GD matches Nesterov

Theorem (general λ)

Assume separability and

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

$$\eta_{\max} = \Theta(1/\lambda^{1/3})$$

Phase transition. GD exists unstable phase in τ steps for

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

$$\tau = \Theta(1/\lambda^{2/3})$$

Stable phase. From τ and onward

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

$$t = \Theta(\ln(1/\epsilon)/\lambda^{2/3})$$

for general λ , large stepsize is faster than small stepsize

$\tilde{O}(1/\lambda^{2/3})$

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

Margin-based generalization

Assume $(x_i, y_i)_{i=1}^n$ are iid copies of (x, y) , where a.s.

- $\|x\| \leq 1, y \in \{\pm 1\}$
- \exists unit vector $w^*, yx^\top w^* \geq \Theta(1)$

Corollary. The best known test error upper bound is $\tilde{O}(1/n)$. To get $\tilde{O}(1/n)$ rate, GD takes

- $O(n)$ steps with $\lambda = 0$ and $\eta = \Theta(1)$
- $O(n)$ steps with $\lambda = 1/n$ and $\eta = 1$
- $\tilde{O}(n^{2/3})$ steps with $\lambda = 1/n$ and $\eta = \Theta(n^{1/3})$

large stepsize accelerates GD without overfitting

Stepsize diagram

