# Large Stepsizes Accelerate Gradient Descent for Regularized Logistic Regression
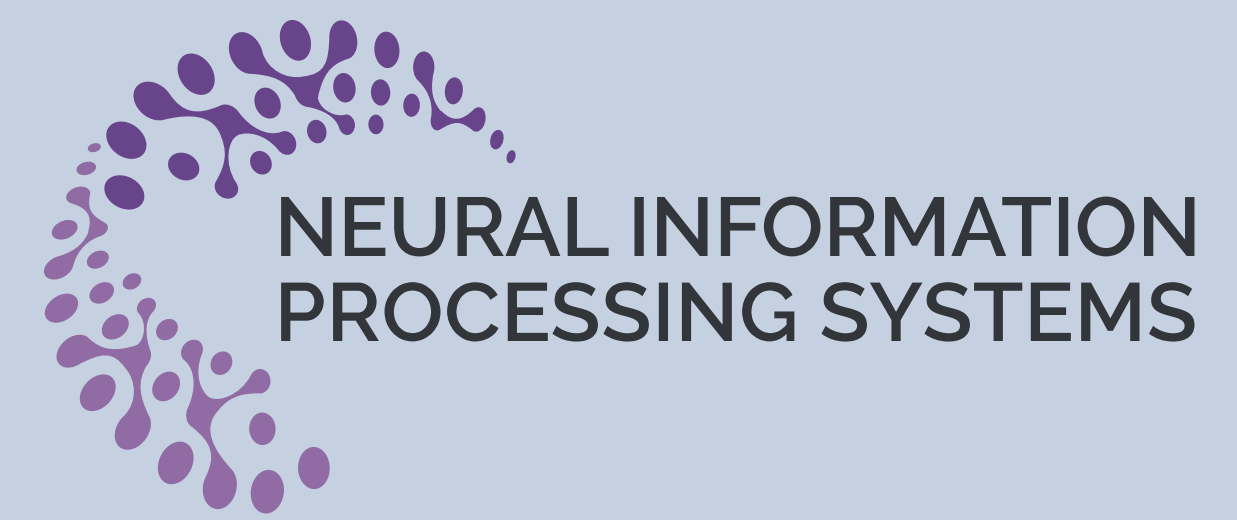
Jingfeng Wu[*1]   Pierre Marion[*3]   Peter Bartlett[12]

[1]UC Berkeley   [2]Google DeepMind
[3]Inria, DI ENS, PSL University

NEURAL INFORMATION PROCESSING SYSTEMS

## Background

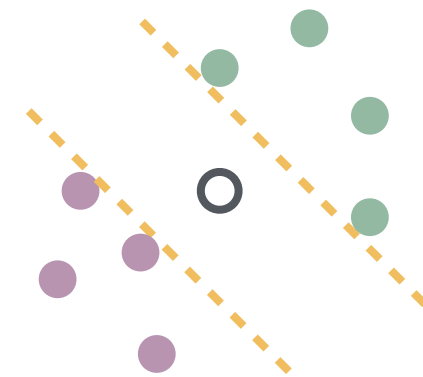$$\tilde{L}(w) = L(w) + \frac{\lambda}{2}\|w\|^2 \quad L(w) = \frac{1}{n}\sum_{i=1}^{n}\ln\big(1 + e^{-y_i x_i^\top w}\big)$$

[Assumption (bounded + separable)]

- $\|x_i\| \leq 1$, $y_i \in \{\pm 1\}$, $i = 1,\ldots,n$

- $\exists$ unit vector $w^*$, $\min\limits_{i} y_i x_i^\top w^* \geq \gamma = \Theta(1)$

> typical case when overparameterized
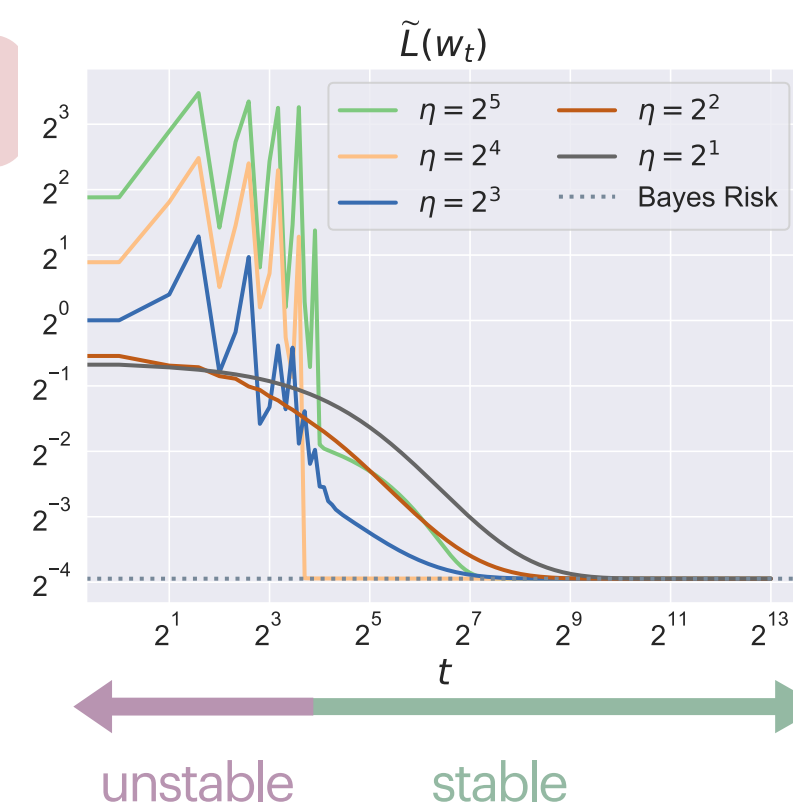
[Basic properties]

- objective $\tilde{L}$ is $\Theta(1)$-smooth and $\lambda$-strongly convex

- condition number is $\kappa = \Theta(1/\lambda)$

- minimizer $w_\lambda = \arg\min\tilde{L}(\,\cdot\,)$ is unique, $\|w_\lambda\| = \Theta(\ln(1/\lambda))$

### Gradient descent

$$w_{t+1} = w_t - \eta\,\nabla\tilde{L}(w_t) \quad w_0 = 0$$

large stepsize is
- unstable
- but faster



$\tilde{L}(w_t)$

unstable — stable

### Prior results

[Without regularization, $\lambda = 0$; Wu, Bartlett, Telgarsky, Yu, 2024]

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

- For $T = \Omega(n)$ and $\eta = \Theta(T)$, we have $L(w_T) \leq \tilde{O}(1/T^2)$

[With regularization; classical optimization theory]

For $\eta = \Theta(1)$, we have $\tilde{L}(w_t) \downarrow$ and $\tilde{L}(w_t) - \min\tilde{L} \leq \epsilon$ for

$$t = O(\ln(1/\epsilon)/\lambda) = \tilde{O}(1/\lambda)$$

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

## Acceleration via large stepsizes

### Small regularization

[Theorem]   for small $\lambda$, large stepsize improves step complexity to $\tilde{O}(1/\lambda^{1/2})$

Assume separability and

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

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

**Unstable phase.** GD is *unstable* for at most $\tau$ steps for

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

> $\tau = \Theta(\eta) \leq \Theta(1/\lambda^{1/2})$

**Stable phase.** From $\tau$ and onward, $\tilde{L}(w_{\tau+t}) \downarrow$ and

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

### General regularization

[Theorem]   large stepsize improves step complexity to $\tilde{O}(1/\lambda^{2/3})$

Assume separability and

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

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

**Unstable phase.** GD is *unstable* for at most $\tau$ steps for

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

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

**Stable phase.** From $\tau$ and onward, $\tilde{L}(w_{\tau+t}) \downarrow$ and

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

### A lower bound

[Theorem]   small stepsize cannot accelerate

Fix $0 < \gamma < 0.1$ and consider a separable dataset

$$x_1 = (\gamma, 0.9) \quad x_2 = (\gamma, -0.5) \quad y_1 = y_2 = 1$$

For all $\lambda \lesssim 1$ and $\epsilon \lesssim \lambda\ln^2(1/\lambda)$, if $\eta$ is such that $\tilde{L}(w_t) \downarrow$ for $t \geq 0$,

then $\tilde{L}(w_t) - \min\tilde{L} \leq \epsilon \;\Rightarrow\; t = \Omega\left(\dfrac{\ln(1/\epsilon)/\ln^2(1/\lambda)}{\lambda}\right)$

## Acceleration without overfitting

[Assumption]  Let $(x_i, y_i)_{i=1}^{n}$ be iid copies of $(x, y)$, where a.s.

- $\|x\| \leq 1$, $y \in \{\pm 1\}$

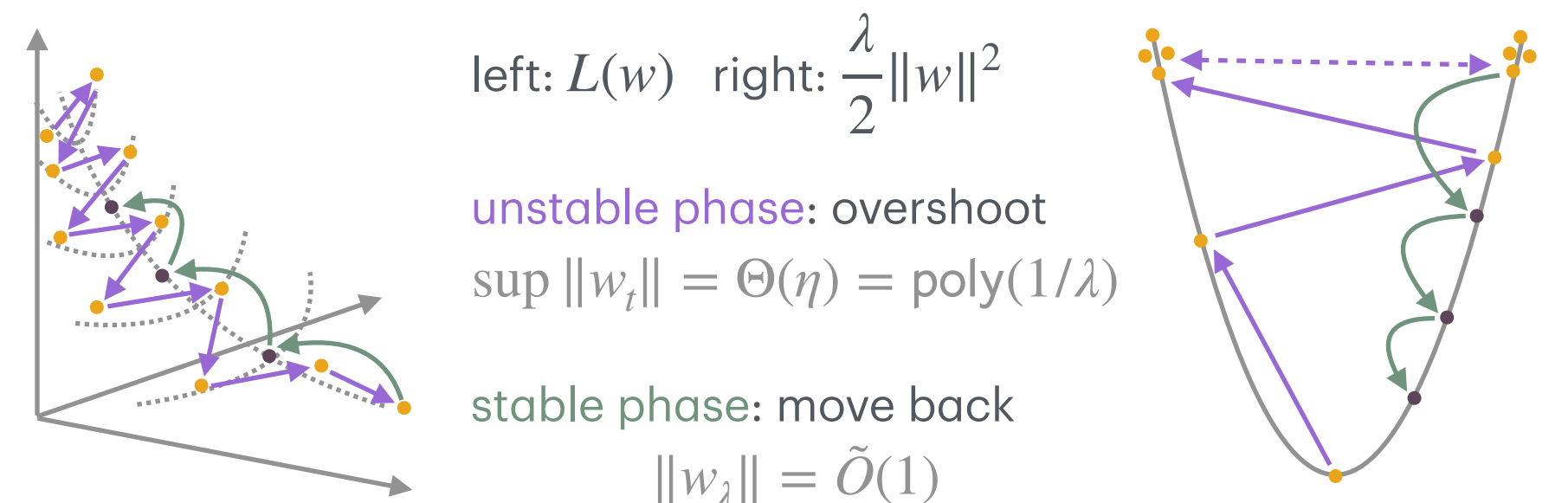- $\exists$ unit vector $w^*$, $yx^\top w^* \geq \gamma = \Theta(1)$

[Classical fast rate]

> data fitting vs estimator norm

For the test error of any estimator $\hat{w}$, w.h.p.

$$L_{\text{test}}(\hat{w}) := \mathbb{E}\ln(1 + e^{-yx^\top\hat{w}}) \lesssim L(\hat{w}) + \tilde{O}(1)\frac{\max\{1, \|\hat{w}\|^2\}}{n}$$

| algorithms | $\lambda$ | $\eta$ | #steps to get 1/n test error |
|---|---|---|---|
| GD | $0$ | $\Theta(1)$ | $O(n)$ |
| | $1/n$ | $1$ | $\tilde{O}(n)$ |
| | $1/n$ | $\Theta(n^{1/3})$ | $\tilde{O}(n^{2/3})$ |
| Nesterov | $1/n$ | $1$ | $\tilde{O}(n^{1/2})$ |



left: $L(w)$   right: $\dfrac{\lambda}{2}\|w\|^2$

unstable phase: overshoot
$\sup\|w_t\| = \Theta(\eta) = \text{poly}(1/\lambda)$

stable phase: move back
$\|w_\lambda\| = \tilde{O}(1)$

## Contribution & open problems



$\eta = \infty$

$\lambda^{-1}$ — divergent

$(\lambda\ln(1/\lambda))^{-1}$ — locally convergent — generic support vectors — unknown global behavior

$\lambda^{-1/2}$ — unstable convergent — match Nesterov

$\lambda^{-1/3}$ — sample size independent

$1$ — stable convergent

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