### **Implicit Bias of Gradient Descent** Jingfeng Wu<sup>1</sup>, Vladimir Braverman<sup>2</sup>, Jason Lee<sup>3</sup> <sup>1</sup>Johns Hopkins University (=> UC Berkeley), for Logistic Regression at the Edge of Stability <sup>2</sup>Rice University, <sup>3</sup>Princeton University



Edge of Stability

# small stepsize works; large stepsize also works

## non-monotonicity against descent lemma



3-layer net + 1,000 samples from MNIST+ GD with const-stepsize

D. risk is bounded by

C.  $L(v_t, \bar{v}_t) \to \infty$ 

## Theory

• binary classification data  $(x_i, y_i = 1)_{i=1}^n$ 

• logistic loss + linear model

$$L(w) := \sum_{i} \log \left( 1 + \exp \left( - \langle w, x_i \rangle \right) \right)$$

constant-stepsize GD

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

• Assumption 1.  $\exists w, \langle w, x_i \rangle > 0, i = 1, ..., n$ 

## Convergence at EoS

For **every** constant stepsize  $\eta > 0$ :

A. in the max-margin subspace,

$$\mathscr{P} \circ w_t \ge \frac{1}{\gamma} \cdot \log(t) + \Theta(1)$$

B. in the non-separable subspace,

 $\|\overline{\mathscr{P}} \circ w_t\|_2 \leq \Theta(1)$ 

C. moreover,

$$G\left(\overline{\mathscr{P}} \circ w_t\right) - \min G(\cdot) \leq \frac{\Theta(1)}{\log(t)}, \text{ where } G(v) := \sum_{x \in \text{supp.}} \exp\left(-\left\langle \overline{\mathscr{P}} \circ x, v\right\rangle\right)$$

$$L(w_t) \le \frac{\Theta(1)}{t}$$

## Negative Results for Exp Loss

Consider exp loss on two 2D samples

$$D = \sum_{i} e^{-\langle w, x_i \rangle} \iff L(v, \bar{v}) = e^{-\gamma v - \bar{v}} + e^{-\gamma v + \bar{v}}$$
  
$$\leq v_0 \leq 2, \quad |\bar{v}_0| \geq 1, \quad 0 < \gamma < 1/4, \quad \eta \geq 4, \text{ then}$$
  
$$t \to \infty$$
  
$$\bar{v}_t | > 2\gamma v_t \text{ and } \bar{v}_t \text{ flips sign every iteration}$$



 $x_2 = (\gamma, -1)$ 



P	Dual form: $\hat{w} = \alpha_1 \cdot x_1 + \ldots + \alpha_s \cdot x_s$
	Orthogonal: $0 = \langle v, \hat{w} \rangle$
	$0 = \alpha_1 \cdot \langle v, x_1 \rangle + \ldots + \alpha_s \cdot \langle v, x_s \rangle$
	<b>Assumption 2.</b> $\alpha_i > 0$ for supp. vectors
	Assumption 3. supp. vectors span the space
ə 死	<b>Lemma.</b> There must exist $\langle v, x_i \rangle > 0$ , $\langle v, x_j \rangle < 0$

# • benefits of logisitic loss

1. Soudry, D., Hoffer, E., Nacson, M. S., Gunasekar, S., & Srebro, N. (2018). The implicit bias of gradient descent on separable data. JMLR 2019. 2. Ji, Z., & Telgarsky, M. (2018). Risk and parameter convergence of logistic regression. COLT 2019.

3. Cohen, J. M., Kaur, S., Li, Y., Kolter, J. Z., & Talwalkar, A. (2021). Gradient descent on neural networks typically occurs at the edge of stability. ICLR 2021