## Implicit Bias of Gradient Descent for Logistic Regression at the *Edge of Stability*

**Jingfeng Wu** (JHU -> Berkeley), 05/2023 Joint work with Vladimir Braverman (Rice) and Jason D. Lee (Princeton)





### What happens in OPT for DL



Risk oscillation caused by large stepsize



4-layer fully connected net 1,000 samples from MNIST gradient descent (GD)



### **Good hyperparameter is often large**



Training imagenet in 1 hour. arXiv preprint arXiv:1706.02677.

Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y. and He, K., 2017. Accurate, large minibatch sgd:



## **Edge of stability (EoS)** $\eta > 2 / ||\nabla^2 L(w)||_2$



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. *arXiv preprint arXiv:2103.00065*.

# Risks converge while oscillating by chance or by math?

#### **Remove "distractions"**

10 classes -> 2 classes (y = "0" or "8")



Logistic regression on "linearly separable" data



#### NN -> linear model (w/o bias)



#### A even more minimal example

#### Logistic regression on four 2D samples



gd 100 train loss  $10^{-1}$ 20 80 60 100 40 ()# steps 0.6 2/eta gd, sharpness 0.5 4.0 sparbness0.30.2 0.1 0.0 60 80 100 20 40 0 # steps



#### **Problem setup** Logistic regression on separable data

• Samples 
$$(x_i, y_i = 1)_{i=1}^n$$

• Risk

$$L(w) := \sum_{i} \log\left(1 + \exp\left(\frac{1}{i}\right)\right)$$

Constant-stepsize GD

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

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

 $\left(-\langle w, x_i \rangle\right)$ 

 $(W_{t-1})$ 



What we know so far Logistic regression on separable data If the stepsize is small (e.g.,  $\eta = 0.01$ )

**Risk minimization** lacksquare

$$L(w_t) \le \tilde{\mathcal{O}}\left(\frac{1}{t}\right)$$
 by des

Margin maximization  $\bullet$ 

$$\left|\frac{w_t}{\|w_t\|_2} - \frac{\hat{w}}{\|\hat{w}\|_2}\right| \leq \tilde{\mathcal{O}}\left(\frac{1}{\log(t)}\right)$$

• Soudry, D., Hoffer, E., Nacson, M. S., Gunasekar, S., & Srebro, N. (2018). The implicit bias of gradient descent on separable data. The Journal of Machine Learning Research, 19(1), 2822-2878. • Ji, Z., & Telgarsky, M. (2018). Risk and parameter convergence of logistic regression. arXiv preprint arXiv:1803.07300.



#### SVM solution

$$\hat{w} = \arg \min \|w\|_2$$
, s.t.  $\forall i, \langle w, x_i \rangle \geq$ 





# Large stepsize works well trust me bro! tried 3 different seeds



# Our Results

(any) const-stepsize + logistic regression + separable data

### Space decomposition



Dual form:  $\hat{w} = \alpha_1 \cdot x_1 + \ldots + \alpha_s \cdot x_s$ Orthogonal:  $0 = \langle v, \hat{w} \rangle$ 

$$0 = \alpha_1 \cdot \langle v, x_1 \rangle + \ldots + \alpha_s \cdot \langle v, x_s \rangle$$

Assumption 2: supp. vectors span the space Assumption 3:  $\alpha_i > 0$  for supp. vectors

Then there must exist  $\langle v, x_i \rangle > 0$ ,  $\langle v, x_j \rangle < 0$ 



#### **Implicit Bias**

For every constant stepsize  $\eta > 0$ :

A. In the max-margin subspace,

B. In the non-separable subspace, C. In the non-separable subspace,  $\Theta(1)$ 

 $\mathscr{P} \circ w_t \ge \frac{1}{\gamma} \cdot \log(t) + \Theta(1)$  $\left\|\overline{\mathscr{P}} \circ w_t\right\|_{\gamma} \leq \Theta(1)$ strongly convex  $G(\overline{\mathscr{P}} \circ w_t) - \min G(\cdot) \le \frac{\Theta(1)}{\log(t)}, \quad \text{where } G(v) := \sum_{v \in \text{current}} \exp\left(-\left\langle \overline{\mathscr{P}} \circ x, v\right\rangle\right)$ *x*∈supp.



#### **Risk minimization**

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

D. Risk is bounded by



theory based!

*large stepsize => risk still converges* 

possibly non-monotonically



Feel free to use large stepsizes?

### GD can diverge under exp loss

Consider exp loss on two 2D samples

$$L(w) = \sum_{i} e^{-\langle w, x_i \rangle} \quad \langle = \rangle \quad L(v, \bar{v})$$

Assume that

 $0 \le v_0 \le 2$ ,  $|\bar{v}_0| \ge 1$ ,  $0 < \gamma < 1/4$ ,  $\eta \ge 4$ . Then

A.  $v_t \rightarrow \infty$ 

B.  $|\bar{v}_t| > 2\gamma v_t$  and  $\bar{v}_t$  flips sign every iteration

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



 $EoS \neq small-stepsize$ 

logistic loss > exp loss

# Feel free to use large stepsizes under logistic loss!

# Techniques Overview



## A new approach for handling EoS

- Existing approaches
- 1. Show risk convergence (by descent lemma)
- 2. Show iterate limiting behaviors

#### descent lemma is broken in EoS

Journal of Machine Learning Research, 19(1), 2822-2878.

Our approach

- 1. Study iterate limiting behaviors
- 2. Show risk convergence

(by iterate limits)

• Soudry, D., Hoffer, E., Nacson, M. S., Gunasekar, S., & Srebro, N. (2018). The implicit bias of gradient descent on separable data. The • Ji, Z., & Telgarsky, M. (2018). Risk and parameter convergence of logistic regression. arXiv preprint arXiv:1803.07300.



#### Work out the two samples

 $L(v, \bar{v}) = \log(1 + e^{-\gamma v - \bar{v}}) + \log(1 + e^{-\gamma v + \bar{v}})$ 

 $v_{t+1} = v_t - \eta \cdot g_t \qquad \quad \bar{v}_{t+1} = \bar{v}_t - \eta \cdot \bar{g}_t$ 

For simplicity, assume  $v_0 = 0$ ,  $|\bar{v}_0| > 0$ .



#### **Step 1:** $(\bar{v}_t)_{t>0}$ are uniformly bounded $|\bar{v}_{t+1}| = ||\bar{v}_t| - \eta \cdot |\bar{g}_t||$ $\leq \max\{|\bar{v}_t|, \eta \cdot |\bar{g}_t|\}$ $\bar{v}_{t}$ and $\bar{g}_{t}$ share the same sign $\leq \max\{|\bar{v}_t|, \eta\}$ $\bar{g}_{t}$ is bounded $\leq \dots$ $\leq \max\{ |\bar{v}_0|, \eta \}$ induction **Step 2:** $g_t \approx -\gamma \cdot e^{-\gamma v_t} \cdot \Theta(1)$ , so $v_t \approx \log(t)/\gamma =>$ margin gets max-ed

**Step 3:**  $\bar{g}_t \approx e^{-\gamma v_t} \cdot \nabla G(\bar{v}_t)$ , so

$$\bar{g}_t := -\left(\frac{1}{1 + e^{\gamma v_t + \bar{v}_t}} - \frac{1}{1 + e^{\gamma v_t}}\right)$$

#### $\bar{v}_{t+1} \approx \bar{v}_t - \eta_t \cdot \nabla G(\bar{v}_t)$ , where $\eta_t = \eta \cdot e^{-\gamma v_t} \approx \Theta(1)/t$





#### Step 4: a modified descent lemma

#### $G(\bar{v}_{t+1}) \leq G(\bar{v}_t) + \langle \nabla G(\bar{v}_t), \bar{v}_{t+1} - \langle \nabla G(\bar{v}_t), \bar{v}_{t+1} -$

 $\approx G(\bar{v}_t) - \eta_t \cdot \|\nabla G(\bar{v}_t)\|_2^2 +$  $\leq G(\bar{v}_t) + \Theta(1) \cdot \eta_t^2 \cdot \|\nabla G$  $\leq G(\bar{v}_t) + \frac{1}{t^2} \cdot \Theta(1)$ 

So for  $T \ge t \ge 1$ ,  $G(\bar{v}_T) \le G(\bar{v}_t) + \frac{1}{t} \cdot \Theta(1)$  For small-stepsize,  $G(\bar{v}_T) < G(\bar{v}_t)$ 

$$\begin{aligned} & -\bar{v}_t \rangle + \frac{\beta}{2} \cdot \|\bar{v}_{t+1} - \bar{v}_t\|_2^2 \\ & + \frac{\beta}{2} \cdot \eta_t^2 \cdot \|\nabla G(\bar{v}_t)\|_2^2 \quad (+\text{higher orders}) \\ & \overline{G}(\bar{v}_t)\|_2^2 \end{aligned}$$

The "increase" of risk must decrease



$$\begin{aligned} & \text{Step 5: convergence of } \bar{v}_t \\ \|\bar{v}_{t+1} - \bar{v}_*\|_2^2 &= \|\bar{v}_t - \bar{v}_*\|_2^2 + 2 \cdot \langle \bar{v}_t - \bar{v}_*, \bar{v}_{t+1} - \bar{v}_t \rangle + \|\bar{v}_{t+1} - \bar{v}_t\|_2^2 \\ &\approx \|\bar{v}_t - \bar{v}_*\|_2^2 - 2\eta_t \cdot \langle \bar{v}_t - \bar{v}_*, \nabla G(\bar{v}_t) \rangle + \eta_t^2 \cdot \|\nabla G(\bar{v}_t)\|_2^2 \\ &\leq \|\bar{v}_t - \bar{v}_*\|_2^2 - 2\eta_t \cdot \left(G(\bar{v}_t) - G(\bar{v}_*)\right) + \frac{1}{t^2} \cdot \Theta(1) \\ &\sum_{t=t_0}^T \eta_t \cdot \left(G(\bar{v}_t) - G(\bar{v}_*)\right) \leq \Theta(1) + \sum_{t=t_0}^T \frac{1}{t^2} \cdot \Theta(1) \leq \Theta(1) \\ &G(\bar{v}_T) - G(\bar{v}_*) \leq \frac{\sum_{t=t_0}^T \eta_t \cdot \left(G(\bar{v}_t) - G(\bar{v}_*)\right) + \sum_{t=t_0}^T \eta_t \cdot \frac{\Theta(1)}{t}}{\sum_{t=t_0}^T \eta_t} \leq \frac{\Theta(1)}{\log(T)} \end{aligned}$$

# $\|_{2}^{2}$

$$\sum_{t=t_0} \eta_t$$



#### Take away

- Convergence/implicit bias theory still holds in the EoS regime
- However, only for logistic loss not for exponential loss
- New analysis ideas

### What's next?

- What happens in poly time?
- Population risk?
- Non-linear model? Other Algos?...

