

Finite-Sample Analysis of Learning High-Dim ReLU Neuron

Jingfeng Wu

with Difan Zou, Zixiang Chen, Vladimir Braverman, Quanquan Gu, Sham Kakade

High-dim ~~ReLU~~ regression

Linear

$$\text{Minimize } R(\mathbf{w}) = \mathbb{E}(\text{ReLU}(\mathbf{x}^\top \mathbf{w}) - y)^2, \quad \mathbf{w} \in \mathbb{R}^d$$

With n samples (iid): $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$

Distributional assumption (simplified)

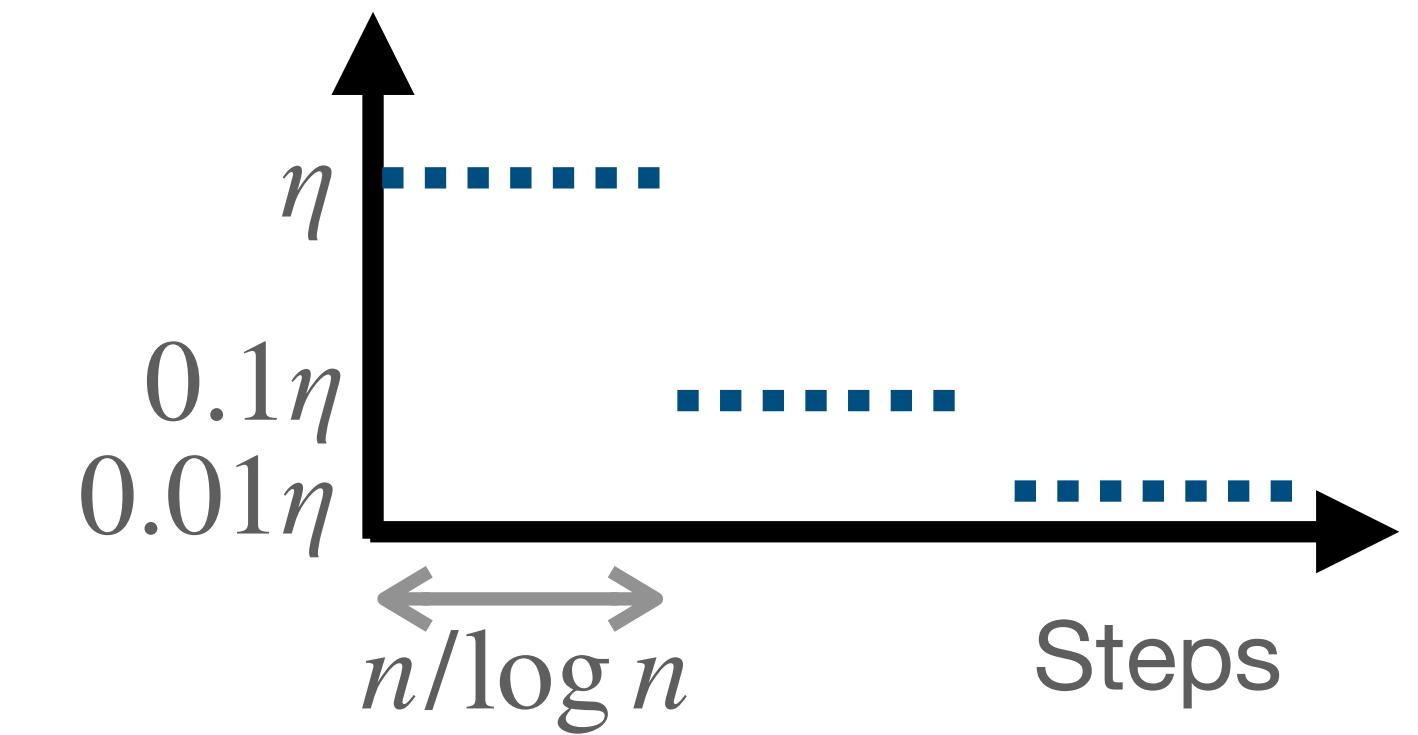
$$y = \text{ReLU}(\mathbf{x}^\top \mathbf{w}_*) + \mathcal{N}(0, 1), \quad \|\mathbf{w}_*\|_2 \leq 1, \quad \mathbf{x} \sim \mathcal{N}(0, \mathbf{H})$$



$(\lambda_i)_{i \geq 1}$ denote eigenvalues of $\mathbf{H} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$

Learning high-dim linear regression by (online) SGD

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \cdot (\mathbf{x}_t^\top \mathbf{w}_{t-1} - y_t) \cdot \mathbf{x}_t, \quad t = 1, \dots, n$$



[ZWBGK'21, WZBGK'22]

$$\frac{\text{DIM}}{n} \lesssim \mathbb{E}R(\mathbf{w}_n) - \text{OPT} \lesssim \frac{\text{DIM}}{n}$$

$$\text{DIM} := \# \left\{ \lambda_i \geq \frac{1}{n\eta_0} \right\} + n^2\eta_0^2 \cdot \sum_{\lambda_i < \frac{1}{n\eta_0}} \lambda_i^2 \ll d$$

Learning high-dim ReLU regression

by (online) SGD

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \cdot (\text{ReLU}(\mathbf{x}_t^\top \mathbf{w}_{t-1}) - y_t) \cdot \mathbf{x}_t \cdot 1_{[\mathbf{x}_t^\top \mathbf{w}_{t-1} > 0]}, \quad t = 1, \dots, n$$

not a good idea

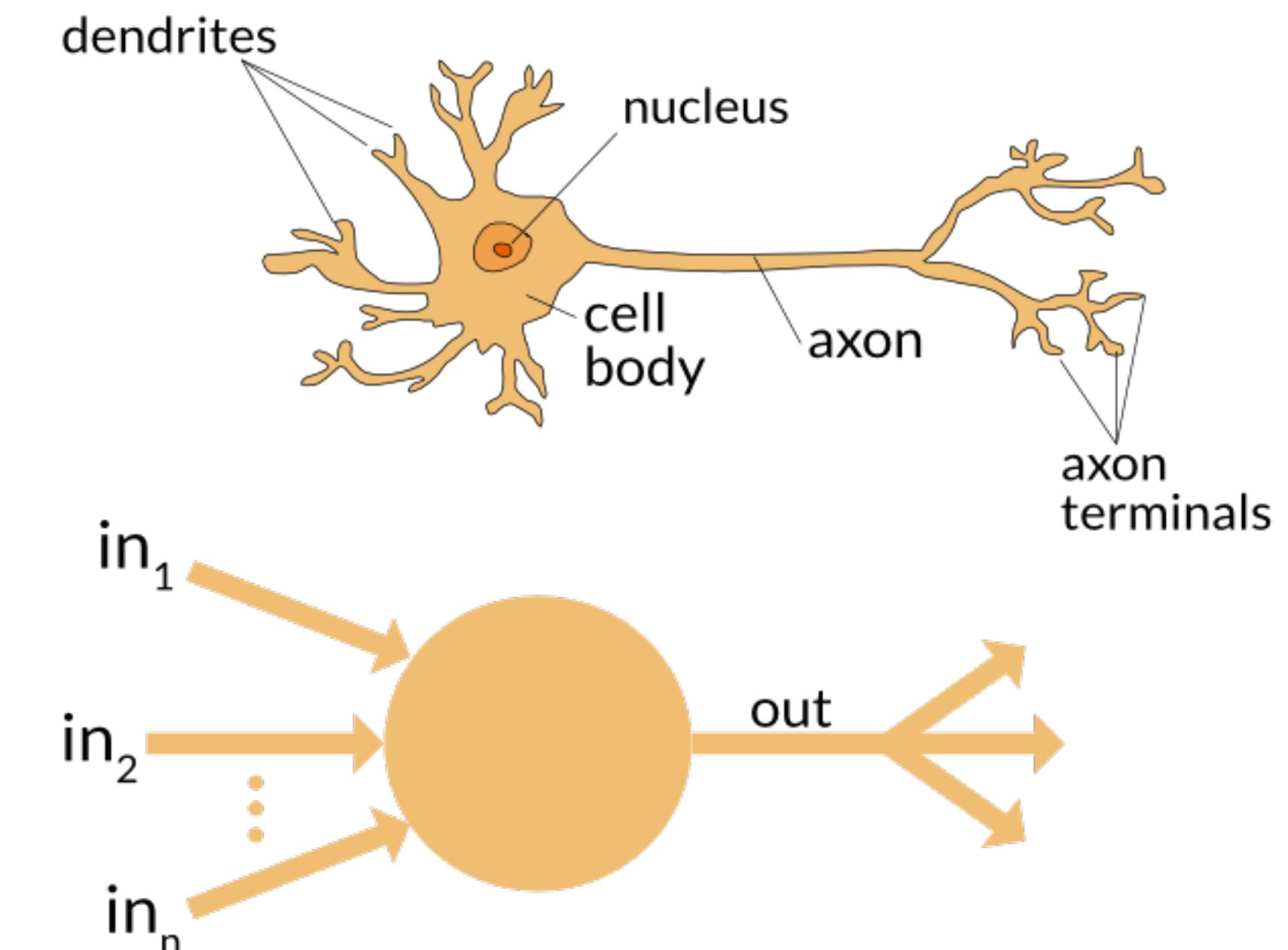


by (online) Perceptron

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \cdot (\text{ReLU}(\mathbf{x}_t^\top \mathbf{w}_{t-1}) - y_t) \cdot \mathbf{x}_t,$$

$$t = 1, \dots, n$$

our main focus



Result 1: well-specified ReLU regression

$$y = \text{ReLU}(\mathbf{x}^\top \mathbf{w}_*) + \mathcal{N}(0, 1), \quad \|\mathbf{w}_*\|_2 \leq 1, \quad \mathbf{x} \sim \mathcal{N}(0, \mathbf{H})$$

[This work]

$$\frac{\text{DIM}}{n} \lesssim \mathbb{E}R(\mathbf{w}_n) - \text{OPT} \lesssim \frac{\text{DIM}}{n}$$

- Linear reg. bounds port over to ReLU reg.
- Point-wise tight, dim-free, ...
- **ReLU isn't harder!**

misspecified Result 2: well-specified ReLU regression

$$y \neq \text{ReLU}(\mathbf{x}^\top \mathbf{w}_*) + \mathcal{N}(0, 1), \quad \|\mathbf{w}_*\|_2 \leq 1, \quad \mathbf{x} \sim \mathcal{N}(0, \mathbf{H})$$

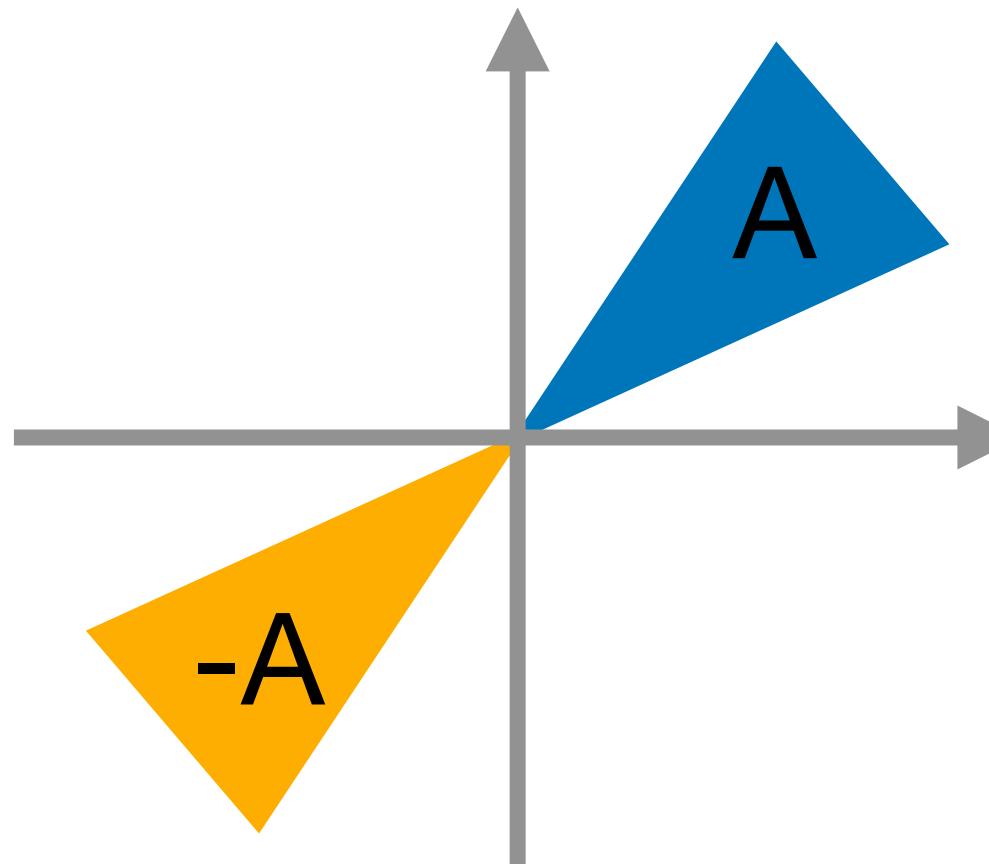
[This work]

$$\mathbb{E}R(\mathbf{w}_n) \lesssim \text{OPT} + \frac{\text{DIM}}{n} \lesssim \mathbb{E}R(\mathbf{w}_n) - \text{OPT} \lesssim \frac{\text{DIM}}{n}$$

- $\mathbb{E}R(\mathbf{w}) \leq \text{OPT} + o(1)$ bound is believed to be “impossible” [GKK’19]
- Previously, best known bound $\mathcal{O}(\text{OPT} + \sqrt{d/n})$ [DGKKS’20]
- A *dim-free*, const. factor approx.

What's more

- In well-specified cases, bounds are ***tight up to const. factor*** as functions of η, n, H, w_* , and noise scale
- $x \sim \mathcal{N}(0, H)$ can be relaxed to
 1. hypercontractivity
 2. moments symmetricity



$$\forall v, \mathbb{E}(x^T v)^4 \leq \alpha (\mathbb{E}(x^T v)^2)^2$$

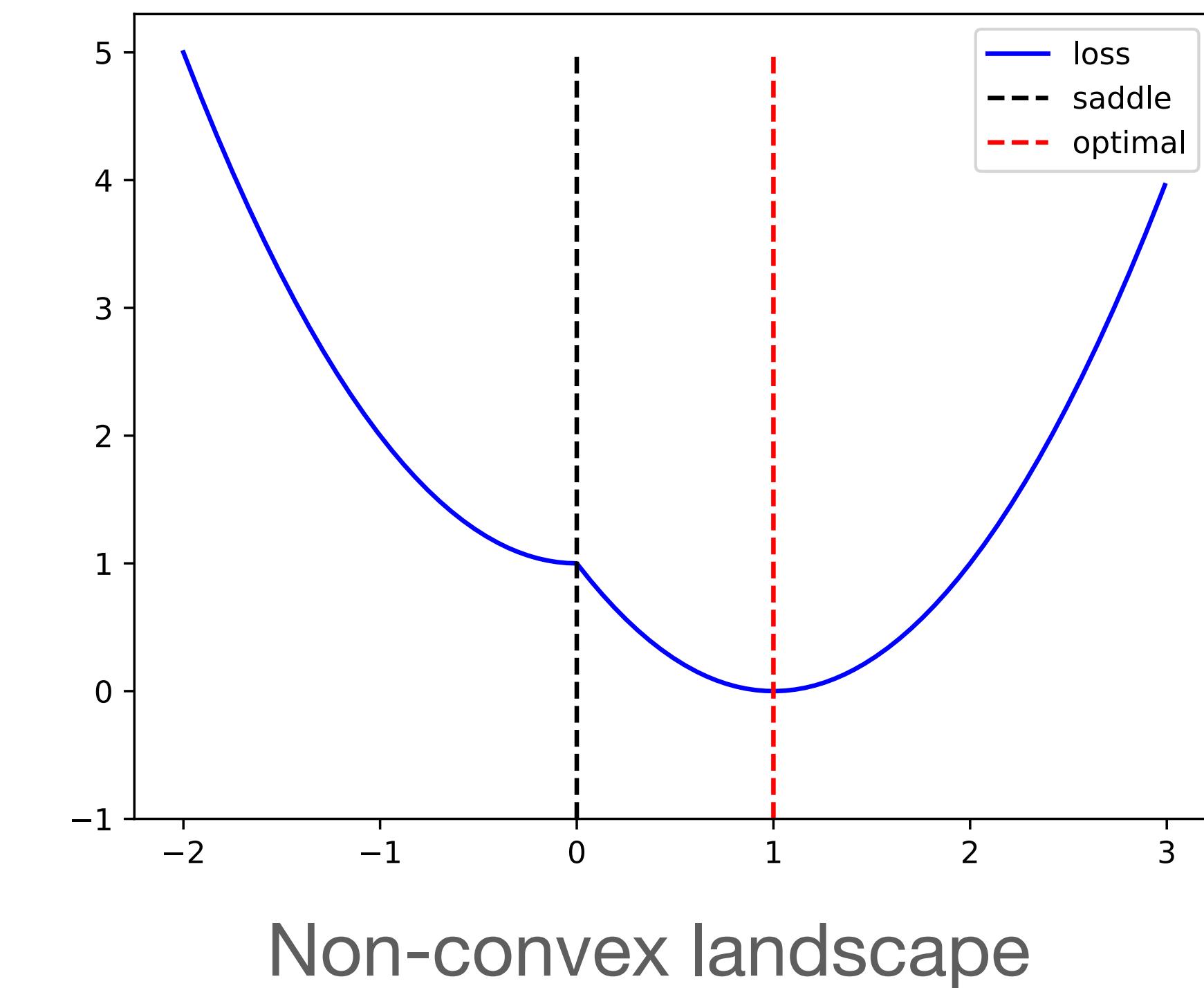
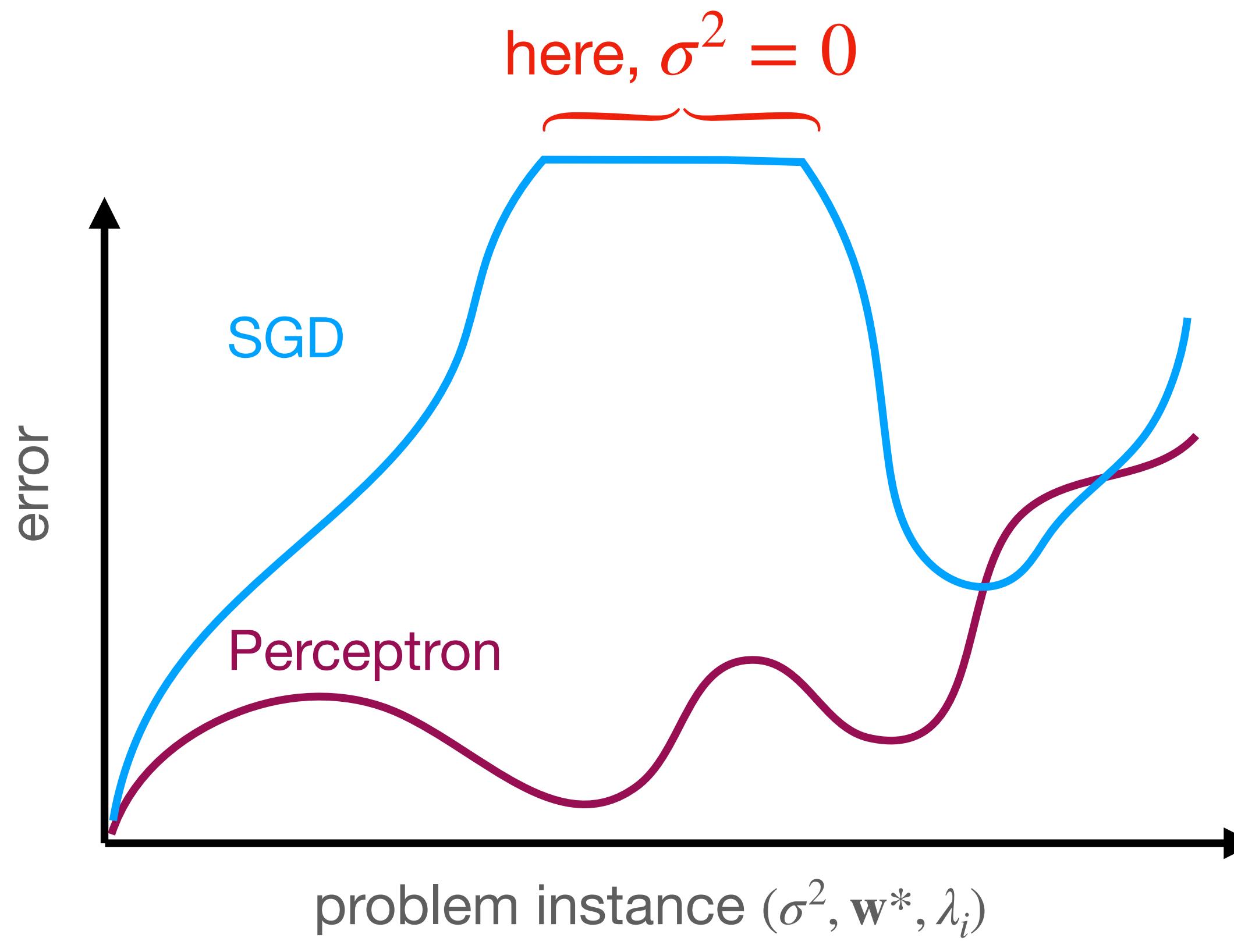
$$\mathbb{E}f(x)\mathbf{1}[x \in A] = \mathbb{E}f(x)\mathbf{1}[x \in -A]$$

$$\text{for } f(x) = x^{\otimes 2} \text{ or } x^{\otimes 4}$$

Result 3: SGD is less effective than Perceptron

$$y = \text{ReLU}(\mathbf{x}^\top \mathbf{w}_*) + \mathcal{N}(0, \sigma^2), \quad \|\mathbf{w}_*\|_2 \leq 1,$$

$$\mathbb{P}\{\mathbf{x} = \mathbf{e}_i\} = \mathbb{P}\{\mathbf{x} = -\mathbf{e}_i\} = \frac{\lambda_i}{2}$$



Take Home: ReLU Neuron

1. Well-specified noise

- ReLU reg. \approx linear reg, $\text{OPT} + \mathcal{O}(\text{DIM}/n)$ bound
- **Perceptron**, allowing $d > n$

2. Misspecified noise

- $\mathcal{O}(\text{OPT} + \text{DIM}/n)$ bound
- **Perceptron**, allowing $d > n$

3. Symmetric Bernoulli data

- SGD no better than Perceptron when noise is well-specified
- SGD suffers const. error if noiseless



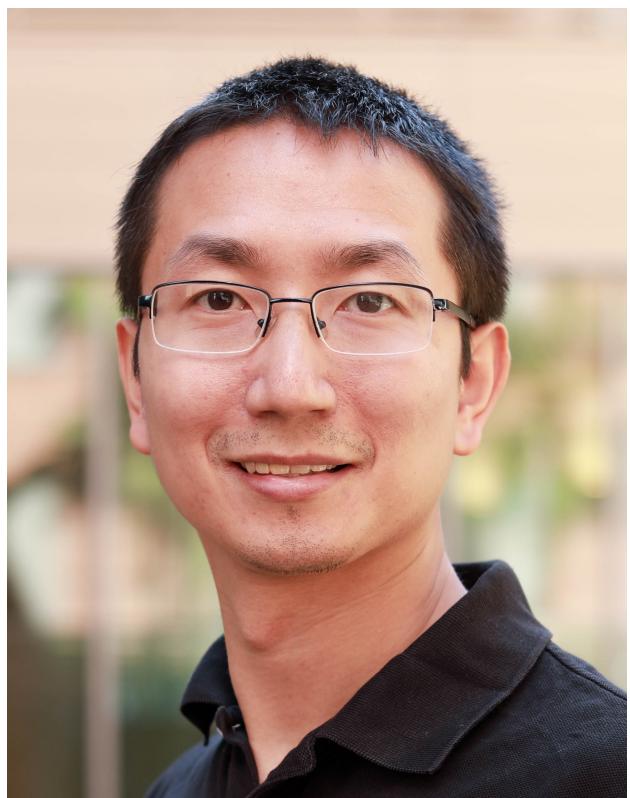
Difan Zou



Zixiang Chen



Vladimir Braverman



Quanquan Gu



Sham Kakade

