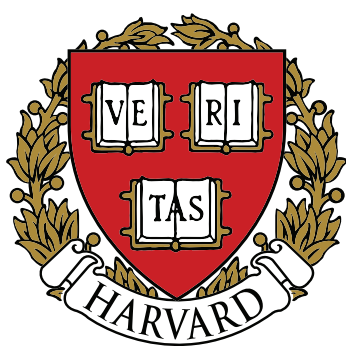


# Last Iterate Risk Bounds of SGD with Decaying Stepsize for Overparameterized Linear Regression

Jingfeng Wu<sup>\*1</sup>, Difan Zou<sup>\*2</sup>, Vladimir Braverman<sup>1</sup>, Quanquan Gu<sup>2</sup>, Sham M. Kakade<sup>3</sup>

<sup>1</sup>Johns Hopkins University, <sup>2</sup>UCLA, <sup>3</sup>Harvard University

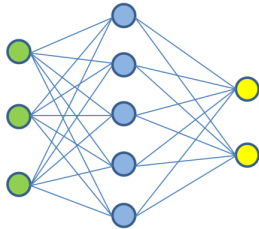


## Population Risk

$$\mathcal{L}(\mathbf{w}) = \mathbb{E}\ell(\mathbf{x}, y; \mathbf{w})$$

## $n$ training samples

$$(\mathbf{x}_1, y_1) \cdots, (\mathbf{x}_n, y_n) \in \mathbb{R}^{d \times 1}$$



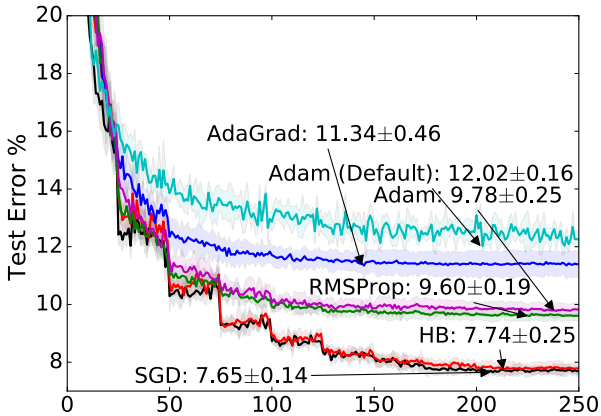
## Large model

$$\mathbf{w} \in \mathbb{R}^d \text{ for large } d$$

## SGD

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \cdot \nabla \ell(\mathbf{x}_i, y_i; \mathbf{w})$$

SGD generalizes well (WRSSR 2017)



## Least Square

True Model  $y = \mathbf{x}^\top \mathbf{w}^* + \mathcal{N}(0, \sigma^2)$

Data Covariance  $\mathbf{H} := \mathbb{E}[\mathbf{x}\mathbf{x}^\top] =: \text{diag}(\lambda_1, \lambda_2, \dots)$ , WOLG

Population Risk  $\mathcal{L}(\mathbf{w}) := \mathbb{E}(y - \mathbf{x}^\top \mathbf{w})^2$

Excess Risk  $\Delta(\mathbf{w}) := \mathcal{L}(\mathbf{w}) - \mathcal{L}(\mathbf{w}^*) = (\mathbf{w} - \mathbf{w}^*)^\top \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$

## [Strongly Contractive Fourth Moment Condition]

Recall that  $\mathbf{H} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$ . Assume that for every PSD matrix  $\mathbf{A}$ ,

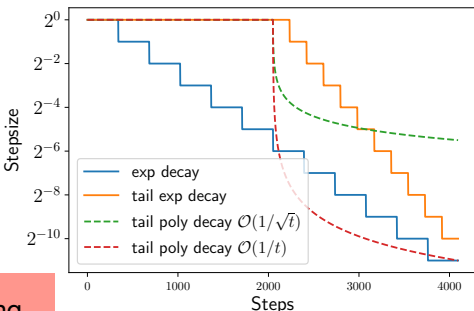
- $\mathbb{E}[\mathbf{x}^\top \mathbf{A} \mathbf{x} \cdot \mathbf{x} \mathbf{x}^\top] \leq \alpha \cdot \text{tr}(\mathbf{H} \mathbf{A}) \cdot \mathbf{H}$  for some constant  $\alpha \geq 1$ ;

- $\mathbb{E}[\mathbf{x}^\top \mathbf{A} \mathbf{x} \cdot \mathbf{x} \mathbf{x}^\top] \geq \beta \cdot \text{tr}(\mathbf{H} \mathbf{A}) \cdot \mathbf{H} + \mathbf{H} \mathbf{A} \mathbf{H}$  for some constant  $\beta > 0$ .

## SGD

$n$  samples  $(\mathbf{x}_t, y_t)_{t=1}^n$

$$\mathbf{w}_t = \mathbf{w}_{t-1} + \eta_t \cdot (y_t - \mathbf{x}_t^\top \mathbf{w}_{t-1}) \cdot \mathbf{x}_t \text{ output} := \mathbf{w}_n$$



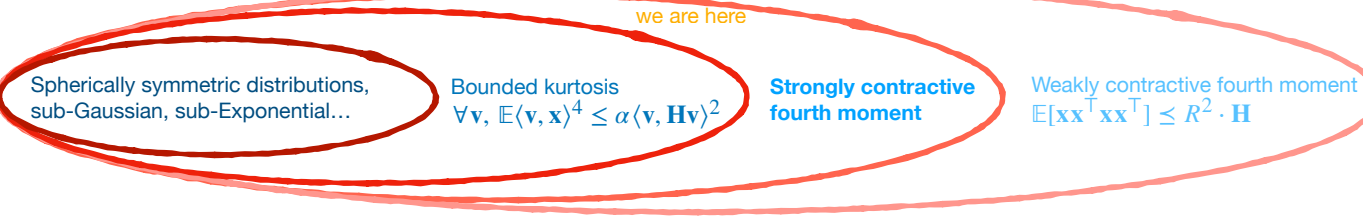
Stepsize Schedulers

Geometrically Decaying

$$\eta_t = \begin{cases} \eta_0, & t \leq s \\ 0.5\eta_{t-1}, & t > s, t \% K = 0 \\ \eta_{t-1}, & \text{otherwise} \end{cases}$$

Polynomially Decaying

$$\eta_t = \begin{cases} \eta_0, & t \leq s \\ \frac{\eta_0}{(t-s)^a}, & t > s \end{cases} \text{ for } 0 \leq a \leq 1$$



## [WZBGK 2022]

Let  $\mathbf{w}_n^{\text{exp}}$  and  $\mathbf{w}_n^{\text{poly}}$  be the SGD outputs with **geometrically** and **polynomially** decaying stepsizes, respectively. Fix same  $s = n/2$ , same  $\mathbf{w}_0$ , same  $\eta_0$ . Then we have

$$\mathbb{E}\Delta(\mathbf{w}_n^{\text{exp}}) \lesssim (1 + \text{SNR} \cdot \log n) \cdot \mathbb{E}\Delta(\mathbf{w}_n^{\text{poly}})$$

where  $\text{SNR} := \|\mathbf{w}_0 - \mathbf{w}^*\|_{\mathbf{H}}^2 / \sigma^2$ .

## [WZBGK 2022]

Consider SGD with **geometrically** decaying stepsizes. Let the stepsize decaying interval be  $K := (n - s) / \log(n - s)$ . For every  $s > 0$ ,  $K > 2$  and every  $\eta_0 < 1 / (4\alpha \text{tr}(\mathbf{H}) \log(n))$ , we have

$$\begin{aligned} \mathbb{E}\Delta(\mathbf{w}_n) \lesssim & \frac{\|(\mathbf{I} - \eta_0 \mathbf{H})^{s+K}(\mathbf{w}_0 - \mathbf{w}^*)\|_{\mathbf{I}_{0:k^*}}^2}{\eta_0 K} + \|(\mathbf{I} - \eta_0 \mathbf{H})^{s+K}(\mathbf{w}_0 - \mathbf{w}^*)\|_{\mathbf{H}_{k^*:\infty}}^2 \\ & + \frac{k^* + \eta_0 K \sum_{k^* < i \leq k^\dagger} \lambda_i + \eta_0^2 K^2 \sum_{i > k^\dagger} \lambda_i^2}{K} \cdot (\sigma^2 + \alpha \cdot \|\mathbf{w}_0 - \mathbf{w}^*\|_{\mathbf{H}}^2 \cdot \log(n)) \end{aligned}$$

Here  $k^*, k^\dagger$  are such that  $\lambda_1 \geq \dots \geq \lambda_{k^*} \geq \frac{1}{\eta_0 K} \geq \lambda_{k^*+1} \geq \dots \geq \lambda_{k^\dagger} \geq \frac{1}{\eta_0(s+K)} \geq \lambda_{k^\dagger+1} \geq \dots$

See the paper for a nearly matching lower bound.

## [GKKN 2019]

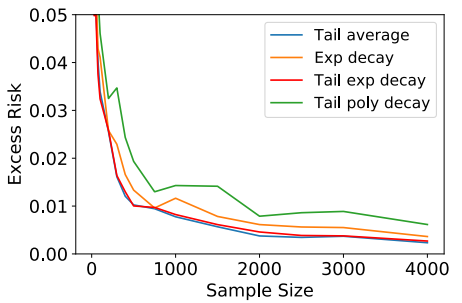
$$\mathbb{E}\Delta(\mathbf{w}_n) \lesssim \left( \frac{d \|\mathbf{w}_0 - \mathbf{w}^*\|_2^2}{\eta_0 n} + \frac{d}{n} \cdot \sigma^2 \right) \cdot \log n$$

## Remarks

- Weakly contractive fourth moment condition
- Variance bound scales with  $d$
- $\ell_2$ -norm implicitly depends on  $d$



Get the Paper!



## Take Home

- Risk of SGD in high-dim  $\approx d_{\text{eff}} / n$
- $d_{\text{eff}}$  determined by  $(\lambda_i)_{i \geq 1}$ ,  $\eta_0$ ,  $n_{\text{eff}}$ , and  $\ll d$  when  $(\lambda_i)_{i \geq 1}$  decay fast
- Geometrical stepsize > polynomially stepsize