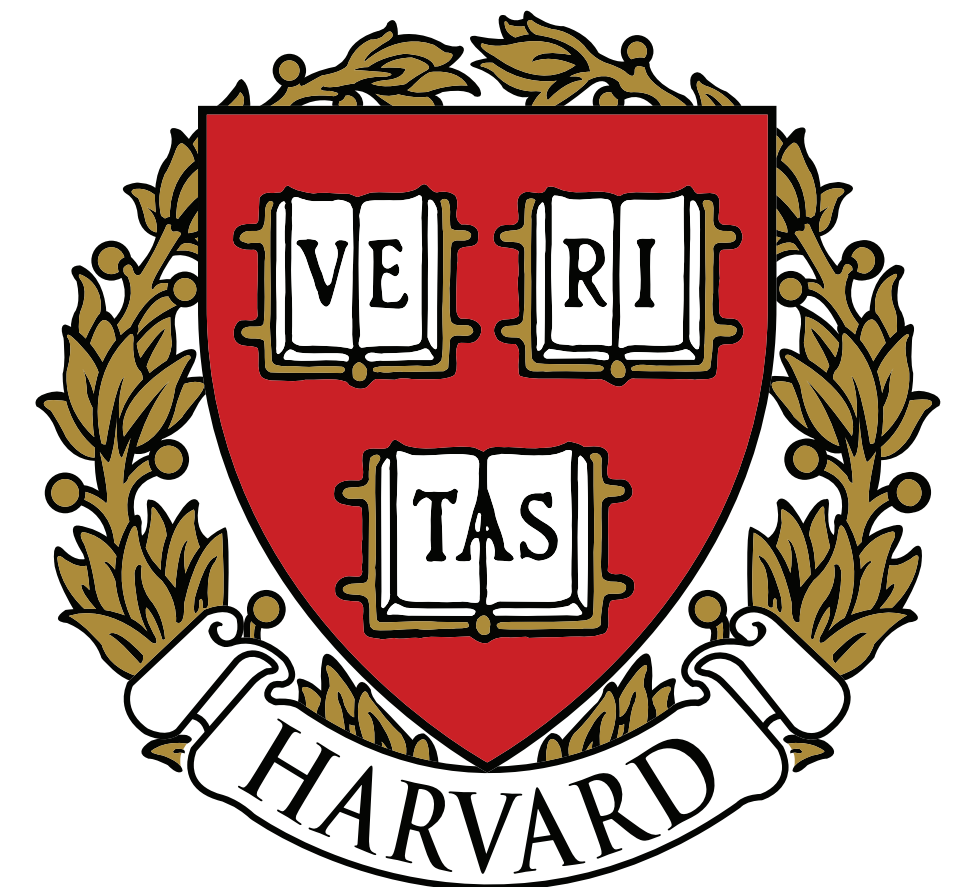


Last Iterate Risk Bounds of SGD with Decaying Stepsize for Overparameterized Linear Regression

Jingfeng Wu

with Difan Zou, Vladimir Braverman, Quanquan Gu, Sham M. Kakade



The Implicit Regularization Effect of SGD

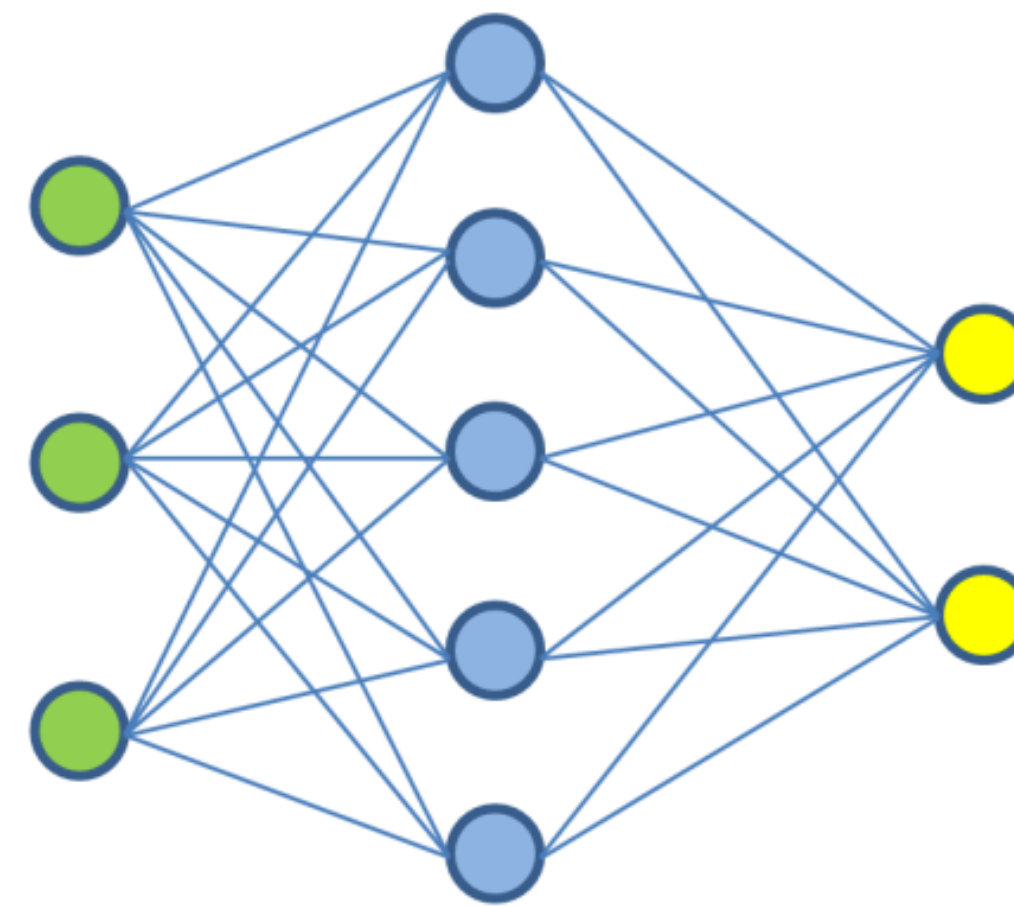
n training samples

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

Population Risk

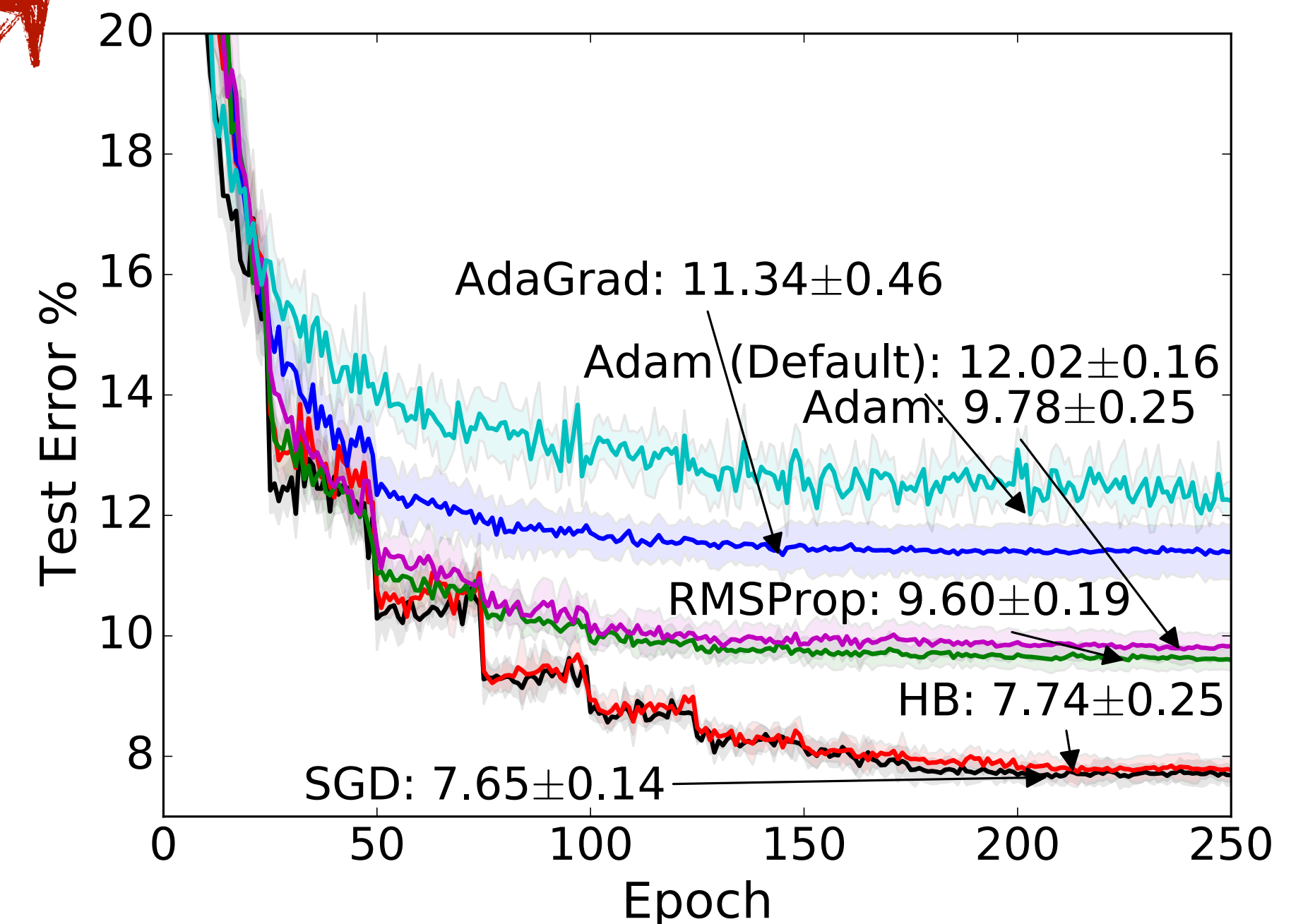
$$\mathcal{L}(\mathbf{w}) = \mathbb{E} \ell(\mathbf{x}, y; \mathbf{w})$$

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



SGD generalizes well for learning high-dim model

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



SGD generalizes well

High Dimensional Linear Regression

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

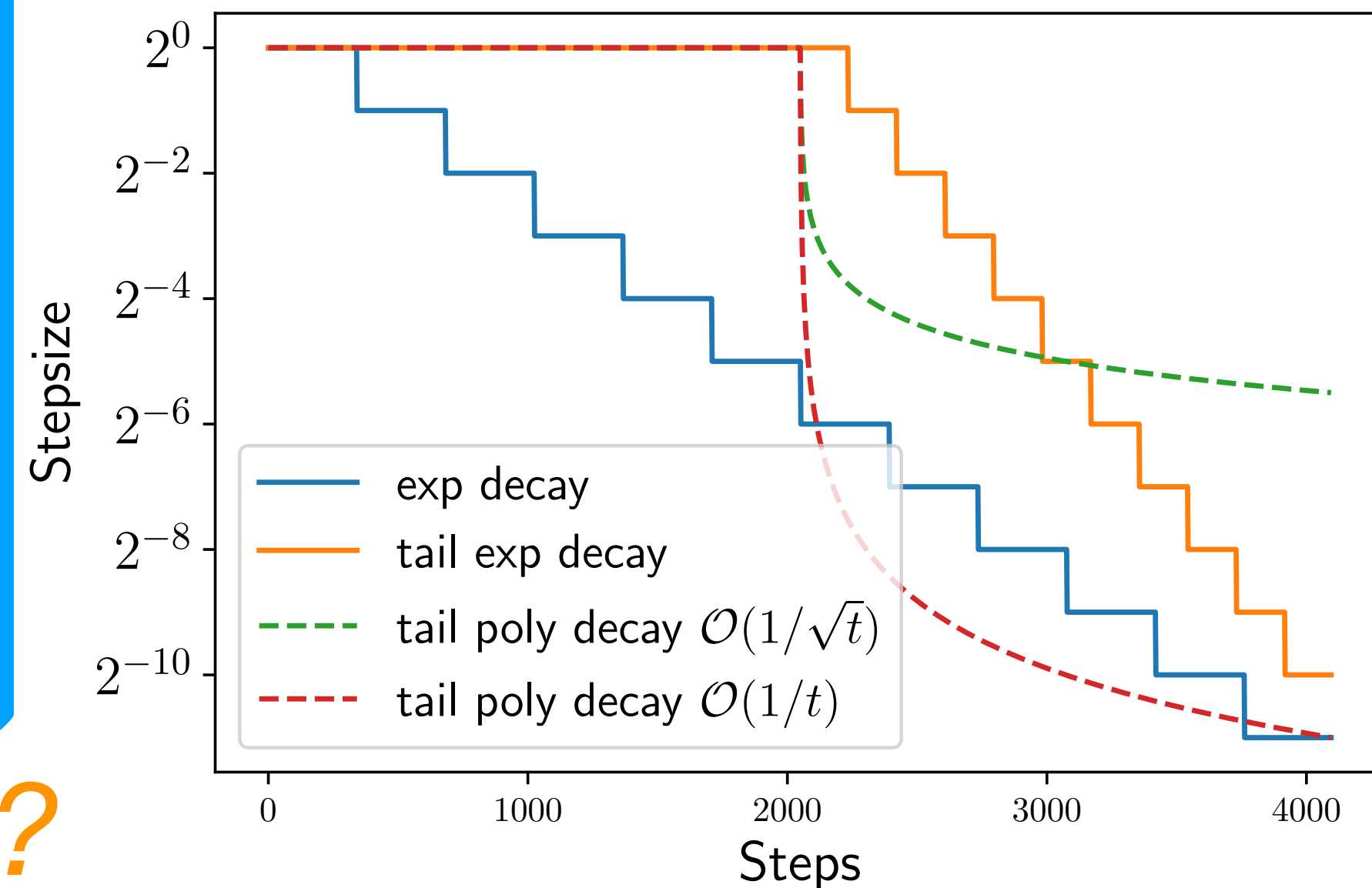
Population Risk $\mathcal{L}(\mathbf{w}) := \mathbb{E}(y - \mathbf{x}^\top \mathbf{w})^2$

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

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

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

output := \mathbf{w}_n



Caveat: One-Pass SGD *Two regimes: $d \lesssim n$?*

Key Assumption: **Strongly Contractive Fourth Moment**

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

- $\mathbb{E}[\mathbf{x}^\top \mathbf{A} \mathbf{x} \cdot \mathbf{x}\mathbf{x}^\top] \preceq \alpha \cdot \text{tr}(\mathbf{H}\mathbf{A}) \cdot \mathbf{H}$ for some constant $\alpha \geq 1$;
- $\mathbb{E}[\mathbf{x}^\top \mathbf{A} \mathbf{x} \cdot \mathbf{x}\mathbf{x}^\top] \succeq \beta \cdot \text{tr}(\mathbf{H}\mathbf{A}) \cdot \mathbf{H} + \mathbf{H}\mathbf{A}\mathbf{H}$ for some constant $\beta > 0$.

we are here

One-hot distributions
(which are easy to analyze)

Spherically symmetric distributions,
sub-Gaussian, sub-Exponential...

e.g., [BLLT 2020]

Bounded kurtosis

$$\forall \mathbf{v}, \mathbb{E}\langle \mathbf{v}, \mathbf{x} \rangle^4 \leq \alpha \langle \mathbf{v}, \mathbf{H}\mathbf{v} \rangle^2$$

**Strongly contractive
fourth moment**

Weakly contractive fourth moment

$$\mathbb{E}[\mathbf{x}\mathbf{x}^\top \mathbf{x}\mathbf{x}^\top] \preceq R^2 \cdot \mathbf{H}$$

e.g., [BM 2013]

Tail Geometrically Decaying Stepsizes

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

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

[GKKN 2019]

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

Useful in practice!

what if $d > n$?

Remarks

1. Weakly contractive fourth moment
2. Variance bound scales with d
3. ℓ_2 -norm or condition number implicitly depends on d

A Fine-Grained Upper Bound

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

$$\mathbb{E}\Delta(\mathbf{w}_n) \lesssim \frac{\|(\mathbf{I} - \eta_0\mathbf{H})^{s+K}(\mathbf{w}_0 - \mathbf{w}^*)\|_{\mathbf{I}_{0:k^*}}^2}{\eta_0 K} + \|(\mathbf{I} - \eta_0\mathbf{H})^{s+K}(\mathbf{w}_0 - \mathbf{w}^*)\|_{\mathbf{H}_{k^*:\infty}}^2$$

$$+ \frac{k^* + \eta_0 K \sum_{k^* < i \leq k^\dagger} \lambda_i + \eta_0^2 K^2 \sum_{i > k^\dagger} \lambda_i^2}{K} \cdot (\sigma^2 + \alpha \cdot \|\mathbf{w}_0 - \mathbf{w}^*\|_{\mathbf{H}}^2 \cdot \log(n))$$

exponentially decaying (pointing to the first term)

effective dimension (pointing to the numerator of the second term)

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

Ambient Dimension d vs.

$$\mathbf{I}_{0:k^*} := \text{diag}(1, \dots, 1, 0, 0, \dots) \quad \mathbf{H}_{k^*:\infty} := \text{diag}(0, \dots, 0, \lambda_{k^*+1}, \lambda_{k^*+2}, \dots)$$

Effective Dimension $k^* + \eta_0 K \sum_{k^* < i \leq k^\dagger} \lambda_i + \eta_0^2 K^2 \sum_{i > k^\dagger} \lambda_i^2$, small when $(\lambda_i)_{i \geq 1}$ decays fast

A Nearly Matching Lower Bound

Let the stepsize decaying interval be $K := (n - s)/\log(n - s)$. For every $s \geq 0$, $K > 10$ and every $\eta_0 < 1/\lambda_1$, we have

$$\mathbb{E} \Delta(\mathbf{w}_n) \gtrsim \|(\mathbf{I} - \eta_0 \mathbf{H})^{s+2K} (\mathbf{w}_0 - \mathbf{w}^*)\|_{\mathbf{H}}^2 +$$

$$\frac{k^* + \eta_0 K \sum_{k^* < i \leq k^\dagger} \lambda_i + \eta_0^2 K^2 \sum_{i > k^*} \lambda_i^2}{K} \cdot \left(\sigma^2 + \beta \cdot \|\mathbf{w}_0 - \mathbf{w}^*\|_{\mathbf{H}_{k^*:\infty}}^2 \right)$$

effective dimension \nearrow

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

Lower bound nearly matches upper bound if SNR is bounded, $\|\mathbf{w}_0 - \mathbf{w}^\|_{\mathbf{H}}^2 \lesssim \sigma^2$*

$$\mathbf{I}_{0:k^*} := \text{diag}(1, \dots, 1, 0, 0, \dots)$$

$$\mathbf{H}_{k^*:\infty} := \text{diag}(0, \dots, 0, \lambda_{k^*+1}, \lambda_{k^*+2}, \dots)$$

Geometrically vs. Polynomially Decaying Stepsize

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

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

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 η_0 . Then we have

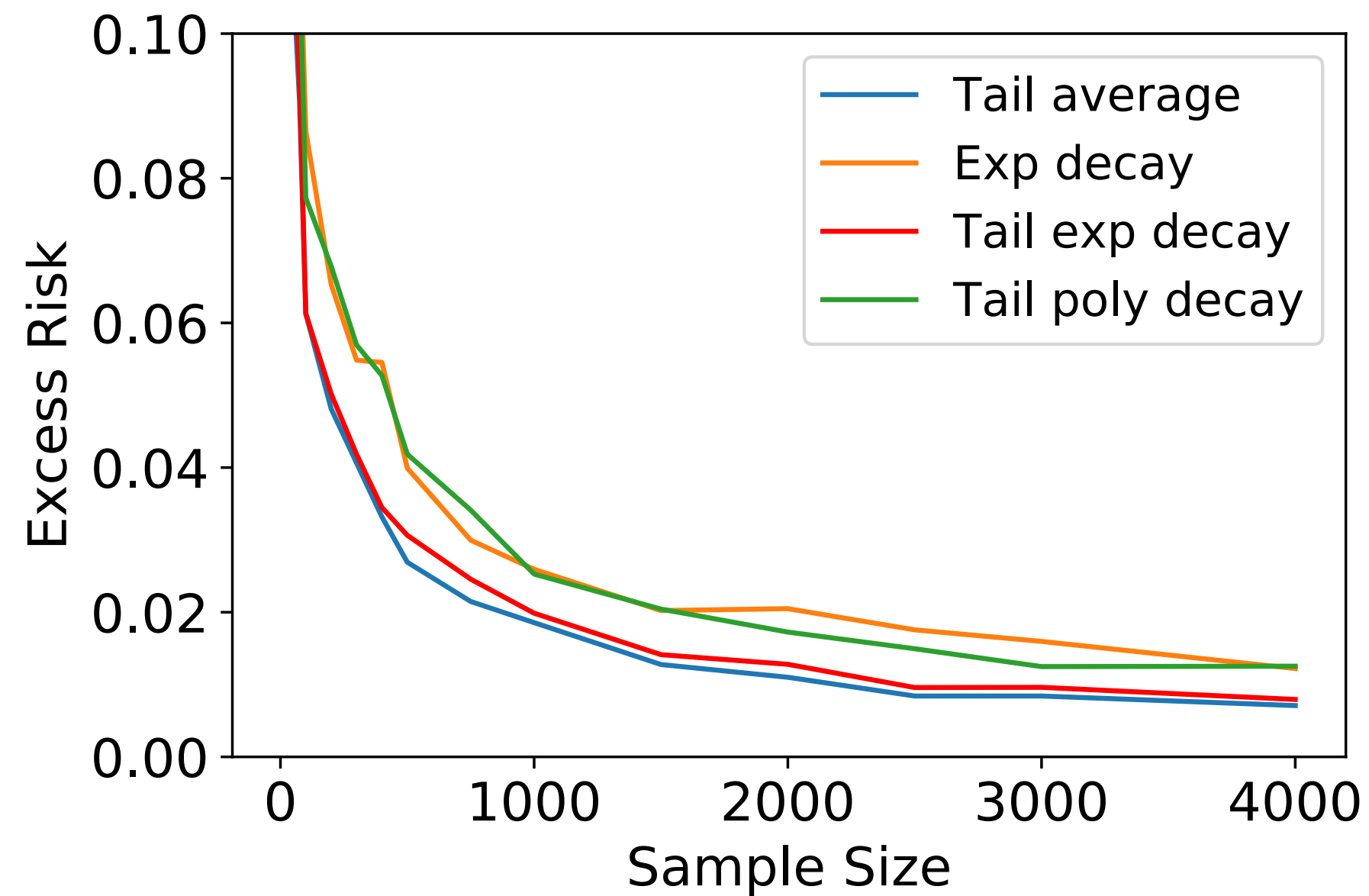
$$\mathbb{E}\Delta(\mathbf{w}_n^{\text{exp}}) \lesssim (1 + \text{SNR} \cdot \log n) \cdot \mathbb{E}\Delta(\mathbf{w}_n^{\text{poly}})$$

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

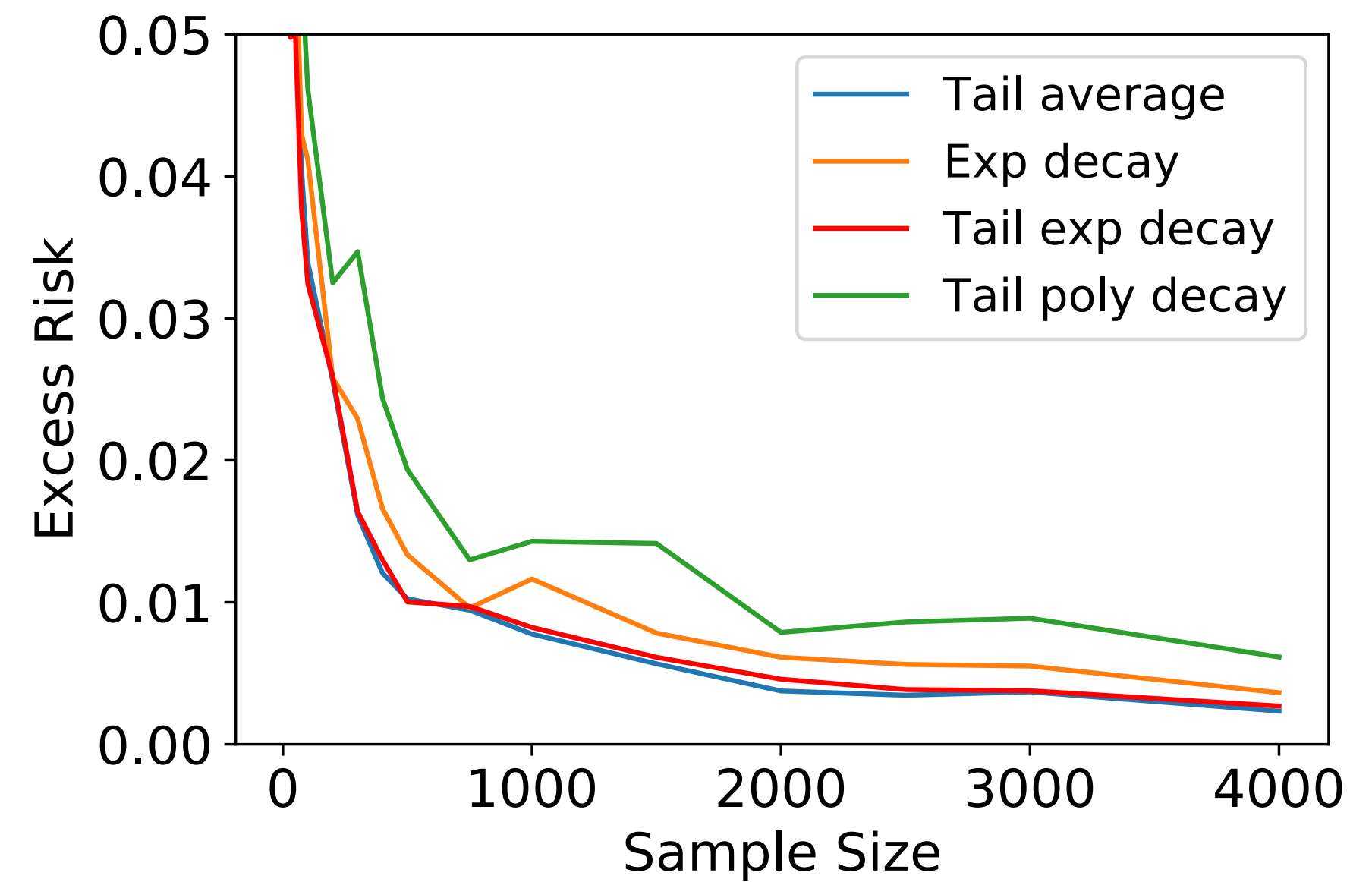
For every least square problem with bounded SNR,

$\mathbf{w}_n^{\text{exp}}$ is always nearly no worse than $\mathbf{w}_n^{\text{poly}}$

Numerical Simulation



$$\lambda_i = i^{-1}, \mathbf{w}^*[i] = i^{-1}$$



$$\lambda_i = i^{-2}, \mathbf{w}^*[i] = i^{-1}$$

Experimental Setting: $\sigma^2 = 1, d = 256, \mathbf{w}_0 = 0, s = n/2, a = 1$

Under each sample size, the initial stepsize is fine-tuned for each algorithm

- SGD can generalize in high-dim least squares
- Geometrically decaying stepsizes > polynomially decaying stepsizes

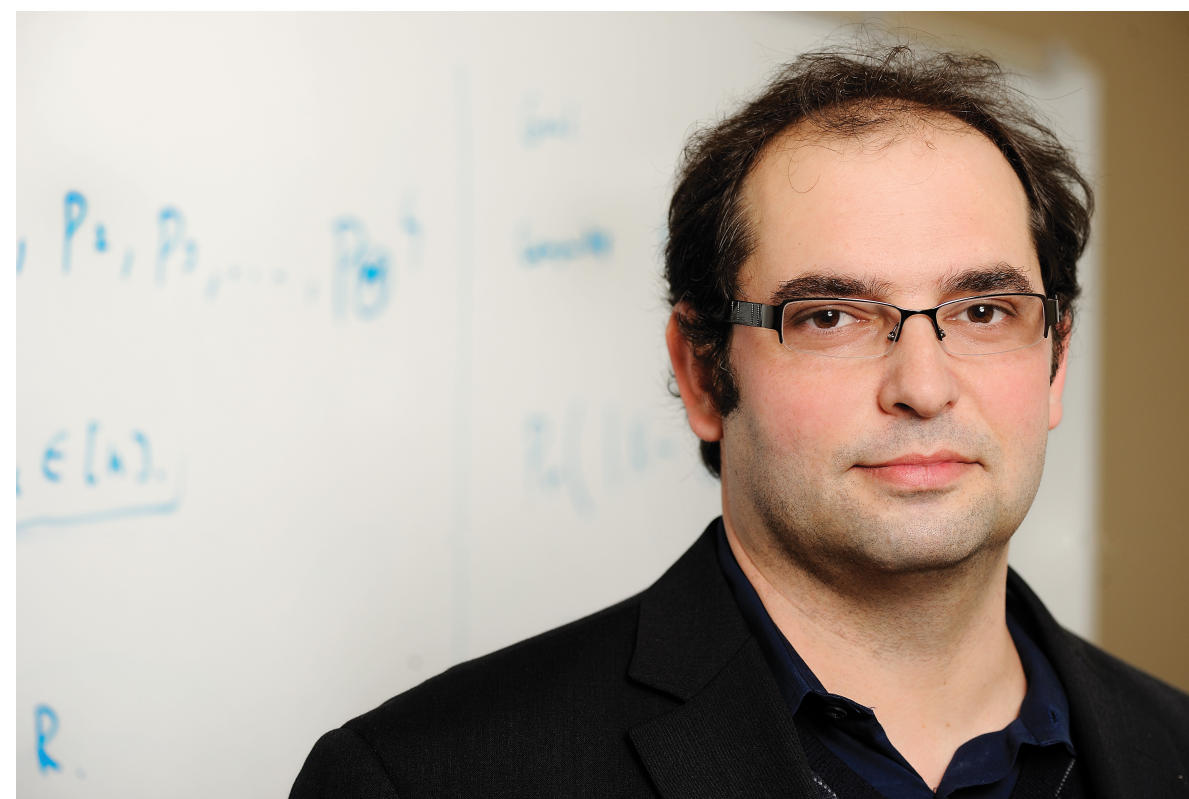
Conclusion

Take Home

- Risk of SGD in high-dim $\approx d_{\text{eff}} / n$
- d_{eff} determined by $(\lambda_i)_{i \geq 1}$, η_0 , n_{eff} ; and $\ll d$ when $(\lambda_i)_{i \geq 1}$ decay fast
- Geometrical stepsize $>$ polynomially stepsize

Limitations

- One-pass SGD
- Linear model
- Strongly contractive fourth moment



Vladimir Braverman @ JHU



Quanquan Gu @ UCLA



Sham M. Kakade @
Harvard



Difan Zou @ UCLA

Get the Paper!

