

## Contributions

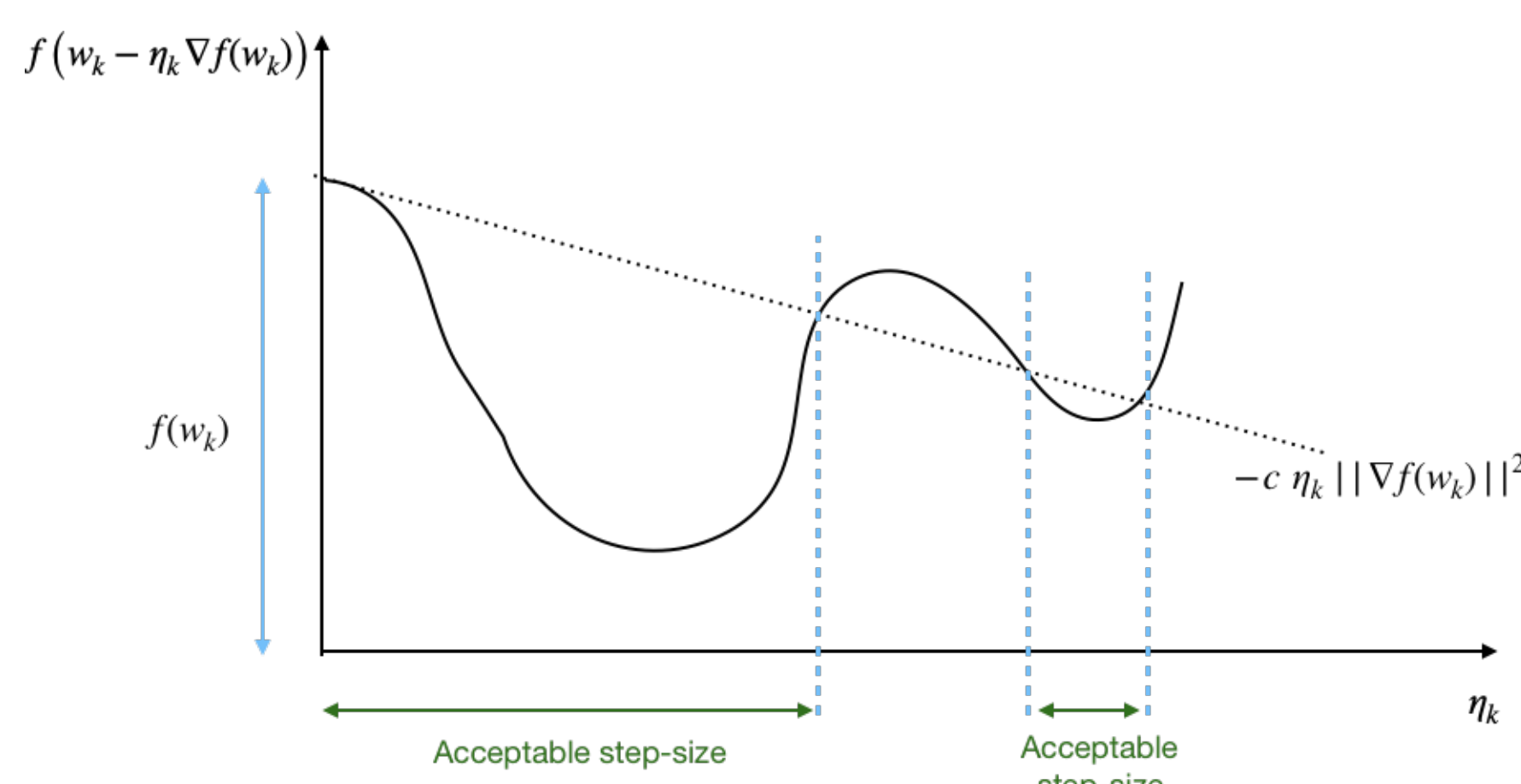
- Use *line-search methods* to automatically set the step-size when training over-parametrized models that can interpolate the data.
- Prove that SGD with the classic Armijo line-search attains the fast convergence rates of full-batch gradient descent in the interpolation setting for convex and strongly-convex functions.
- Prove that a stochastic extra-gradient method with a Lipschitz line-search attains fast convergence rates for an important class of non-convex functions and saddle-point problems satisfying interpolation.
- Compare against numerous optimization methods for standard classification tasks using both kernel methods and deep networks.

## General Setup

- **Objective:** Find  $w^* = f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$ .
- Technical assumptions: Lower bounded by  $f^*$ ,  $L$ -smoothness,  $\mu$ -strong convexity.
- **Interpolation:** If  $\nabla f(w^*) = 0$ , then for all  $f_i$ ,  $\nabla f_i(w^*) = 0$ .
- **Strong growth condition ( $\rho$ -SGC):**  $\mathbb{E}_i \|\nabla f_i(w)\|^2 \leq \rho \|\nabla f(w)\|^2$

## SGD + Armijo line-search

- **SGD update:**  $w_{k+1} = w_k - \eta_k \nabla f_{ik}(w_k)$  where  $\mathbb{E}_i [\nabla f_{ik}(w)] = \nabla f(w)$  for all  $x$ .
- **Stochastic Armijo condition:** Choose an  $\eta_k$  such that it satisfies the condition:  $f_{ik}(w_k - \eta_k \nabla f_{ik}(w_k)) \leq f_{ik}(w_k) - c \cdot \eta_k \|\nabla f_{ik}(w_k)\|^2$ .
- Use **back-tracking line-search** to choose a “large” step-size that satisfies the above condition.



## Lemma

The step-size  $\eta_k$  returned by the Armijo line-search and constrained to lie in the  $(0, \eta_{\max}]$  range satisfies the following inequality,

$$\eta_k \geq \min \left\{ \frac{2(1-c)}{L_{ik}}, \eta_{\max} \right\}.$$

Here  $L_{ik}$  is the Lipschitz constant of  $\nabla f_{ik}$ .

## Convergence of SGD with Armijo line-search

## Theorem (Strongly-Convex)

Assuming (a) interpolation, (b)  $L_i$ -smoothness and (c) convexity of  $f_i$ 's and (d)  $\mu$  strong-convexity of  $f$ , SGD with Armijo line-search with  $c = 1/2$  achieves the rate:

$$\mathbb{E} [\|w_T - w^*\|^2] \leq \max \left\{ \left(1 - \frac{\bar{\mu}}{L_{\max}}\right), (1 - \bar{\mu} \eta_{\max}) \right\}^T \|w_0 - w^*\|^2.$$

Here  $\bar{\mu} = \sum_{i=1}^n \mu_i / n$  is the average strong-convexity of the finite sum and  $L_{\max} = \max_i L_i$  is the maximum smoothness constant in the  $f_i$ 's.

## Convergence of SGD with Armijo line-search

## Theorem (Convex)

Assuming (a) interpolation, (b)  $L_i$ -smoothness and (c) convexity of  $f_i$ 's, SGD with Armijo line-search for all  $c \geq 1/2$  and iterate averaging achieves the rate:

$$\mathbb{E} [f(\bar{w}_T) - f(w^*)] \leq \frac{c \cdot \max \left\{ \frac{L_{\max}}{2(1-c)}, \frac{1}{\eta_{\max}} \right\}}{(2c-1)T} \|w_0 - w^*\|^2.$$

Here,  $\bar{w}_T = \frac{\sum_{i=1}^T w_i}{T}$  is the averaged iterate after  $T$  iterations and  $L_{\max} = \max_i L_i$ .

## Theorem (Non-convex)

Assuming (a) the SGC with constant  $\rho$  and (b)  $L_i$ -smoothness of  $f_i$ 's, SGD with Armijo line-search with  $c = 1/2$  and setting  $\eta_{\max} = \frac{3}{2\rho L}$  achieves the rate:

$$\min_{k=0, \dots, T-1} \mathbb{E} \|\nabla f(w_k)\|^2 \leq \frac{4L_{\max}}{T} \left( \frac{2\rho}{3} + 1 \right) (f(w_0) - f(w^*))$$

## Practical considerations and algorithm

- For allowing the step-size to increase, we (i) reset the step-size to a larger value for initializing the line-search procedure. (Algorithm 2). (ii) Use an alternative Goldstein condition that checks the curvature and allows for an increase in the step-size.
- Consider accelerated variants with both Polyak and Nesterov momentum.

Algorithm 1 SGD+Armijo( $f, w_0, \eta_{\max}, b, c, \beta, \gamma, \text{opt}$ )

```

1: for  $k = 0, \dots, T$  do
2:    $i_k \leftarrow$  sample mini-batch of size  $b$ 
3:    $\eta \leftarrow \text{reset}(\eta, \eta_{\max}, \gamma, b, k, \text{opt})/\beta$ 
4:   repeat
5:      $\eta \leftarrow \beta \cdot \eta$ 
6:      $\tilde{w}_k \leftarrow w_k - \eta \nabla f_{ik}(w_k)$ 
7:   until  $f_{ik}(\tilde{w}_k) \leq f_{ik}(w_k) - c \cdot \eta \|\nabla f_{ik}(w_k)\|^2$ 
8:    $w_{k+1} \leftarrow \tilde{w}_k$ 
9: end for
10: return  $w_{k+1}$ 

```

Algorithm 2 reset( $\eta, \eta_{\max}, \gamma, b, k, \text{opt}$ )

```

1: if  $k = 1$  then
2:   return  $\eta_{\max}$ 
3: else if  $\text{opt} = 0$  then
4:    $\eta \leftarrow \eta$ 
5: else if  $\text{opt} = 1$  then
6:    $\eta \leftarrow \eta_{\max}$ 
7: else if  $\text{opt} = 2$  then
8:    $\eta \leftarrow \eta \cdot \gamma^{b/n}$ 
9: end if
10: return  $\eta$ 

```

## Stochastic Extra-Gradient + Lipschitz line-search

- **SEG update:**  $w'_k = w_k - \eta_k \nabla f_{ik}(w_k)$ ,  $w_{k+1} = w_k - \eta_k \nabla f_{ik}(w'_k)$
- **Lipschitz line-search condition:** Choose  $\eta_k$  such that it satisfies the condition:  $\|\nabla f_{ik}(w_k - \eta_k \nabla f_{ik}(w_k)) - \nabla f_{ik}(w_k)\| \leq c \|\nabla f_{ik}(w_k)\|$ .
- The step-size returned by the Lipschitz line-search satisfies  $\eta_k \geq \min \left\{ \frac{c}{L_{ik}}, \eta_{\max} \right\}$ .

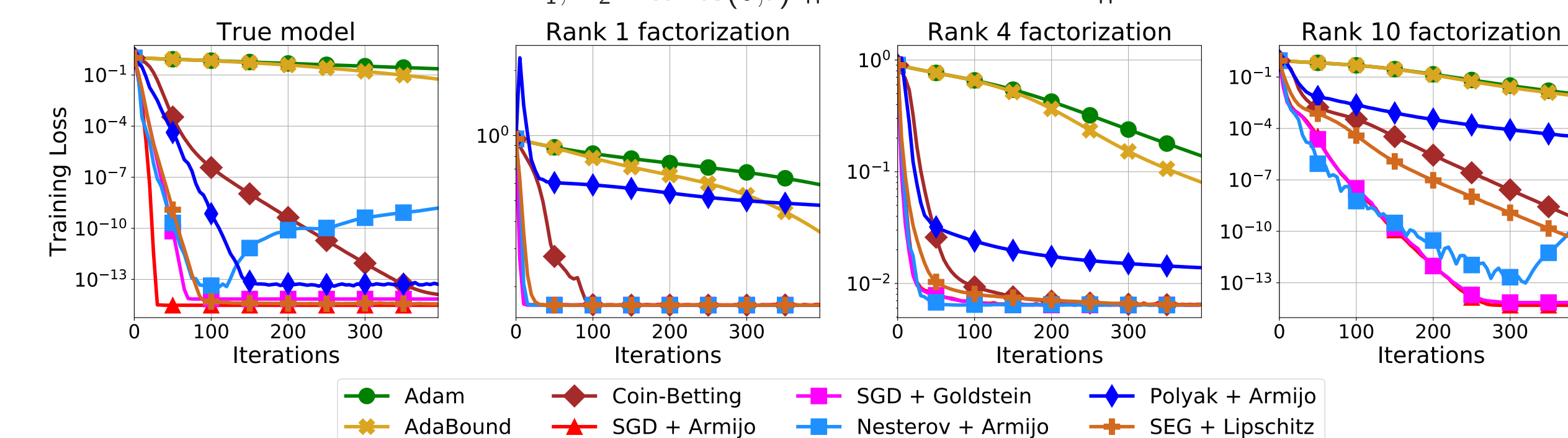
## Theorem

Assuming (a) interpolation, (b)  $L_i$ -smoothness, and (c)  $\mu_i$ -RSI of  $f_i$ 's, SEG with Lipschitz line-search with  $c = 14$  and  $\eta_{\max} \leq \min_i 1/4\mu_i$  achieves the rate:

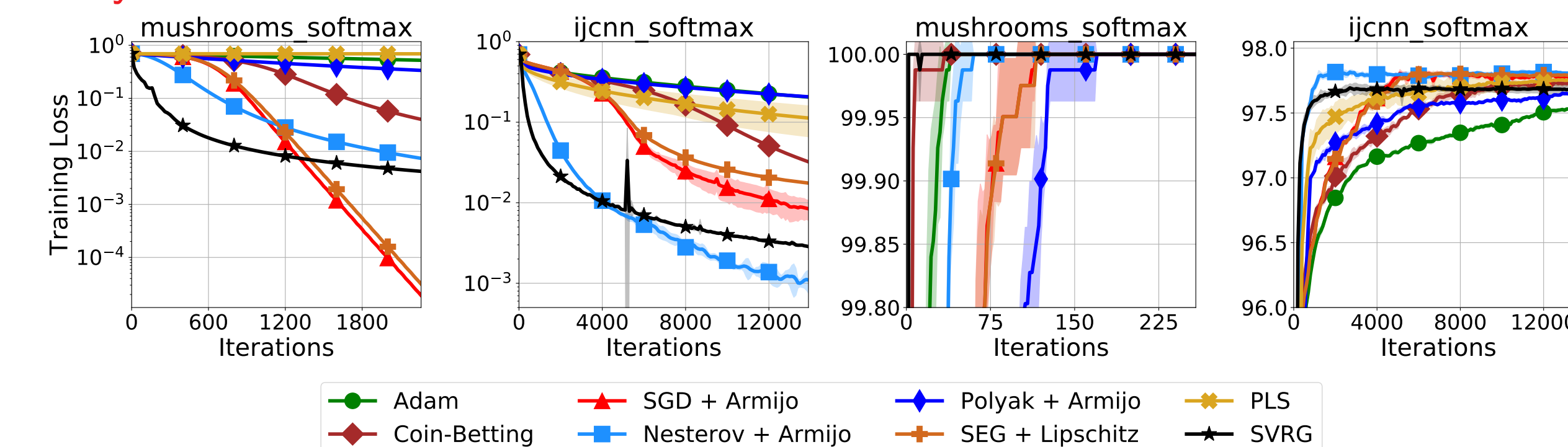
$$\mathbb{E} [\|w_T - \mathcal{P}_{\mathcal{X}^*}[w_T]\|^2] \leq \max \left\{ \left(1 - \frac{\bar{\mu}}{4L_{\max}}\right), (1 - \eta_{\max} \bar{\mu}) \right\}^T \|w_0 - \mathcal{P}_{\mathcal{X}^*}[w_0]\|^2,$$

where  $\bar{\mu} = \frac{\sum_{i=1}^n \mu_i}{n}$  is the average RSI constant of the finite sum and  $\mathcal{X}^*$  is the non-empty set of optimal solutions.  $\mathcal{P}_{\mathcal{X}^*}[w]$  denotes the projection of  $w$  onto  $\mathcal{X}^*$ .

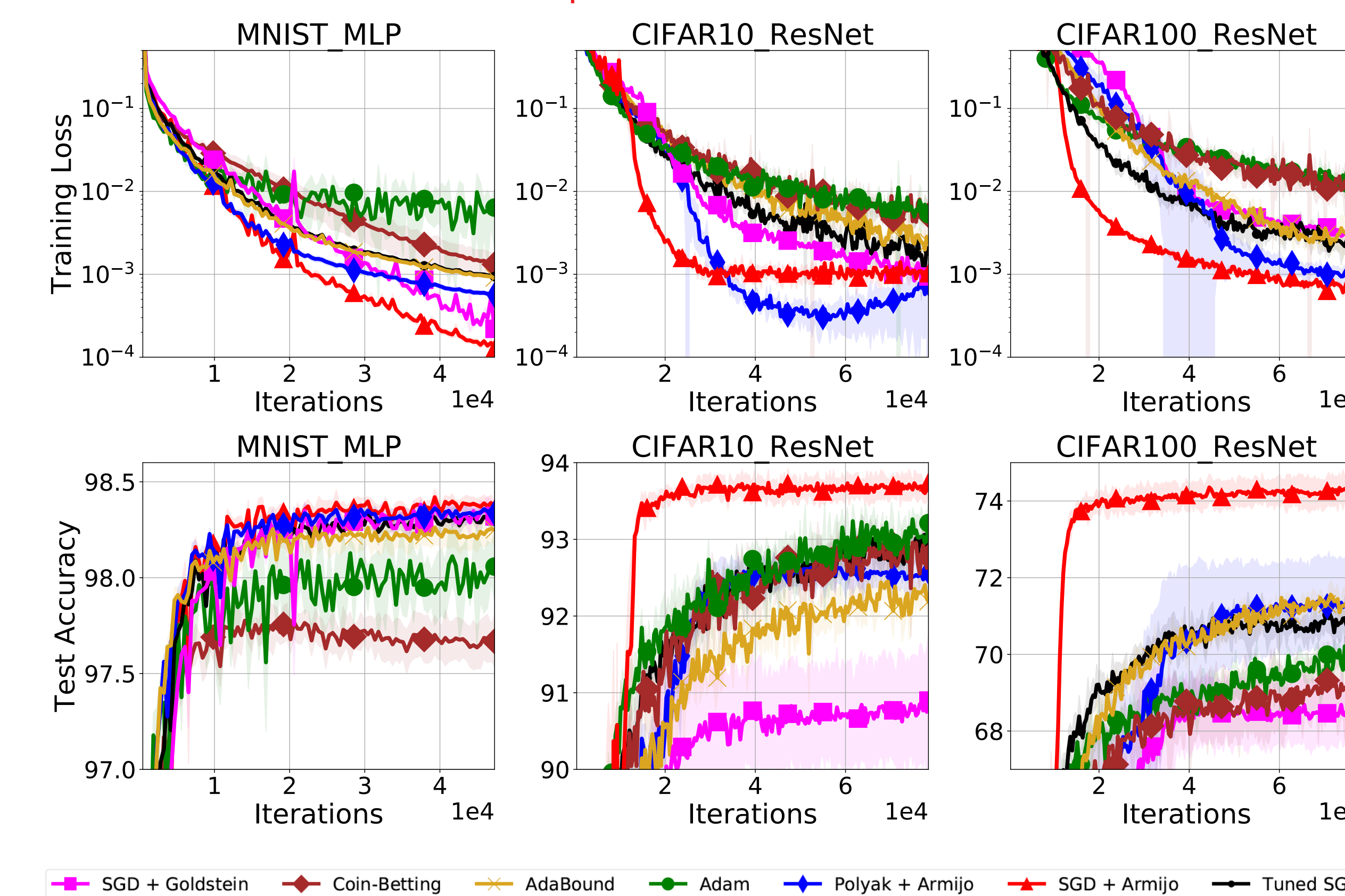
- Derive an  $O(1/T)$  rate when minimizing convex functions.
- Derive linear convergence rates for strongly-convex strongly-concave and bilinear saddle point problems satisfying interpolation.

Experiments <https://github.com/IssamLaradji/sls>► Matrix Factorization:  $\min_{W_1, W_2} \mathbb{E}_{x \sim \mathcal{N}(0, I)} \|W_2 W_1 x - Ax\|^2$ 

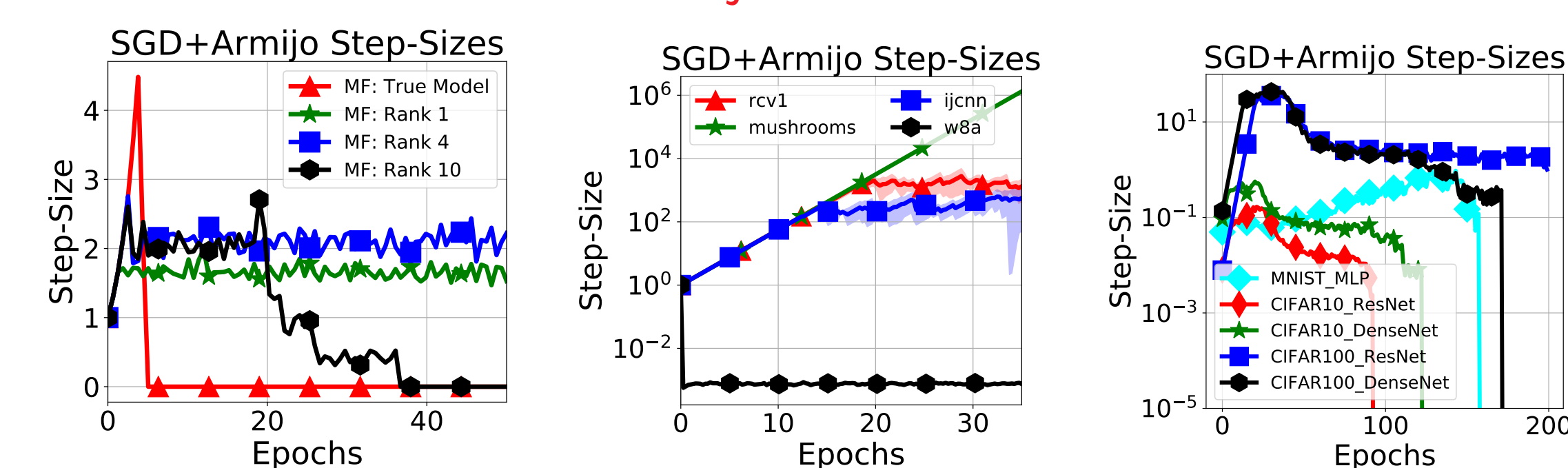
## ► Binary classification with RBF kernels



## ► Multi-class classification with deep networks



## ► Step-size variation for SGD with Armijo line-search



## ► Runtimes

