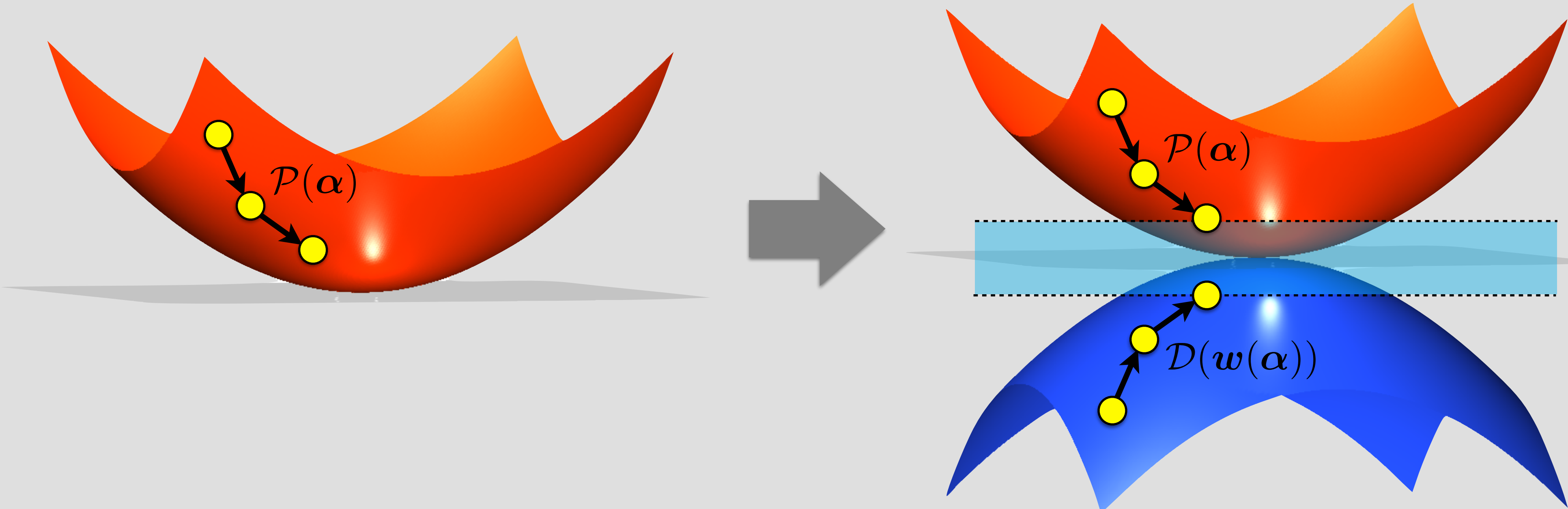


## overview

**Challenge:** most existing optimization algorithms lack accuracy certificates.



### Contributions

- ✓ equip existing optimizers with accuracy **certificates**
- ✓ algorithm-independent new **primal-dual convergence rates** for a larger problem class

## setup

Problem formulation

$$\min_{\alpha \in \mathbb{R}^n} f(A\alpha) + g(\alpha)$$

$$\min_{\mathbf{w} \in \mathbb{R}^d} f^*(\mathbf{w}) + g^*(-A^\top \mathbf{w})$$

gap

$f$  convex, smooth     $g$  convex

**correspondence mapping**  
 $\mathbf{w} = \mathbf{w}(\alpha) := \nabla f(A\alpha)$

Optimality conditions

$$\mathbf{w} \in \partial f(A\alpha) \quad -A^\top \mathbf{w} \in \partial g(\alpha),$$

$$A\alpha \in \partial f^*(\mathbf{w}) \quad \alpha \in \partial g^*(-A^\top \mathbf{w})$$

Convex conjugate:  $h^*(\mathbf{v}) := \sup_{\mathbf{u} \in \mathbb{R}^d} \mathbf{v}^\top \mathbf{u} - h(\mathbf{u})$

## main results

**existing rates  $\Rightarrow$  primal-dual rates**

- ▶ algorithm agnostic

$g$  strongly convex

- ✦ linear rate  $\Rightarrow$  linear primal-dual rate  
 $f = L2, g$  separable: see SDCA

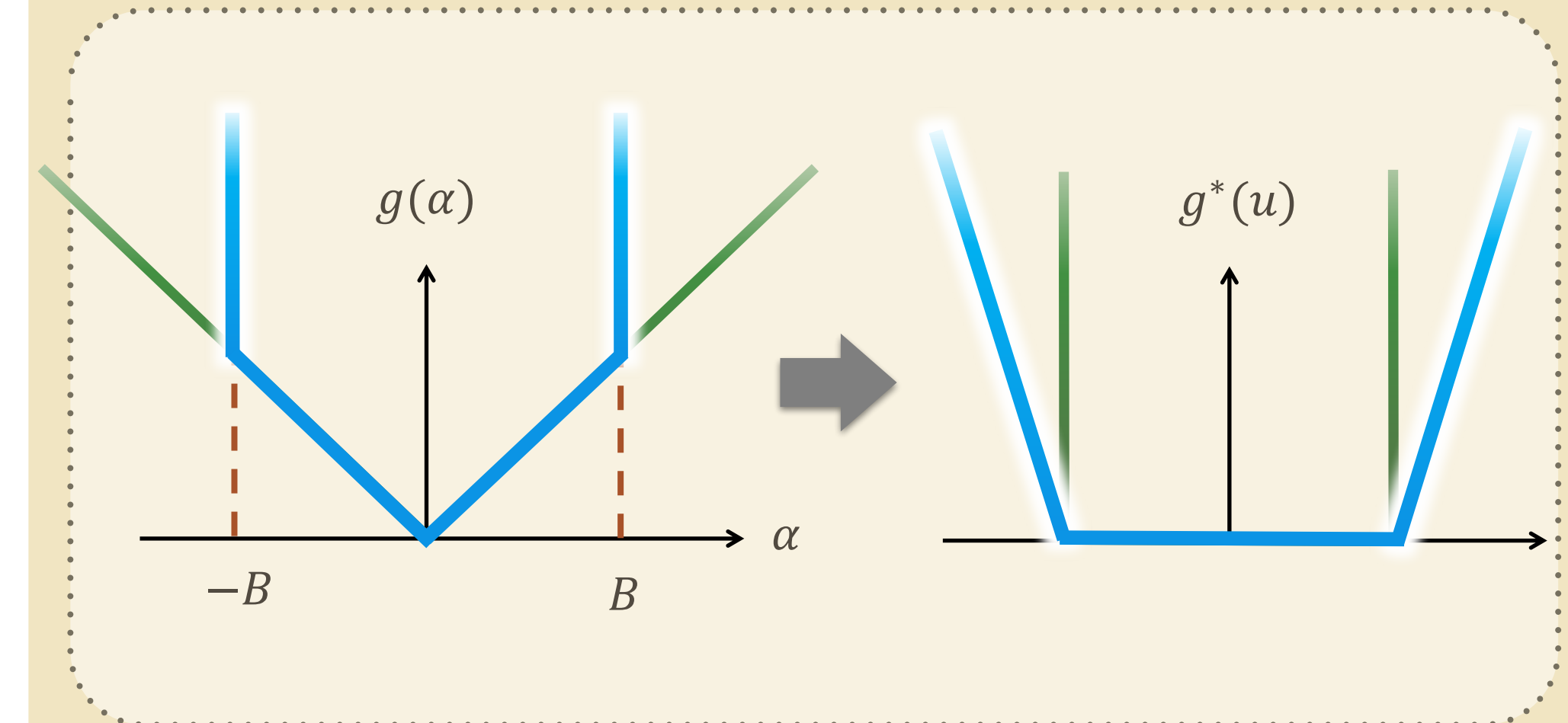
$g$  bounded support

- ✦ linear rate  $\Rightarrow$  linear primal-dual rate  
new: SVM
- ✦  $1/T$  rate  $\Rightarrow \sqrt{1/T}$  primal-dual rate

- $g$  general convex?
- ✦ same! using trick, see next

examples:  
L1, elastic-n,  
group lasso,  
TV, fused L1,  
structured

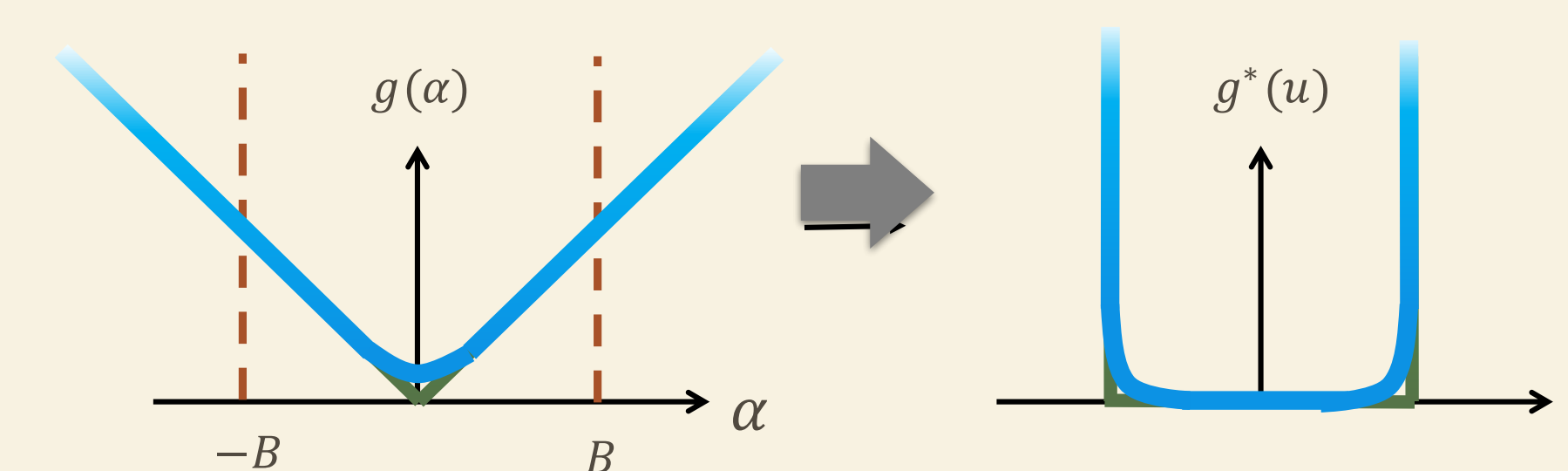
## lipschitzing trick



- ▶ makes  $g^*$  globally Lipschitz
- ▶ gives duality gap defined on entire region of interest  
 $B$  easy to choose for norm-reg. problems
- ▶ problem and algorithms unaffected!

▶ can re-use all existing algorithms!

## Comparison to Nesterov Smoothing



- ▶ makes  $g^*$  strongly convex
- ▶ changes iterates

## proof details

**Lemma 1.** Consider an optimization problem of the form (A). Let  $f$  be  $1/\beta$ -smooth w.r.t. a norm  $\|\cdot\|_f$  and let  $g$  be  $\mu$ -strongly convex with convexity parameter  $\mu \geq 0$  w.r.t. a norm  $\|\cdot\|_g$ . The general convex case  $\mu = 0$  is explicitly allowed, but only if  $g$  has bounded support.

Then, for any  $\alpha \in \text{dom}(\mathcal{D})$  and any  $s \in [0, 1]$ , it holds that

$$\mathcal{D}(\alpha) - \mathcal{D}(\alpha^*) \geq sG(\alpha) \quad (5)$$

$$+ \frac{s^2}{2} \left( \frac{\mu(1-s)}{s} \|\mathbf{u} - \alpha\|_g^2 - \frac{1}{\beta} \|A(\mathbf{u} - \alpha)\|_f^2 \right)$$

where  $G(\alpha)$  is the gap function defined in (4) and

$$\mathbf{u} \in \partial g^*(-A^\top \mathbf{w}(\alpha)). \quad (6)$$

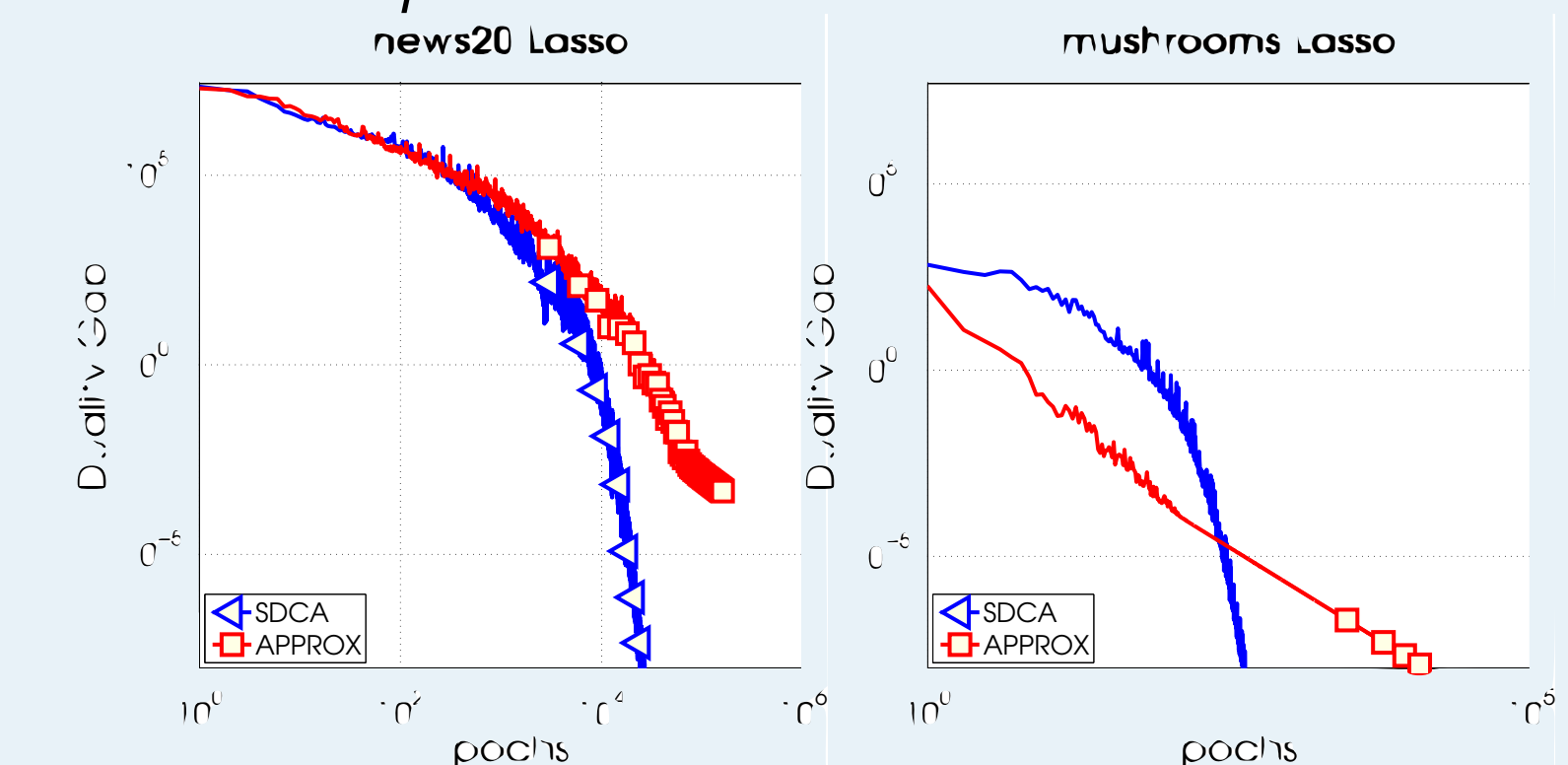
## coordinate descent

$$f(A\alpha) + \sum_i g_i(\alpha_i)$$

$$f^*(\mathbf{w}) + \sum_i g_i^*(-A_{:,i}^\top \mathbf{w})$$

- ▶ new primal-dual rates for CD on *L1, group-lasso, elastic-net*, etc
- ▶ certificates
- ▶ generalize SDCA. no square root

illustrative experiments for L1

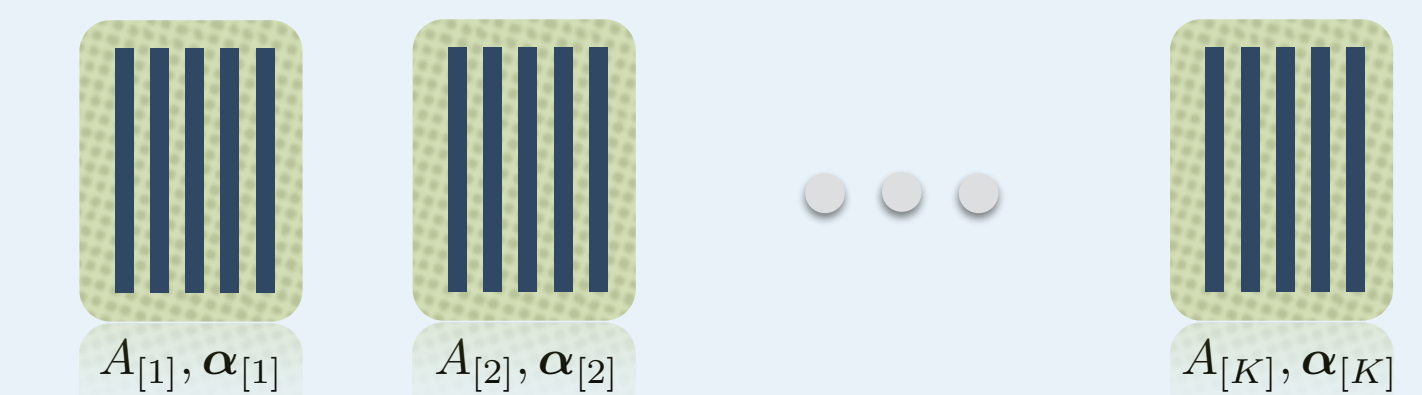


## distributed version

ML Systems Workshop on Friday

**L1-Regularized Distributed Optimization: A Communication-Efficient Primal-Dual Framework**

Smith, V., Forte, S., Jordan, M. I., & Jaggi, M.



arXiv 1512.04011

## references

- [1] Shalev-Shwartz and Zhang. *Stochastic dual coordinate ascent methods for regularized loss minimization*. JMLR, 14:567–599, 2013
- [2] Necoara, I. (2015). Linear convergence of first order methods under weak nondegeneracy assumptions for convex programming. arXiv/1504.06298