

Recent Advances in Primal-Dual Coordinate Methods for Empirical Risk Minimization

Qi Lei *

Based on joint work with

Ian E.H. Yen[†], Chao-yuan Wu*, Inderjit S. Dhillon* and Pradeep Ravikumar[†].

* University of Texas at Austin

† Carnegie Mellon University

Problem Setup

$$\begin{aligned}\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) &= \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}) \\ &\stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \phi_i(\mathbf{a}_i^\top \mathbf{x}) + g(\mathbf{x})\end{aligned}$$

- Large-scale supervised machine learning: **very large feature dimension d , and very large data size n**
- Assumptions: μ -strongly convex, L -smooth. Condition number $\kappa = \frac{L}{\mu}$.

Problem Setup

$$\begin{aligned}\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) &= \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}) \\ &\stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \phi_i(\mathbf{a}_i^\top \mathbf{x}) + g(\mathbf{x})\end{aligned}$$

- Large-scale supervised machine learning: **very large feature dimension d , and very large data size n**
- Assumptions: μ -strongly convex, L -smooth. Condition number $\kappa = \frac{L}{\mu}$.
e.g.,

squared loss $\phi \Rightarrow$ linear regression problem

sigmoid loss $\phi \Rightarrow$ logistic regression problem

smooth hinge loss $\phi \Rightarrow$ support vector machine

First-Order Optimization

$$\begin{aligned}\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) &= \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}) \\ &\stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \phi_i(\mathbf{a}_i^\top \mathbf{x}) + g(\mathbf{x})\end{aligned}$$

- Classical methods comparison

	method	per Iteration	#Iteration
$\min_{\mathbf{x}} f(\mathbf{x})$	GD	$\mathcal{O}(nd)$	$\mathcal{O}(\kappa \log^{-1} \epsilon)$
	RCD	$\mathcal{O}(n)$	$\mathcal{O}(d\kappa \log^{-1} \epsilon)$
	SGD	$\mathcal{O}(d)$	$\mathcal{O}(\kappa \epsilon^{-1})$

method	per Iteration	# Iteration
GD	$\mathcal{O}(nd)$	$\mathcal{O}(\kappa \log^{-1} \epsilon)$
RCD	$\mathcal{O}(n)$	$\mathcal{O}(d\kappa \log^{-1} \epsilon)$
SGD	$\mathcal{O}(d)$	$\mathcal{O}(\kappa \epsilon^{-1})$
• Expansions		
MISO		$\mathcal{O}((n + \kappa) \log^{-1} \epsilon)$
SVRG		$\mathcal{O}((n + \sqrt{n\kappa}) \log^{-1} \epsilon)$
SAG	$\mathcal{O}(d)$	
SAGA		$\mathcal{O}((n + \sqrt{n\kappa}) \log^{-1} \epsilon)$ if accelerated
SDCA	$\mathcal{O}(d)$	$\mathcal{O}((n + \kappa) \log^{-1} \epsilon)$
SPDC	$\mathcal{O}(d)$	$\mathcal{O}((n + \sqrt{n\kappa}) \log^{-1} \epsilon)$

method	per Iteration	# Iteration
GD	$\mathcal{O}(nd)$	$\mathcal{O}(\kappa \log^{-1} \epsilon)$
RCD	$\mathcal{O}(n)$	$\mathcal{O}(d\kappa \log^{-1} \epsilon)$
SGD	$\mathcal{O}(d)$	$\mathcal{O}(\kappa \epsilon^{-1})$
MISO		
SVRG		$\mathcal{O}((n + \kappa) \log^{-1} \epsilon)$
SAG	$\mathcal{O}(d)$	$\mathcal{O}((n + \sqrt{n\kappa}) \log^{-1} \epsilon)$
SAGA		if accelerated
SDCA	$\mathcal{O}(d)$	$\mathcal{O}((n + \kappa) \log^{-1} \epsilon)$
SPDC	$\mathcal{O}(d)$	$\mathcal{O}((n + \sqrt{n\kappa}) \log^{-1} \epsilon)$

from lazy gradient evaluations to variance reduction:

MISO (Mairal, 2015)

SVRG (Johnson & Zhang, 2013)

SAG(A) (Defazio et al., 2014)

method	per Iteration	# Iteration
GD	$\mathcal{O}(nd)$	$\mathcal{O}(\kappa \log^{-1} \epsilon)$
RCD	$\mathcal{O}(n)$	$\mathcal{O}(d\kappa \log^{-1} \epsilon)$
SGD	$\mathcal{O}(d)$	$\mathcal{O}(\kappa \epsilon^{-1})$
• Expansions		
MISO		$\mathcal{O}((n + \kappa) \log^{-1} \epsilon)$
SVRG		$\mathcal{O}((n + \sqrt{n\kappa}) \log^{-1} \epsilon)$
SAG	$\mathcal{O}(d)$	
SAGA	$\mathcal{O}(d)$	$\mathcal{O}((n + \sqrt{n\kappa}) \log^{-1} \epsilon)$ if accelerated
SDCA	$\mathcal{O}(d)$	$\mathcal{O}((n + \kappa) \log^{-1} \epsilon)$
SPDC	$\mathcal{O}(d)$	$\mathcal{O}((n + \sqrt{n\kappa}) \log^{-1} \epsilon)$

similar update rule but improved convergence rate
 SDCA (Shalev-Shwartz and Zhang, 2013)

method	per Iteration	# Iteration
GD	$\mathcal{O}(nd)$	$\mathcal{O}(\kappa \log^{-1} \epsilon)$
RCD	$\mathcal{O}(n)$	$\mathcal{O}(d\kappa \log^{-1} \epsilon)$
SGD	$\mathcal{O}(d)$	$\mathcal{O}(\kappa \epsilon^{-1})$
• Expansions		
MISO		$\mathcal{O}((n + \kappa) \log^{-1} \epsilon)$
SVRG		$\mathcal{O}((n + \sqrt{n\kappa}) \log^{-1} \epsilon)$
SAG	$\mathcal{O}(d)$	
SAGA		if accelerated
SDCA	$\mathcal{O}(d)$	$\mathcal{O}((n + \kappa) \log^{-1} \epsilon)$
SPDC	$\mathcal{O}(d)$	$\mathcal{O}((n + \sqrt{n\kappa}) \log^{-1} \epsilon)$

added extrapolation terms
 SPDC (Zhang and Xiao, 2015)

Primal, Dual, Lagrangian forms

Primal form

$$\min_{\mathbf{x} \in \mathbb{R}^d} \left\{ P(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \phi_i(\mathbf{a}_i^\top \mathbf{x}) + g(\mathbf{x}) \right\}$$

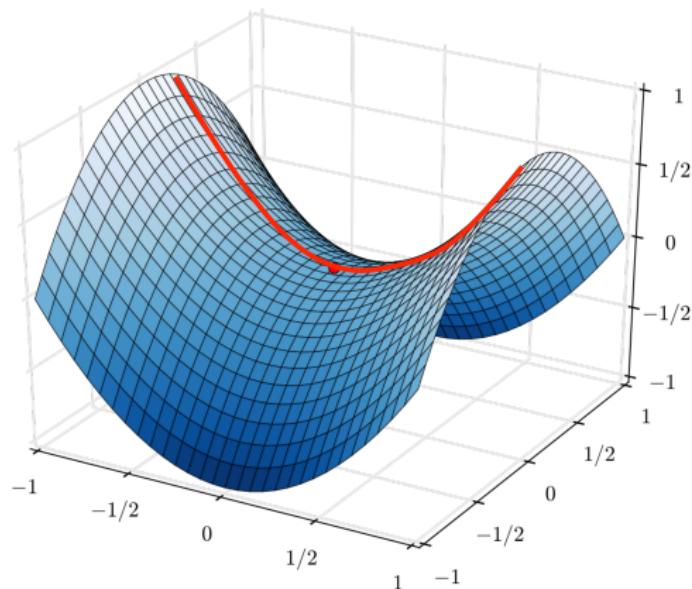
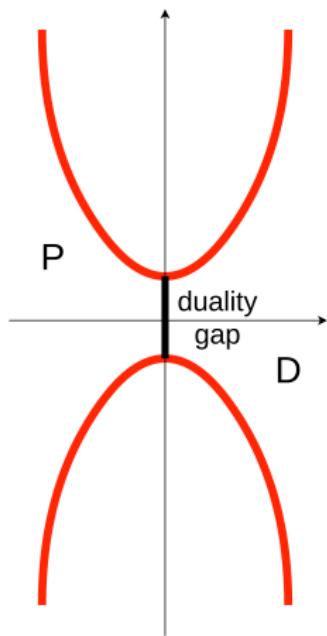
⇒ Dual form

$$\max_{\mathbf{y} \in \mathbb{R}^n} \left\{ D(\mathbf{y}) \stackrel{\text{def}}{=} -g^*\left(-\frac{\mathbf{A}^\top \mathbf{y}}{n}\right) - \frac{1}{n} \sum_{i=1}^n \phi_i^*(y_i) \right\}$$

⇒ Lagrangian (convex-concave saddle point problem):

$$\max_{\mathbf{y} \in \mathbb{R}^n} \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ \mathcal{L}(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} g(\mathbf{x}) + \frac{1}{n} \mathbf{y}^\top \mathbf{A} \mathbf{x} - \frac{1}{n} \sum_{i=1}^n \phi_i^*(y_i) \right\}$$

Illustration



Building up

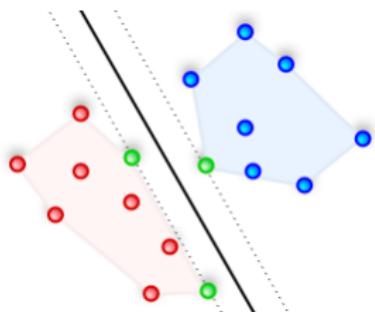
Why this convex-concave formulations?

$$\min_{\mathbf{x} \in \mathbb{R}^d} \max_{\mathbf{y} \in \mathbb{R}^n} \{ L(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top A \mathbf{x} - g(\mathbf{y}) \}$$

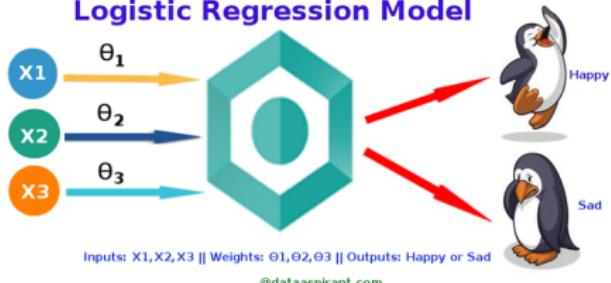
- ① many machine learning applications

Machine Learning Applications

Empirical Risk Minimization



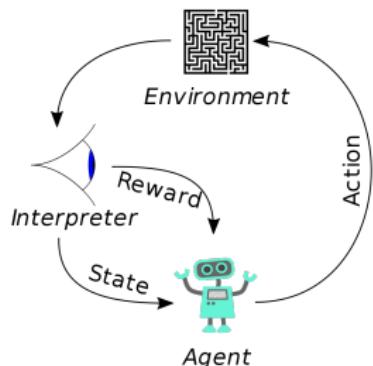
Logistic Regression Model



Inputs: X_1, X_2, X_3 || Weights: $\theta_1, \theta_2, \theta_3$ || Outputs: Happy or Sad

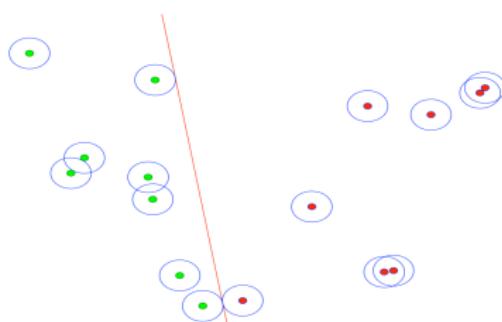
@dataaspirant.com

Reinforcement Learning



(Du et al. 2017)

Robust Optimization



(Ben-Tal et al., 2009)

Building up

Why this convex-concave formulations?

$$\min_{\mathbf{x} \in \mathbb{R}^d} \max_{\mathbf{y} \in \mathbb{R}^n} \{ L(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top A \mathbf{x} - g(\mathbf{y}) \}$$

- ① many machine learning applications
- ② more suitable for greedy coordinate methods

A theoretical vignette

Q

What are some advantages of using greedy coordinate methods?

- Given a function $f(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^d$ that is μ -strongly convex, and L -smooth,
 - CD with random selection/permutation: $\mathcal{O}\left(\frac{dL}{\mu} \log \frac{1}{\epsilon}\right)$ iterations to achieve ϵ error.
 - CD with greedy selection: $\mathcal{O}\left(\frac{L}{\mu_1} \log \frac{1}{\epsilon}\right)$ iterations, if loss function is μ_1 ℓ_1 -strongly convex (Nutini et al. 2015)

A theoretical vignette

Q

What are some advantages of using greedy coordinate methods?

- Given a function $f(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^d$ that is μ -strongly convex, and L -smooth,
 - CD with random selection/permutation: $\mathcal{O}\left(\frac{dL}{\mu} \log \frac{1}{\epsilon}\right)$ iterations to achieve ϵ error.
 - CD with greedy selection: $\mathcal{O}\left(\frac{L}{\mu_1} \log \frac{1}{\epsilon}\right)$ iterations, if loss function is μ_1 ℓ_1 -strongly convex ([Nutini et al. 2015](#))
 - μ_1 could be as worse as μ/d
 - A special case, when iterates and optimal values are s -sparse, it only requires $\mathcal{O}\left(\frac{sL}{\mu} \log \frac{1}{\epsilon}\right)$ iterations to achieve ϵ error

A theoretical vignette

Q

What are some advantages of using greedy coordinate methods?

- Given a function $f(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^d$ that is μ -strongly convex, and L -smooth,
 - 1) CD with random selection/permutation: $\mathcal{O}\left(\frac{dL}{\mu} \log \frac{1}{\epsilon}\right)$ iterations to achieve ϵ error.
 - 2) CD with greedy selection: $\mathcal{O}\left(\frac{L}{\mu_1} \log \frac{1}{\epsilon}\right)$ iterations, if loss function is μ_1 ℓ_1 -strongly convex ([Nutini et al. 2015](#))
 - ① μ_1 could be as worse as μ/d
 - ② A special case, when iterates and optimal values are s -sparse, it only requires $\mathcal{O}\left(\frac{sL}{\mu} \log \frac{1}{\epsilon}\right)$ iterations to achieve ϵ error
 - ③ Many applications have such primal sparsity:
 - Large-scale multiclass/multilabel
 - Low-degree polynomial data mapping
 - N-gram feature mapping
 - Random feature kernel machines

Problem with GCD on Primal

- Selection of greedy coordinate is expensive!!
- Naive implementation per iteration: $\mathcal{O}(nd)$
- Some heuristics NN ([Dhillon et al. 2015](#)), MIPS ([Karimireddy et al. 2018](#)) try to do the selection with time sublinear to d
- The only exception: gradient could be maintained: i.e. gradient is linear in x , $g(x) = Ax + \mathbf{b}$. Or $P(x) = \frac{1}{2}x^\top Ax + \mathbf{b}^\top x$
- since when $x^+ \leftarrow x + \delta_i \mathbf{e}_i$, $g^+ \leftarrow g + \delta_i \mathbf{a}_i$. Time complexity $\mathcal{O}(d)$

Primal-Dual formulation to the Rescue!

$$\min_{\mathbf{x} \in \mathbb{R}^d} \max_{\mathbf{y} \in \mathbb{R}^n} \left\{ L(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top A \mathbf{x} - g(\mathbf{y}) \right\}$$

Maintain $\mathbf{w} = A\mathbf{x}$ and $\mathbf{z} = A^\top \mathbf{y}$.

For each iteration,

Operation	cost
Compute full gradient $\partial_y L = A\mathbf{x} - g'(\mathbf{y})$	$\mathcal{O}(nd)$
Choose a greedy coordinate j for \mathbf{y}	$\mathcal{O}(n)$
$\mathbf{y}^+ \leftarrow \mathbf{y} + \delta_y \mathbf{e}_j$	$\mathcal{O}(1)$
$\mathbf{z}^+ \leftarrow \mathbf{z} + \delta_y A^\top \mathbf{e}_j$	$\mathcal{O}(d)$
vice versa for \mathbf{x} and \mathbf{w}	$\mathcal{O}(n + d)$

Primal-Dual formulation to the Rescue!

$$\min_{\mathbf{x} \in \mathbb{R}^d} \max_{\mathbf{y} \in \mathbb{R}^n} \left\{ L(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top A \mathbf{x} - g(\mathbf{y}) \right\}$$

Maintain $\mathbf{w} = A\mathbf{x}$ and $\mathbf{z} = A^\top \mathbf{y}$.

For each iteration,

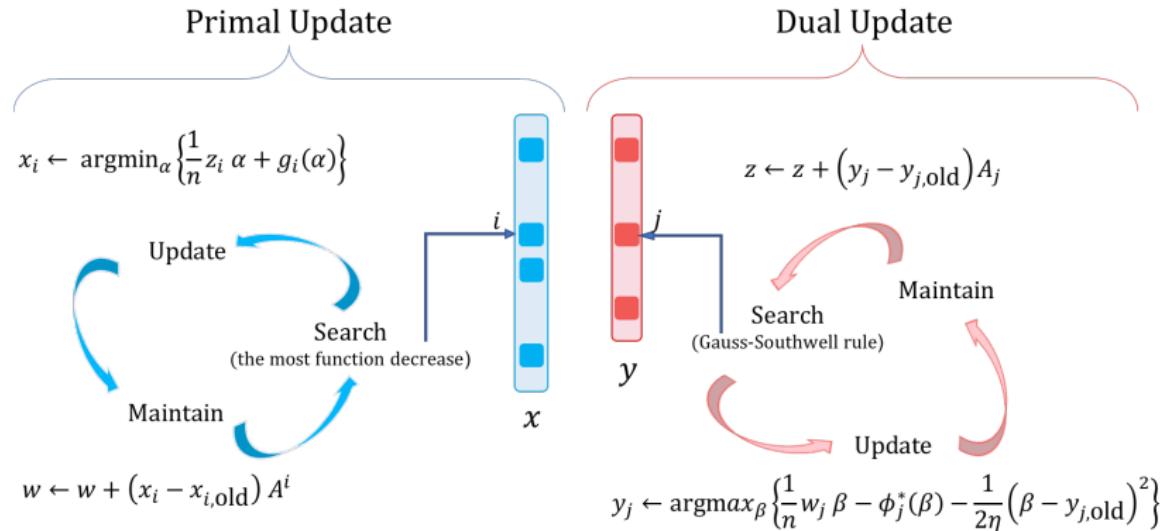
Operation	cost
Compute full gradient $\partial_{\mathbf{y}} L = \mathbf{w} - g'(\mathbf{y})$	$\mathcal{O}(n)$
Choose a greedy coordinate j for \mathbf{y}	$\mathcal{O}(n)$
$\mathbf{y}^+ \leftarrow \mathbf{y} + \delta_{\mathbf{y}} \mathbf{e}_j$	$\mathcal{O}(1)$
$\mathbf{z}^+ \leftarrow \mathbf{z} + \delta_{\mathbf{y}} A^\top \mathbf{e}_j$	$\mathcal{O}(d)$
vice versa for \mathbf{x} and \mathbf{w}	$\mathcal{O}(n + d)$

Our method: Doubly Greedy Primal-Dual Coordinate methods

- Recall we want to achieve

$$\max_{\mathbf{y} \in \mathbb{R}^n} \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ \mathcal{L}(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} g(\mathbf{x}) + \frac{1}{n} \mathbf{y}^\top A \mathbf{x} - \frac{1}{n} \sum_{i=1}^n \phi_i^*(y_i) \right\}$$

- Start from $\mathbf{x} \leftarrow 0, \mathbf{y} \leftarrow 0, \mathbf{w} (\equiv A\mathbf{x}) \leftarrow 0 \in \mathbb{R}^n, \mathbf{z} (\equiv A^\top \mathbf{y}) \leftarrow 0 \in \mathbb{R}^d$.
- $\mathcal{O}(d + n)$ operations per iteration



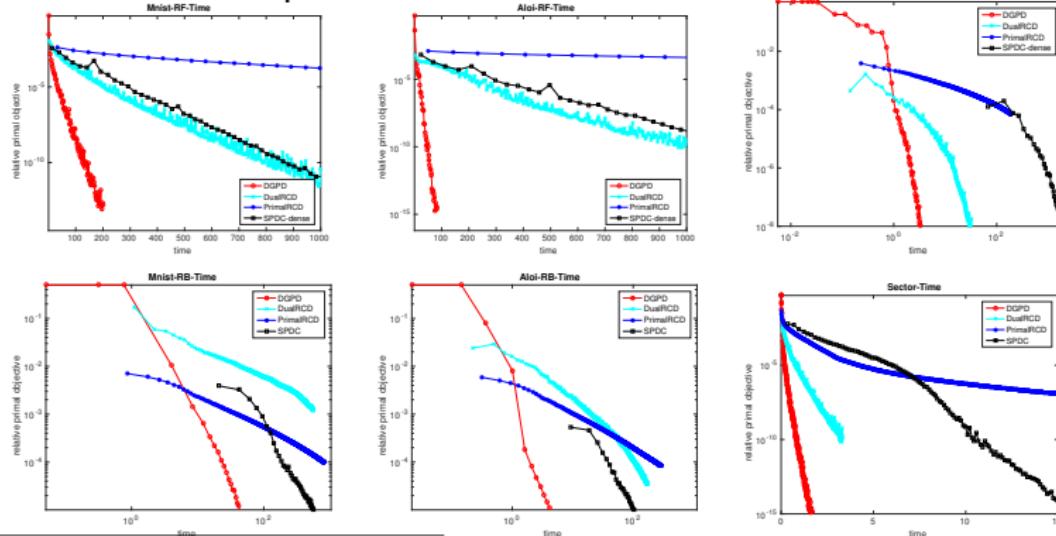
Time complexity comparisons

method	per Iteration	# Iteration
GD	$\mathcal{O}(nd)$	$\mathcal{O}(\kappa \log^{-1} \epsilon)$
RCD	$\mathcal{O}(n)$	$\mathcal{O}(d\kappa \log^{-1} \epsilon)$
SGD	$\mathcal{O}(d)$	$\mathcal{O}(\kappa \epsilon^{-1})$
MISO		$\mathcal{O}((n + \kappa) \log^{-1} \epsilon)$
SVRG		$\mathcal{O}((n + \sqrt{n\kappa}) \log^{-1} \epsilon)$
SAG		$\mathcal{O}((n + \sqrt{n\kappa}) \log^{-1} \epsilon)$ if accelerated
SAGA		
SDCA	$\mathcal{O}(d)$	$\mathcal{O}((n + \kappa) \log^{-1} \epsilon)$
SPDC	$\mathcal{O}(d)$	$\mathcal{O}((n + \sqrt{n\kappa}) \log^{-1} \epsilon)$
DGPDC	$\mathcal{O}(d + n)$	$\mathcal{O}\left(\frac{s}{n}(n + \kappa) \log^{-1} \epsilon\right)$

- DGPDC (Lei et al. 2017)

Experiments

- Compared methods:
 - Stochastic Dual Coordinate Ascent (SDCA, or DualRCD)
 - Primal Randomize Coordinate Descent (PrimalRCD)
 - Stochastic Primal-Dual Coordinate method (SPDC)
- Performance comparisons:¹



¹Relative objective versus time for $\lambda = 0.1, \mu = 0.01$. All datasets are downloaded from LIBSVM <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

Other trial on the Primal-Dual formulation

$$\min_{\mathbf{x} \in \mathbb{R}^d} \max_{\mathbf{y} \in \mathbb{R}^n} \left\{ L(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top A \mathbf{x} - g(\mathbf{y}) \right\}$$

method	per Iteration	# Iteration
SDCA	$\mathcal{O}(d)$	$\mathcal{O}((n + \kappa) \log^{-1} \epsilon)$
SPDC	$\mathcal{O}(d)$	$\mathcal{O}((n + \sqrt{n\kappa}) \log^{-1} \epsilon)$
DGPDC (ours)	$\mathcal{O}(d + n)$	$\mathcal{O}((s + \frac{s}{n}\kappa) \log^{-1} \epsilon)$
DSPDC	$\mathcal{O}(d)$	$\mathcal{O}((n + k\sqrt{n\kappa_1}) \log \frac{1}{\epsilon})$

- DSPDC (Yu et al. 2015), assuming data is factorized with rank k .

$$\kappa_1 \stackrel{\text{def}}{=} \frac{\|\mathbf{a}_i\|_\infty^2}{\lambda\mu} \in \left[\frac{\kappa}{d}, \kappa\right]$$

Q&A

Thank you!

Datasets

Table: Data statistics and number of non-zero primal & dual variables from DGPD (w/ $\lambda = 0.1$, $\mu = 0.01$).

Data set	#features	#nz/sample	#train samples	#test samples	#classes	#nz-primal	#nz-dual
Mnist-RF	10,000	10,000	58,000	2,000	10	1,730	2,000
Alois-RF	10,000	10,000	90,000	8,000	1,000	891	1,428
Mnist-RB	1,572,556	1,000	58,000	2,000	10	1,733	1,208
Alois-RB	636,910	200	90,000	8,000	1,000	1,032	782
RCV1-Regions Sector	47,236 55,197	68.38 162.94	199,328 7,793	23,149 961	225 105	1,123 610	1,447 655