

Primal-Dual Block Generalized Frank-Wolfe

Qi Lei*, Jiacheng Zhuo*, Constantine Caramanis*, Inderjit S. Dhillon*,[†]
and Alexandros G. Dimakis*.

* University of Texas at Austin

[†] Amazon

Convex-concave saddle point Problem (with constraints):

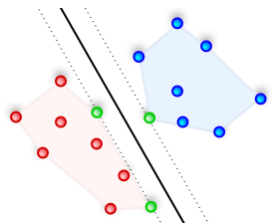
$$\min_{\mathbf{x} \in C \subset \mathbb{R}^d} \max_{\mathbf{y} \in \mathbb{R}^n} \{L(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top A \mathbf{x} - g(\mathbf{y})\}$$

Why is this formulation important?

- 1 Many machine learning applications

Machine Learning Applications with Convex-Concave Formulations

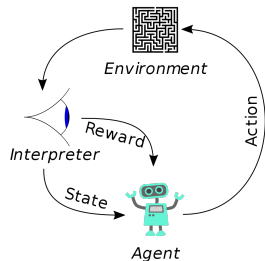
Empirical Risk Minimization



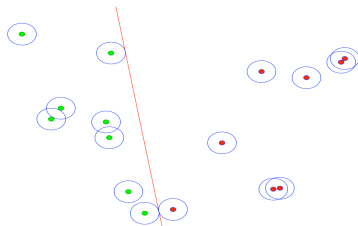
Logistic Regression Model



Reinforcement Learning



Robust Optimization



Convex-Concave Saddle Point Problem:

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

Why is this formulation important?

- 1 Many machine learning applications
- 2 To exploit special structure induced by the constraints

Observations and Challenges On Frank-Wolfe algorithm

Lessons from simple constrained minimization problems:

Observations.

Frank-Wolfe conducts *partial updates*:

1. For ℓ_1 ball constraint, FW conducts **1-sparse** update
2. For nuclear norm ball constraint, FW conducts **rank-1** update

Challenges to get full benefits from FW and the partial updates.

1. FW yield **sublinear convergence** even for strongly convex problems
2. Even with partial updates, FW requires to compute the full gradient.
(For big data setting, **per iteration complexity is the same with projected gradient descent.**)

Tackle challenge 1: To achieve linear convergence

- Continue to look at simple minimization problems:

$$\min_{\mathbf{x} \in \mathbb{R}^d, \|\mathbf{x}\|_1 \leq \tau} \{f(\mathbf{x})\}$$

Method	iteration complexity	# update per iteration
Projected GD	$\kappa \log \frac{1}{\epsilon}$	d (feature dimension)
Frank-Wolfe	$\frac{1}{\epsilon}$	1
Ours	$\kappa \log \frac{1}{\epsilon}$	s (optimal sparsity)

Tackle challenge 1: block Frank-Wolfe

1: **Input:** Data matrix $A \in \mathbb{R}^{n \times d}$, label matrix b , iteration T .

2: **Initialize:** $x_1 \leftarrow 0$.

3: **for** $t = 1, 2, \dots, T - 1$ **do**

4:

$$\text{ProjectedGD: } \Delta x_t \leftarrow \underset{\|\Delta x\|_1 \leq \tau}{\operatorname{argmin}} \left\{ \langle \nabla f(x_t), \Delta x \rangle + \frac{\beta}{2} \eta \|\Delta x - x_t\|_2^2 \right\}$$

5:

$$x_{t+1} \leftarrow (1 - \eta)x_t + \eta \Delta x_t$$

6: **end for**

7: **Output:** x_T

Tackle challenge 1: block Frank-Wolfe

1: **Input:** Data matrix $A \in \mathbb{R}^{n \times d}$, label matrix b , iteration T .

2: **Initialize:** $x_1 \leftarrow 0$.

3: **for** $t = 1, 2, \dots, T - 1$ **do**

4:

$$\text{ProjectedGD: } \Delta x_t \leftarrow \underset{\|\Delta x\|_1 \leq \tau}{\operatorname{argmin}} \left\{ \langle \nabla f(x_t), \Delta x \rangle + \frac{\beta}{2} \eta \|\Delta x - x_t\|_2^2 \right\}$$

$$\text{FW: } \Delta x_t \leftarrow \underset{\|\Delta x\|_1 \leq \tau}{\operatorname{argmin}} \left\{ \langle \nabla f(x_t), \Delta x \rangle \right\}$$

5:

$$x_{t+1} \leftarrow (1 - \eta)x_t + \eta \Delta x_t$$

6: **end for**

7: **Output:** x_T

Tackle challenge 1: block Frank-Wolfe

1: **Input:** Data matrix $A \in \mathbb{R}^{n \times d}$, label matrix b , iteration T .

2: **Initialize:** $x_1 \leftarrow 0$.

3: **for** $t = 1, 2, \dots, T - 1$ **do**

4:

$$\text{ProjectedGD: } \Delta x_t \leftarrow \underset{\|\Delta x\|_1 \leq \tau}{\operatorname{argmin}} \left\{ \langle \nabla f(x_t), \Delta x \rangle + \frac{\beta}{2} \eta \|\Delta x - x_t\|_2^2 \right\}$$

$$\text{FW: } \Delta x_t \leftarrow \underset{\|\Delta x\|_1 \leq \tau}{\operatorname{argmin}} \left\{ \langle \nabla f(x_t), \Delta x \rangle \right\}$$

$$\text{ours: } \Delta x_t \leftarrow \underset{\|\Delta x\|_1 \leq \tau, \|\Delta x\|_0 \leq s}{\operatorname{argmin}} \left\{ \langle \nabla f(x_t), \Delta x \rangle + \frac{\beta}{2} \eta \|\Delta x - x_t\|_2^2 \right\}$$

5:

$$x_{t+1} \leftarrow (1 - \eta)x_t + \eta \Delta x_t$$

6: **end for**

7: **Output:** x_T

Tackle challenge 2: reduce iteration complexity from partial updates

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

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

For each iteration,

Operation	cost
Compute full gradient $\partial_{\mathbf{x}} L = A^\top \mathbf{y} + f'(\mathbf{x})$	$\mathcal{O}(nd)$
Conduct BlockFW on \mathbf{x} to find s -sparse update $\Delta \mathbf{x}$	$\mathcal{O}(d)$
$\mathbf{x}^+ \leftarrow (1 - \eta)\mathbf{x} + \eta \Delta \mathbf{x}$	$\mathcal{O}(d)$
Greedy block- k coordinate ascent for \mathbf{y}	$\mathcal{O}(nd)$

Remark 1: take $k = ns/d$ the iteration complexity is $\mathcal{O}(sn)$.

Remark 2: the advantage comes from the fact that gradient could be maintained with the bilinear form

Tackle challenge 2: reduce iteration complexity from partial updates

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

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

For each iteration,

Operation	cost
Compute full gradient $\partial_{\mathbf{x}} L = \mathbf{z} + f'(\mathbf{x})$	$\mathcal{O}(d)$
Conduct BlockFW on \mathbf{x} to find s -sparse update $\Delta \mathbf{x}$	$\mathcal{O}(d)$
$\mathbf{x}^+ \leftarrow (1 - \eta) \mathbf{x} + \eta \Delta \mathbf{x}$	$\mathcal{O}(d)$
$\mathbf{w}^+ \leftarrow (1 - \eta) \mathbf{w} + \eta \mathbf{A} \Delta \mathbf{x}$	$\mathcal{O}(sn)$
Greedy block- k coordinate ascent for \mathbf{y} and \mathbf{z}	$\mathcal{O}(kd)$

Remark 1: take $k = ns/d$ the iteration complexity is $\mathcal{O}(sn)$.

Remark 2: the advantage comes from the fact that gradient could be maintained with the bilinear form

Time complexity comparisons

Algorithm	Per Iteration Cost	Iteration Complexity
Frank Wolfe	$\mathcal{O}(nd)$	$\mathcal{O}(\frac{1}{\epsilon})$
Accelerated PGD (Nesterov et al. 2013)	$\mathcal{O}(nd)$	$\mathcal{O}(\sqrt{\kappa} \log \frac{1}{\epsilon})$
SVRG (Rie et al. 2013)	$\mathcal{O}(nd)$	$\mathcal{O}((1 + \kappa/n) \log \frac{1}{\epsilon})$
SCGS (Lan et al. 2016)	$\mathcal{O}(\kappa^2 \frac{\#\text{iter}^3}{\epsilon^2} d)$	$\mathcal{O}(\frac{1}{\epsilon})$
STORC (Hazan et al. 2016)	$\mathcal{O}(\kappa^2 d + nd)$	$\mathcal{O}(\log \frac{1}{\epsilon})$
Primal Dual FW (ours)	$\mathcal{O}(ns)$	$\mathcal{O}((1 + \kappa/n) \log \frac{1}{\epsilon})$

Remark 1: s is the sparsity of primal optimal induced by ℓ_1 constraint.

Remark 2: for algorithm and complexity for nuclear norm constraints, refer to our paper to details.

Experiments

Compared methods: (1) Accelerated Projected Gradient Descent (Acc PG) (2) Frank-Wolfe algorithm (FW) (3) Stochastic Variance Reduced Gradient (SVRG) (4) Stochastic Conditional Gradient Sliding (SCGS) and (5) Stochastic Variance-Reduced Conditional Gradient Sliding (STORC)

