

Efficient Learning with Forward-Backward Splitting

John Duchi^{1,2} Yoram Singer²

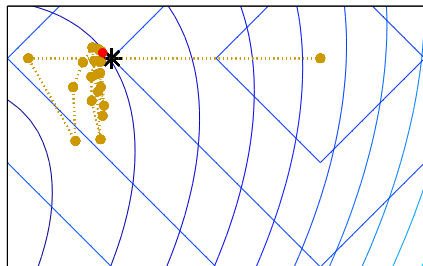
¹University of California, Berkeley

²Google

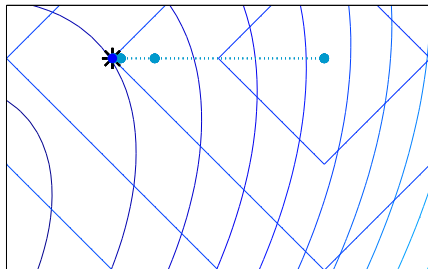
Neural Information Processing Systems, 2009

Motivating Example

Minimize $\frac{1}{2}\mathbf{w}^\top A\mathbf{w} + \mathbf{c}^\top \mathbf{w} + \lambda \|\mathbf{w}\|_1$. True solution: $\mathbf{w}^* = [-1 \ 0]^\top$.



Subgradient

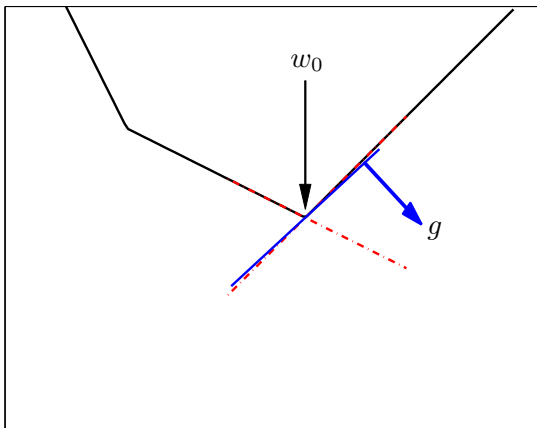


Fobos

Subgradients

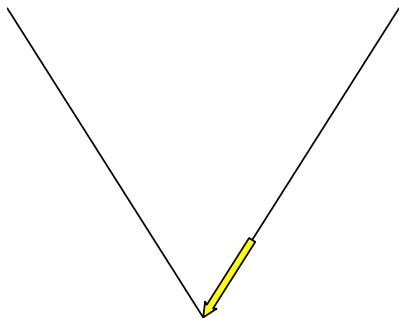
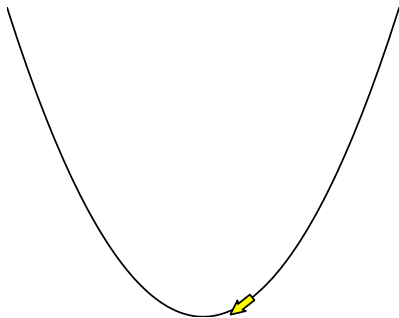
- ▶ Subgradient set of a function f

$$\partial f(\mathbf{w}_0) = \{ \mathbf{g} \in \mathbb{R}^n \mid f(\mathbf{w}) \geq f(\mathbf{w}_0) + \mathbf{g}^\top (\mathbf{w} - \mathbf{w}_0) \}$$



What is the problem?

- ▶ Subgradient set is large at singularities
- ▶ Subgradients are non-informative at singularities



Outline

Algorithmic Framework

Convergence and Regret

Derived Algorithms

Experimental Results

Conclusions and Related Work

The Fobos Algorithm

$$\text{Goal: } \min_{\mathbf{w}} L(\mathbf{w}) + R(\mathbf{w}).$$

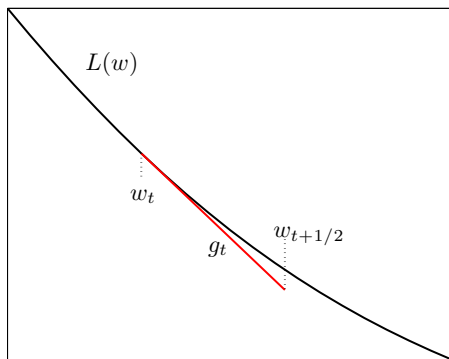
- ▶ Repeat
 - I. Unconstrained (stochastic sub) gradient of loss
 - II. Incorporate regularization
- ▶ Similar to forward-backward splitting (Lions and Mercier 79), composite gradient methods (Wright et al. 09, Nesterov 07), dual averaging with regularization (Xiao 09).

Fobos: Step I

Goal: $\min_{\mathbf{w}} L(\mathbf{w}) + R(\mathbf{w})$

- Unconstrained (stochastic sub) gradient of loss

$$\mathbf{w}_{t+\frac{1}{2}} = \mathbf{w}_t - \eta_t \mathbf{g}_t \quad \text{where} \quad \mathbb{E} \mathbf{g}_t \in \partial L(\mathbf{w}_t)$$

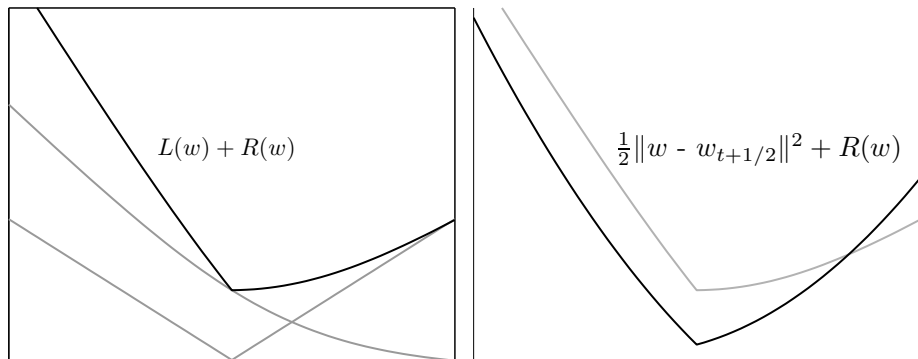


Fobos: Step II

Goal: $\min_{\mathbf{w}} L(\mathbf{w}) + R(\mathbf{w})$

- Incorporate regularization

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w}} \left\{ \frac{1}{2} \left\| \mathbf{w} - \mathbf{w}_{t+\frac{1}{2}} \right\|^2 + \eta_t R(\mathbf{w}) \right\}.$$



Forward Looking Property

- ▶ The optimum \mathbf{w}_{t+1} satisfies

$$\mathbf{0} \in \mathbf{w}_{t+1} - \mathbf{w}_t + \eta_t \partial L(\mathbf{w}_t) + \eta_t \partial R(\mathbf{w}_{t+1})$$

- ▶ Pick $\mathbf{g}_t^L \in \partial L(\mathbf{w}_t)$ and $\mathbf{g}_{t+1}^R \in \partial R(\mathbf{w}_{t+1})$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{g}_t^L - \eta_t \mathbf{g}_{t+1}^R$$

current loss

forward regularization

- ▶ *Current* subgradient of loss, *forward* subgradient of regularization

Batch Convergence and Online Regret

- ▶ Set $\eta_t \propto \frac{1}{\sqrt{T}}$ or $\frac{1}{\sqrt{t}}$ to obtain batch convergence

$$L(\mathbf{w}_t) + R(\mathbf{w}_t) - (L(\mathbf{w}^*) + R(\mathbf{w}^*)) = O\left(\frac{1}{\sqrt{T}}\right).$$

- ▶ Online (average) regret bounds

$$\text{Regret}(T) \triangleq \frac{1}{T} \left[\sum_{t=1}^T L_t(\mathbf{w}_t) + R(\mathbf{w}_t) - \sum_{t=1}^T L_t(\mathbf{w}^*) + R(\mathbf{w}^*) \right]$$

$$\eta_t \propto \frac{1}{\sqrt{t}} \Rightarrow \text{Regret}(T) = O\left(\frac{1}{\sqrt{T}}\right)$$

$$\eta_t \propto \frac{1}{t} \Rightarrow \text{Regret}(T) = O\left(\frac{\log T}{T}\right) \text{ (strong convexity)}$$

Derived Algorithms

We show step II for

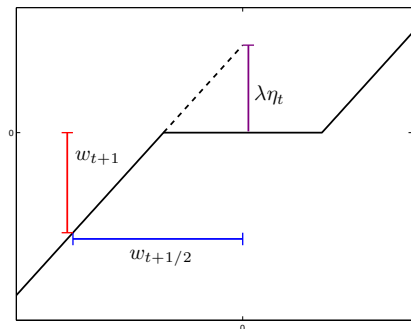
- ▶ FOBOS with ℓ_1 -regularization
- ▶ FOBOS with ℓ_2 -regularization
- ▶ FOBOS with mixed norms (ℓ_1/ℓ_2 or ℓ_1/ℓ_∞)

FOBOS with ℓ_1

$$\min \frac{1}{2} \left\| \mathbf{w} - \mathbf{w}_{t+\frac{1}{2}} \right\|^2 + \lambda \|\mathbf{w}\|_1$$

- ▶ Separable: minimize $\frac{1}{2} (w - w_{t+\frac{1}{2},j})^2 + \lambda |w|$.
- ▶ Coordinate-wise update yields sparsity:

$$w_{t+1,j} = \text{sign} (w_{t+\frac{1}{2},j}) \max \left\{ |w_{t+\frac{1}{2},j}| - \lambda \eta_t, 0 \right\}$$



Truncated gradient

(Langford et al. 08)

Iterative shrinkage and
thresholding

(Donoho 95, Daubechies et al. 04)

FOBOS with ℓ_2

- ▶ When $R(\mathbf{w}) = \frac{\lambda}{2} \|\mathbf{w}\|_2^2$, gradient descent & geometric shrinkage

$$\mathbf{w}_{t+1} = \frac{\mathbf{w}_{t+\frac{1}{2}}}{1 + \lambda\eta_t} = \frac{\mathbf{w}_t - \eta_t \mathbf{g}_t}{1 + \lambda\eta_t}$$

- ▶ When $R(\mathbf{w}) = \lambda \|\mathbf{w}\|_2$, all or nothing update

$$\mathbf{w}_{t+1} = \left[1 - \frac{\lambda\eta_t}{\|\mathbf{w}_{t+\frac{1}{2}}\|_2} \right]_+ \mathbf{w}_{t+\frac{1}{2}}$$

FOBOS with mixed norms

$$r(W) = \|W\|_{\ell_1/\ell_q} = \sum_{j=1}^d \|\bar{\mathbf{w}}_j\|_q$$

$$W = \begin{bmatrix} \bar{\mathbf{w}}_1 \\ \bar{\mathbf{w}}_2 \\ \vdots \\ \bar{\mathbf{w}}_d \end{bmatrix} \Rightarrow \begin{matrix} \|\bar{\mathbf{w}}_1\|_q \\ \|\bar{\mathbf{w}}_2\|_q \\ \vdots \\ \|\bar{\mathbf{w}}_d\|_q \end{matrix}$$

- ▶ Separable and solvable using previous methods
- ▶ Multitask and multiclass learning
 - ▶ $\bar{\mathbf{w}}_j$ associated with feature j
 - ▶ Penalize $\bar{\mathbf{w}}_j$ once

Sparse Gradients

	g				
$t = 1$	[1	3	0]
$t = 2$	[2	0	.5]
$t = 3$	[1	0	.5]
$t = 4$	[.1	0	-.25]
$t = 5$	[-.5	0	.25]

High Dimensional Efficiency

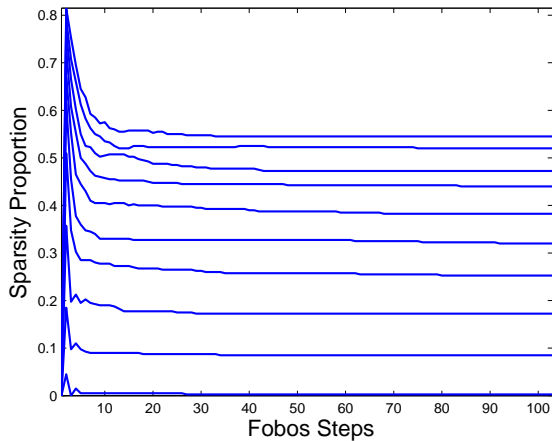
- ▶ Input space is sparse but huge
- ▶ Need to perform lazy updates to \mathbf{w}
- ▶ **Proposition:** The following are equivalent:

$$\mathbf{w}_t = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w} - \mathbf{w}_{t-1}\|^2 + \eta_t \lambda \|\mathbf{w}\|_q \quad \text{for } t = 1 \text{ to } T$$

$$\mathbf{w}_T = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w} - \mathbf{w}_0\|^2 + \left(\sum_{t=1}^{T-1} \eta_t \lambda \right) \|\mathbf{w}\|_q$$

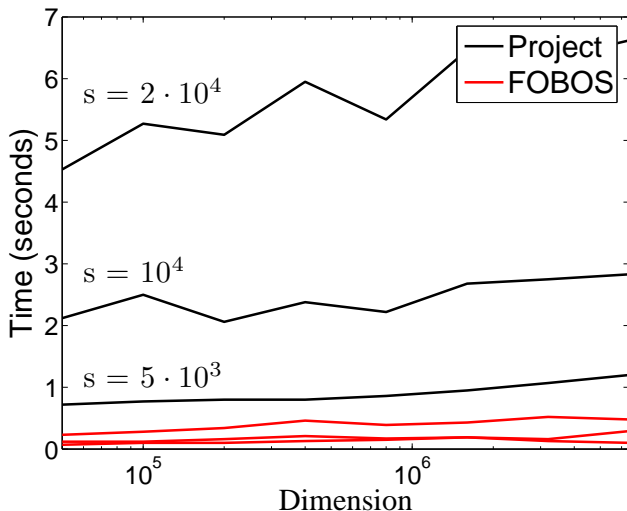
Experimental Results

Sparsity



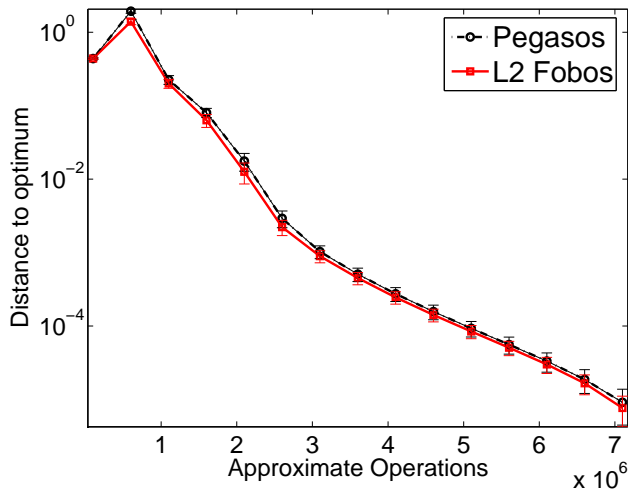
Sparsity as function of F_{OBOS} steps

Sparse timing experiments



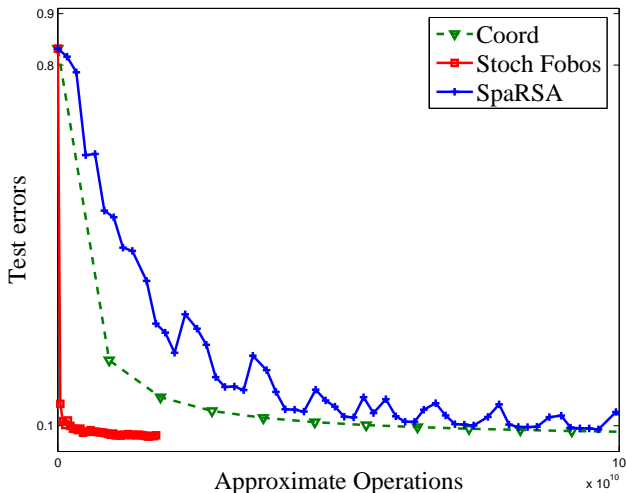
Comparison of ℓ_1 -projection to FOBOS update

ℓ_2^2 regularized experiments



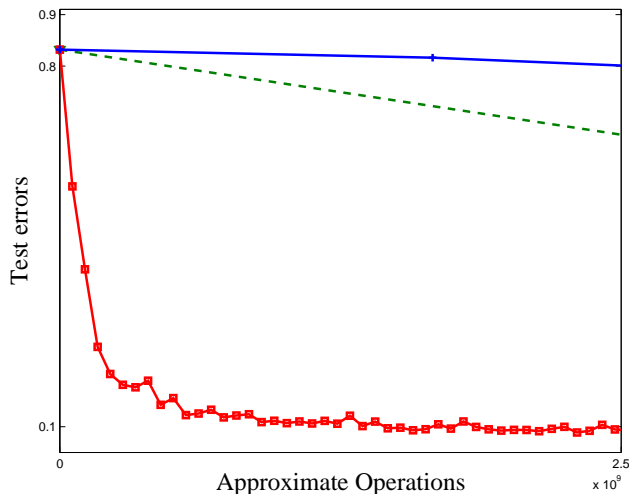
Convergence of FOBOS versus Pegasus on ℓ_2^2 regularized problem

MNIST experiments



Comparison of test error rate of FOBOS, Sparsa (Wright et al. 2009), coordinate descent (Tseng 2007).

MNIST experiments



Comparison of test error rate of FOBOS, Sparsa (Wright et al. 2009), coordinate descent (Tseng 2007).

Conclusions

- ▶ General framework for stochastic gradient with arbitrary regularization
- ▶ Regret bounds and convergence rates
- ▶ Extensions for mixed-norm regularization
- ▶ Future: hierarchical and structured models, faster algorithms.
- ▶ Long version to appear in *Journal of Machine Learning Research*

Thanks!