

Stochastic Methods for L1 Regularized Loss Minimization

S. Shalev-Shwartz and A. *Tewari*

Toyota Technological Institute at Chicago

June 16, 2009

L1 regularized loss minimization

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \underbrace{\frac{1}{m} \sum_{i=1}^m L(\langle \mathbf{w}, \mathbf{x}_i \rangle, y_i)}_{\mathcal{L}_m(\mathbf{w})} + \underbrace{\lambda \|\mathbf{w}\|_1}_{L_1 \text{ penalty}}$$

Examples:

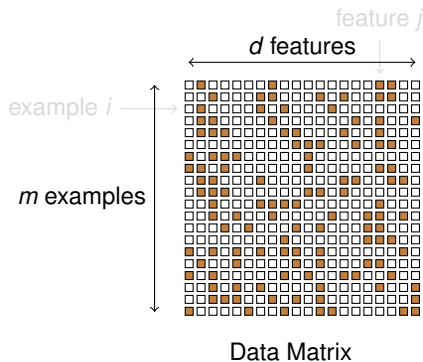
- **LASSO**: $L(a, y) = (a - y)^2$
- L_1 regularized **logistic regression**:
 $L(a, y) = \log(1 + \exp(-ay))$
- L_1 regularized **hinge-loss**: $L(a, y) = \max\{0, 1 - ya\}$

Our Contribution

- L_1 penalty encourages sparsity
- With $L(\cdot, y)$ convex, results in a convex optimization problem: problem solved?
- Run-time of IP methods scales typically as $\max\{d^3, m^3\}$: this is bad for large datasets
- Can live with worse dependence on $\frac{1}{\epsilon}$
- Our contribution: Two practical methods for solving L_1 regularized problem for large datasets

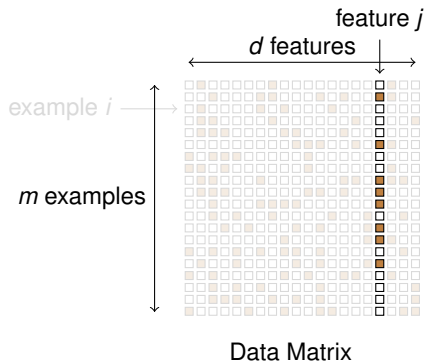
Main Theme – Stochastic Updates

Idea: Many simple updates are more efficient than few sophisticated updates and randomness helps!



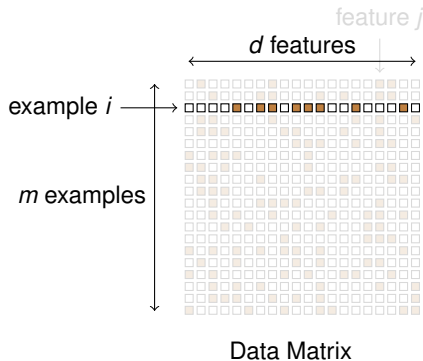
Main Theme – Stochastic Updates

Idea: Many simple updates are more efficient than few sophisticated updates and randomness helps!



Main Theme – Stochastic Updates

Idea: Many simple updates are more efficient than few sophisticated updates and randomness helps!



Outline

- 1 **Stochastic Coordinate Descent**
 - Coordinate Descent for L1
 - Stochastic Coordinate Descent
 - Efficient Implementation
 - Run-time Bound
- 2 **SMIDAS**
 - Stochastic Gradient Descent
 - SMIDAS
 - Efficient Implementation
 - Run-time Bound
- 3 **Experiments**
- 4 **Summary**

Outline

- 1 **Stochastic Coordinate Descent**
 - Coordinate Descent for L1
 - Stochastic Coordinate Descent
 - Efficient Implementation
 - Run-time Bound
- 2 **SMIDAS**
 - Stochastic Gradient Descent
 - SMIDAS
 - Efficient Implementation
 - Run-time Bound
- 3 Experiments
- 4 Summary

Outline

- 1 **Stochastic Coordinate Descent**
 - Coordinate Descent for L1
 - Stochastic Coordinate Descent
 - Efficient Implementation
 - Run-time Bound
- 2 **SMIDAS**
 - Stochastic Gradient Descent
 - SMIDAS
 - Efficient Implementation
 - Run-time Bound
- 3 **Experiments**
- 4 **Summary**

Outline

- 1 Stochastic Coordinate Descent
 - Coordinate Descent for L1
 - Stochastic Coordinate Descent
 - Efficient Implementation
 - Run-time Bound
- 2 SMIDAS
 - Stochastic Gradient Descent
 - SMIDAS
 - Efficient Implementation
 - Run-time Bound
- 3 Experiments
- 4 Summary

Outline

- 1 Stochastic Coordinate Descent
 - Coordinate Descent for L1
 - Stochastic Coordinate Descent
 - Efficient Implementation
 - Run-time Bound
- 2 SMIDAS
 - Stochastic Gradient Descent
 - SMIDAS
 - Efficient Implementation
 - Run-time Bound
- 3 Experiments
- 4 Summary

Coordinate Descent

Goal: To minimize $R(\mathbf{w}) = R(w_1, \dots, w_d)$

Generic CD Algorithm

FOR $t = 1, 2, \dots$

[SELECTION STEP]

Choose coordinate $j \in \{1, \dots, d\}$

[UPDATE STEP]

Update w_j so that $R(\mathbf{w})$ decreases (descent)

Coordinate Descent

Goal: To minimize $R(\mathbf{w}) = R(w_1, \dots, w_d)$

Generic CD Algorithm

FOR $t = 1, 2, \dots$

[SELECTION STEP]

Choose coordinate $j \in \{1, \dots, d\}$

[UPDATE STEP]

Update w_j so that $R(\mathbf{w})$ decreases (descent)

Coordinate Descent

Goal: To minimize $R(\mathbf{w}) = R(w_1, \dots, w_d)$

Generic CD Algorithm

FOR $t = 1, 2, \dots$

[SELECTION STEP]

Choose coordinate $j \in \{1, \dots, d\}$

[UPDATE STEP]

Update w_j so that $R(\mathbf{w})$ decreases (descent)

Coordinate Descent

Goal: To minimize $R(\mathbf{w}) = R(w_1, \dots, w_d)$

Generic CD Algorithm

FOR $t = 1, 2, \dots$

[SELECTION STEP]

Choose coordinate $j \in \{1, \dots, d\}$

[UPDATE STEP]

Update w_j so that $R(\mathbf{w})$ decreases (descent)

Coordinate Descent

Goal: To minimize $R(\mathbf{w}) = R(w_1, \dots, w_d)$

Generic CD Algorithm

FOR $t = 1, 2, \dots$

[SELECTION STEP]

Choose coordinate $j \in \{1, \dots, d\}$

[UPDATE STEP]

Update w_j so that $R(\mathbf{w})$ decreases (descent)

Coordinate Descent

Goal: To minimize $R(\mathbf{w}) = R(w_1, \dots, w_d)$

Generic CD Algorithm

FOR $t = 1, 2, \dots$

[SELECTION STEP]

Choose **coordinate** $j \in \{1, \dots, d\}$

[UPDATE STEP]

Update w_j so that $R(\mathbf{w})$ decreases (**descent**)

Coordinate Descent for L1

Coordinate-wise descent algorithms deserve more attention in convex optimization. They are simple and well-suited to large problems.

– Friedman et al., 2007

Our reasons for liking coordinate descent boil down to simplicity, speed and stability.

– Wu & Lange, 2008

Lots of Recent Activity

- Genkin, Lewis, Madigan (2007)* Minimize along 1 coordinate
(using quadratic approximation)
- Friedman et al. (2007)* show empirically that CD is competitive
cycle through coordinates
- Wu and Lange (2008)* cyclic & greedy coordinate selection rules
proof of convergence
- Tseng and Yun (2009)* theoretical analysis
sophisticated coordinate selection rule
rule takes $O(d)$ time

A Simple Transformation

Recall we're minimizing

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m L(\langle \mathbf{w}, \mathbf{x}_i \rangle, y_i) + \lambda \|\mathbf{w}\|_1$$

Doubling dimension by letting $\hat{\mathbf{x}}_i = [\mathbf{x}_i; -\mathbf{x}_i]$, we get rid of non-differentiability of objective

Equivalent Problem

$$\min_{\mathbf{w} \in \mathbb{R}^{2d}} R(\mathbf{w}) \quad \text{s.t.} \quad \mathbf{w} \geq \mathbf{0}$$

where

$$R(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m L(\langle \mathbf{w}, \hat{\mathbf{x}}_i \rangle, y_i) + \lambda \sum_{j=1}^{2d} w_j$$

Stochastic Coordinate Descent

SCD Algorithm

$\beta =$ upper bound on $L''(\cdot, y)$ (e.g. $\beta = 1$ for LASSO)

FOR $t = 1, 2, \dots$

[SELECTION STEP]

Choose j randomly from $\{1, \dots, 2d\}$ simple rule; takes $O(1)$ time!

[UPDATE STEP]

$$g_j = (\nabla R(\mathbf{w}))_j$$

$$w_j = \max\{0, w_j - g_j/\beta\}$$

j th entry in the gradient
to keep $w_j \geq 0$

Stochastic Coordinate Descent

SCD Algorithm

$\beta =$ upper bound on $L''(\cdot, y)$ (e.g. $\beta = 1$ for LASSO)

FOR $t = 1, 2, \dots$

[SELECTION STEP]

Choose j randomly from $\{1, \dots, 2d\}$ simple rule; takes $O(1)$ time!

[UPDATE STEP]

$$g_j = (\nabla R(\mathbf{w}))_j$$

$$w_j = \max\{0, w_j - g_j/\beta\}$$

j th entry in the gradient
to keep $w_j \geq 0$

Stochastic Coordinate Descent

SCD Algorithm

$\beta =$ upper bound on $L''(\cdot, y)$ (e.g. $\beta = 1$ for LASSO)

FOR $t = 1, 2, \dots$

[SELECTION STEP]

Choose j randomly from $\{1, \dots, 2d\}$ simple rule; takes $O(1)$ time!

[UPDATE STEP]

$$g_j = (\nabla R(\mathbf{w}))_j$$

$$w_j = \max\{0, w_j - g_j/\beta\}$$

j th entry in the gradient
to keep $w_j \geq 0$

Stochastic Coordinate Descent

SCD Algorithm

$\beta =$ upper bound on $L''(\cdot, y)$ (e.g. $\beta = 1$ for LASSO)

FOR $t = 1, 2, \dots$

[SELECTION STEP]

Choose j randomly from $\{1, \dots, 2d\}$ simple rule; takes $O(1)$ time!

[UPDATE STEP]

$g_j = (\nabla R(\mathbf{w}))_j$ j th entry in the gradient
 $w_j = \max\{0, w_j - g_j/\beta\}$ to keep $w_j \geq 0$

Stochastic Coordinate Descent

SCD Algorithm

$\beta =$ upper bound on $L''(\cdot, y)$ (e.g. $\beta = 1$ for LASSO)

FOR $t = 1, 2, \dots$

[SELECTION STEP]

Choose j randomly from $\{1, \dots, 2d\}$ simple rule; takes $O(1)$ time!

[UPDATE STEP]

$g_j = (\nabla R(\mathbf{w}))_j$ j th entry in the gradient
 $w_j = \max\{0, w_j - g_j/\beta\}$ to keep $w_j \geq 0$

Stochastic Coordinate Descent

SCD Algorithm

$\beta =$ upper bound on $L''(\cdot, y)$ (e.g. $\beta = 1$ for LASSO)

FOR $t = 1, 2, \dots$

[SELECTION STEP]

Choose j randomly from $\{1, \dots, 2d\}$ simple rule; takes $O(1)$ time!

[UPDATE STEP]

$$g_j = (\nabla R(\mathbf{w}))_j$$

$$w_j = \max\{0, w_j - g_j/\beta\}$$

j th entry in the gradient
to keep $w_j \geq 0$

Stochastic Coordinate Descent

SCD Algorithm

$\beta =$ upper bound on $L''(\cdot, y)$ (e.g. $\beta = 1$ for LASSO)

FOR $t = 1, 2, \dots$

[SELECTION STEP]

Choose j randomly from $\{1, \dots, 2d\}$ simple rule; takes $O(1)$ time!

[UPDATE STEP]

$$g_j = (\nabla R(\mathbf{w}))_j$$

$$w_j = \max\{0, w_j - g_j/\beta\}$$

j th entry in the gradient
to keep $w_j \geq 0$

Stochastic Coordinate Descent

SCD Algorithm

$\beta =$ upper bound on $L''(\cdot, y)$ (e.g. $\beta = 1$ for LASSO)

FOR $t = 1, 2, \dots$

[SELECTION STEP]

Choose j randomly from $\{1, \dots, 2d\}$ simple rule; takes $O(1)$ time!

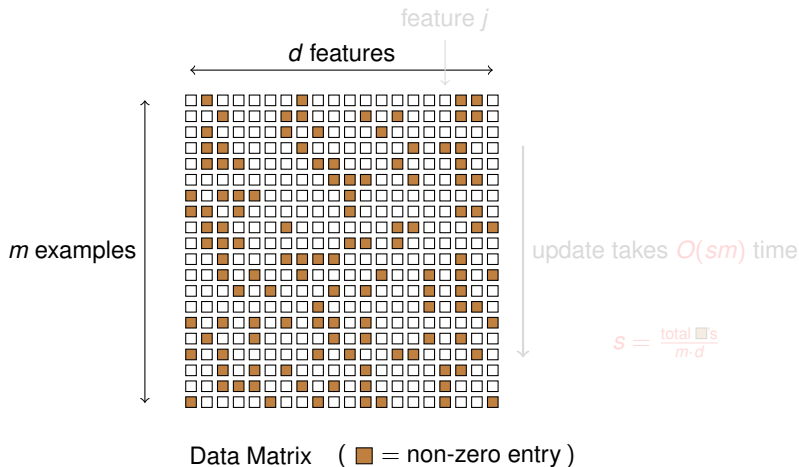
[UPDATE STEP]

$$g_j = (\nabla R(\mathbf{w}))_j$$

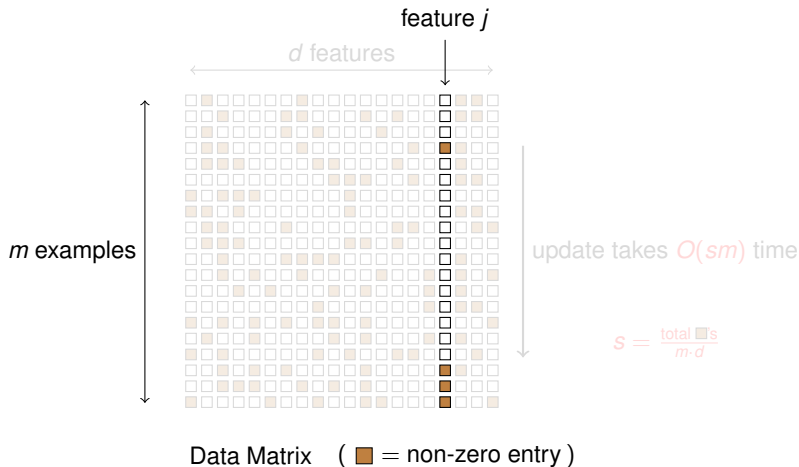
$$w_j = \max\{0, w_j - g_j/\beta\}$$

j th entry in the gradient
to keep $w_j \geq 0$

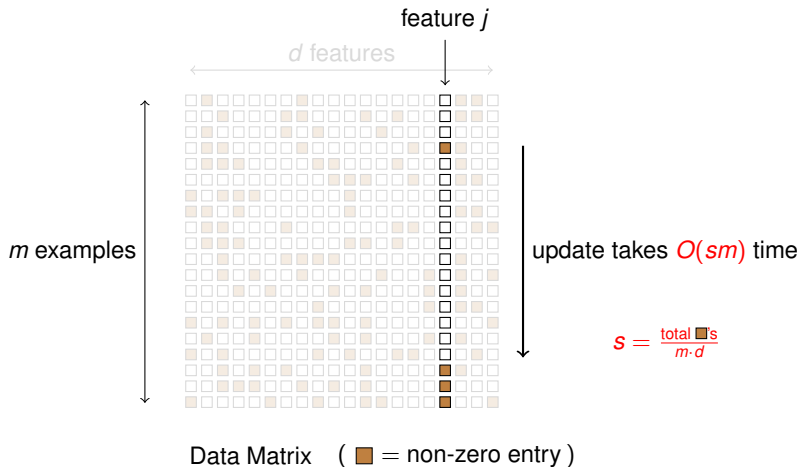
Implementing Updates Efficiently



Implementing Updates Efficiently



Implementing Updates Efficiently



Run-time Guarantee

- Simple selection rule
- Parameter free
- No Hessian computations; no line searches
- Yet, best convergence guarantee: time to achieve (expected) ϵ -accuracy is

$$O\left(\frac{m d \beta \|\mathbf{w}^*\|_2^2}{\epsilon}\right)$$

- We save a factor of d compared to Tseng-Yun's (2009) bound

$$O\left(\frac{m d^2 \beta \|\mathbf{w}^*\|_2^2}{\epsilon}\right)$$

Run-time Guarantee

- Simple selection rule
- Parameter free
- No Hessian computations; no line searches
- Yet, best convergence guarantee: time to achieve (expected) ϵ -accuracy is

$$O\left(\frac{m d \beta \|\mathbf{w}^*\|_2^2}{\epsilon}\right)$$

- We save a factor of d compared to Tseng-Yun's (2009) bound

$$O\left(\frac{m d^2 \beta \|\mathbf{w}^*\|_2^2}{\epsilon}\right)$$

Run-time Guarantee

- Simple selection rule
- Parameter free
- No Hessian computations; no line searches
- Yet, best convergence guarantee: time to achieve (expected) ϵ -accuracy is

$$O\left(\frac{m d \beta \|\mathbf{w}^*\|_2^2}{\epsilon}\right)$$

- We save a factor of d compared to Tseng-Yun's (2009) bound

$$O\left(\frac{m d^2 \beta \|\mathbf{w}^*\|_2^2}{\epsilon}\right)$$

Proof Idea

- A common proof strategy is to find a *potential* that decreases along the way
- Common potentials
 - Distance from the minimum $\|\mathbf{w} - \mathbf{w}^*\|_2^2$
 - Objective function itself $R(\mathbf{w})$
- Analysis of SCD simplifies when we use the “double potential”

$$\frac{\beta}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 + R(\mathbf{w})$$

Proof Idea

- A common proof strategy is to find a *potential* that decreases along the way
- Common potentials
 - Distance from the minimum $\|\mathbf{w} - \mathbf{w}^*\|_2^2$
 - Objective function itself $R(\mathbf{w})$
- Analysis of SCD simplifies when we use the “double potential”

$$\frac{\beta}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 + R(\mathbf{w})$$

Outline

- 1 Stochastic Coordinate Descent
 - Coordinate Descent for L1
 - Stochastic Coordinate Descent
 - Efficient Implementation
 - Run-time Bound
- 2 SMIDAS
 - Stochastic Gradient Descent
 - SMIDAS
 - Efficient Implementation
 - Run-time Bound
- 3 Experiments
- 4 Summary

Stochastic Gradient Descent

$$\nabla \mathcal{L}_m(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m L'(\langle \mathbf{w}, \mathbf{x}_i \rangle, y_i) \mathbf{x}_i$$

- Get unbiased estimate of gradient using one *random* example
- Use it to perform a gradient descent step
- Avoids incurring m dependence
- Shown to be state-of-the-art for large-dataset regime

SGD for L1

$$\nabla R(\mathbf{w}) = \nabla \mathcal{L}_m(\mathbf{w}) + \lambda \text{sign}(\mathbf{w})$$

- SGD applied directly fails to give sparsity
- One fix: work with constrained formulation (Duchi et al., 2008)

$$\min_{\mathbf{w}} \mathcal{L}_m(\mathbf{w}) \quad \text{s.t.} \quad \|\mathbf{w}\|_1 \leq B$$

- Second fix: truncate w_i to 0 if it crosses 0 during SGD update (Langford et al., 2009)
- SMIDAS combines Mirror Descent with Langford et al.'s *truncated gradient* idea

SGD for L1

$$\nabla R(\mathbf{w}) = \nabla \mathcal{L}_m(\mathbf{w}) + \lambda \text{sign}(\mathbf{w})$$

- SGD applied directly fails to give sparsity
- One fix: work with constrained formulation (Duchi et al., 2008)

$$\min_{\mathbf{w}} \mathcal{L}_m(\mathbf{w}) \quad \text{s.t.} \quad \|\mathbf{w}\|_1 \leq B$$

- Second fix: truncate w_i to 0 if it crosses 0 during SGD update (Langford et al., 2009)
- SMIDAS combines Mirror Descent with Langford et al.'s *truncated gradient* idea

SGD for L1

$$\nabla R(\mathbf{w}) = \nabla \mathcal{L}_m(\mathbf{w}) + \lambda \text{sign}(\mathbf{w})$$

- SGD applied directly fails to give sparsity
- One fix: work with constrained formulation (Duchi et al., 2008)

$$\min_{\mathbf{w}} \mathcal{L}_m(\mathbf{w}) \quad \text{s.t.} \quad \|\mathbf{w}\|_1 \leq B$$

- Second fix: truncate w_i to 0 if it crosses 0 during SGD update (Langford et al., 2009)
- SMIDAS combines Mirror Descent with Langford et al.'s *truncated gradient* idea

SGD for L1

$$\nabla R(\mathbf{w}) = \nabla \mathcal{L}_m(\mathbf{w}) + \lambda \text{sign}(\mathbf{w})$$

- SGD applied directly fails to give sparsity
- One fix: work with constrained formulation (Duchi et al., 2008)

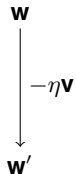
$$\min_{\mathbf{w}} \mathcal{L}_m(\mathbf{w}) \quad \text{s.t.} \quad \|\mathbf{w}\|_1 \leq B$$

- Second fix: truncate w_i to 0 if it crosses 0 during SGD update (Langford et al., 2009)
- SMIDAS combines Mirror Descent with Langford et al.'s *truncated gradient* idea

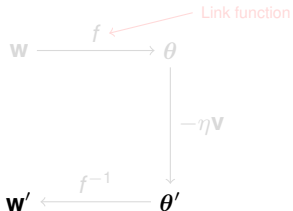
Mirror Descent

\mathbf{v} = (estimate of) gradient

Gradient Descent



Mirror Descent

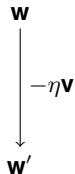


We'll use the q -norm link function, i.e. $f(\mathbf{w}) = \nabla \|\mathbf{w}\|_q^2$

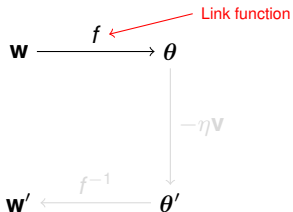
Mirror Descent

\mathbf{v} = (estimate of) gradient

Gradient Descent



Mirror Descent

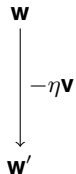


We'll use the q -norm link function, i.e. $f(\mathbf{w}) = \nabla \|\mathbf{w}\|_q^2$

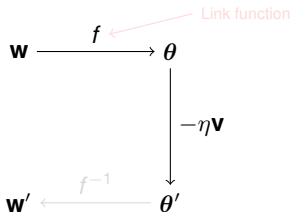
Mirror Descent

\mathbf{v} = (estimate of) gradient

Gradient Descent



Mirror Descent

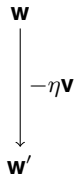


We'll use the q -norm link function, i.e. $f(\mathbf{w}) = \nabla \|\mathbf{w}\|_q^2$

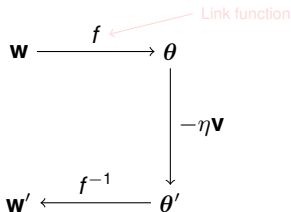
Mirror Descent

\mathbf{v} = (estimate of) gradient

Gradient Descent



Mirror Descent

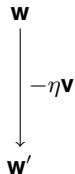


We'll use the q -norm link function, i.e. $f(\mathbf{w}) = \nabla \|\mathbf{w}\|_q^2$

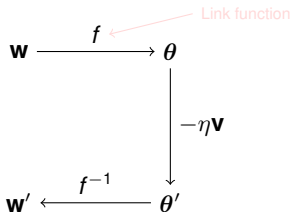
Mirror Descent

\mathbf{v} = (estimate of) gradient

Gradient Descent



Mirror Descent



We'll use the q -norm link function, i.e. $f(\mathbf{w}) = \nabla \|\mathbf{w}\|_q^2$

Stochastic **M**irror **D**escent **A**lgorithm made **S**pars

SMIDAS

Parameter: $\eta > 0$

$$\theta = \mathbf{0}, \mathbf{w} = f^{-1}(\theta)$$

FOR $t = 1, 2, \dots$

[GRADIENT ESTIMATION STEP]

Choose i randomly from $\{1, \dots, m\}$

$$\mathbf{v} = L'(\langle \mathbf{w}, \mathbf{x}_i \rangle, y_i) \mathbf{x}_i$$

[TRUNCATED GRADIENT STEP]

$$\tilde{\theta} = \theta - \eta \mathbf{v}$$

$$\forall j, \theta_j = \text{sign}(\tilde{\theta}_j) \max\{0, |\tilde{\theta}_j| - \eta \lambda\}$$

truncate if 0 is crossed

[MIRROR STEP]

$$\mathbf{w} = f^{-1}(\theta)$$

preserves sparsity pattern!

Stochastic **M**irror **D**escent **A**lgorithm made **S**pars

SMIDAS

Parameter: $\eta > 0$

$$\theta = \mathbf{0}, \mathbf{w} = f^{-1}(\theta)$$

FOR $t = 1, 2, \dots$

[GRADIENT ESTIMATION STEP]

Choose i randomly from $\{1, \dots, m\}$

$$\mathbf{v} = L'(\langle \mathbf{w}, \mathbf{x}_i \rangle, y_i) \mathbf{x}_i$$

[TRUNCATED GRADIENT STEP]

$$\tilde{\theta} = \theta - \eta \mathbf{v}$$

$$\forall j, \theta_j = \text{sign}(\tilde{\theta}_j) \max\{0, |\tilde{\theta}_j| - \eta \lambda\} \quad \text{truncate if 0 is crossed}$$

[MIRROR STEP]

$$\mathbf{w} = f^{-1}(\theta) \quad \text{preserves sparsity pattern!}$$

Stochastic **M**irror **D**escent **A**lgorithm made **S**pars

SMIDAS

Parameter: $\eta > 0$

$$\theta = \mathbf{0}, \mathbf{w} = f^{-1}(\theta)$$

FOR $t = 1, 2, \dots$

[GRADIENT ESTIMATION STEP]

Choose i randomly from $\{1, \dots, m\}$

$$\mathbf{v} = L'(\langle \mathbf{w}, \mathbf{x}_i \rangle, y_i) \mathbf{x}_i$$

[TRUNCATED GRADIENT STEP]

$$\tilde{\theta} = \theta - \eta \mathbf{v}$$

$$\forall j, \theta_j = \text{sign}(\tilde{\theta}_j) \max\{0, |\tilde{\theta}_j| - \eta \lambda\}$$

truncate if 0 is crossed

[MIRROR STEP]

$$\mathbf{w} = f^{-1}(\theta)$$

preserves sparsity pattern!

Stochastic **M**irror **D**escent **A**lgorithm made **S**pars

SMIDAS

Parameter: $\eta > 0$

$$\theta = \mathbf{0}, \mathbf{w} = f^{-1}(\theta)$$

FOR $t = 1, 2, \dots$

[GRADIENT ESTIMATION STEP]

Choose i randomly from $\{1, \dots, m\}$

$$\mathbf{v} = L'(\langle \mathbf{w}, \mathbf{x}_i \rangle, y_i) \mathbf{x}_i$$

[TRUNCATED GRADIENT STEP]

$$\tilde{\theta} = \theta - \eta \mathbf{v}$$

$$\forall j, \theta_j = \text{sign}(\tilde{\theta}_j) \max\{0, |\tilde{\theta}_j| - \eta \lambda\}$$

truncate if 0 is crossed

[MIRROR STEP]

$$\mathbf{w} = f^{-1}(\theta)$$

preserves sparsity pattern!

Stochastic **M**irror **D**escent **A**lgorithm made **S**pars

SMIDAS

Parameter: $\eta > 0$

$$\theta = \mathbf{0}, \mathbf{w} = f^{-1}(\theta)$$

FOR $t = 1, 2, \dots$

[GRADIENT ESTIMATION STEP]

Choose i randomly from $\{1, \dots, m\}$

$$\mathbf{v} = L'(\langle \mathbf{w}, \mathbf{x}_i \rangle, y_i) \mathbf{x}_i$$

[TRUNCATED GRADIENT STEP]

$$\tilde{\theta} = \theta - \eta \mathbf{v}$$

$$\forall j, \theta_j = \text{sign}(\tilde{\theta}_j) \max\{0, |\tilde{\theta}_j| - \eta \lambda\}$$

truncate if 0 is crossed

[MIRROR STEP]

$$\mathbf{w} = f^{-1}(\theta)$$

preserves sparsity pattern!

Stochastic **M**irror **D**escent **A**lgorithm made **S**pars

SMIDAS

Parameter: $\eta > 0$

$$\theta = \mathbf{0}, \mathbf{w} = f^{-1}(\theta)$$

FOR $t = 1, 2, \dots$

[GRADIENT ESTIMATION STEP]

Choose i randomly from $\{1, \dots, m\}$

$$\mathbf{v} = L'(\langle \mathbf{w}, \mathbf{x}_i \rangle, y_i) \mathbf{x}_i$$

[TRUNCATED GRADIENT STEP]

$$\tilde{\theta} = \theta - \eta \mathbf{v}$$

$$\forall j, \theta_j = \text{sign}(\tilde{\theta}_j) \max\{0, |\tilde{\theta}_j| - \eta\lambda\}$$

truncate if 0 is crossed

[MIRROR STEP]

$$\mathbf{w} = f^{-1}(\theta)$$

preserves sparsity pattern!

Stochastic **M**irror **D**escent **A**lgorithm made **S**pars

SMIDAS

Parameter: $\eta > 0$

$$\theta = \mathbf{0}, \mathbf{w} = f^{-1}(\theta)$$

FOR $t = 1, 2, \dots$

[GRADIENT ESTIMATION STEP]

Choose i randomly from $\{1, \dots, m\}$

$$\mathbf{v} = L'(\langle \mathbf{w}, \mathbf{x}_i \rangle, y_i) \mathbf{x}_i$$

[TRUNCATED GRADIENT STEP]

$$\tilde{\theta} = \theta - \eta \mathbf{v}$$

$$\forall j, \theta_j = \text{sign}(\tilde{\theta}_j) \max\{0, |\tilde{\theta}_j| - \eta \lambda\}$$

truncate if 0 is crossed

[MIRROR STEP]

$$\mathbf{w} = f^{-1}(\theta)$$

preserves sparsity pattern!

Stochastic **M**irror **D**escent **A**lgorithm made **S**pars

SMIDAS

Parameter: $\eta > 0$

$$\theta = \mathbf{0}, \mathbf{w} = f^{-1}(\theta)$$

FOR $t = 1, 2, \dots$

[GRADIENT ESTIMATION STEP]

Choose i randomly from $\{1, \dots, m\}$

$$\mathbf{v} = L'(\langle \mathbf{w}, \mathbf{x}_i \rangle, y_i) \mathbf{x}_i$$

[TRUNCATED GRADIENT STEP]

$$\tilde{\theta} = \theta - \eta \mathbf{v}$$

$$\forall j, \theta_j = \text{sign}(\tilde{\theta}_j) \max\{0, |\tilde{\theta}_j| - \eta \lambda\}$$

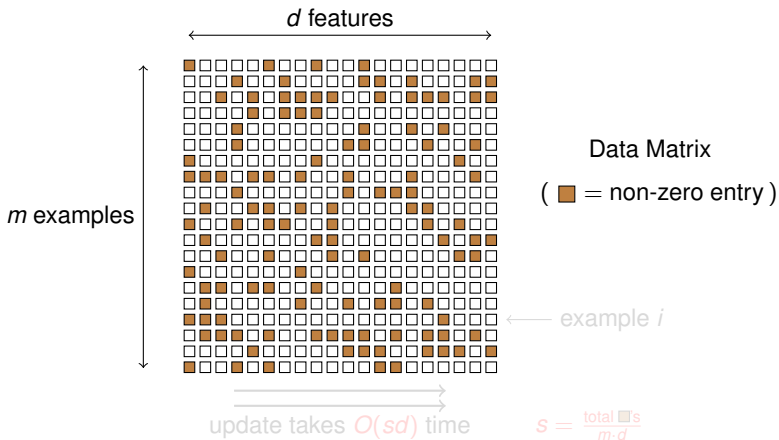
truncate if 0 is crossed

[MIRROR STEP]

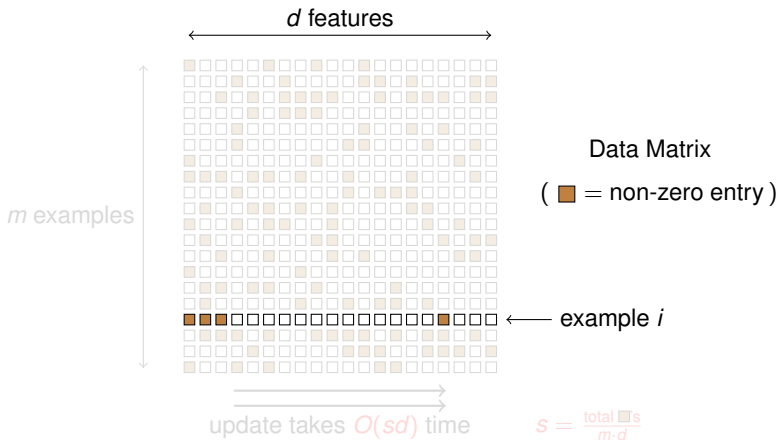
$$\mathbf{w} = f^{-1}(\theta)$$

preserves sparsity pattern!

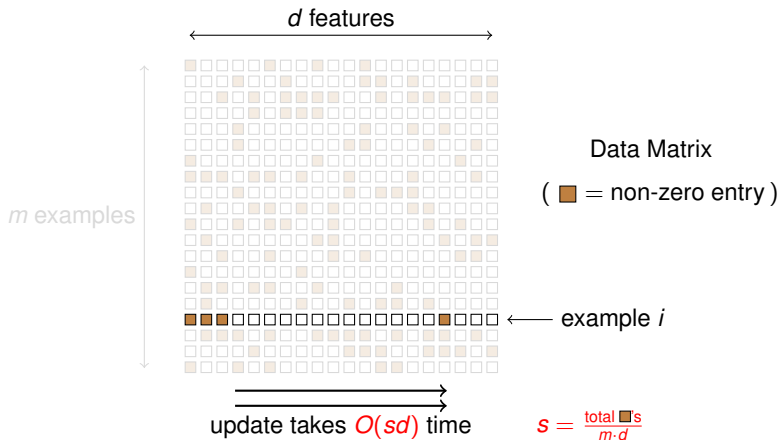
Implementing Updates Efficiently



Implementing Updates Efficiently



Implementing Updates Efficiently



Run-time Guarantee

- Time to achieve (expected) ϵ -accuracy is ($\|x_i\|_\infty \leq X_\infty$)

$$O\left(\frac{d \log(d) X_\infty^2 \|\mathbf{w}^*\|_1^2}{\epsilon^2}\right) \quad (B1)$$

- For Truncated Gradient (Langford et al., 2009), bound is ($X_2^2 =$ average squared 2-norm of x_i 's)

$$O\left(\frac{d X_2^2 \|\mathbf{w}^*\|_2^2}{\epsilon^2}\right) \quad (B2)$$

- \mathbf{w}^* sparse, \mathbf{x}_i 's dense: (B1) is better
- \mathbf{x}_i 's sparse, \mathbf{w}^* dense: (B2) is better

Run-time Guarantee

- Time to achieve (expected) ϵ -accuracy is ($\|x_i\|_\infty \leq X_\infty$)

$$O\left(\frac{d \log(d) X_\infty^2 \|\mathbf{w}^*\|_1^2}{\epsilon^2}\right) \quad (B1)$$

- For Truncated Gradient (Langford et al., 2009), bound is ($X_2^2 =$ average squared 2-norm of x_i 's)

$$O\left(\frac{d X_2^2 \|\mathbf{w}^*\|_2^2}{\epsilon^2}\right) \quad (B2)$$

- \mathbf{w}^* sparse, \mathbf{x}_i 's dense: (B1) is better
- \mathbf{x}_i 's sparse, \mathbf{w}^* dense: (B2) is better

Run-time Guarantee

- Time to achieve (expected) ϵ -accuracy is ($\|x_i\|_\infty \leq X_\infty$)

$$O\left(\frac{d \log(d) X_\infty^2 \|\mathbf{w}^*\|_1^2}{\epsilon^2}\right) \quad (B1)$$

- For Truncated Gradient (Langford et al., 2009), bound is ($X_2^2 =$ average squared 2-norm of x_i 's)

$$O\left(\frac{d X_2^2 \|\mathbf{w}^*\|_2^2}{\epsilon^2}\right) \quad (B2)$$

- \mathbf{w}^* sparse, \mathbf{x}_i 's dense: (B1) is better
- \mathbf{x}_i 's sparse, \mathbf{w}^* dense: (B2) is better

Outline

- 1 Stochastic Coordinate Descent
 - Coordinate Descent for L1
 - Stochastic Coordinate Descent
 - Efficient Implementation
 - Run-time Bound
- 2 SMIDAS
 - Stochastic Gradient Descent
 - SMIDAS
 - Efficient Implementation
 - Run-time Bound
- 3 Experiments
- 4 Summary

Algorithms for Experiments

- SCD: Stochastic Coordinate Descent
- DETCD (Wu & Lange): Deterministic version that chooses coordinate based on guaranteed decrease of objective
- SMIDAS
- TruncGrad (Langford et al.): Truncated Gradient (with $g = \theta = \infty$)

Last two need a parameter η . We report results for best η (chosen over an exponentially spaced grid)

Algorithms for Experiments

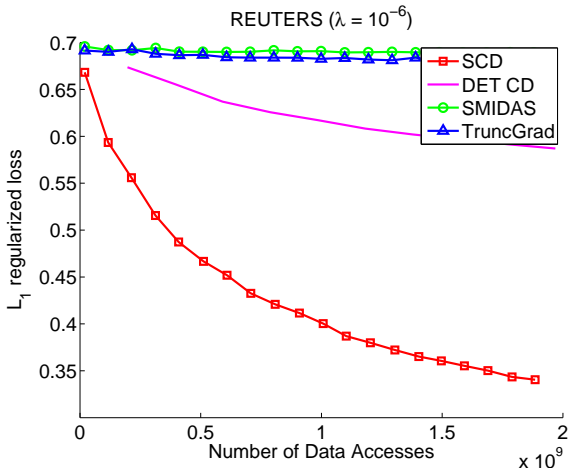
- SCD: Stochastic Coordinate Descent
- DETCD (Wu & Lange): Deterministic version that chooses coordinate based on guaranteed decrease of objective
- SMIDAS
- TruncGrad (Langford et al.): Truncated Gradient (with $g = \theta = \infty$)

Last two need a parameter η . We report results for best η (chosen over an exponentially spaced grid)

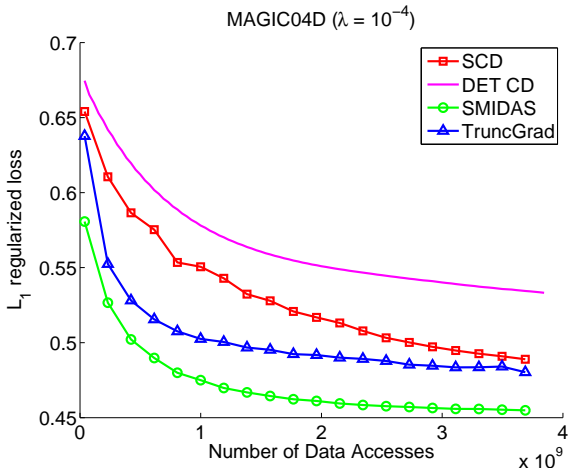
Datasets

- Reuters
 - Dataset derived from Reuters RCV1 collection
 - 806,791 examples
 - 378,452 features
 - Sparsity level of instances is 0.03%
- Magic04D
 - Magic04 dataset from UCI + 1000 random features (± 1 with prob. 0.5)
 - 19,020 examples
 - $10+1000 = 1010$ features
 - Sparsity level of instances is 100% (i.e. not sparse at all)

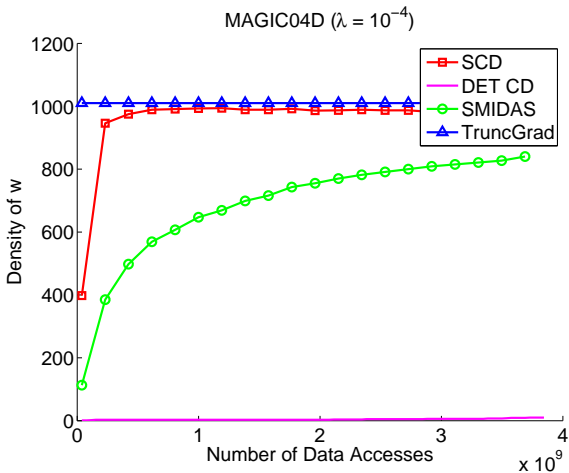
Reuters



Magic04D



Magic04D



Outline

- 1 Stochastic Coordinate Descent
 - Coordinate Descent for L1
 - Stochastic Coordinate Descent
 - Efficient Implementation
 - Run-time Bound
- 2 SMIDAS
 - Stochastic Gradient Descent
 - SMIDAS
 - Efficient Implementation
 - Run-time Bound
- 3 Experiments
- 4 Summary

Summary

- SCD & SMIDAS both perform lots of simple updates + randomization
- Contrast with IP methods: few but sophisticated updates
- Despite simplicity, get (provably) good convergence rates in terms of m, d
- Experiments suggest SCD & SMIDAS have potential to be successful in practice