

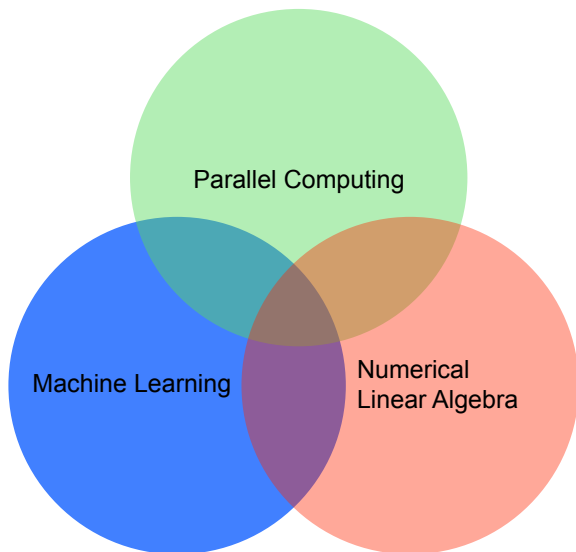
Communication-avoiding coordinate descent methods for linear systems

Aditya Devarakonda

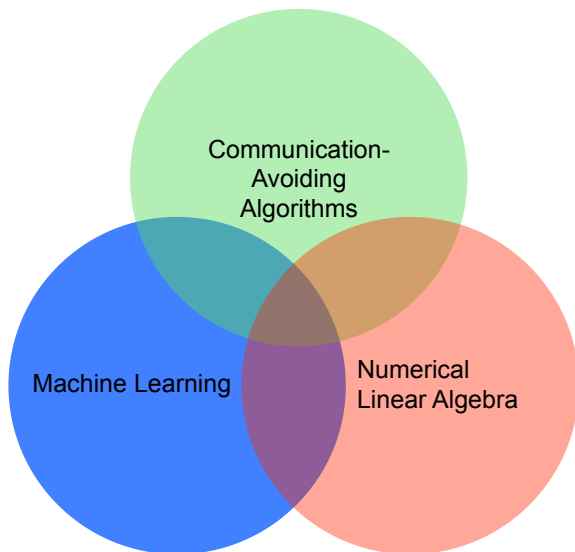
with: Kimon Fountoulakis, Jim Demmel and Michael Mahoney

UC Berkeley

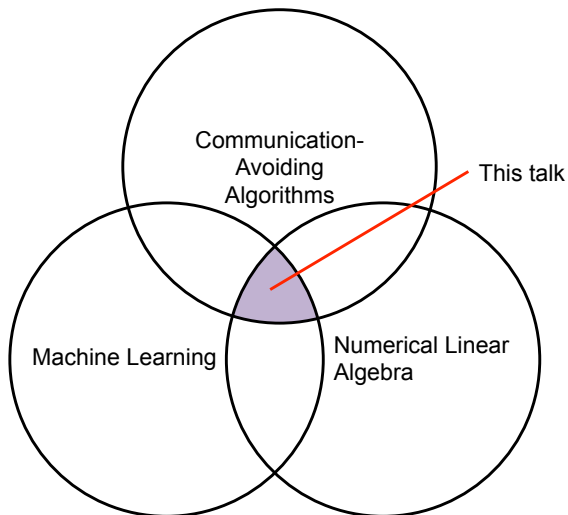
Introduction



Introduction



Introduction



Preliminaries: Cost Model

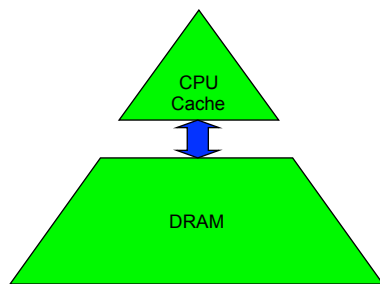


Figure 1: Sequential Case

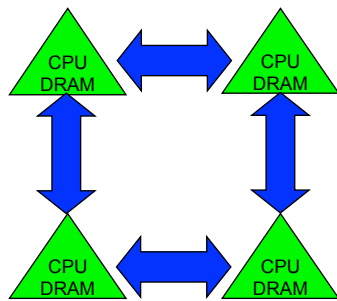


Figure 2: Parallel Case

Algorithms have two costs: **Computation** and **Communication**.

Preliminaries: Cost Model



Computation: flops

Communication: bandwidth
latency

Running Time Model:

$(\gamma \times \text{flops}) + (\beta \times \text{amounts of data moved}) + (\alpha \times \# \text{ of messages})$.

γ = time per floating point operation

β = time per unit of data moved

α = time per message



$$\gamma \ll \beta \ll \alpha$$

Gap is growing.



$$\gamma \ll \beta \ll \alpha$$

Gap is growing.

Avoid communication to save time.



$$\gamma \ll \frac{1}{\beta} \ll \alpha$$

Gap is growing.

Avoid communication to save time **and energy**.



Processors can communicate several ways:



Processors can communicate several ways:
point to point.



Processors can communicate several ways:

point to point.

one to many.



Processors can communicate several ways:

point to point.

one to many.

many to one.



Processors can communicate several ways:

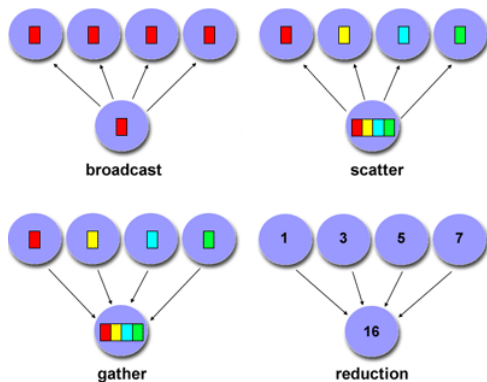
- point to point.

- one to many.

- many to one.

- all to all.

Preliminaries: Communication Costs



Preliminaries: Communication Costs



Communication Primitive	Routing Algorithm	Latency (L)	Bandwidth (W)
Point-to-Point	–	$O(1)$	$O(d)$
Broadcast, Reduction and Scatter	Virtual N -ary Tree	$O(\log P)$	$O(d \log P)$
Gather	Virtual N -ary Tree	$O(\log P)$	$O(dP \log P)$

$P = \#$ of processors

$d =$ amount of data moved

Preliminaries: Communication Costs



Communication Primitive	Routing Algorithm	Latency (L)	Words (W)
Point-to-Point	–	$O(1)$	$O(d)$
Broadcast, Reduction and Scatter	Virtual N -ary Tree	$O(\log P)$	$O(d \log P)$
Gather	Virtual N -ary Tree	$O(\log P)$	$O(dP \log P)$
All-to-All	Recursive Doubling (Short Message)	$O(\log P)$	$O(dP \log P)$
	Virtual Ring (Long Message)	$O(P)$	$O(dP)$



Basic kernel is a matrix-vector product (no CA version).



Basic kernel is a matrix-vector product (no CA version).

Strategy: Unroll Krylov method iteration loop and apply matrix powers kernel.

Factor of s savings in communication (Hoemmen, Carson, Demmel and others).

Regularized Least Squares



Given $X \in \mathbb{R}^{m \times n}$ ($m < n$) and $y \in \mathbb{R}^n$, find $w \in \mathbb{R}^m$ such that

$$\min_{w \in \mathbb{R}^m} \frac{1}{2n} \sum_{i=1}^n (w^T x_i - y_i)^2 + \frac{\lambda}{2} \|w\|_2^2$$

Closed form solution: $w = (XX^T + \lambda n I_m)^{-1} Xy$.

Regularized Least Squares



Many ways to solve using direct and iterative methods.

Regularized Least Squares



Many ways to solve using direct and iterative methods.

We solve it using stochastic dual coordinate ascent (SDCA).

Comparable to (or better than) stochastic gradient descent (SGD).



Primal problem:

$$\min_{w \in \mathbb{R}^m} \frac{1}{2n} \sum_{i=1}^n (w^T x_i - y_i)^2 + \frac{\lambda}{2} \|w\|_2^2$$

Dual problem¹:

$$\max_{\alpha \in \mathbb{R}^n} -\frac{1}{2n} \sum_{i=1}^n \alpha_i^2 + \alpha_i y_i - \frac{\lambda}{2} \left\| \frac{1}{\lambda n} X \alpha \right\|_2^2$$
$$\min_{\alpha \in \mathbb{R}^n} \frac{1}{2n} \|\alpha + y\|^2 + \frac{\lambda}{2} \left\| \frac{1}{\lambda n} X \alpha \right\|_2^2.$$

¹By definition of conjugate functions (Legendre Transform)



Primal-Dual Relationship:

$$w = -\frac{1}{\lambda n} X\alpha$$



Dual problem (not an algorithm):

$$\min_{\alpha \in \mathbb{R}^n} \frac{1}{2n} \|\alpha + y\|^2 + \frac{\lambda}{2} \left\| \frac{1}{\lambda n} X \alpha \right\|_2^2.$$



Dual problem (not an algorithm):

$$\min_{\alpha \in \mathbb{R}^n} \frac{1}{2n} \|\alpha + y\|_2^2 + \frac{\lambda}{2} \left\| \frac{1}{\lambda n} X \alpha \right\|_2^2.$$

Coordinate-wise update method:

Let h be an iteration index.

Select a coordinate $i \in [n]$.

$$\min_{\Delta \alpha \in \mathbb{R}} \frac{1}{2n} \left\| \alpha^{(h-1)} + \Delta \alpha e_i + y \right\|_2^2 + \frac{\lambda}{2} \left\| \frac{1}{\lambda n} X \alpha^{(h-1)} + \frac{\Delta \alpha}{\lambda n} x_i \right\|_2^2$$

$$\alpha^{(h)} = \alpha^{(h-1)} + \Delta \alpha e_i.$$



$$\min_{\Delta\alpha \in \mathbb{R}} \frac{1}{2n} \left\| \alpha^{(h-1)} + \Delta\alpha \mathbf{e}_i + y \right\|_2^2 + \frac{\lambda}{2} \left\| \frac{1}{\lambda n} X \alpha^{(h-1)} + \frac{\Delta\alpha}{\lambda n} \mathbf{x}_i \right\|_2^2$$

Solve for $\Delta\alpha$:

$$\Delta\alpha = \frac{-1}{\frac{1}{\lambda n^2} \mathbf{x}_i^T \mathbf{x}_i + \frac{1}{n}} \mathbf{e}_i^T \left\{ \left(\frac{1}{\lambda n^2} X^T X + \frac{1}{n} I \right) \alpha^{(h-1)} + \frac{1}{n} y \right\}.$$



$$\Delta\alpha = \frac{-1}{\frac{1}{\lambda n^2} x_i^T x_i + \frac{1}{n}} e_i^T \left\{ \left(\frac{1}{\lambda n^2} X^T X + \frac{1}{n} I \right) \alpha^{(h-1)} + \frac{1}{n} y \right\}.$$

Let²:

$$\begin{aligned}\gamma_i &= \frac{-1}{\frac{1}{\lambda n^2} x_i^T x_i + \frac{1}{n}} \\ A &= \frac{1}{\lambda n^2} X^T X + \frac{1}{n} I \\ r^{(h-1)} &= A\alpha^{(h-1)} + \frac{1}{n} y.\end{aligned}$$

²A is a notational change, **not explicitly formed**



Before:

$$\Delta\alpha = \frac{-1}{\frac{1}{\lambda n^2} \mathbf{x}_i^T \mathbf{x}_i + \frac{1}{n}} \mathbf{e}_i^T \left\{ \left(\frac{1}{\lambda n^2} \mathbf{X}^T \mathbf{X} + \frac{1}{n} \mathbf{I} \right) \alpha^{(h-1)} + \frac{1}{n} \mathbf{y} \right\}$$

$$\alpha^{(h)} = \alpha^{(h-1)} + \Delta\alpha \mathbf{e}_i$$

$$\mathbf{w}^{(h)} = \mathbf{w}^{(h-1)} - \frac{1}{\lambda n} \Delta\alpha \mathbf{X} \mathbf{e}_i.$$

After:

$$\Delta\alpha = \gamma_i \mathbf{e}_i^T \mathbf{r}^{(h-1)}$$

$$\alpha^{(h)} = \alpha^{(h-1)} + \gamma_i \mathbf{e}_i \mathbf{e}_i^T \mathbf{r}^{(h-1)}$$

$$\mathbf{w}^{(h)} = \mathbf{w}^{(h-1)} - \frac{1}{\lambda n} \Delta\alpha \mathbf{X} \mathbf{e}_i.$$



Since $r^{(h-1)} = A\alpha^{(h-1)} + \frac{1}{n}y$
and $\alpha^{(h)} = \alpha^{(h-1)} + \gamma_i e_i e_i^T r^{(h-1)}$

$$r^{(h)} = (I + \gamma_i A e_i e_i^T) r^{(h-1)}$$

This residual update helps us avoid communication.



Let $i_j \in [n]$ (an index chosen from 1 to n)
and $r^{(0)}$ be the starting residual.

$$\mathcal{R}_1 := \left[r^{(0)}, \left(I + \gamma_{i_1} A e_{i_1} e_{i_1}^T \right) r^{(0)}, \right. \\ \left. \left(I + \gamma_{i_2} A e_{i_2} e_{i_2}^T \right) \left(I + \gamma_{i_1} A e_{i_1} e_{i_1}^T \right) r^{(0)}, \right. \\ \left. \dots, \prod_{j=s-1}^1 \left(I + \gamma_{i_j} A e_{i_j} e_{i_j}^T \right) r^{(0)} \right] \in \mathbb{R}^{n \times s}.$$

This “unrolls” the residual by s iterations.



Let $i_j \in [n]$ (an index chosen from 1 to n)
and $r^{(0)}$ be the starting residual.

$$\mathcal{A}_{sh} := \left[\begin{array}{l} \gamma_{i_1} \mathbf{e}_{i_1}^T r^{(0)}, \gamma_{i_2} \mathbf{e}_{i_2}^T \left(I + \gamma_{i_1} \mathbf{A} \mathbf{e}_{i_1} \mathbf{e}_{i_1}^T \right) r^{(0)}, \\ \gamma_{i_3} \mathbf{e}_{i_3}^T \left(I + \gamma_{i_2} \mathbf{A} \mathbf{e}_{i_2} \mathbf{e}_{i_2}^T \right) \left(I + \gamma_{i_1} \mathbf{A} \mathbf{e}_{i_1} \mathbf{e}_{i_1}^T \right) r^{(0)}, \\ \dots, \gamma_{i_s} \mathbf{e}_{i_s}^T \prod_{j=s-1}^1 \left(I + \gamma_{i_j} \mathbf{A} \mathbf{e}_{i_j} \mathbf{e}_{i_j}^T \right) r^{(0)} \end{array} \right] \in \mathbb{R}^{1 \times s}.$$

Due to unrolling the residual.



The hard part is computing: $\gamma_{i_s} e_{i_s}^T \prod_{j=s-1}^1 \left(I + \gamma_{i_j} A e_{i_j} e_{i_j}^T \right) r^{(0)}$

However, recall that

$$\gamma_{i_j} = \frac{-1}{\frac{1}{\lambda n^2} x_{i_j}^T x_{i_j} + \frac{1}{n}} \text{ requires a dot-product.}$$



The hard part is computing: $\gamma_{i_s} e_{i_s}^T \prod_{j=s-1}^1 \left(I + \gamma_{i_j} A e_{i_j} e_{i_j}^T \right) r^{(0)}$

However, recall that

$$\gamma_{i_j} = \frac{-1}{\frac{1}{\lambda n^2} x_{i_j}^T x_{i_j} + \frac{1}{n}} \text{ requires a dot-product.}$$

Easy to avoid communication. Requires one reduction.



Recall that: $A := \frac{1}{\lambda n^2} X^T X + \frac{1}{n} I$

Notice the $e_{i_j}^T$ terms: $\gamma_{i_s} e_{i_s}^T \prod_{j=s-1}^1 \left(I + \gamma_{i_j} A e_{i_j} e_{i_j}^T \right) r^{(0)}$

Expanding the product results in dot-products.



Let $Y = [x_{i_1}, x_{i_2}, \dots, x_{i_s}] \in \mathbb{R}^{m \times s}$

then compute the Gram-like matrix $G = \frac{1}{\lambda n^2} Y^T Y + \frac{1}{n} I$.

SDCA Krylov-like Basis



Let $Y = [x_{i_1}, x_{i_2}, \dots, x_{i_s}] \in \mathbb{R}^{m \times s}$

then compute the Gram-like matrix $G = \frac{1}{\lambda n^2} Y^T Y + \frac{1}{n} I$.

G contains all of the coefficients required to compute \mathcal{A}_1 .



Let $Y = [x_{i_1}, x_{i_2}, \dots, x_{i_s}] \in \mathbb{R}^{m \times s}$

then compute the Gram-like matrix $G = \frac{1}{\lambda n^2} Y^T Y + \frac{1}{n} I$.

G contains all of the coefficients required to compute \mathcal{A}_1 .

Computing G requires a single reduction.

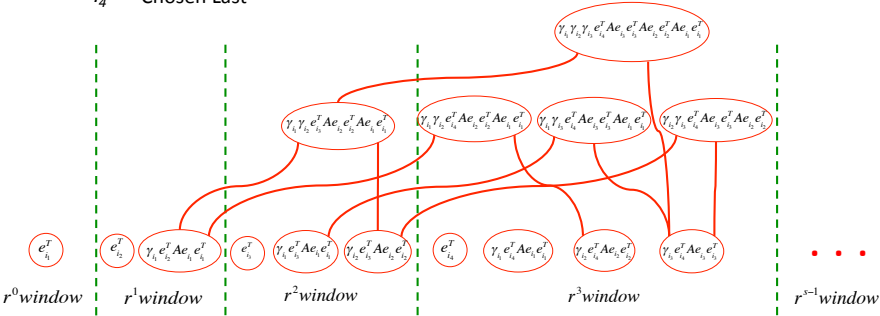
CA-SDCA Dependency Graph



G contains the coefficients but still need to compute the expansion.

Index Order

i_1 Chosen First
 i_2 ↓
 i_3 ↓
 i_4 Chosen Last



SDCA Polynomial Expansion



$$\begin{array}{l}
 \text{Step 1:} \\
 \left[\begin{array}{cccc} \gamma_{i_1} & 0 & 0 & 0 \\ \text{...} & \text{...} & \text{...} & \text{...} \\ \text{...} & \text{...} & \text{...} & \text{...} \end{array} \right] \times \left[\begin{array}{cccc} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{array} \right] \\
 \hline
 \text{Step 2:} \\
 \left[\begin{array}{cccc} \text{...} & \text{...} & \text{...} & \text{...} \\ \gamma_{i_2} e_{i_2}^T A e_{i_1} & \gamma_{i_2} & 0 & 0 \\ \text{...} & \text{...} & \text{...} & \text{...} \end{array} \right] \times \left[\begin{array}{cccc} \gamma_{i_1} & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{array} \right] \\
 \hline
 \text{Step 3:} \\
 \left[\begin{array}{cccc} \text{...} & \text{...} & \text{...} & \text{...} \\ \gamma_{i_3} e_{i_3}^T A e_{i_1} & \gamma_{i_3} e_{i_3}^T A e_{i_2} & \gamma_{i_3} & 0 \\ \text{...} & \text{...} & \text{...} & \text{...} \end{array} \right] \times \left[\begin{array}{cccc} \gamma_{i_1} & 0 & 0 & 0 \\ \gamma_{i_1} \gamma_{i_2} e_{i_2}^T A e_{i_1} & \gamma_{i_2} & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{array} \right] \\
 \hline
 \text{Step 4:} \\
 \left[\begin{array}{cccc} \text{...} & \text{...} & \text{...} & \text{...} \\ \gamma_{i_4} e_{i_4}^T A e_{i_1} & \gamma_{i_4} e_{i_4}^T A e_{i_2} & \gamma_{i_4} e_{i_4}^T A e_{i_3} & \gamma_{i_4} \\ \text{...} & \text{...} & \text{...} & \text{...} \end{array} \right] \times \left[\begin{array}{cccc} \gamma_{i_1} & 0 & 0 & 0 \\ \gamma_{i_1} \gamma_{i_2} e_{i_2}^T A e_{i_1} & \gamma_{i_2} & 0 & 0 \\ \gamma_{i_1} \gamma_{i_2} \gamma_{i_3} e_{i_3}^T A e_{i_2} e_{i_2}^T A e_{i_1} + \gamma_{i_1} \gamma_{i_3} e_{i_3}^T A e_{i_1} & \gamma_{i_2} \gamma_{i_3} e_{i_3}^T A e_{i_2} & \gamma_{i_3} & 0 \\ 0 & 0 & 0 & 1 \end{array} \right]
 \end{array}$$

=

$$\left[\begin{array}{cccc} \gamma_{i_1} & 0 & 0 & 0 \\ \gamma_{i_1} \gamma_{i_2} e_{i_2}^T A e_{i_1} & \gamma_{i_2} & 0 & 0 \\ \gamma_{i_1} \gamma_{i_2} \gamma_{i_3} e_{i_3}^T A e_{i_2} e_{i_2}^T A e_{i_1} + \gamma_{i_1} \gamma_{i_3} e_{i_3}^T A e_{i_1} & \gamma_{i_2} \gamma_{i_3} e_{i_3}^T A e_{i_2} & \gamma_{i_3} & 0 \\ \gamma_{i_1} \gamma_{i_2} \gamma_{i_4} e_{i_4}^T A e_{i_3} e_{i_3}^T A e_{i_2} e_{i_2}^T A e_{i_1} + \gamma_{i_1} \gamma_{i_2} \gamma_{i_4} e_{i_4}^T A e_{i_3} e_{i_3}^T A e_{i_1} + \gamma_{i_1} \gamma_{i_4} e_{i_4}^T A e_{i_1} & \gamma_{i_2} \gamma_{i_3} \gamma_{i_4} e_{i_4}^T A e_{i_3} e_{i_3}^T A e_{i_2} + \gamma_{i_2} \gamma_{i_4} e_{i_4}^T A e_{i_2} & \gamma_{i_3} \gamma_{i_4} e_{i_4}^T A e_{i_3} & \gamma_{i_4} \end{array} \right]$$

CA-SDCA Step Sizes



$$\begin{bmatrix} \Delta\alpha_{i_1} \\ \Delta\alpha_{i_2} \\ \Delta\alpha_{i_3} \\ \Delta\alpha_{i_4} \end{bmatrix} = \begin{bmatrix} \gamma_{i_1} & 0 & 0 & 0 \\ \gamma_{i_1}\gamma_{i_2}e_{i_2}^T Ae_{i_1} & \gamma_{i_2} & 0 & 0 \\ \gamma_{i_1}\gamma_{i_2}\gamma_{i_3}e_{i_3}^T Ae_{i_2}e_{i_2}^T Ae_{i_1} + \gamma_{i_1}\gamma_{i_3}e_{i_3}^T Ae_{i_1} & \gamma_{i_2}\gamma_{i_3}e_{i_3}^T Ae_{i_2} & \gamma_{i_3} & 0 \\ \gamma_{i_1}\gamma_{i_2}\gamma_{i_4}e_{i_4}^T Ae_{i_2}e_{i_2}^T Ae_{i_1} + \gamma_{i_4}e_{i_4}^T Ae_{i_3}(\gamma_{i_1}\gamma_{i_2}\gamma_{i_3}e_{i_3}^T Ae_{i_2}e_{i_2}^T Ae_{i_1} + \gamma_{i_1}\gamma_{i_3}e_{i_3}^T Ae_{i_1}) & \gamma_{i_2}\gamma_{i_3}\gamma_{i_4}e_{i_4}^T Ae_{i_3}e_{i_3}^T Ae_{i_2} + \gamma_{i_2}\gamma_{i_4}e_{i_4}^T Ae_{i_2} & \gamma_{i_3}\gamma_{i_4}e_{i_4}^T Ae_{i_3} & \gamma_{i_4} \end{bmatrix} \times \begin{bmatrix} r_{i_1}^{(0)} \\ r_{i_2}^{(0)} \\ r_{i_3}^{(0)} \\ r_{i_4}^{(0)} \end{bmatrix}$$

No communication required.



The subroutine packs s iterations into a Gram-matrix and matrix-vector computations.

Algorithm	Latency (L)	Bandwidth (W)
SDCA	$O(H \log P)$	$O(Hm \log P)$
CA-SDCA	$O(\frac{H}{s} \log P)$	$O(Hs \log P + \frac{Hm}{s} \log P)$

Reduces latency cost.

Reduces bandwidth cost if $s^2 < m$.

CA-SDCA Results

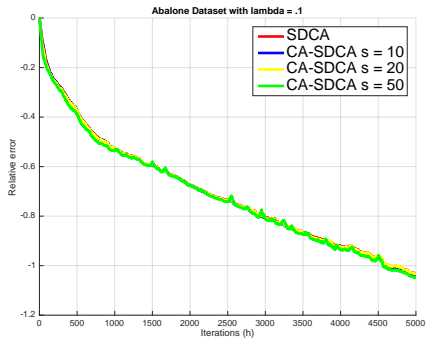
Matlab experiments using datasets from libsvm.

Measuring relative error: $\frac{\|w_{alg} - w_{opt}\|}{\|w_{opt}\|}$

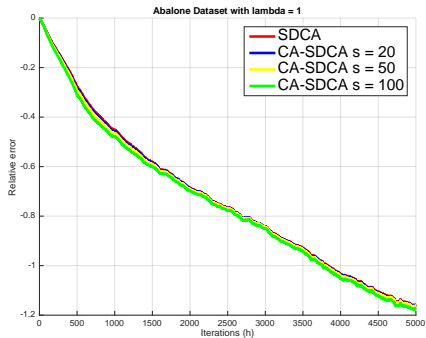
w_{opt} obtained from conjugate gradients.

Dataset name	Dimensions	Problem type
Abalone	4177×8	regression
w1a	2477×300	classification
a9a	32561×123	classification

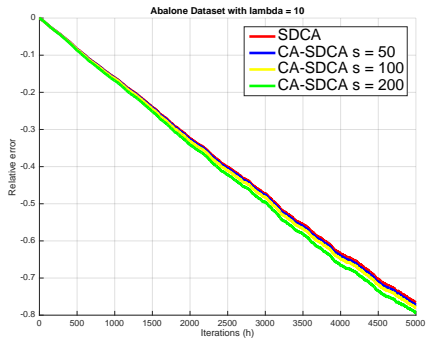
CA-SDCA Results: Abalone



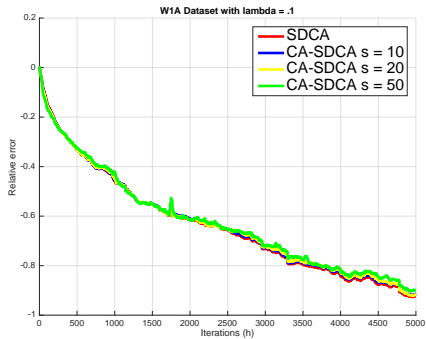
CA-SDCA Results: Abalone



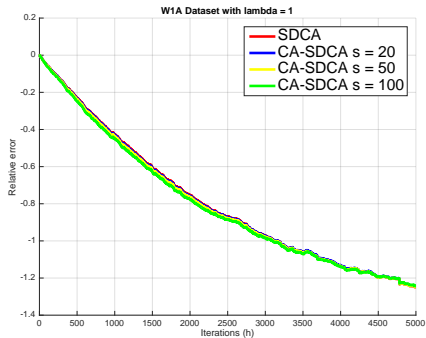
CA-SDCA Results: Abalone



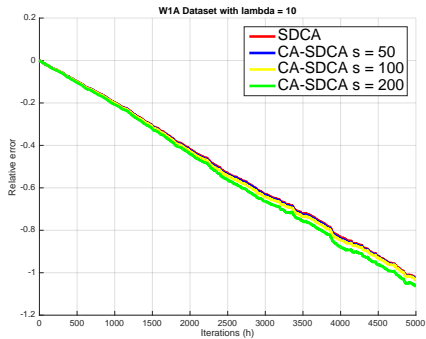
CA-SDCA Results: w1a



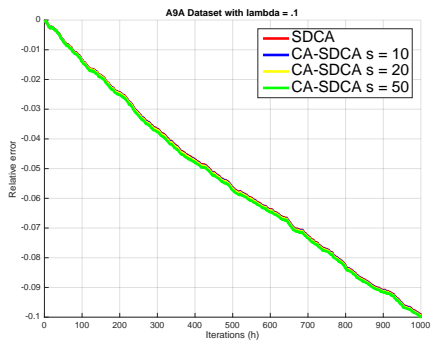
CA-SDCA Results: w1a



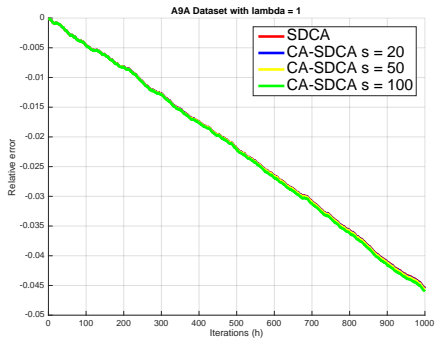
CA-SDCA Results: w1a



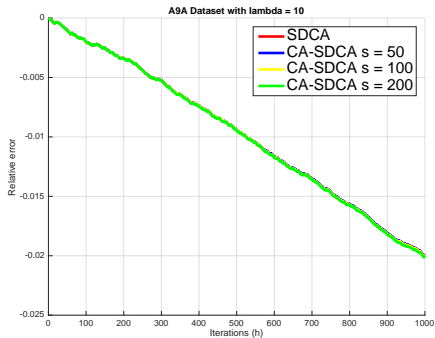
CA-SDCA Results: a9a



CA-SDCA Results: a9a



CA-SDCA Results: a9a



CA-SDCA Summary

Maximum attainable s depends on conditioning of X and λ .

CA-SDCA Summary

Maximum attainable s depends on conditioning of X and λ .

In practice, s should be chosen to balance:

- numerical error

- computation vs. communication

CA-mini-batch SDCA Derivation



Updates batches of b coordinates independently.



Recall for SDCA:

$$r^{(h)} = \left(I + \gamma_i A e_i e_i^T \right) r^{(h-1)}$$

residual update for mbSDCA:

$$r^{(h)} = \left(I + \sum_{k=1}^b \gamma_{i_k} A e_{i_k} e_{i_k}^T \right) r^{(h-1)}$$



Let:

$$\mathbb{I}_h = [e_{i_1} \quad e_{i_2} \quad \cdots \quad e_{i_k} \quad \cdots \quad e_{i_b}],$$

and define the diagonal matrix $\Gamma_h \in \mathbb{R}^{b \times b}$ such that $\Gamma_h(i_k, i_k) = \gamma_{i_k}$.

then residual update for mbSDCA:

$$r^{(h)} = \left(I + A \mathbb{I}_h \Gamma_h \mathbb{I}_h^T \right) r^{(h-1)}.$$



All updates are now block updates:

$$\mathbb{A}_h = \begin{bmatrix} \Delta\alpha_{i_1} \\ \vdots \\ \Delta\alpha_{i_k} \\ \vdots \\ \Delta\alpha_{i_b} \end{bmatrix} = \Gamma_h \mathbb{I}_h^T r^{(h-1)}$$
$$\alpha^{(h)} = \alpha^{(h-1)} + \mathbb{I}_h \mathbb{A}_h$$
$$w^{(h)} = w^{(h-1)} - \frac{1}{\lambda n} X \mathbb{I}_h \mathbb{A}_h.$$



$$\mathcal{R}_1 := \left[r^{(0)}, \left(I + A\mathbb{I}_1\Gamma_1\mathbb{I}_1^T \right) r^{(0)}, \right. \\ \left. \left(I + A\mathbb{I}_2\Gamma_2\mathbb{I}_2^T \right) \left(I + A\mathbb{I}_1\Gamma_1\mathbb{I}_1^T \right) r^{(0)}, \right. \\ \left. \dots, \prod_{j=s-1}^1 \left(I + A\mathbb{I}_j\Gamma_j\mathbb{I}_j^T \right) r^{(0)} \right] \in \mathbb{R}^{n \times s}.$$



$$\mathcal{A}_1 := \left[\Gamma_1 \mathbb{I}_1^T r^{(0)}, \Gamma_2 \mathbb{I}_2^T \left(I + A \mathbb{I}_1 \Gamma_1 \mathbb{I}_1^T \right) r^{(0)}, \right. \\ \Gamma_3 \mathbb{I}_3^T \left(I + \mathbb{I}_2 A \Gamma_2 \mathbb{I}_2^T \right) \left(I + A \mathbb{I}_1 \Gamma_1 \mathbb{I}_1^T \right) r^{(0)}, \\ \left. \dots, \Gamma_s \mathbb{I}_s^T \prod_{j=s-1}^1 \left(I + A \mathbb{I}_j \Gamma_j \mathbb{I}_j^T \right) r^{(0)} \right] \in \mathbb{R}^{b \times s}.$$

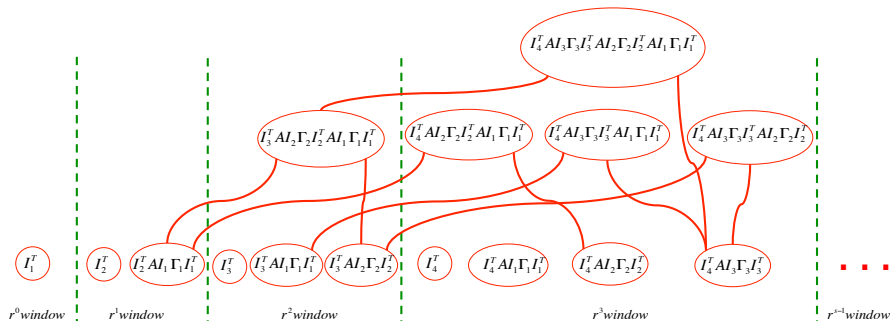
CA-mini-batch SDCA Dependency Graph



Index Order

I_1 Chosen First
 I_2
 I_3
 I_4 Chosen Last

↓



CA-mini-batch SDCA Poly. Expansion



$$\begin{array}{l}
 \text{Step 1:} \\
 \text{Step 2:} \\
 \text{Step 3:} \\
 \text{Step 4:}
 \end{array}
 \left[\begin{array}{cccc}
 \Gamma_1 & 0 & 0 & 0 \\
 \text{...} & \text{...} & \text{...} & \text{...} \\
 \Gamma_2 \mathbb{I}_2^T A \mathbb{I}_1 & \Gamma_2 & 0 & 0 \\
 \text{...} & \text{...} & \text{...} & \text{...} \\
 \Gamma_3 \mathbb{I}_3^T A \mathbb{I}_1 & \Gamma_3 \mathbb{I}_3^T A \mathbb{I}_2 & \Gamma_3 & 0 \\
 \text{...} & \text{...} & \text{...} & \text{...} \\
 \Gamma_4 \mathbb{I}_4^T A \mathbb{I}_1 & \Gamma_4 \mathbb{I}_4^T A \mathbb{I}_2 & \Gamma_4 \mathbb{I}_4^T A \mathbb{I}_3 & \Gamma_4
 \end{array} \right]
 \times
 \left[\begin{array}{cccc}
 I & 0 & 0 & 0 \\
 0 & I & 0 & 0 \\
 0 & 0 & I & 0 \\
 0 & 0 & 0 & I
 \end{array} \right]$$

$$= \left[\begin{array}{cccc}
 \Gamma_1 & 0 & 0 & 0 \\
 \Gamma_2 \mathbb{I}_2^T A \mathbb{I}_1 \Gamma_1 & \Gamma_2 & 0 & 0 \\
 \Gamma_3 \mathbb{I}_3^T A \mathbb{I}_2 \Gamma_2 \mathbb{I}_2^T A \mathbb{I}_1 \Gamma_1 + \Gamma_3 \mathbb{I}_3^T A \mathbb{I}_1 \Gamma_1 & \Gamma_2 \Gamma_3 \mathbb{I}_3^T A \mathbb{I}_2 & \Gamma_3 & 0 \\
 \Gamma_4 \mathbb{I}_4^T A \mathbb{I}_2 \Gamma_2 \mathbb{I}_2^T A \mathbb{I}_1 \Gamma_1 + \Gamma_4 \mathbb{I}_4^T A \mathbb{I}_3 (\Gamma_3 \mathbb{I}_3^T A \mathbb{I}_2 \Gamma_2 \mathbb{I}_2^T A \mathbb{I}_1 \Gamma_1 + \Gamma_3 \mathbb{I}_3^T A \mathbb{I}_1 \Gamma_1) + \Gamma_4 \mathbb{I}_4^T A \mathbb{I}_1 \Gamma_1 & \Gamma_4 \mathbb{I}_4^T A \mathbb{I}_3 \Gamma_3 \mathbb{I}_3^T A \mathbb{I}_2 \Gamma_2 + \Gamma_4 \mathbb{I}_4^T A \mathbb{I}_2 \Gamma_2 & \Gamma_4 \mathbb{I}_4^T A \mathbb{I}_3 \Gamma_3 & \Gamma_4
 \end{array} \right]$$

CA-mini-batch SDCA Step Sizes



$$\begin{bmatrix} A_1 \\ A_2 \\ A_3 \\ A_4 \end{bmatrix} = \begin{bmatrix} \Gamma_1 & 0 & 0 & 0 \\ \Gamma_2 \mathbb{I}_2^T A \mathbb{I}_1 \Gamma_1 & \Gamma_2 & 0 & 0 \\ \Gamma_3 \mathbb{I}_3^T A \mathbb{I}_2 \Gamma_2 \mathbb{I}_2^T A \mathbb{I}_1 \Gamma_1 + \Gamma_3 \mathbb{I}_3^T A \mathbb{I}_1 \Gamma_1 & \Gamma_3 \mathbb{I}_3^T A \mathbb{I}_2 \Gamma_2 & \Gamma_3 & 0 \\ \Gamma_4 \mathbb{I}_4^T A \mathbb{I}_2 \Gamma_2 \mathbb{I}_2^T A \mathbb{I}_1 \Gamma_1 + \Gamma_4 \mathbb{I}_4^T A \mathbb{I}_3 (\Gamma_3 \mathbb{I}_3^T A \mathbb{I}_2 \Gamma_2 \mathbb{I}_2^T A \mathbb{I}_1 \Gamma_1 + \Gamma_3 \mathbb{I}_3^T A \mathbb{I}_1 \Gamma_1) + \Gamma_4 \mathbb{I}_4^T A \mathbb{I}_1 \Gamma_1 & \Gamma_4 \mathbb{I}_4^T A \mathbb{I}_3 \Gamma_3 \mathbb{I}_3^T A \mathbb{I}_2 \Gamma_2 + \Gamma_4 \mathbb{I}_4^T A \mathbb{I}_2 \Gamma_2 & \Gamma_4 \mathbb{I}_4^T A \mathbb{I}_3 \Gamma_3 & \Gamma_4 \end{bmatrix} \times \begin{bmatrix} \mathbb{I}_1^T r^{(0)} \\ \mathbb{I}_2^T r^{(0)} \\ \mathbb{I}_3^T r^{(0)} \\ \mathbb{I}_4^T r^{(0)} \end{bmatrix}$$

CA-mini-batch SDCA Communication Cost



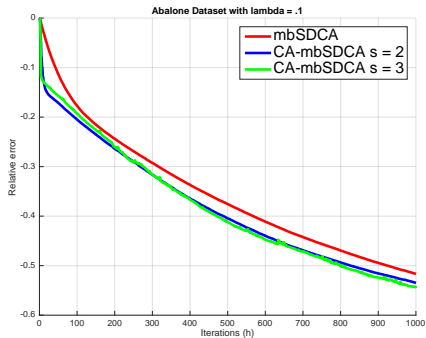
The subroutine packs s iterations into Gram-matrix and matrix-vector computations.

Algorithm	Latency (L)	Bandwidth (W)
mbSDCA	$O(H \log P)$	$O(Hb \log P + Hm \log P)$
CA-mbSDCA	$O(\frac{H}{s} \log P)$	$O(Hsm^2 \log P + \frac{Hm}{s} \log P)$

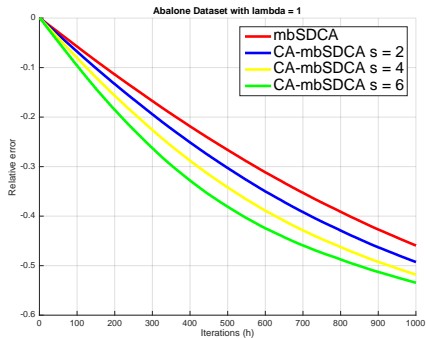
Reduces latency cost.

Reduces bandwidth cost if $(sb)^2 < m$.

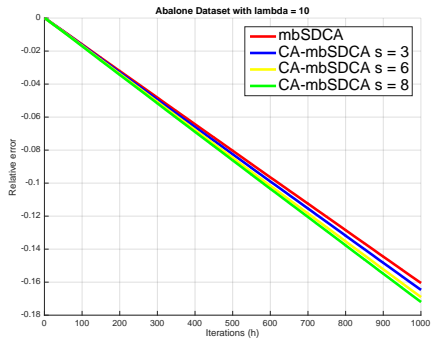
CA-mbSDCA Results: Abalone



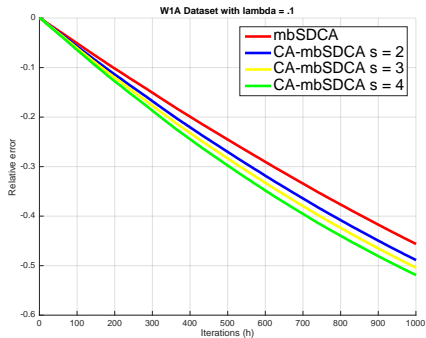
CA-mbSDCA Results: Abalone



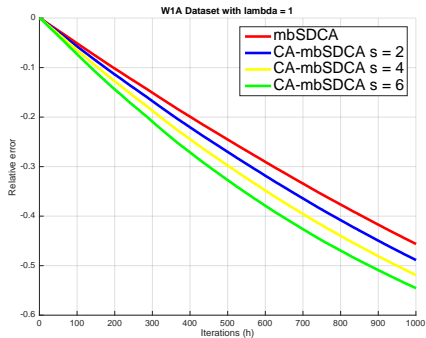
CA-mbSDCA Results: Abalone



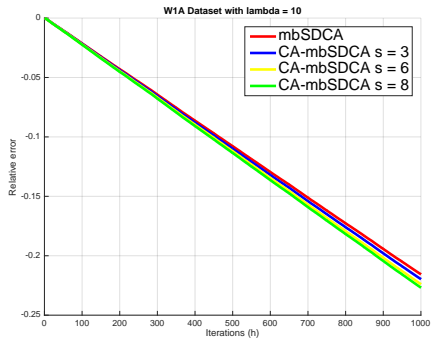
CA-mbSDCA Results: w1a



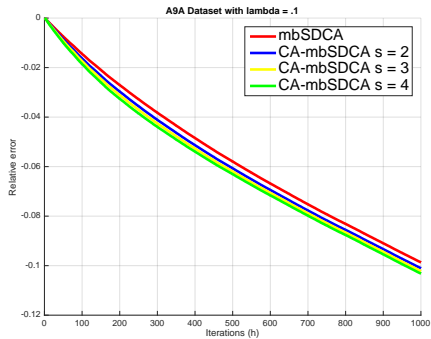
CA-mbSDCA Results: w1a



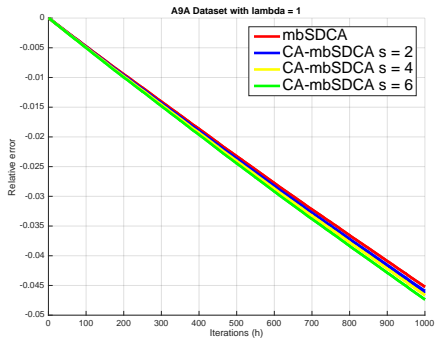
CA-mbSDCA Results: w1a



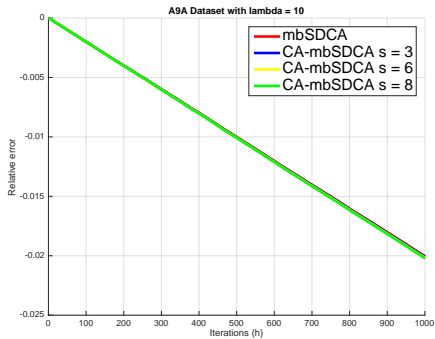
CA-mbSDCA Results: a9a



CA-mbSDCA Results: a9a



CA-mbSDCA Results: a9a



CA-SDCA Summary

Maximum attainable s depends on conditioning of X and λ .

CA-SDCA Summary

Maximum attainable s depends on conditioning of X and λ .

In practice, s and b should be chosen to balance:

- numerical error

- computation vs. communication

CA-SDCA Summary

Maximum attainable s depends on conditioning of X and λ .

In practice, s and b should be chosen to balance:

- numerical error

- computation vs. communication

Interesting to explore tradeoff between s and b .

Conclusion and Future Work

Derived CA versions of SDCA, mbSDCA and block-coordinate descent algorithms.

Experiments show these methods are numerically viable.

More work to be done!

Parallel implementations in progress.

Need to extend to L1-regularization and non-linear loss functions.