
Communication-Avoiding Machine Learning

— Aditya Devarakonda —

ASPIRE Summer Retreat

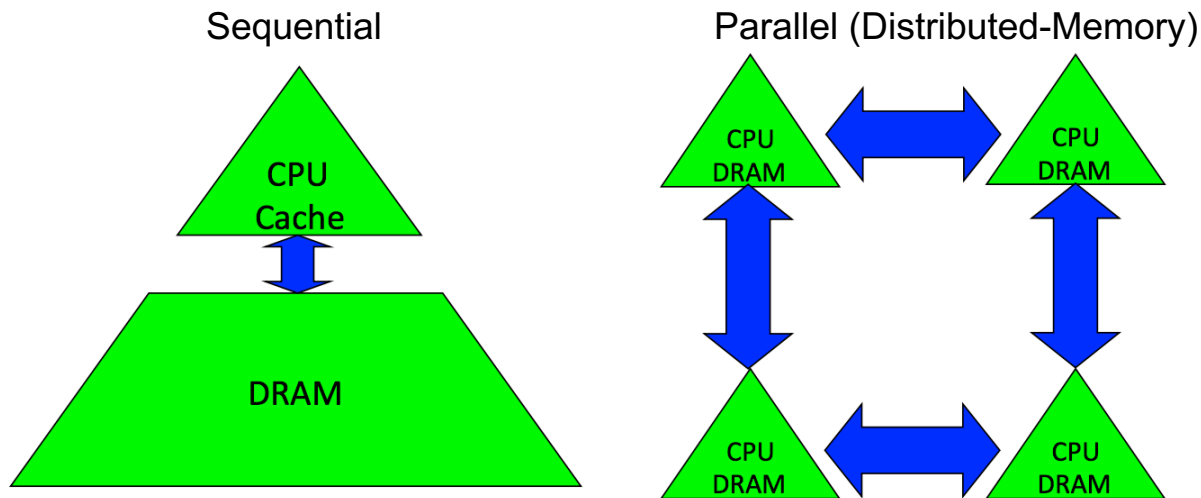
June 14th, 2017

Collaborators:

Jim Demmel (Adviser), Kimon Fountoulakis (Post-Doc), and Michael W. Mahoney (Co-adviser)

Definition

Communication is data movement.



Courtesy:
Demmel

Least-Squares (Linear Regression)

Many ways to solve.

Direct

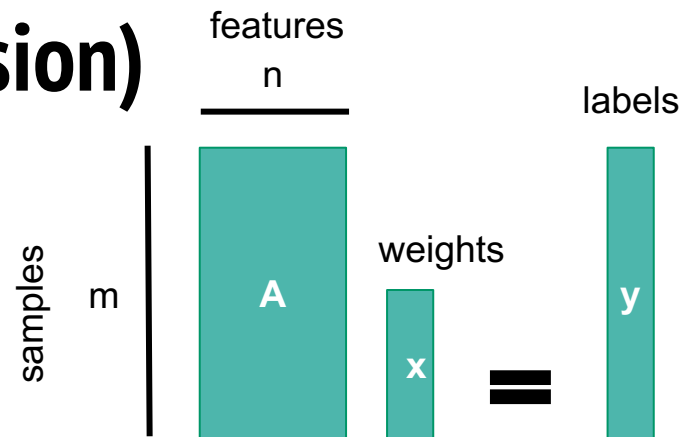
Explicitly solve normal equation.

Implicitly through matrix factorizations.

Iterative

Krylov methods (e.g. Conjugate Gradients).

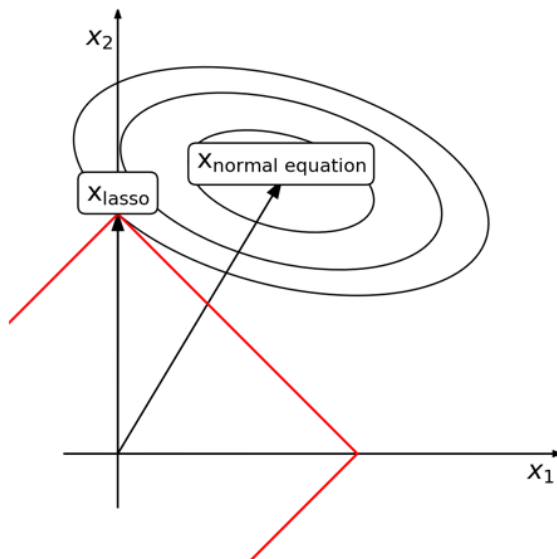
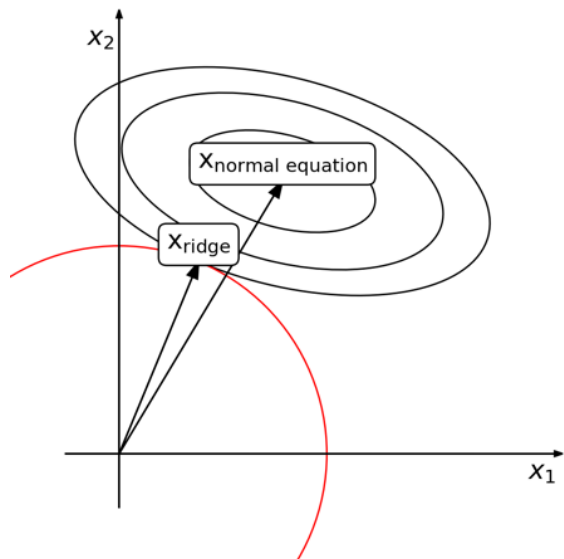
(Block) Coordinate Descent.



$$\text{Least-Squares: } \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - y\|_2^2$$

$$\text{Normal Equation: } x = (A^T A)^{-1} A^T y$$

Regularized Least-Squares (Regression)



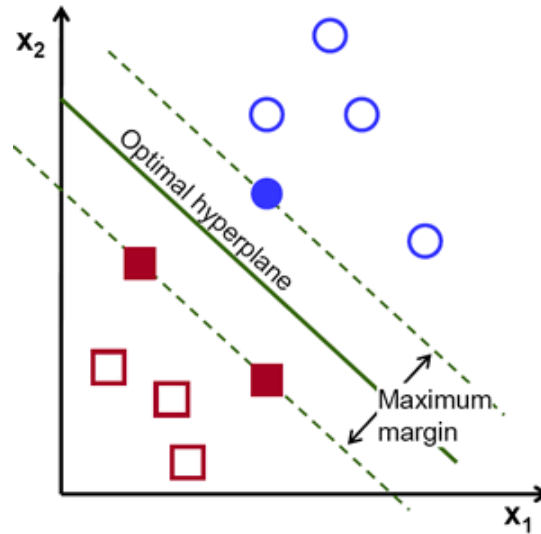
$$\text{Ridge: } \operatorname{argmin}_{x \in \mathbb{R}^n} \underbrace{\frac{1}{2} \|Ax - y\|_2^2}_{\text{Loss function.}} + \underbrace{\frac{\lambda}{2} \|x\|_2^2}_{\text{Regularization function.}}$$

$$\text{LASSO: } \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - y\|_2^2 + \boxed{\lambda \|x\|_1}$$

Loss function. Regularization function.

Binary Classification

Support Vector Machines



Ridge Regression with Block Coordinate Descent

Ridge solution:
(closed-form)

$$\begin{matrix} n \\ \left| \begin{matrix} x \end{matrix} \right. \end{matrix} = \left(\underbrace{\begin{matrix} m & \\ \hline A^T \end{matrix}}_{m} + \underbrace{\begin{matrix} n \\ \hline A \end{matrix}}_n \right)^{-1} \begin{matrix} & \\ A^T & \end{matrix} \begin{matrix} y \end{matrix}$$

λ = regularization parameter
 I_n = $n \times n$ Identity matrix
 y = labels

Similar to normal equation.

Ridge Regression with Block Coordinate Descent

Ridge solution:
(closed-form)

$$\mathbf{x} = \left(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_n \right)^{-1} \mathbf{A}^T \mathbf{y}$$

Block Coordinate Descent

λ = regularization parameter
 $\mathbf{I}_b = b \times b$ Identity matrix
 \mathbf{r} = residual

$\mathbf{A}' = \text{subsampled columns.}$

b | $\Delta \mathbf{x}_h$

Solution to sub-problem.

$$\Delta \mathbf{x}_h = \left(\mathbf{A}'^T \mathbf{A}' + \lambda \mathbf{I}_b \right)^{-1} \mathbf{A}'^T \mathbf{r}$$

Block Coordinate Descent in Parallel

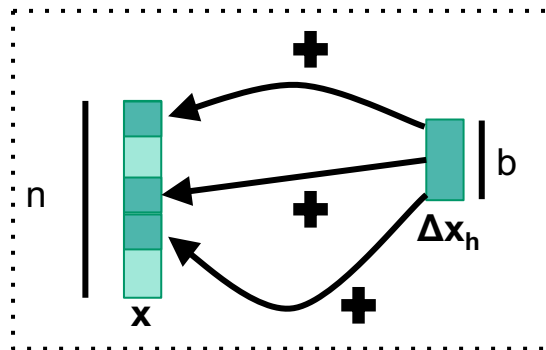
The diagram illustrates the `MPI_Allreduce` operation. It shows a local computation on the left: a vector $\Delta \mathbf{x}_h$ is calculated as $\mathbf{b} - \mathbf{A}'^T (\mathbf{A}' \Delta \mathbf{x}_h + \lambda \mathbf{I}_b \mathbf{r}_h)$. The terms $\mathbf{A}'^T \mathbf{A}' \Delta \mathbf{x}_h$ and $\lambda \mathbf{I}_b \mathbf{r}_h$ are then combined and reduced across all processes using `MPI_Allreduce` to produce the global residual \mathbf{r}_h .

Block Coordinate Descent in Parallel

$$b \mid \Delta x_h = \left(\begin{bmatrix} A'^T & A' \end{bmatrix} + \lambda_b \right)^{-1} \begin{bmatrix} A'^T & r_h \end{bmatrix}$$

MPI_Allreduce

Update solution

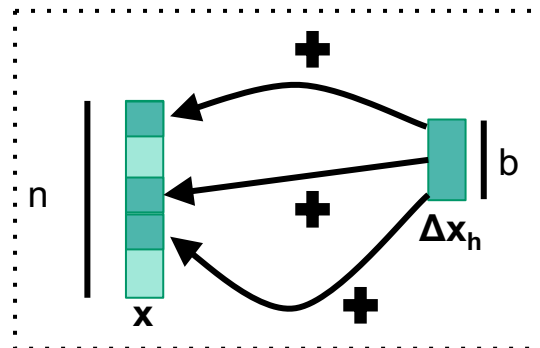


Block Coordinate Descent in Parallel

$$b \mid \Delta x_h = \left(\begin{bmatrix} A'^T & A' \end{bmatrix} + \lambda_b \begin{bmatrix} A'^T & 0 \end{bmatrix} \right)^{-1} \begin{bmatrix} A'^T & 0 \end{bmatrix} r_h$$

MPI_Allreduce

Update solution



Select new coordinates.
Repeat until convergence.

Communication-Avoiding Block Coordinate Descent

$$b \mid \begin{array}{|c|} \hline \Delta x_h \\ \hline \end{array} = \left(\begin{array}{|c|} \hline \begin{array}{|c|} \hline \vdots \vdots \vdots \vdots \vdots \\ \hline \end{array} A'^T \begin{array}{|c|} \hline \vdots \vdots \vdots \vdots \vdots \\ \hline \end{array} A' \begin{array}{|c|} \hline \vdots \vdots \vdots \vdots \vdots \\ \hline \end{array} \\ \hline \end{array} + \lambda I_b \right)^{-1} \begin{array}{|c|} \hline \begin{array}{|c|} \hline \vdots \vdots \vdots \vdots \vdots \\ \hline \end{array} A'^T \begin{array}{|c|} \hline \vdots \vdots \vdots \vdots \vdots \\ \hline \end{array} r_h \begin{array}{|c|} \hline \vdots \vdots \vdots \vdots \vdots \\ \hline \end{array} \\ \hline \end{array}$$

Block Coordinate Descent

Communication-Avoiding Block Coordinate Descent

$$b \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} \Delta \mathbf{x}_h = \left(\begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \mathbf{A}'^T \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} \mathbf{A}' + \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \lambda \mathbf{I}_b \right)^{-1} \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \mathbf{A}'^T \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} \mathbf{r}_h$$

s = CA parameter
 \mathbf{I}_{sb} = $sb \times sb$ Identity
 \mathbf{r} = residual

$$sb \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \mathbf{A}'^T \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \mathbf{A}' + \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \lambda \mathbf{I}_{sb}$$

Compute larger Gram matrix, \mathbf{G} .
 (Do not invert!)

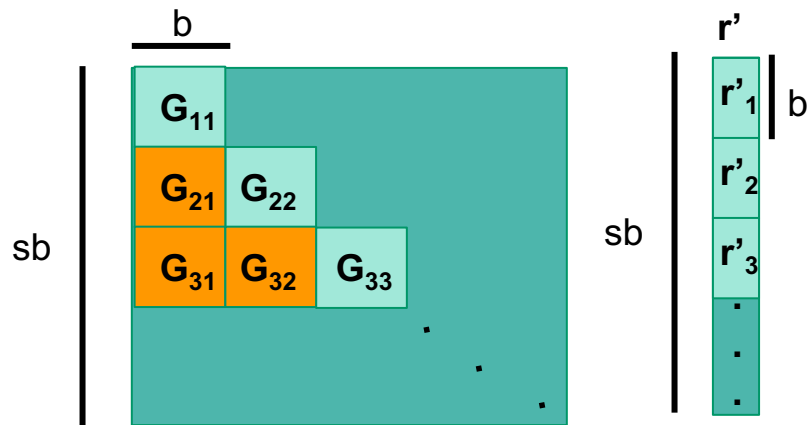
$$sb \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \mathbf{A}'^T \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} \mathbf{r}$$

Larger residual contribution.
 \mathbf{r} is from this iteration 12

Redundantly on all processors.

Dimensions exaggerated for clarity.

Communication-Avoiding Block Coordinate Descent



Partition G and r' .

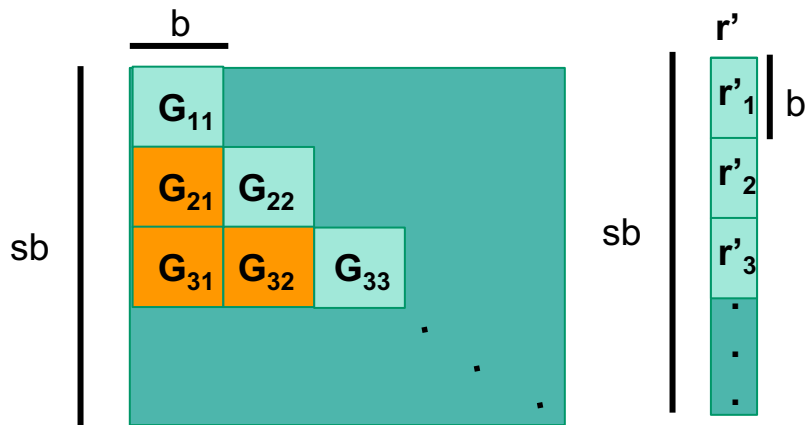
(G is symmetric)

Redundantly on all processors.

No communication for inner iterations.

Communication-Avoiding Block Coordinate Descent

Solution to 1st sub-problem.



$$\Delta x_1 = \begin{bmatrix} G_{11} \end{bmatrix}^{-1} r'_1$$

Inner iteration 1

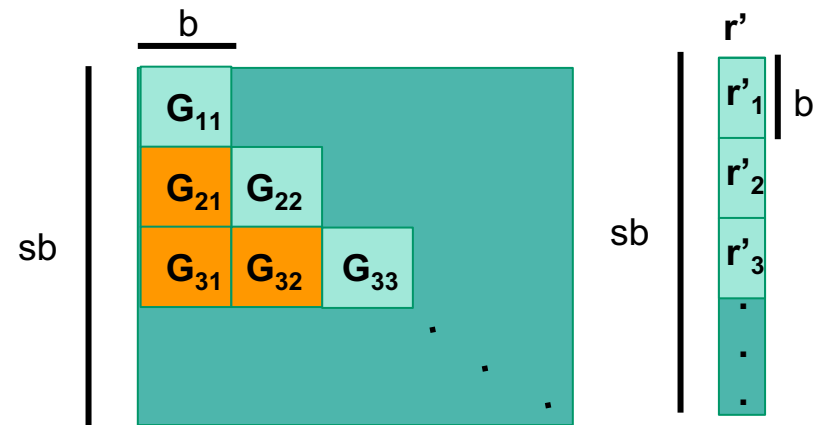
r'_2, r'_3, \dots become stale as soon as Δx_1 is computed.

Orange blocks are the updates

Redundantly on all processors.

No communication for inner iterations.

Communication-Avoiding Block Coordinate Descent



$$\Delta x_1 = \begin{bmatrix} G_{11} \end{bmatrix}^{-1} r'_1$$

Inner iteration 1

$$\Delta x_2 = \begin{bmatrix} G_{22} \end{bmatrix}^{-1} r'_2 - G_{21} \Delta x_1$$

Inner iteration 2

Solution to 2nd sub-problem.

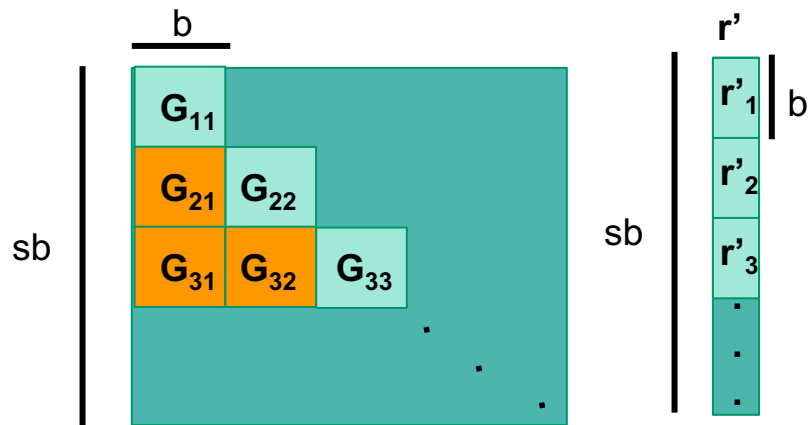
r'_2, r'_3, \dots become stale as soon as Δx_1 is computed.

Orange blocks are the updates

Redundantly on all processors.

No communication for inner iterations.

Communication-Avoiding Block Coordinate Descent



Inner iteration 1

$$\Delta x_1 = \left(G_{11} \right)^{-1} r'_1$$

...

Inner iteration 2

$$\Delta x_2 = \left(G_{22} \right)^{-1} r'_2 - G_{21} \Delta x_1$$

...

Inner iteration 3

$$\Delta x_3 = \left(G_{33} \right)^{-1} r'_3 - G_{31} \Delta x_1 - G_{32} \Delta x_2$$

Solution to 3rd sub-problem.

Subtraction terms are needed for CA and non-CA solutions to agree.

Algorithm 1

Block Coordinate Descent (BCD) Algorithm

- 1: **Input:** $X \in \mathbb{R}^{d \times n}, y \in \mathbb{R}^n, H > 1, w_0 \in \mathbb{R}^d, b \in \mathbb{Z}_+ \text{ s.t. } b \leq d$
- 2: **for** $h = 1, 2, \dots, H$ **do**
- 3: choose $\{i_m \in [d] | m = 1, 2, \dots, b\}$ uniformly at random without replacement
- 4: $\mathbb{I}_h = [e_{i_1}, e_{i_2}, \dots, e_{i_b}]$
- 5: $\Gamma_h = \frac{1}{n} \mathbb{I}_h^T X X^T \mathbb{I}_h + \lambda \mathbb{I}_h^T \mathbb{I}_h$
- 6: $\Delta w_h = \Gamma_h^{-1} \left(-\lambda \mathbb{I}_h^T w_{h-1} - \frac{1}{n} \mathbb{I}_h^T X z_{h-1} + \frac{1}{n} \mathbb{I}_h^T X y \right)$
- 7: $w_h = w_{h-1} + \mathbb{I}_h \Delta w_h$
- 8: $z_h = z_{h-1} + X^T \mathbb{I}_h \Delta w_h$
- 9: **Output** w_H

Communication
(every outer iteration)

No communication

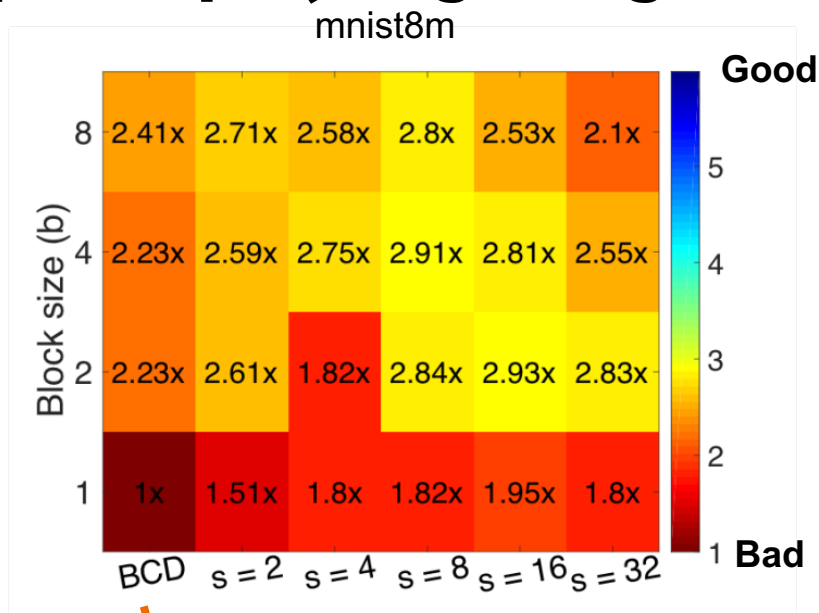
} Communication
(every iteration)

Algorithm 2

Communication-Avoiding Block Coordinate Descent (CA-BCD) Algorithm

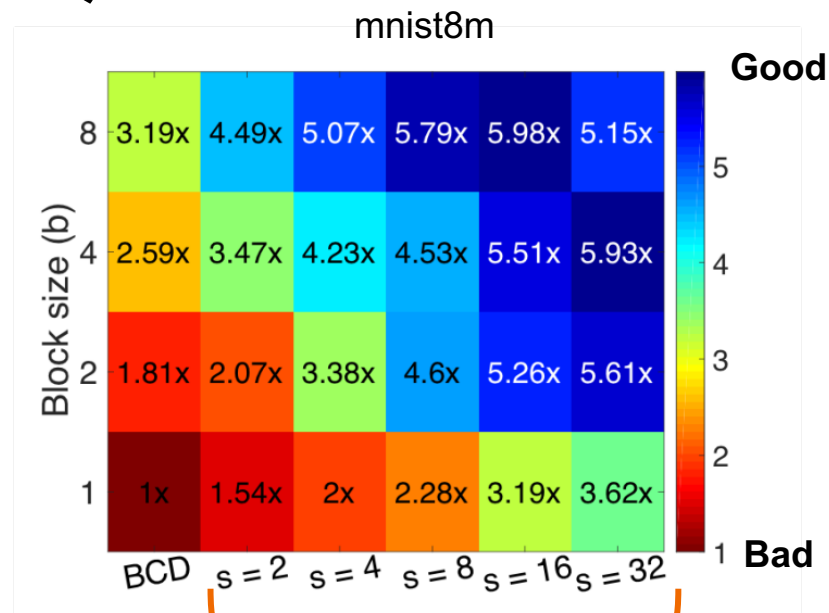
- 1: **Input:** $X \in \mathbb{R}^{d \times n}, y \in \mathbb{R}^n, H > 1, w_0 \in \mathbb{R}^d, b \in \mathbb{Z}_+ \text{ s.t. } b \leq d$
- 2: **for** $k = 0, 1, \dots, \frac{H}{s}$ **do**
- 3: **for** $j = 1, 2, \dots, s$ **do**
- 4: choose $\{i_m \in [d] | m = 1, 2, \dots, b\}$ uniformly at random without replacement
- 5: $\mathbb{I}_{sk+j} = [e_{i_1}, e_{i_2}, \dots, e_{i_b}]$
- 6: let $Y = [\mathbb{I}_{sk+1}, \mathbb{I}_{sk+2}, \dots, \mathbb{I}_{sk+s}]^T X$.
- 7: compute the Gram matrix, $G = \frac{1}{n} Y Y^T + \lambda I$.
- 8: **for** $j = 1, 2, \dots, s$ **do**
- 9: Γ_{sk+j} are the $b \times b$ diagonal blocks of G .
- 10:
$$\Delta w_{sk+j} = \Gamma_{sk+j}^{-1} \left(-\lambda \mathbb{I}_{sk+j}^T w_{sk} - \lambda \sum_{t=1}^{j-1} \left(\mathbb{I}_{sk+j}^T \mathbb{I}_{sk+t} \Delta w_{sk+t} \right) - \frac{1}{n} \mathbb{I}_{sk+j}^T X z_{sk} \right. \\ \left. - \frac{1}{n} \sum_{t=1}^{j-1} \left(\mathbb{I}_{sk+j}^T X X^T \mathbb{I}_{sk+t} \Delta w_{sk+t} \right) + \frac{1}{n} \mathbb{I}_{sk+j}^T X y \right)$$
- 11: $w_{sk+j} = w_{sk+j-1} + \mathbb{I}_{sk+j} \Delta w_{sk+j}$
- 12: $z_{sk+j} = z_{sk+j-1} + X^T \mathbb{I}_{sk+j} \Delta w_{sk+j}$
- 13: **Output** w_H

Speedups (Ridge Regression)



64 nodes of Edison

Standard Algorithm



1K nodes of Edison

Communication-Avoiding

Not Just Ridge Regression

Applies to **proximal methods** (e.g. LASSO, Group LASSO, elastic-net).

Recall that, Inner loop = no communication.

If nonlinearity is in inner loop, then **CA possible**.

Also applies to **Support Vector Machines** (current work).

Nonlinearity is in inner loop.

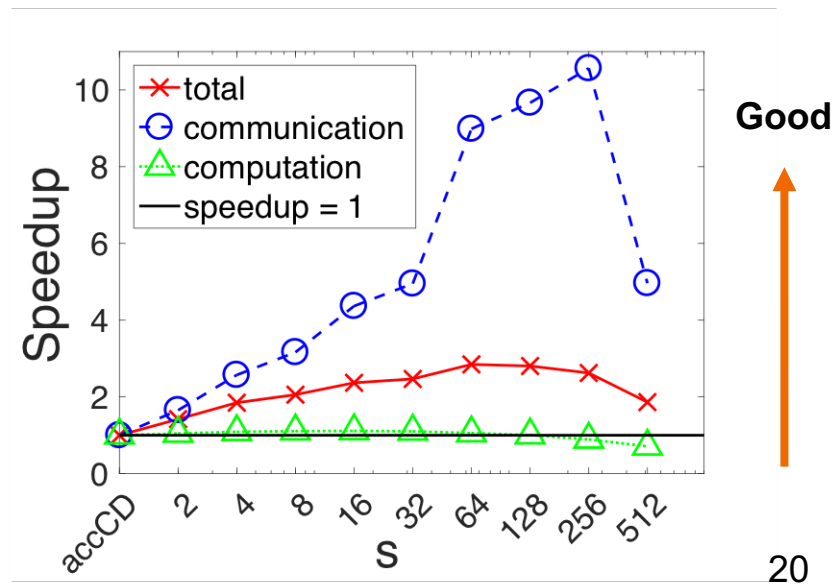
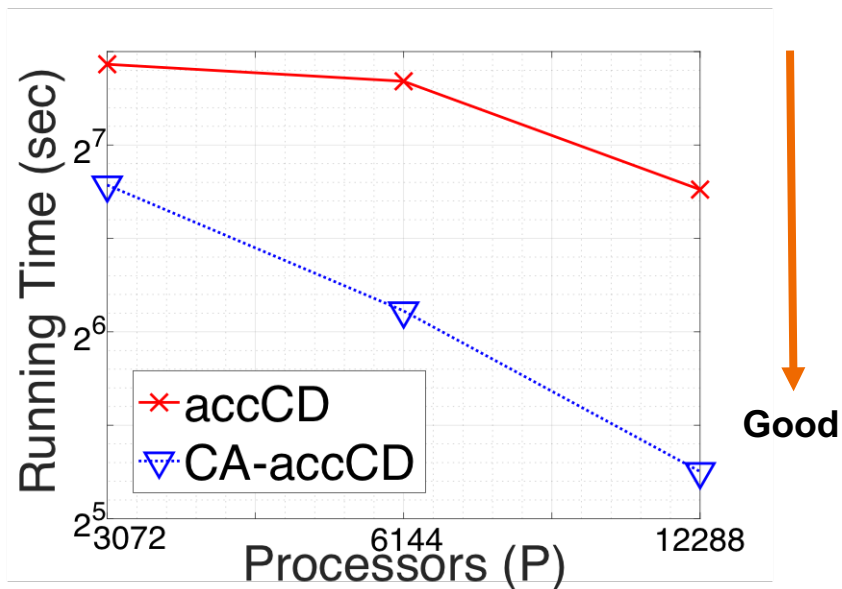
CALASSOS: Scalable Proximal Methods

Strong scaling and speedups on **Url dataset (2M by 3M)**.

accCD = accelerated
Coordinate Descent



CA-technique
applies to accelerated
methods.



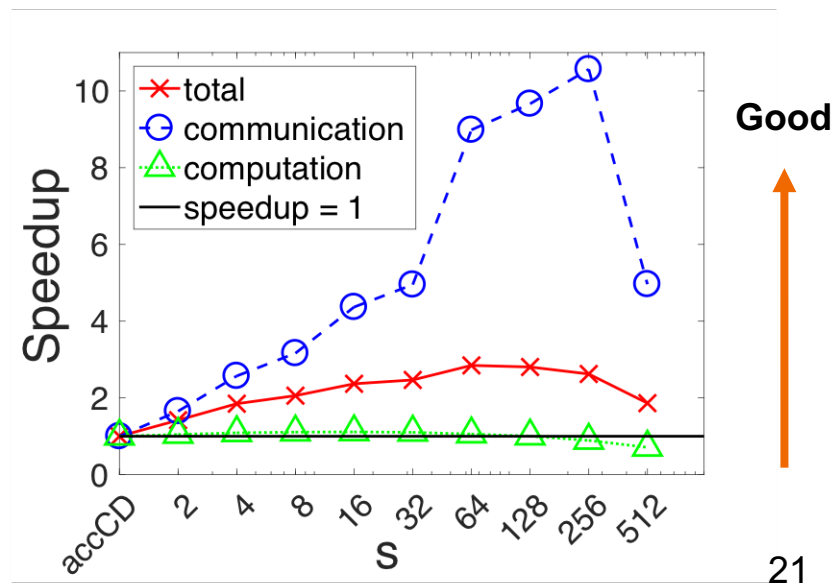
CALASSOS: Scalable Proximal Methods

CA-method perform s^2 more flops.

But still get **computation speedup**.

Due to BLAS-3 calls instead of BLAS-1 in non-CA.

BLAS-3 = Cache-efficient computation + higher flops rate.



CA-SVM: Preliminary Results

Based on Dual Coordinate Descent for Linear SVM (Hsieh, et. al.)

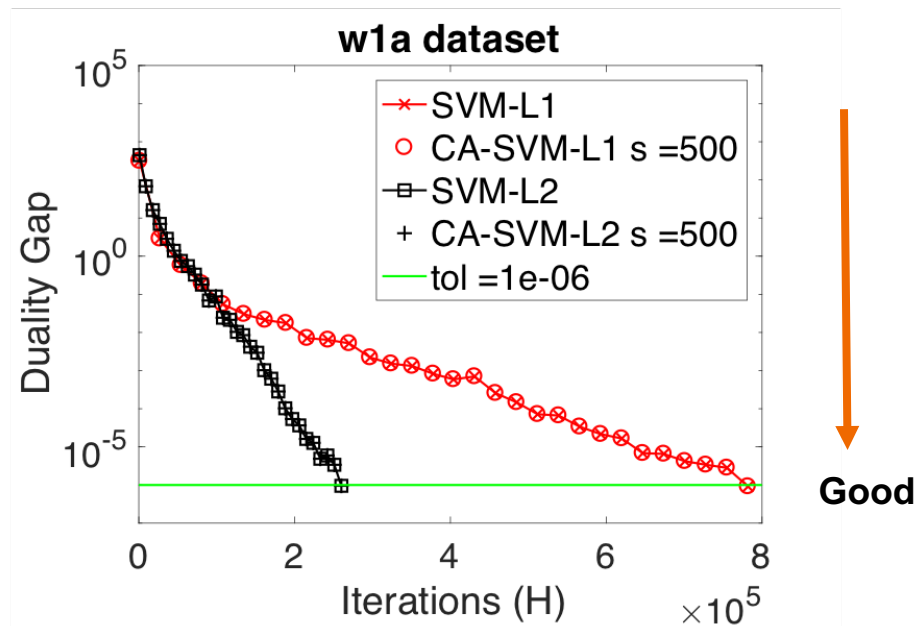
SVM-L1 = Hinge loss = $\max(0, 1 - A^T_i x y_i)$

SVM-L2 = (Hinge loss)²

easier so, converges quickly.

Numerically stable (**unlike CA-Krylov**)!

Ditto for Ridge and Proximal.



Summary and Future Work

Large speedups when **latency dominates**.

Provably communication-avoiding.

CA-technique applies to **non-linear optimization**.

How far can we go (e.g. Logistic regression)?

Speedups on **other platforms and frameworks**?

Example: Cloud + Spark is latency dominated.

Expect greater speedups!

Problem	MPI Speedup
Ridge Regression	Up to 6.1x
Proximal Least-Squares	Up to 5.1x
SVM	Similar expected

Questions?

Thanks to co-authors, collaborators, and sponsors!