

Sparse Linear Programming via Primal and Dual Augmented Coordinate Descent

Ian E.H. Yen¹, Kai Zhong¹, Cho-Jui Hsieh², Pradeep Ravikumar¹ and Inderjit S. Dhillon¹

¹University of Texas at Austin. ²University of California at Davis

Abstract

- Linear Programming (LP) has been widely used in different areas.
- State-of-the-art algorithms, e.g. interior-point method and primal, dual simplex methods, have complexity at least quadratic in the number of variables or constraints.
- We investigate a general LP algorithm based on the combination of Augmented Lagrangian and Coordinate Descent (AL-CD), giving an iteration complexity of $O((\log(1/\epsilon))^2)$ with $O(\text{nnz}(A))$ cost per iteration.
- Experiment on large-scale LP shows that AL-CD is orders of magnitude faster than state-of-the-art implementations of IPM and Simplex Methods.

Sparse Linear Program

Given vectors $c \in \mathbb{R}^n$, $b \in \mathbb{R}^m$ and $m \times n$ matrix

$$A = \begin{bmatrix} A_I \\ A_E \end{bmatrix} = [A_b \ A_f],$$

the *primal* and *dual* forms of Linear Program (LP) are

$$\begin{aligned} \min_{x \in \mathbb{R}^n} c^T x & & \min_{y \in \mathbb{R}^m} b^T y \\ \text{s.t. } A_I x &\leq b_I & \text{s.t. } -A_b^T y &\leq c_b \\ A_E x &= b_E & -A_f^T y &= c_f \\ x_j &\geq 0, j \in [n_b] & y_i &\geq 0, i \in [m]. \end{aligned}$$

We say a LP is sparse in the sense that

- (i) $\text{nnz}(A) \ll m \times n$.
- (ii) $\text{nnz}(y^*) \ll m$ (dual sparsity) or $\text{nnz}(x^*) \ll n$ (primal sparsity).

Sparse Linear Program: Examples

- L1-regularized Multiclass SVM ($\{x_i, y_i\}_{i=1,2,\dots,l}$ is the given dataset)

$$\begin{aligned} \min_{w_m, \xi_i} \lambda \sum_{m=1}^k \|w_m\|_1 + \sum_{i=1}^l \xi_i \\ \text{s.t. } w_{y_i}^T x_i - w_m^T x_i \geq e_i^m - \xi_i, \forall (i, m) \end{aligned}$$

- Sparse Inverse Covariance Estimation (S is the covariance matrix)

$$\begin{aligned} \min_{\Omega \in \mathbb{R}^{d \times d}} \|\Omega\|_1 \\ \text{s.t. } \|\Omega - I_d\|_{\max} \leq \lambda \end{aligned}$$

- Non-negative Matrix Factorization (M is the given matrix and p is any positive vector of size n)

$$\begin{aligned} \min_{X \in \mathbb{R}^{n \times n}} p^T \text{diag}(X) \\ \text{s.t. } \|M - MX\|_1 \leq \epsilon, X(i, i) \leq 1, X(i, j) \leq X(j, i), \forall i, j \end{aligned}$$

Algorithms for Linear Program

- Simplex Method [1]:** Moving between corner points requires solutions of two $m \times m$ linear system ($O(m^3)$ for m iterations). #iter $\approx O(m)$ typically.
- Interior Point Method (IPM) [1]:** Reduce LP to series of (asymptotically ill-conditioned) *unconstrained problems* via Barrier function. Newton step requires solution of $m \times m$ linear system ($O(m^3)$). #iter= $O(\log(1/\epsilon))$.
- Augmented Lagrangian Method (ALM):** Reduce LP to series of *bound-constrained Quadratic Problem*. #iter= $O(\log(1/\epsilon))$. Previous works employ Proj-Newton / Proj-GD for sub-problems. However, Proj-Newton has no guarantee on #iter, while Proj-GD converges very slow.
- This paper** \Rightarrow Randomized Coordinate Descent (RCD) with ALM gives $O(\text{nnz}(A) \log^2(1/\epsilon))$ overall complexity, and more efficient for primal, dual-sparse problems via an Active-Set strategy.

Augmented Lagrangian Method (ALM)

Let $g(y)$ denote the objective of dual LP (taking ∞ for infeasible points). Then a primal ALM is equivalent to a dual Proximal-Point iterates

$$y^{t+1} = \underset{y}{\text{argmin}} g(y) + \frac{1}{2\eta_t} \|y - y^t\|^2, \quad (1)$$

where we find y^{t+1} by solving the dual of (1)

$$\begin{aligned} \min_{x, \xi} c^T x + \frac{\eta_t}{2} \left\| \begin{bmatrix} A_I x - b_I + \xi \\ A_E x - b_E \end{bmatrix} + \frac{1}{\eta_t} \begin{bmatrix} y_I^t \\ y_E^t \end{bmatrix} \right\|^2 \\ \text{s.t. } x_b \geq 0, \xi \geq 0 \end{aligned} \quad (2)$$

to obtain

$$y^{t+1} = y(x^*, \xi^*) = \eta_t \begin{bmatrix} A_I x^* - b_I + \xi^* \\ A_E x^* - b_E \end{bmatrix} + \begin{bmatrix} y_I^t \\ y_E^t \end{bmatrix}.$$

ALM with Coordinate Descent (AL-CD)

The quadratic program (with ξ eliminated)

$$\begin{aligned} \min_x F(x) = c^T x + \frac{\eta_t}{2} \left\| \begin{bmatrix} A_I x - b_I + y_I^t / \eta_t \\ A_E x - b_E + y_E^t / \eta_t \end{bmatrix} \right\|^2 \\ \text{s.t. } x_b \geq 0. \end{aligned}$$

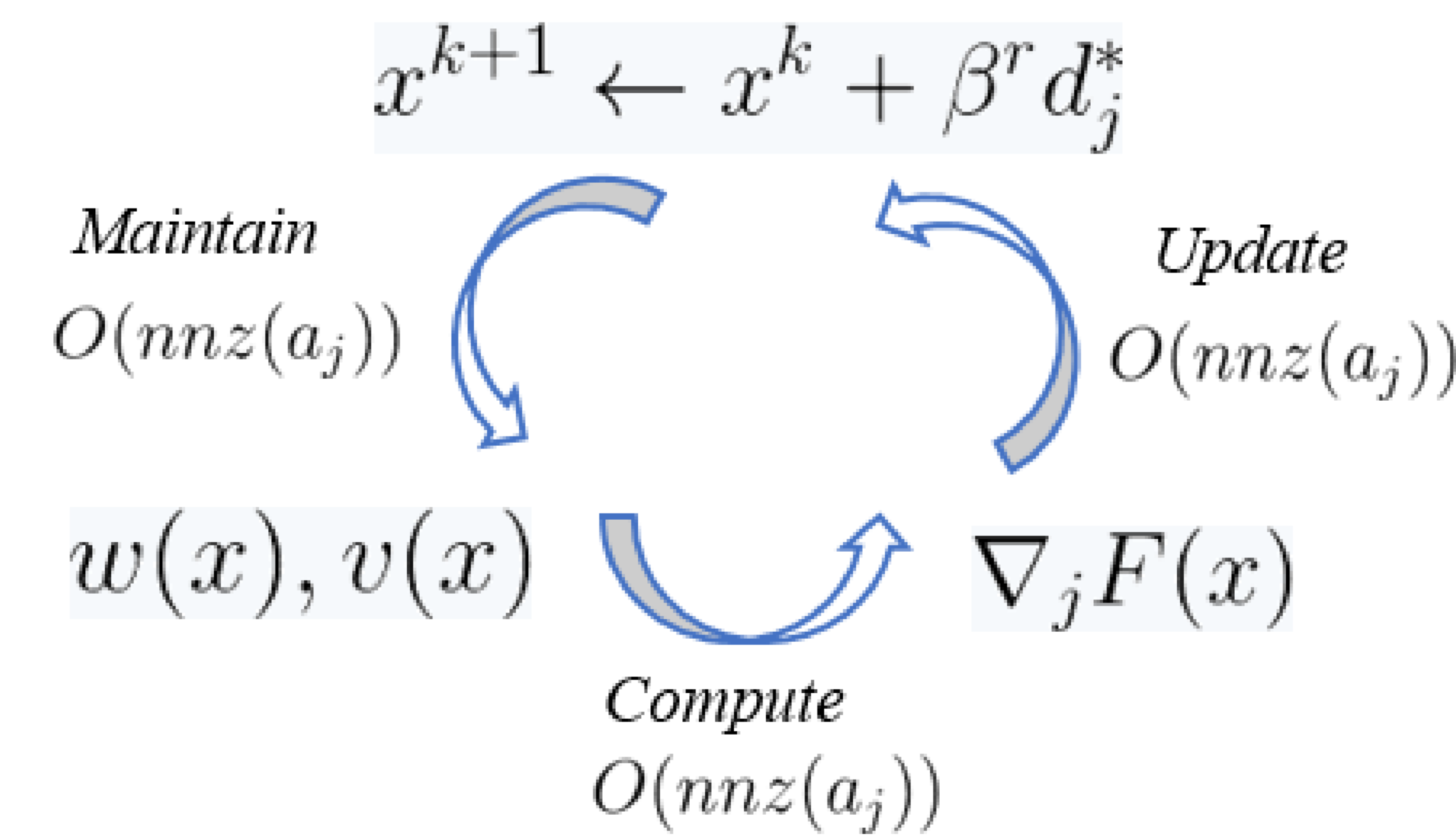
has

$$\nabla F(x) = c + \eta_t A_I^T [w(x)]_+ + \eta_t A_E^T v(x)$$

where

$$\begin{aligned} w(x) &= A_I x - b_I + y_I^t / \eta_t \\ v(x) &= A_E x - b_E + y_E^t / \eta_t \end{aligned}$$

- Given $w(x)$, $v(x)$, gradient $\nabla_j F(x)$ of coordinate can be evaluated in $O(\text{nnz}(a_j))$.
- Maintaining $w(x)$, $v(x)$ after each coordinate j update needs $O(\text{nnz}(a_j))$.



Randomized Coordinate Descent for Subproblem

INPUT: $\eta_t > 0$ and $(x^{t,0}, w^{t,0}, v^{t,0})$

OUTPUT: $(x^{t,k}, w^{t,k}, v^{t,k})$

repeat

- Pick a coordinate j uniformly at random
- Compute $\nabla_j F(x)$, $\nabla_j^2 F(x)$.
- Obtain Newton direction d_j^* , by solving a single-variable QP:

$$d_j^* = [x_j - \nabla_j F(x) / \nabla_j^2 F(x)]_+ - x_j,$$

- Do line search to find step size β^r .

$$F(x^{t,k} + \beta^r d_j^* e_j) - F(x^{t,k}) \leq \sigma \beta^r (\nabla_j F(x^{t,k}) d_j^*). \quad (3)$$

- Update $x^{t,k+1} \leftarrow x^{t,k} + \beta^r d_j^*$.

- Maintain relation $w(x)$, $v(x)$.

- $k \leftarrow k + 1$.

until $\|d^*(x)\|_\infty \leq \epsilon_t$.

Convergence Analysis

- Subproblem is **strongly convex when restricted to $N(\bar{A})^\perp$** (CNSC[2]), where $\bar{A} = [A_I \ I; A_E \ O]$. \Rightarrow We show RCD converges to $F(x) - F(x^*) \leq \epsilon$ in $O(n \log(1/\epsilon))$ number of iterations w.h.p., and the cost for sub-problem is $O(\text{nnz}(A) \log(1/\epsilon))$
- The ALM does not amplify error when solving subproblem inexactly. \Rightarrow If each subproblem is approximated to ϵ then after t iterations we have approximation error at most $t\epsilon$.
- Overall #iterations= $O(tn \log(t/\epsilon)) = O(n \log^2(1/\epsilon))$.

Implementation Details

- Maintain an active set \mathcal{A}^t when solving subproblem. Remove coordinate $j \in [n_b]$ from \mathcal{A}^t if $x_j = 0$ and $\nabla_j F(x) > \epsilon_A$.
- Two-Phase strategy:
Phase-I: Spend constant cost on each sub-problem. ($|\mathcal{A}^t|$ is large)
Phase-II: Solve each sub-problem to precision ϵ_t , with η_t, ϵ_t dynamically adjusted. ($|\mathcal{A}^t|$ is small)
- For ill-conditioned problem, use Projected-Newton-CG to achieve asymptotic fast convergence (when active set is almost fixed).

Experiments on L1 Multiclass SVM

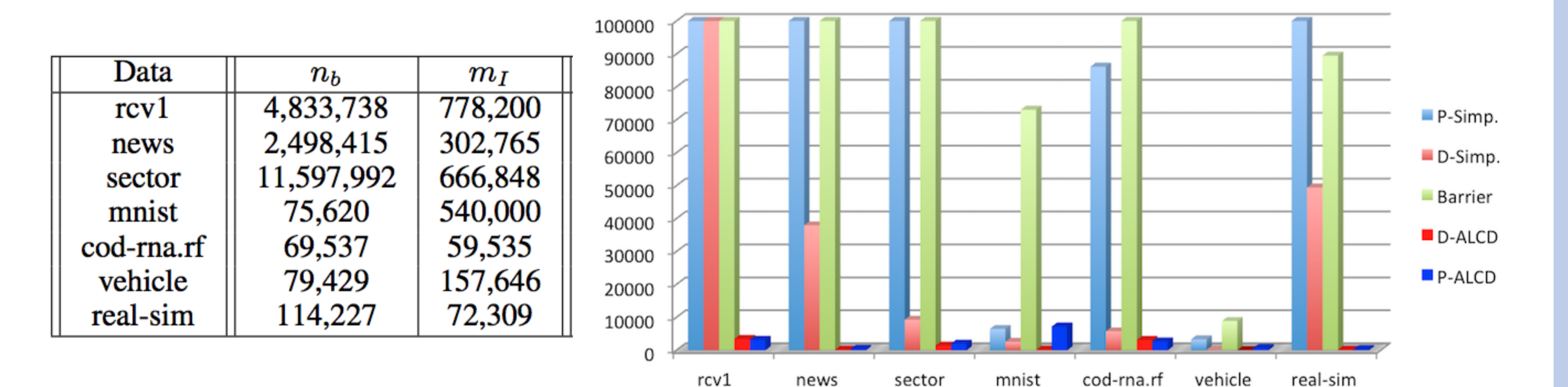


Figure: Timing Results (in sec., 100,000s denotes >100,000s) on L1-SVM

L1-regularized Multiclass SVM: $\text{nnz}(A) \ll mn$, $\text{nnz}(y^*) \ll m$, and $\text{nnz}(x^*) \ll n$.

Experiments on Sparse Inverse Covariance Estimation

Data	n_b	m_I	m_E	n_f	P-Simp	D-Simp	Barrier	D-ALCD	P-ALCD
textmine	60,876	60,876	43,038	43,038	> 48hr	> 48hr	> 48hr	43,096	18,507
E2006	55,834	55,834	32,174	32,174	> 48hr	> 48hr	94623	> 48hr	4,207
dorothea	47,232	47,232	1,600	1,600	3,980	103	82	47	38

Figure: Timing Results (in sec.) on Sparse Inverse Covariance Estimation

- Sparse Inverse Covariance Estimation:

$$\min_{\Omega \in \mathbb{R}^{d \times d}} \|\Omega\|_1, \text{ s.t. } \|\Omega - I_n\|_{\max} \leq \lambda$$

- $S = Z^T Z$, Z is $l \times n \ll n^2 \Rightarrow$ Transform to:

$$\min_{\Omega \in \mathbb{R}^{n \times n}, Y \in \mathbb{R}^{l \times n}} \|\Omega\|_1, \text{ s.t. } \|Z^T Y - I_n\|_{\max} \leq \lambda, Y = Z \Omega$$

- $\text{nnz}(A) \ll mn$ and $\text{nnz}(\Omega^*) \ll n$.

Experiments on Non-negative Matrix Factorization

Data	n_b	m_I	P-Simp	D-Simp	Barrier	D-ALCD	P-ALCD
micromass	2,896,770	4,107,438	> 96hr	> 96hr	280,230	12,966	12,119
ocr	6,639,433	13,262,864	> 96hr	> 96hr	284,530	40,242	> 96hr

Figure: Timing Results (in sec.) on non-negative matrix factorization

NMF: $\text{nnz}(A) \ll mn$.

References

- Jorge Nocedal and Stephen Wright. Numerical optimization. Springer Science & Business Media, 2006.
- I. E.H. Yen, C.J. Hsieh, P. Ravikumar, and I. S. Dhillon. "Constant nullspace strong convexity and fast convergence of proximal methods under high-dimensional settings." NIPS 2014.