

Faster Coordinate Descent via Adaptive Importance Sampling

Dmytro Perekrestenko, Volkan Cevher, Martin Jaggi

pdmytro@nari.ee.ethz.ch, volkan.cevher@epfl.ch, martin.jaggi@epfl.ch



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Abstract

Coordinate descent methods employ random partial updates of decision variables in order to solve huge-scale convex optimization problems. In this work, we introduce new *adaptive* rules for the random selection of their updates. By adaptive, we mean that our selection rules are based on the dual residual or the primal-dual gap estimates and can change at each iteration. We theoretically characterize the performance of our selection rules and demonstrate improvements over the state-of-the-art, and extend our theory and algorithms to general convex objectives. Numerical evidence with hinge-loss support vector machines and Lasso confirm that the practice follows the theory.

Problem

We solve the problem of the following template:

$$\min_{\alpha \in \mathbb{R}^n} f(A\alpha) + \sum_i g_i(\alpha_i), \quad (1)$$

where A is the data matrix, f is a smooth convex function, and each g_i is a general convex function.

Primal-dual framework

We are working in the following primal-dual optimization framework ([Dünner et al., 2016]):

$$\min_{\alpha \in \mathbb{R}^n} [\mathcal{O}_A(\alpha) := f(A\alpha) + \sum_i g_i(\alpha_i)], \quad (A)$$

$$\min_{w \in \mathbb{R}^d} [\mathcal{O}_B(w) := f^*(w) + \sum_i g_i^*(-\mathbf{a}_i^\top w)], \quad (B)$$

where we have $A = [\mathbf{a}_1, \dots, \mathbf{a}_n]$.

Duality gap

The duality gap is the difference between primal and dual solutions:

$$G(\alpha, w) := \mathcal{O}_A(\alpha) - (-\mathcal{O}_B(w)), \quad (2)$$

which can be written as a sum of coordinate-wise gaps:

$$G(\alpha) = \sum_i G_i(\alpha_i) := \sum_i (g_i^*(-\mathbf{a}_i^\top w) + g_i(\alpha_i) + \alpha_i \mathbf{a}_i^\top w). \quad (3)$$

Preliminaries

Definition (Dual Residual. A generalization of [Csiba et al., 2015]). i -th dual residue on iteration t is given by:

$$\kappa_i^{(t)} := \min_{u \in \partial g_i^*(-\mathbf{a}_i^\top w^{(t)})} |u - \alpha_i^{(t)}|.$$

Definition (Nonuniformity measure, [Osokin et al., 2016]). The nonuniformity measure of a vector $\mathbf{x} \in \mathbb{R}^n$:

$$\chi(\mathbf{x}) := \sqrt{1 + n^2 \text{Var}[\mathbf{p}]},$$

Definition (Averaged residual). Averaged residual on iteration t is given by:

$$F^{(t)} := \frac{1}{n^2 \beta} \sum_{i \in I_t} \left(\frac{|\kappa_i^{(t)}|^2 \|\mathbf{a}_i\|^2}{p_i^{(t)}} \right). \quad (4)$$

Algorithm 1 Coordinate Descent

- 1: Let $\alpha^{(0)} := \mathbf{0} \in \mathbb{R}^n$, $w^{(0)} := w(\alpha^{(0)})$
- 2: **for** $t = 0, 1, \dots, T$ **do**
- 3: Sample $i \in [n]$ randomly according to $\mathbf{p}^{(t)}$
- 4: Find $\Delta \alpha_i$ minimizing $\mathcal{O}_A(\alpha^{(t)} + \mathbf{e}_i \Delta \alpha_i)$
- 5: $\alpha^{(t+1)} := \alpha^{(t)} + \mathbf{e}_i \Delta \alpha_i$
- 6: $w^{(t+1)} := w(\alpha^{(t+1)})$
- 7: **end for**

Adaptive Sampling residual-based CD

Theorem. If f is a $\frac{1}{\beta}$ -smooth and g_i^* is L_i -Lipschitz for each i , then the residual-based CD iterates satisfy

$$\mathbb{E}[\varepsilon_A^{(t)}] \leq \frac{2F^\circ n^2 + \frac{2\varepsilon_A^{(0)}}{p_{\min}}}{\frac{2}{p_{\min}} + t}. \quad (5)$$

A duality gap $G(\bar{\alpha}) \leq \varepsilon$ is reached after an overall number of iterations T whenever

$$T \geq \max \left\{ 0, \frac{1}{p_{\min}} \log \left(\frac{2\varepsilon_A^{(0)}}{n^2 p_{\min} F^\circ} \right) \right\} + \frac{5F^\circ n^2}{\varepsilon} - \frac{1}{p_{\min}}. \quad (6)$$

Here F° is an upper bound on $\mathbb{E}[F^{(t)}]$.

Adaptive Sampling gap-based CD

Theorem. If f be a $\frac{1}{\beta}$ -smooth and g_i^* is L_i -Lipschitz for each i , then the gap-based CD iterates satisfy

$$\mathbb{E}[\varepsilon_A^{(t)}] \leq \frac{2F_g^\circ n^2 + 2n\varepsilon_A^{(0)}}{t + 2n}, \quad (7)$$

where F_g° is an upper bound on $\mathbb{E}[F_g^{(t)}]$. The \vec{G} and \vec{F} are defined as follows:

$$\vec{G} := (G_i(\alpha^{(t)}))_{i=1}^n, \quad \vec{F} := (\|\mathbf{a}_i\|^2 |\kappa_i^{(t)}|^2)_{i=1}^n,$$

and $F_g^{(t)}$ is:

$$F_g^{(t)} := \frac{\chi(\vec{F})}{n\beta(\chi(\vec{G}))^3} \sum_i \|\mathbf{a}_i\|^2 |\kappa_i^{(t)}|^2. \quad (8)$$

Variations along the theme

- *uniform* - sample uniformly at random.
- *supportSet uniform* - $p_i^{(t)} := \begin{cases} \frac{1}{m_i}, & \text{if } \kappa_i^{(t)} \neq 0 \\ 0, & \text{otherwise.} \end{cases}$
- *adaptive* - $p_i^{(t)} := \frac{|\kappa_i^{(t)}| \|\mathbf{a}_i\|}{\sum_j |\kappa_j^{(t)}| \|\mathbf{a}_j\|}$.

- *ada-uniform* - $p_i^{(t)} := \begin{cases} \frac{0.5}{m_i} + 0.5 \frac{|\kappa_i^{(t)}| \|\mathbf{a}_i\|}{\sum_j |\kappa_j^{(t)}| \|\mathbf{a}_j\|}, & \text{if } \kappa_i^{(t)} \neq 0 \\ 0, & \text{otherwise} \end{cases}$
- *importance* - sample with a fixed non-uniform variant of *adaptive* obtained by bounding $\kappa_i^{(t)}$ with $2L_i$: $p_i := \frac{L_i \|\mathbf{a}_i\|}{\sum_j L_j \|\mathbf{a}_j\|}$.
- *ada-gap* - $p_i^{(t)} := \frac{G_i(\alpha^{(t)})}{G(\alpha^{(t)})}$.
- *gap-per-epoch* - Use *ada-gap* but with updates per-epoch. The gap-based distribution is only recomputed at the beginning of each epoch and stays fixed during each epoch.

Experimental results

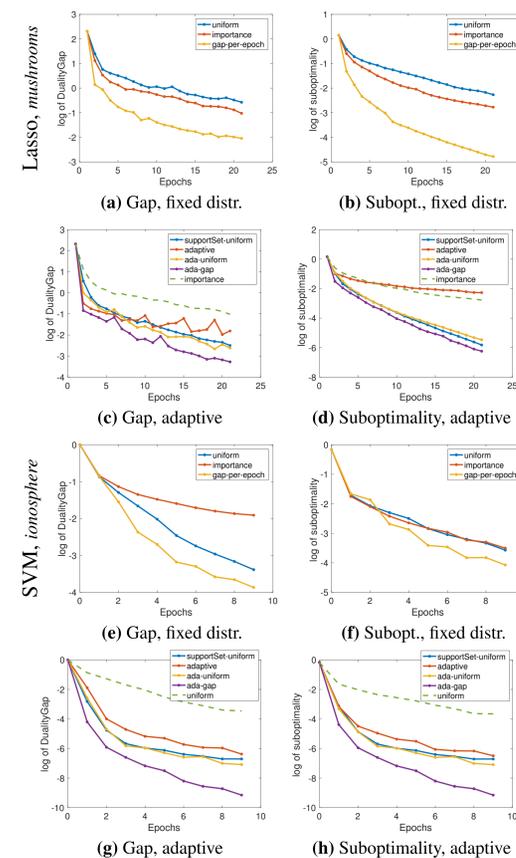


Figure 1: Lasso (two top rows) and SVM (two bottom rows). Comparison of different fixed and adaptive variants of CD, reporting duality gap and suboptimality measures vs. epochs - mushrooms and ionosphere datasets.

Dataset	d	n	$\text{nnz}/(nd)$	$c_v = \frac{\mu(\ \mathbf{a}_i\)}{\sigma(\ \mathbf{a}_i\)}$
mushrooms	112	8124	18.8%	1.34
ionosphere	351	33	88%	3.07

Table 1: Datasets

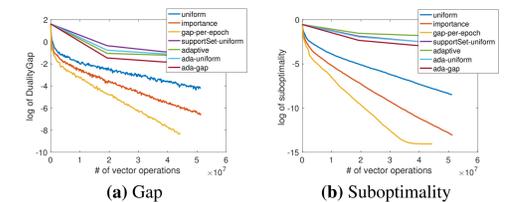


Figure 2: Lasso on the mushrooms dataset. Performance in terms of duality gap and suboptimality, plotted against the total number of vector operations.

Algorithm	Cost per Epoch
uniform	$\mathcal{O}(\text{nnz})$
importance	$\mathcal{O}(\text{nnz} + n \log(n))$
gap-per-epoch	$\mathcal{O}(\text{nnz} + n \log(n))$
supportSet-uniform	$\mathcal{O}(n \cdot \text{nnz})$
adaptive	$\mathcal{O}(n \cdot \text{nnz})$
ada-uniform	$\mathcal{O}(n \cdot \text{nnz})$
ada-gap	$\mathcal{O}(n \cdot \text{nnz})$

Table 2: A summary of computational costs

Conclusion

- We investigated *adaptive* rules for adjusting the sampling probabilities in coordinate descent.
- Our theoretical results provide improved convergence rates for a more general class of algorithm schemes on one hand, and optimization problems on the other hand, where we are able to directly analyze CD on general convex objectives (as opposed to strongly convex regularizers in previous works).
- Our results are particularly useful for L1 problems and (original) hinge-loss objectives.
- We advocate the use of the computationally efficient gap-per-epoch sampling scheme in practice. While the scheme is close to the ones supported by our theory, an explicit primal-dual convergence analysis remains a future research question.

References

- [Csiba et al., 2015] Csiba, D., Qu, Z., and Richtárik, P. (2015). Stochastic Dual Coordinate Ascent with Adaptive Probabilities. In *ICML 2015 - Proceedings of the 32th International Conference on Machine Learning*.
- [Dünner et al., 2016] Dünner, C., Forte, S., Takáč, M., and Jaggi, M. (2016). Primal-Dual Rates and Certificates. In *ICML 2016 - Proceedings of the 33th International Conference on Machine Learning*.
- [Osokin et al., 2016] Osokin, A., Alayrac, J.-B., Lukaszewicz, I., Dokania, P., and Lacoste-Julien, S. (2016). Minding the Gaps for Block Frank-Wolfe Optimization of Structured SVMs. In *ICML 2016 - Proceedings of the 33th International Conference on Machine Learning*, pages 593–602.

Acknowledgements

We thank Dominik Csiba and Peter Richtárik for helpful discussions. VC was supported in part by the European Commission under Grant ERC Future Proof, SNF 200021-146750, and SNF CRSII2-147633.