Faster optimization through adaptive importance sampling [Perekrestenko et al., 2017]

Student: Dmytro Perekrestenko

Supervisors:

Prof. Martin Jaggi Prof. Volkan Cevher

May 25, 2018

Outline

1 Introduction and motivation

- 2 Core lemma
- 3 Convergence theorem for arbitrary sampling distributions
- Gap-wise sampling
- 5 Numerical experiments lasso
- 6 Conclusion and references

Primal-dual setting [Dünner et al., 2016]

The following pair of dual to each other optimization problems is considered:

$$\begin{split} \min_{\boldsymbol{\alpha} \in \mathbb{R}^n} \Big[D(\boldsymbol{\alpha}) &:= f(A\boldsymbol{\alpha}) + \sum_i g_i(\alpha_i) \Big],\\ \min_{\boldsymbol{w} \in \mathbb{R}^d} \Big[P(\boldsymbol{w}) &:= f^*(\boldsymbol{w}) + \sum_i g_i^*(-\boldsymbol{a}_i^\top \boldsymbol{w}) \Big]. \end{split}$$

Examples:

• SVM:
$$f^*(\boldsymbol{w}) = \frac{\lambda}{2} \|\boldsymbol{w}\|_2^2$$
, $g_i^*(-\boldsymbol{a}_i^\top \boldsymbol{w}) = \max(0, 1 - y_i \boldsymbol{a}_i^\top \boldsymbol{w})$.
• LASSO: $f(\boldsymbol{A}\boldsymbol{\alpha}) = \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{\alpha}\|_2^2$, $g_i(\alpha_i) = \lambda |\alpha_i|$.

Goal

Find a ϵ_P -suboptimal \boldsymbol{w} or ϵ_D -suboptimal α , i.e., $P(\boldsymbol{w}) - P(\boldsymbol{w}^*) \leq \epsilon_P$ or $D(\alpha) - D(\alpha^*) \leq \epsilon_D$.

Optimality conditions and classic algorithm

Optimality Conditions

Under assumption of strong duality: $P(oldsymbol{w}^{\star}) = -D(lpha^{\star})$ and

$$\boldsymbol{w}^{\star} \in \partial f(\boldsymbol{A} \boldsymbol{\alpha}^{\star}) \qquad \alpha_{i}^{\star} \in \partial g_{i}^{\star}(-\boldsymbol{a}_{i}^{\top} \boldsymbol{w}^{\star}) \text{ for all } i \in [n]$$

Duality gap is defined as: $G(\alpha, \boldsymbol{w}) := P(\boldsymbol{w}(\alpha)) - (-D(\alpha))$. It can act as a stopping criterion $G(\alpha) \ge P(\boldsymbol{w}(\alpha)) - P(\boldsymbol{w}^*)$.

Algorithm 1 Stochastic Coordinate Descent

1: let
$$\boldsymbol{\alpha}^{(0)} = \mathbf{0} \in \mathbb{R}^n$$
, $\boldsymbol{w}^{(0)} = \boldsymbol{w}(\boldsymbol{\alpha}^{(0)})$

2: for
$$t = 0, 1, ... T$$
 do

3: Sample $i \in [n]$ randomly according to distribution \boldsymbol{p}

4: Find
$$\Delta \alpha_i$$
 minimizing $D(\alpha^{(t)} + \boldsymbol{e}_i \Delta \alpha_i)$

5: $\boldsymbol{\alpha}^{(t+1)} = \boldsymbol{\alpha}^{(t)} + \boldsymbol{e}_i \Delta \alpha_i, \ \boldsymbol{w}^{(t+1)} = \boldsymbol{w}(\boldsymbol{\alpha}^{(t+1)})$

6: end for

Motivation

The Limitations of Existing Algorithms

- sampling of the active datapoint uniformly at random in each iteration
- the convergence rate is negatively affected

Stochastic Optimization with Adaptive Importance Sampling

- adaptively changes the sampling probability distribution according to the data and values of the dual variables
- improved practical and theoretical convergence rates

Definition (Dual Residual, ([Csiba et al., 2015], Def. 1))

Consider our primal-dual setting. Given dual variable α , the *i*-th dual residue on iteration *t* is given by:

$$\kappa_i^{(t)} = u_i^{(t)} - \alpha_i^{(t)},$$

where $u_i^{(t)} = \nabla g_i^* (-\boldsymbol{a}_i^\top \boldsymbol{w}^{(t)})$.

Definition (Coherent probability vector, ([Csiba et al., 2015], Def. 2))

We say that probability vector $\boldsymbol{p}^{(t)} \in \mathbb{R}^n$ is coherent with the dual residue vector $\boldsymbol{\kappa}^{(t)}$ if for all $i \in [n]$, we have $\kappa_i^{(t)} \neq 0 \quad \rightarrow \quad p_i^{(t)} > 0$.

Definition (*t*-support set)

We call set
$$I_t$$
: $I_t = \{i \in [n] : \kappa_i^{(t)} \neq 0\} \subseteq [n]$ a *t*-support set.

Dmytro Perekrestenko

Outline

Introduction and motivation

2 Core lemma

- 3) Convergence theorem for arbitrary sampling distributions
- 4 Gap-wise sampling
- 5 Numerical experiments lasso
- 6 Conclusion and references

Lemma (A generalization of ([Csiba et al., 2015], Lemma 3))

Consider Stochastic Coordinate Descent. Let f be $1/\beta$ -smooth and g_i be μ_i -strongly convex with convexity parameter $\mu_i \ge 0 \ \forall i \in [n]$. For the case $\mu_i = 0$ we require g_i to have a bounded support. Then for any iteration t, any sampling distribution $\mathbf{p}^{(t)}$ coherent with $\kappa^{(t)}$ and any $\theta \in [0, \min_{i \in I_t} p_i^{(t)}]$ it holds that:

$$\mathbb{E}[D(\boldsymbol{\alpha}^{(t+1)})|\boldsymbol{\alpha}^{(t)}] \leq D(\boldsymbol{\alpha}^{(t)}) - \theta G(\boldsymbol{\alpha}^{(t)}) + \frac{\theta^2 n^2}{2} F^{(t)}$$

where

$$F^{(t)} = \frac{1}{n^2 \beta \theta} \sum_{i \in I_t} \left(\frac{\theta(\mu_i \beta + \|\boldsymbol{a}_i\|^2)}{p_i^{(t)}} - \mu_i \beta \right) |\kappa_i^{(t)}|^2.$$

Outline

Introduction and motivation

2 Core lemma

3 Convergence theorem for arbitrary sampling distributions

- 4 Gap-wise sampling
- 5 Numerical experiments lasso
- 6 Conclusion and references

Theorem

Consider Stochastic Coordinate Descent. Assume f is $\frac{1}{\beta}$ -smooth function. Then, if g_i^* is L_i -Lipschitz for each i and $\mathbf{p}^{(t)}$ is coherent with $\kappa^{(t)}$, it suffices to have a total number of iterations of

$$T \geq \max\left\{0, \frac{1}{p_{\min}}\log\left(\frac{2\epsilon_D^0}{n^2 p_{\min} F^\circ}\right)\right\} + \frac{5F^\circ n^2}{\epsilon} - \frac{1}{p_{\min}}$$

or alternatively

$$T \geq rac{5F^{\circ}n^2}{\epsilon} + rac{5\epsilon_D^{(0)}}{\epsilon
ho_{\mathsf{min}}} - rac{1}{
ho_{\mathsf{min}}}$$

to obtain a duality gap $G(\bar{\alpha}) \leq \epsilon$, where ϵ_D^0 is the initial dual suboptimality and F° is an upper bound on $\mathbb{E}[F^{(t)}]$ taken over all coordinates at $1, \ldots, T$ algorithm iterations.

From this theorem we can recover as a special case:

1 ([Dünner et al., 2016], Theorem 9) by setting $p_i = \frac{1}{n}$ and using R instead of $||\mathbf{a}_i||$ ($R \ge ||\mathbf{a}_i|| \forall i \in [n]$).

2 ([Shalev-Shwartz and Zhang, 2013], Theorem 2) is a special case of ([Dünner et al., 2016], Theorem 9) for quadratic regularizer.

3 ([Zhao and Zhang, 2014], Theorem 5) by setting $p_i = \frac{L_i}{\sum_i L_j}$.

Maximizing rate in case of infinitesimal ϵ

Recall that the number of iterations sufficient to achieve ϵ -accuracy is:

$$T \geq \max\left\{0, \frac{1}{p_{\min}}\log\left(\frac{2\epsilon_D^0}{n^2 p_{\min}F^\circ}\right)\right\} + \frac{5F^\circ n^2}{\epsilon} - \frac{1}{p_{\min}}$$

In limit $\epsilon \to 0$, the only significant term is $\frac{5F^{\circ}n^2}{\epsilon}$, therefore we want to minimize $F^{(t)}$, which consequently minimizes F° . We solve the following optimization problem:

$$\boldsymbol{p}^{(t)} = \arg\min_{\boldsymbol{p}^{(t)}} \boldsymbol{F}^{(t)} := \arg\min_{\boldsymbol{p}^{(t)}} \frac{1}{n^2\beta} \sum_i \Big(\frac{|\kappa_i^{(t)}|^2 \|\boldsymbol{a}_i\|^2}{p_i^t} \Big).$$

Optimal sampling distributions

1) For uniform sampling $\forall t$:

$$p_i := \frac{1}{n}; \quad F_{\text{unif}}^\circ = \mathbb{E}\left[\frac{1}{\beta}\sum_i \left(\frac{|\kappa_i^{(t)}|^2 \|\boldsymbol{a}_i\|^2}{n}\right)\right];$$

2) For importance sampling $\forall t$:

$$p_i := \frac{L_i \|\boldsymbol{a}_i\|}{\sum_j L_j \|\boldsymbol{a}_j\|}; \quad F_{imp}^{\circ} = \mathbb{E}\left[\frac{1}{\beta} \sum_i \left(\frac{|\kappa_i^{(t)}|^2 \|\boldsymbol{a}_i\|}{L_i n}\right) \sum_j \left(\frac{L_j \|\boldsymbol{a}_j\|}{n}\right)\right];$$

3) For adaptive sampling $\forall t$:

$$p_i^{(t)} := \frac{|\kappa_i^{(t)}| \|\boldsymbol{a}_i\|}{\sum_j |\kappa_j^{(t)}| \|\boldsymbol{a}_j\|}; \quad F_{\mathsf{ada}}^\circ = \mathbb{E}\left[\frac{1}{\beta}\left(\sum_i \frac{|\kappa_i^{(t)}| \|\boldsymbol{a}_i\|}{n}\right)^2\right];$$

Support set uniform sampling

Let's assume that the size of the *t*-support set never exceeds some $m \in [1, n]$ and compare two sampling methods:

• Uniform sampling: $p_i^{(t)} = \frac{1}{n}$; $p_{\min} = \frac{1}{n}$, $F^{(t)} = \frac{1}{n^{2\beta}} \sum_{i \in I_t} \left(\frac{|\kappa_i^{(t)}|^2 \|\boldsymbol{a}_i\|^2}{n^{(t)}} \right) = \frac{1}{n\beta} \sum_{i \in I_t} |\kappa_i^{(t)}|^2 \|\boldsymbol{a}_i\|^2 \le F_{\text{unif}}$ Number of iterations: $T \ge \max\left\{0, n \log\left(\frac{2\epsilon_D^0}{nF_{\text{unif}}}\right)\right\} + \frac{5n^2F_{\text{unif}}}{\epsilon}$ • Support set uniform: $\begin{cases} p_i^{(t)} = \frac{1}{m}, & \text{if } \kappa_i^{(t)} \neq 0\\ p_i^{(t)} = 0, & \text{otherwise} \end{cases}; \ p_{\min} = \frac{1}{m}, \end{cases}$ $F^{(t)} = rac{1}{n^{2eta}} \sum_{i \in I_t} \left(rac{|\kappa_i^{(t)}|^2 \|m{a}_i\|^2}{n^{(t)}}
ight) = rac{m}{n^{2eta}} \sum_{i \in I_t} |\kappa_i^{(t)}|^2 \|m{a}_i\|^2 \le rac{m}{n} F_{ ext{unif}}$ Number of iterations: $T \ge \max \left\{ 0, m \log \left(\frac{2\epsilon_D^{(0)}}{nF_{\text{unif}}} \right) \right\} + \frac{5nmF_{\text{unif}}}{\epsilon}$

Maximizing rate in case of constant ϵ

The number of iterations T is directly proportional to F° and $1/p_{\min}$, the optimal distribution p should minimize F° and $1/p_{\min}$ at the same time. We define mixed distribution as:

$$\begin{cases} p_i^{(t)} = \frac{\sigma}{m} + (1 - \sigma) \frac{|\kappa_i^{(t)}| \|\boldsymbol{a}_i\|}{\sum_j |\kappa_j^{(t)}| \|\boldsymbol{a}_j\|}, & \text{if } \kappa_i^{(t)} \neq 0\\ p_i^{(t)} = 0, & \text{otherwise} \end{cases}$$

where $\sigma \in [0, 1]$. This distribution gives us the following bounds on F° and $1/p_{\min}$:

$$F_{\min}^{\circ} \leq rac{F_{ada}^{\circ}}{1-\sigma} \qquad rac{1}{p_{\min}} \leq rac{m}{\sigma},$$

and bound on the number of iterations:

$$T \geq rac{5F_{\mathsf{ada}}^\circ n^2}{\epsilon(1-\sigma)} + rac{\epsilon_D^{(0)}m}{\epsilon\sigma}.$$

Outline

Introduction and motivation

- 2 Core lemma
- 3 Convergence theorem for arbitrary sampling distributions
- Gap-wise sampling
 - 5 Numerical experiments lasso
 - 6 Conclusion and references

Remark

The duality gap can be decomposed as a sum of coordinate-wise duality gaps:

$$G(\boldsymbol{\alpha}) = \sum_{i} G_{i}(\alpha_{i}, \boldsymbol{w}) = \sum_{i} \left(g_{i}^{*}(-\boldsymbol{a}_{i}^{\top}\boldsymbol{w}) + g_{i}(\alpha_{i}) + \alpha_{i}\boldsymbol{a}_{i}^{\top}\boldsymbol{w} \right)$$

Definition (Nonuniformity measure, [Osokin et al., 2016])

The nonuniformity measure $\chi(\mathbf{x})$ of a vector $\mathbf{x} \in \mathbb{R}^n$, is defined as:

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

where $\boldsymbol{p} := \frac{\boldsymbol{x}}{\|\boldsymbol{x}\|_1}$ is the probability vector obtained by normalizing \boldsymbol{x} .

Theorem

Consider Stochastic Coordinate Descent. Assume f is $\frac{1}{\beta}$ -smooth function. Then, if g_i^* is L_i -Lipschitz for each i and $p_i^{(t)} := \frac{G_i(\alpha^{(t)})}{G(\alpha^{(t)})}$ then on each iteration it holds that

$$\mathbb{E}[\epsilon_D^{(t)}] \leq \frac{C + 2n\epsilon_D^{(0)}}{t + 2n},$$

where *C* is an upper bound on $\mathbb{E}\left[\frac{2n\chi(\vec{F})\sum_{i} ||\boldsymbol{a}_{i}||^{2}|\kappa_{i}^{(t)}|^{2}}{(\chi(\vec{G}))^{3}\beta}\right]$, where the expectation is taken over the random choice of the sampled coordinate at iterations $1, \ldots, t$ of the algorithm. Here \vec{G} and \vec{F} are defined as:

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

Gap-wise vs uniform

From intermediate result in the proof of Theorem 1, uniform distribution $(p_i = 1/n)$ has the following convergence rate:

$$\mathbb{E}[\epsilon_D^{(t)}] \leq \frac{\frac{2n}{\beta} \mathbb{E}\left[\sum_i |\kappa_i^{(t)}|^2 \|\boldsymbol{a}_i\|^2\right] + 2n\epsilon_D^{(0)}}{2n+t}.$$

The rate of gap-wise sampling depends on non-uniformity measures $\chi(\vec{G})$ and $\chi(\vec{F})$:

$$\mathbb{E}[\epsilon_D^t] \le \frac{\frac{2n}{\beta} \mathbb{E}\left[\frac{\chi(\vec{F})}{(\chi(\vec{G}))^3} \sum_i |\kappa_i^{(t)}|^2 \|\boldsymbol{a}_i\|^2\right] + 2n\epsilon_D^{(0)}}{2n+t}$$

In the worst case scenario when variance is maximal in $(|\kappa_i^{(t)}|^2 \|\mathbf{a}_i\|^2)_{i=1}^n$, $\chi(\vec{F}) \approx \sqrt{n}$, the rate of gap-wise sampling is better than of uniform only when the gaps are non-uniform enough i.e., $\chi(\vec{G}) \ge n^{\frac{1}{6}}$.

Outline

Introduction and motivation

- 2 Core lemma
- 3 Convergence theorem for arbitrary sampling distributions
- 4 Gap-wise sampling
- 5 Numerical experiments lasso
 - 6 Conclusion and references

Lasso problem

$$egin{aligned} \min_{oldsymbol{lpha} \in \mathbb{R}^n} \Big[D(oldsymbol{lpha}) &:= f(Aoldsymbol{lpha}) + \sum_i g_i(lpha_i) \Big], \ \min_{oldsymbol{w} \in \mathbb{R}^d} \Big[P(oldsymbol{w}) &:= f^*(oldsymbol{w}) + \sum_i g_i^*(-oldsymbol{a}_i^\top oldsymbol{w}) \Big]. \end{aligned}$$

Here $f(A\alpha) = \frac{1}{2n} ||A\alpha - \mathbf{y}||_2^2$ and $g_i(\alpha_i) = \lambda |\alpha_i|$. To satisfy requirements of the theorem we use "Lipschitzing trick" [Dünner et al., 2016] on $g_i(\alpha_i)$ to make g_i^* a Lipschitz function:

$$ar{g}_i(lpha_i) = egin{cases} \lambda |lpha_i|, & ext{if } |lpha_i| \leq B \ +\infty, & ext{otherwise} \end{cases}$$

The \bar{g}_i -conjugate will be:

$$\bar{g}_i^*(u_i) = \max_{\alpha_i:|\alpha_i| \le B} u_i \alpha_i - \lambda |\alpha_i| = B \big[|u_i| - \lambda \big]_+$$

Algorithm 2 Stochastic Coordinate Descent

1: let
$$\alpha^{(0)} = 0$$
, $\boldsymbol{w}^{(0)} = \nabla f(A \alpha^{(0)})$

- $2: \ \text{for} \ t=0,1,... \ \text{do}$
- 3: sample j from [d] according to distribution p

4: let
$$z_j = (\nabla f(\boldsymbol{\alpha}^{(t)}))_j$$

5:
$$\alpha_j^{(t+1)} = s_\lambda (\alpha_j^{(t)} - z_j)$$

6:
$$\boldsymbol{w}^{(t+1)} =
abla f(A \boldsymbol{\alpha}^{(t+1)})$$

7: end for

• uniform
$$p_i = \frac{1}{d}$$

• importance
$$p_i = \frac{L_i \|\boldsymbol{a}_i\|}{\sum_i L_j \|\boldsymbol{a}_j\|}$$

• (heuristic) gap-init
$$p_i = rac{G_i^{(0)}}{\sum_j G_j^{(0)}}$$

(0)

Algorithm 3 Stochastic Coordinate Descent (adaptive)

1: let
$$\boldsymbol{lpha}^{(0)}=0,~ \boldsymbol{w}^{(0)}=
abla f(A \boldsymbol{lpha}^{(0)})$$

- 2: for t = 0, 1, ... do
- calculate absolute values of dual residuals $|\kappa_i^{(t)}|$ for all $j \in [d]$ 3.
- generate adapted probabilities distribution $\boldsymbol{p}^{(t)}$: 4:

$$p_i^{(t)} = \frac{|\kappa_i^{(t)}| \|\boldsymbol{a}_i\|}{\sum_j |\kappa_j^{(t)}| \|\boldsymbol{a}_j\|}$$

sample j from [d] according to $p^{(t)}$ 5:

6: let
$$z_j = (\nabla f(\boldsymbol{\alpha}^{(t)}))_j$$

7:
$$\alpha_j^{(t+1)} = s_\lambda(\alpha_j^{(t)} - z_j)$$

8: $\mathbf{w}^{(t+1)} = \nabla f(A \alpha^{(t+1)})$

8:

9: end for

Algorithm 4 Stochastic Coordinate Descent (supportSet-uniform)

1: let
$$\boldsymbol{lpha}^{(0)}=0,~ \boldsymbol{w}^{(0)}=
abla f(A \boldsymbol{lpha}^{(0)})$$

- $2: \ \text{for} \ t=0,1,... \ \text{do}$
- 3: calculate absolute values of dual residuals $|\kappa_j^{(t)}|$ for all $j \in [d]$

4: find *t*-support set
$$I_t = \{i \in [d] : \kappa_i^{(t)} \neq 0\} \subseteq [d]$$

5: generate adapted probabilities distribution $\boldsymbol{p}^{(t)}$:

$$\begin{cases} p_i^{(t)} = \frac{1}{|l_t|}, & \text{if } \kappa_i^{(t)} \neq 0\\ p_i^{(t)} = 0, & \text{otherwise} \end{cases}$$

6: sample *j* from [*d*] according to $\boldsymbol{p}^{(t)}$

7: let
$$z_j = (\nabla f(\alpha^{(t)}))_j$$

8: $\alpha_j^{(t+1)} = s_\lambda(\alpha_j^{(t)} - z_j), \ \boldsymbol{w}^{(t+1)} = \nabla f(A\alpha^{(t+1)})$
9: end for

Algorithm 5 Stochastic Coordinate Descent (ada-uniform)

1: let
$$\alpha^{(0)} = 0$$
, $\boldsymbol{w}^{(0)} =
abla f(A \alpha^{(0)})$

- $2: \ \text{for} \ t=0,1,... \ \text{do}$
- 3: calculate absolute values of dual residuals $|\kappa_j^{(t)}|$ for all $j \in [d]$
- 4: find *t*-support set $I_t = \{i \in [d] : \kappa_i^{(t)} \neq 0\} \subseteq [d]$
- 5: generate adapted probabilities distribution $\boldsymbol{p}^{(t)}$:

$$\begin{cases} p_i^{(t)} = \frac{1}{2|I_t|} + \frac{|\kappa_i^{(t)}| \| \mathbf{a}_i \|}{2\sum_j |\kappa_j^{(t)}| \| \mathbf{a}_j \|}, & \text{if } \kappa_i^{(t)} \neq 0\\ p_i^{(t)} = 0, & \text{otherwise} \end{cases}$$

6: sample *j* from [*d*] according to $\boldsymbol{p}^{(t)}$

7: let
$$z_j = (\nabla f(\alpha^{(t)}))_j$$

8:
$$\alpha_{j}^{(t+1)} = s_{\lambda}(\alpha_{j}^{(t)} - z_{j}), \ \boldsymbol{w}^{(t+1)} = \nabla f(A\alpha^{(t+1)})$$

9: end for

Algorithm 6 Stochastic Coordinate Descent (ada-gap)

1: let
$$\boldsymbol{lpha}^{(0)}=0,~ \boldsymbol{w}^{(0)}=
abla f(A \boldsymbol{lpha}^{(0)})$$

- 2: for $t = 0,1,\dots$ do
- 3: calculate feature-wise duality gaps $G_i^{(t)}$ for all $j \in [d]$
- 4: generate adapted probabilities distribution $\boldsymbol{p}^{(t)}$: $p_i^{(t)} = \frac{G_i^{(t)}}{\sum G_i^{(t)}}$
- 5: sample *j* from [*d*] according to $\boldsymbol{p}^{(t)}$

6: let
$$z_j = (\nabla f(\boldsymbol{\alpha}^{(t)}))_j$$

7: $\alpha_i^{(t+1)} = s_\lambda(\alpha_i^{(t)} - z_i)$

7: $\alpha_j^{(t+1)} = s_\lambda(\alpha_j^{(t)} - z_j)$ 8: $\mathbf{w}^{(t+1)} = \nabla f(A\alpha^{(t+1)})$

9: end for

Dataset	Features	Points	nnz/(nd)	mean of $\ \boldsymbol{a}_i\ $	$Var(\ \boldsymbol{a}_i\)$
mushrooms	112	8124	18.8%	31.35	545
rcv1*	809	7438	0.3%	2.58	17.3

Algorithm	Cost of an Epoch	Mode
Lasso uniform	O(nnz)	uniform
Lasso importance	$O(nnz + n \log(n))$	fixed non-uniform
Lasso gap-init	$O(nnz + n \log(n))$	fixed non-uniform
Lasso supportSet-uniform	$O(n \cdot nnz)$	adaptive
Lasso adaptive	$O(n \cdot nnz)$	adaptive
Lasso ada-uniform	$O(n \cdot nnz)$	adaptive
Lasso ada-division	$O(nnz + n \log(n))$	adaptive
Lasso ada-gap	$O(n \cdot nnz)$	adaptive

Numerical experiment - fixed



Figure: Lasso. Comparison of different fixed distribution versions of Stochastic Coordinate Descent Algorithm based on duality gap(left) and suboptimality(right) measure - rcv1 dataset

Numerical experiment - fixed



Figure: Lasso. Comparison of different fixed distribution versions of Stochastic Coordinate Descent Algorithm based on duality gap(left) and suboptimality(right) measure - mushrooms dataset

Numerical experiment - adaptive



Figure: Lasso. Comparison of different adaptive versions of Stochastic Coordinate Descent Algorithm based on duality gap(left) and suboptimality(right) measure - rcv1 dataset

Numerical experiment - adaptive



Figure: Lasso. Comparison of different adaptive versions of Stochastic Coordinate Descent Algorithm based on duality gap(left) and suboptimality(right) measure - mushrooms dataset

Outline

Introduction and motivation

- 2 Core lemma
- 3 Convergence theorem for arbitrary sampling distributions
- Gap-wise sampling
- 5 Numerical experiments lasso
- 6 Conclusion and references

Summary and future work

Summary

- studied SCD with adaptive sampling
- proposed new adaptive and fixed nonuniform sampling schemes
- analyzed their theoretical convergence rates
- showed in practice that they outperform the conventional sampling schemes

Future work

- to find a better solution for the constant ϵ optimization problem
- to find a relation between dual residuals sampling and gap-wise sampling
- to apply the theory to hinge loss SVM

References I

- 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, 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, A., Alayrac, J.-B., Lukasewitz, 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.

Image: Image:

Perekrestenko, D., Cevher, V., and Jaggi, M. (2017). Faster coordinate descent via adaptive importance sampling. *Proc. of the 20th International Conference on Artificial Intelligence and Statistics*, 54.

Shalev-Shwartz, S. and Zhang, T. (2013). Stochastic Dual Coordinate Ascent Methods for Regularized Loss Minimization.

JMLR, 14:567-599.

Zhao, P. and Zhang, T. (2014).

Stochastic Optimization with Importance Sampling. *arXiv*.