

Modern Optimization Methods for Big Data Problems

MATH11146
The University of Edinburgh

Peter Richtárik

Week 3
Randomized Coordinate Descent With Arbitrary Sampling
January 27, 2016



1 / 30

The Problem

We first consider the following problem:

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{subject to} & x = (x_1, \dots, x_n) \in \mathbb{R}^n \end{array} \quad (1)$$

We will assume that f is:

- ▶ **“smooth”** (will be made precise later)
- ▶ **strongly convex** (will be made precise later)

So, this is unconstrained minimization of a smooth convex function.



2 / 30

Randomized Coordinate Descent with Arbitrary Sampling

NSync Algorithm (R. and Takáč 2014, [4])

Input: initial point $x_0 \in \mathbb{R}^n$

subset probabilities $\{p_S\}$ for each $S \subseteq [n] \stackrel{\text{def}}{=} \{1, 2, \dots, n\}$

stepsize parameters $v_1, \dots, v_n > 0$

for $k = 0, 1, 2, \dots$ **do**

a) **Select a random set of coordinates** $S_k \subseteq [n]$ following the law

$$\mathbf{P}(S_k = S) = p_S, \quad S \subseteq [n]$$

b) **Update (possibly in parallel) selected coordinates:**

$$x_{k+1} = x_k - \sum_{i \in S_k} \frac{1}{v_i} (e_i^T \nabla f(x_k)) e_i$$

(e_i is the i th unit coordinate vector)

end for

Remark: The **NSync algorithm** was introduced in 2013. The first coordinate descent algorithm using **arbitrary sampling**.



3 / 30

Two More Ways of Writing the Update Step

1. **Coordinate-by-coordinate:**

$$x_i^{k+1} = \begin{cases} x_i^k, & i \notin S_k, \\ x_i^k - \frac{1}{v_i} (\nabla f(x_k))_i, & i \in S_k. \end{cases}$$

2. **Via projection to a subset of blocks:** If for $h \in \mathbb{R}^n$ and $S \subseteq [n]$ we write

$$h_S \stackrel{\text{def}}{=} \sum_{i \in S} h_i e_i, \quad (2)$$

then

$$x^{k+1} = x^k + h_{S_k} \quad \text{for} \quad h = -(\text{Diag}(v))^{-1} \nabla f(x_k). \quad (3)$$

Depending on context, we shall interchangeably denote the i th partial derivative of f at x by

$$\nabla_i f(x) = e_i^T \nabla f(x) = (\nabla f(x))_i$$



4 / 30

Samplings

Definition 1 (Sampling)

By the name **sampling** we refer to a set valued random mapping with values being subsets of $[n] = \{1, 2, \dots, n\}$. For sampling \hat{S} we define the **probability vector** $p = (p_1, \dots, p_n)^T$ by

$$p_i = \mathbf{P}(i \in \hat{S}) \quad (4)$$

We say that \hat{S} is **proper**, if $p_i > 0$ for all i .

- ▶ A sampling \hat{S} is uniquely characterized by the **probability mass function**

$$p_S \stackrel{\text{def}}{=} \mathbf{P}(\hat{S} = S), \quad S \subseteq [n]; \quad (5)$$

that is, by assigning probabilities to all subsets of $[n]$.

- ▶ Later on it will be useful to also consider the **probability matrix** $P = P(\hat{S}) = (p_{ij})$ given by

$$p_{ij} \stackrel{\text{def}}{=} \mathbf{P}(i \in \hat{S}, j \in \hat{S}) = \sum_{S: \{i,j\} \subseteq S} p_S. \quad (6)$$



5 / 30

Samplings: A Basic Identity

Lemma 2 ([3])

For any sampling \hat{S} we have

$$\sum_{i=1}^n p_i = \mathbf{E}[|\hat{S}|]. \quad (7)$$

Proof.

$$\sum_{i=1}^n p_i \stackrel{(4)+(5)}{=} \sum_{i=1}^n \sum_{S \subseteq [n]: i \in S} p_S = \sum_{S \subseteq [n]} \sum_{i: i \in S} p_S = \sum_{S \subseteq [n]} p_S |S| = \mathbf{E}[|\hat{S}|].$$

□



6 / 30

Sampling Zoo - Part I

Why consider different samplings?

1. **Basic Considerations.** It is important that each block i has a positive probability of being chosen, otherwise NSync will not be able to update some blocks and hence will not converge to optimum. For technical/sanity reasons, we define:
 - ▶ **Proper sampling.** $p_i = \mathbf{P}(i \in \hat{S}) > 0$ for all $i \in [n]$
 - ▶ **Nil sampling:** $\mathbf{P}(\hat{S} = \emptyset) = 1$
 - ▶ **Vacuous sampling:** $\mathbf{P}(\hat{S} = \emptyset) > 0$
2. **Parallelism.** Choice of sampling affects the level of parallelism:
 - ▶ $\mathbf{E}[|\hat{S}|]$ is the average number of updates performed in parallel in one iteration; and is hence closely related to the number of iterations.
 - ▶ **serial sampling:** picks one block:

$$\mathbf{P}(|\hat{S}| = 1) = 1$$

We call this sampling serial although nothing prevents us from computing the actual update to the block, and/or to apply the update in parallel.



7 / 30

Sampling Zoo - Part II

- ▶ **fully parallel sampling:** always picks all blocks:

$$\mathbf{P}(\hat{S} = \{1, 2, \dots, n\}) = 1$$

3. **Processor reliability.** Sampling may be induced/informed by the computing environment:
 - ▶ **Reliable/dedicated processors.** If one has reliable processors, it is sensible to choose sampling \hat{S} such that $\mathbf{P}(|\hat{S}| = \tau) = 1$ for some τ related to the number of processors.
 - ▶ **Unreliable processors.** If processors given a computing task are busy or unreliable, they return answer later or not at all - it is then sensible to ignore such updates and move on. This then means that $|\hat{S}|$ varies from iteration to iteration.
4. **Distributed computing.** In a distributed computing environment it is sensible:
 - ▶ to allow each compute node as much autonomy as possible so as to **minimize communication cost**,
 - ▶ to make sure **all nodes are busy** at all times



8 / 30

Sampling Zoo - Part III

This suggests a strategy where the set of blocks is partitioned, with each node owning a partition, and independently picking a “chunky” subset of blocks at each iteration it will update, ideally from local information.

5. **Uniformity.** It may or may not make sense to update some blocks more often than others:

- ▶ **uniform samplings:**

$$\mathbf{P}(i \in \hat{S}) = \mathbf{P}(j \in \hat{S}) \quad \text{for all } i, j \in [n]$$

- ▶ **doubly uniform (DU):** These are samplings characterized by:

$$|S'| = |S''| \Rightarrow \mathbf{P}(\hat{S} = S') = \mathbf{P}(\hat{S} = S'') \quad \text{for all } S', S'' \subseteq [n]$$

- ▶ **τ -nice:** DU sampling with the additional property that

$$\mathbf{P}(|\hat{S}| = \tau) = 1$$

- ▶ **distributed τ -nice:** will define later

- ▶ **independent sampling:** union of independent uniform serial samplings

- ▶ **nonuniform samplings**



9 / 30

Sampling Zoo - Part IV

6. **Complexity of generating a sampling.** Some samplings are computationally more efficient to generate than others: the potential benefits of a sampling may be completely ruined by the difficulty to generate sets according to the sampling's distribution.

- ▶ a τ -nice sampling can be well approximated by an independent sampling, which is easy to generate. . .
- ▶ in general, many samplings will be hard to generate



10 / 30

Assumption: Strong convexity

Assumption 1 (Strong convexity)

Function f is differentiable and λ -strongly convex (with $\lambda > 0$) with respect to the standard Euclidean norm

$$\|h\| \stackrel{\text{def}}{=} \left(\sum_{i=1}^n h_i^2 \right)^{1/2}.$$

That is, we assume that for all $x, h \in \mathbb{R}^n$,

$$f(x+h) \geq f(x) + \langle \nabla f(x), h \rangle + \frac{\lambda}{2} \|h\|^2. \quad (8)$$



11 / 30

Assumption: Expected Separable Overapproximation

Assumption 2 (ESO)

Assume \hat{S} is proper and that for some vector of positive weights $v = (v_1, \dots, v_n)$ and all $x, h \in \mathbb{R}^n$,

$$\mathbf{E}[f(x + h_{\hat{S}})] \leq f(x) + \langle \nabla f(x), h \rangle_p + \frac{1}{2} \|h\|_{p \bullet v}^2. \quad (9)$$

Note that **the ESO parameters v, p depend on both f and \hat{S}** . For simplicity, we will often instead of (9) use the compact notation

$$(f, \hat{S}) \sim \text{ESO}(v).$$

Notation used above:

$$h_S \stackrel{\text{def}}{=} \sum_{i \in S} h_i e_i \in \mathbb{R}^n \quad (\text{projection of } h \in \mathbb{R}^n \text{ onto coordinates } i \in S)$$

$$\langle g, h \rangle_p \stackrel{\text{def}}{=} \sum_{i=1}^n p_i g_i h_i \in \mathbb{R} \quad (\text{weighted inner product})$$

$$p \bullet v \stackrel{\text{def}}{=} (p_1 v_1, \dots, p_n v_n) \in \mathbb{R}^n \quad (\text{Hadamard product})$$



12 / 30

Assumption: Expected Separable Overapproximation

Here the ESO inequality again, now without the simplifying notation:

$$\underbrace{\mathbf{E} \left[f \left(x + \sum_{i \in \hat{S}} h_i e_i \right) \right]}_{\text{complicated}} \leq f(x) + \underbrace{\sum_{i=1}^n p_i \nabla_i f(x) h_i}_{\text{linear in } h} + \underbrace{\frac{1}{2} \sum_{i=1}^n p_i v_i h_i^2}_{\text{quadratic and separable in } h}$$



13 / 30

Complexity of NSync

Theorem 3 (R. and Takáč 2013, [4])

Let x^* be a minimizer of f . Let Assumptions 1 and 2 be satisfied for a proper sampling \hat{S} (that is, $(f, \hat{S}) \sim \text{ESO}(v)$). Choose

- ▶ starting point $x^0 \in \mathbb{R}^n$,
- ▶ error tolerance $0 < \epsilon < f(x^0) - f(x^*)$ and
- ▶ confidence level $0 < \rho < 1$.

If $\{x^k\}$ are the random iterates generated by **NSync**, where the random sets S_k are iid following the distribution of \hat{S} , then

$$\mathbf{K} \geq \frac{\Omega}{\lambda} \log \left(\frac{f(x^0) - f(x^*)}{\epsilon \rho} \right) \Rightarrow \mathbf{P}(f(x^{\mathbf{K}}) - f(x^*) \leq \epsilon) \geq 1 - \rho, \quad (10)$$

where

$$\Omega \stackrel{\text{def}}{=} \max_{i=1, \dots, n} \frac{v_i}{p_i} \geq \frac{\sum_{i=1}^n v_i}{\mathbf{E}[|\hat{S}|]}. \quad (11)$$



14 / 30

What does this mean?

- ▶ **Linear convergence.** NSync converges linearly (i.e., logarithmic dependence on ϵ)
- ▶ **High confidence is not a problem.** ρ appears inside the logarithm, so it is easy to achieve high confidence (by running the method longer; there is no need to restart)
- ▶ **Focus on the leading term.** The leading term is Ω ; and we have a closed-form expression for it in terms of
 - ▶ parameters v_1, \dots, v_n (which depend on f and \hat{S})
 - ▶ parameters p_1, \dots, p_n (which depend on \hat{S})
- ▶ **Parallelization speedup.** The lower bound suggests that *if it was the case that* the parameters v_i did not grow with increasing $\tau \stackrel{\text{def}}{=} \mathbf{E}[|\hat{S}|]$, then we could potentially be getting linear speedup in τ (average number of updates per iteration).
 - ▶ So we shall **study the dependence of v_i on τ** (this will depend on f and \hat{S})
 - ▶ As we shall see, speedup is often guaranteed for **sparse or well-conditioned problems**.

Question: How to **design** sampling \hat{S} so that Ω is minimized?



15 / 30

Analysis of the Algorithm (Proof of Theorem 3)



16 / 30

Tool: Markov's Inequality

Theorem 4 (Markov's Inequality)

Let X be a nonnegative random variable. Then for any $\epsilon > 0$,

$$\mathbf{P}(X \geq \epsilon) \leq \frac{\mathbf{E}[X]}{\epsilon}.$$

Proof.

Let $1_{X \geq \epsilon}$ be the random variable which is equal to 1 if $X \geq \epsilon$ and 0 otherwise. Then

$$1_{X \geq \epsilon} \leq \frac{X}{\epsilon}.$$

By taking expectations of all terms, we obtain

$$\mathbf{P}(X \geq \epsilon) = \mathbf{E}[1_{X \geq \epsilon}] \leq \mathbf{E}\left[\frac{X}{\epsilon}\right] = \frac{\mathbf{E}[X]}{\epsilon}.$$



17 / 30

Tool: Tower Property of Expectations (Motivation)

Example 5

Consider discrete random variables \mathbf{X} and \mathbf{Y} :

- \mathbf{X} has 2 outcomes: \mathbf{x}_1 and \mathbf{x}_2
- \mathbf{Y} has 3 outcomes: \mathbf{y}_1 , \mathbf{y}_2 and \mathbf{y}_3

Their joint probability mass function is given in this table:

	\mathbf{y}_1	\mathbf{y}_2	\mathbf{y}_3	
\mathbf{x}_1	0.05	0.20	0.03	0.28
\mathbf{x}_2	0.25	0.30	0.17	0.72
	0.30	0.50	0.20	1

Obviously, $\mathbf{E}[\mathbf{X}] = 0.28\mathbf{x}_1 + 0.72\mathbf{x}_2$. But we can also write:

$$\begin{aligned}\mathbf{E}[\mathbf{X}] &= (0.05\mathbf{x}_1 + 0.25\mathbf{x}_2) + (0.20\mathbf{x}_1 + 0.30\mathbf{x}_2) + (0.03\mathbf{x}_1 + 0.17\mathbf{x}_2) \\ &= \underbrace{0.30}_{\mathbf{P}(\mathbf{Y}=\mathbf{y}_1)} \underbrace{\left(\frac{0.05}{0.30}\mathbf{x}_1 + \frac{0.25}{0.30}\mathbf{x}_2\right)}_{\mathbf{E}[\mathbf{X} | \mathbf{Y}=\mathbf{y}_1]} + 0.50 \left(\frac{0.20}{0.50}\mathbf{x}_1 + \frac{0.30}{0.50}\mathbf{x}_2\right) + 0.20 \left(\frac{0.03}{0.20}\mathbf{x}_1 + \frac{0.17}{0.20}\mathbf{x}_2\right) \\ &= \mathbf{E}[\mathbf{E}[\mathbf{X} | \mathbf{Y}]].\end{aligned}$$



18 / 30

Tower Property

Lemma 6 (Tower Property / Iterated Expectation)

For any random variables X and Y , we have $\mathbf{E}[X] = \mathbf{E}[\mathbf{E}[X | Y]]$.

Proof.

We shall only prove this for discrete random variables; the proof is more technical in the continuous case.

$$\begin{aligned}\mathbf{E}[X] &= \sum_x x \mathbf{P}(X = x) = \sum_x x \sum_y \mathbf{P}(X = x \& Y = y) \\ &= \sum_y \sum_x x \mathbf{P}(X = x \& Y = y) \\ &= \sum_y \sum_x \mathbf{P}(Y = y) x \frac{\mathbf{P}(X = x \& Y = y)}{\mathbf{P}(Y = y)} \\ &= \sum_y \mathbf{P}(Y = y) \underbrace{\sum_x x \mathbf{P}(X = x | Y = y)}_{\mathbf{E}[X | Y=y]} \\ &= \mathbf{E}[\mathbf{E}[X | Y]].\end{aligned}$$



19 / 30

Proof of Theorem 3 - Part I

- If we let $\mu \stackrel{\text{def}}{=} \lambda/\Omega$, then

$$\begin{aligned}f(x + h) &\stackrel{(8)}{\geq} f(x) + \langle \nabla f(x), h \rangle + \frac{\lambda}{2} \|h\|^2 \\ &\geq f(x) + \langle \nabla f(x), h \rangle + \frac{\mu}{2} \|h\|_{v \bullet p-1}^2.\end{aligned}\quad (12)$$

Indeed, one can easily verify that μ is defined to be the largest number for which

$$\lambda \|h\|^2 \geq \mu \|h\|_{v \bullet p-1}^2$$

holds for all h . Hence, f is μ -strongly convex with respect to the norm $\|\cdot\|_{v \bullet p-1}$.

- Let x^* be a minimizer of f , i.e., an optimal solution of (1). Minimizing both sides of (12) in h , we get

$$\begin{aligned}f(x^*) - f(x) &\stackrel{(12)}{\geq} \min_{h \in \mathbb{R}^n} \langle \nabla f(x), h \rangle + \frac{\mu}{2} \|h\|_{v \bullet p-1}^2 \\ &= -\frac{1}{2\mu} \|\nabla f(x)\|_{p \bullet v-1}^2.\end{aligned}\quad (13)$$



20 / 30

Proof of Theorem 3 - Part II

- Let $h^k \stackrel{\text{def}}{=} -v^{-1} \bullet \nabla f(x^k)$. Then in view of (3), we have $x^{k+1} = x^k + h_{S_k}^k$. Utilizing Assumption 2, we get

$$\begin{aligned}
 \mathbf{E}[f(x^{k+1}) \mid x^k] &= \mathbf{E}[f(x^k + h_{S_k}^k) \mid x^k] \\
 &\stackrel{(9)}{\leq} f(x^k) + \langle \nabla f(x^k), h^k \rangle_p + \frac{1}{2} \|h^k\|_{p \bullet v}^2 \\
 &= f(x^k) - \frac{1}{2} \|\nabla f(x^k)\|_{p \bullet v^{-1}}^2 \\
 &\stackrel{(13)}{\leq} f(x^k) - \mu(f(x^k) - f(x^*)).
 \end{aligned}$$

- Taking expectations in the last inequality, using the **tower property**, and subtracting $f(x^*)$ from both sides of the inequality, we get

$$\mathbf{E}[f(x^{k+1}) - f(x^*)] \leq (1 - \mu)\mathbf{E}[f(x^k) - f(x^*)].$$

Unrolling the recurrence, we get

$$\mathbf{E}[f(x^k) - f(x^*)] \leq (1 - \mu)^k (f(x^0) - f(x^*)). \quad (14)$$



21 / 30

Proof of Theorem 3 - Part III

- Using **Markov's inequality**, (14) and the definition of K , we get

$$\begin{aligned}
 \mathbf{P}(f(x^K) - f(x^*) \geq \epsilon) &\leq \mathbf{E}[f(x^K) - f(x^*)] / \epsilon \\
 &\stackrel{(14)}{\leq} (1 - \mu)^K (f(x^0) - f(x^*)) / \epsilon \stackrel{(10)}{\leq} \rho.
 \end{aligned}$$

- Finally, let us now establish the lower bound on Ω . Letting

$$\Delta \stackrel{\text{def}}{=} \left\{ p' \in \mathbb{R}^n : p' \geq 0, \sum_{i=1}^n p'_i = \mathbf{E}[|\hat{S}|] \right\},$$

we have

$$\Omega \stackrel{(11)}{=} \max_i \frac{v_i}{p_i} \stackrel{(7)}{\geq} \min_{p' \in \Delta} \max_i \frac{v_i}{p'_i} = \frac{1}{\mathbf{E}[|\hat{S}|]} \sum_{i=1}^n v_i,$$

where the last equality follows since optimal p'_i is proportional to v_i .



22 / 30

How to compute the ESO “stepsize” parameters
 v_1, \dots, v_n ?



23 / 30

$C^1(A)$ Functions

By definition, the ESO parameters v depend on both f and \hat{S} . This is also highlighted by the notation we use:

$$(f, \hat{S}) \sim ESO(v).$$

Hence, in order to compute these parameters, we need to specify f and \hat{S} . The following definition describes a wide class of functions, often appearing in computational practice, for which we will be able to do this.

Definition 7 ($C^1(A)$ functions)

Let $A \in \mathbb{R}^{m \times n}$. By $C^1(A)$ we denote the set of continuously differentiable functions $f : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying the following inequality for all $x, h \in \mathbb{R}^n$:

$$f(x + h) \leq f(x) + (\nabla f(x))^T h + \frac{1}{2} h^T A^T A h. \quad (15)$$

Example 8 (Least Squares)

The function $f(x) = \frac{1}{2} \|Ax - b\|^2$ satisfies (15) with an equality. Hence, $f \in C^1(A)$.



24 / 30

Are There More $C^1(A)$ Functions?

- Functions we wish to minimize in **machine learning** often have a “**finite sum**” structure:

$$f(x) \stackrel{\text{def}}{=} \sum_j \phi_j(M_j x), \quad (16)$$

where $M_j \in \mathbb{R}^{m \times n}$ are **data matrices** and $\phi_j : \mathbb{R}^m \rightarrow \mathbb{R}$ are **loss functions**.

- The next result says that under certain assumptions on the loss functions, $f \in C^1(A)$ for some A , and describes what this A is.

Theorem 9 ([5])

Assume that for each j , function ϕ_j is γ_j -smooth:

$$\|\nabla \phi_j(s) - \nabla \phi_j(s')\| \leq \gamma_j \|s - s'\|, \quad \text{for all } s, s' \in \mathbb{R}^m.$$

Then $f \in C^1(A)$, where A satisfies

$$A^T A = \sum_{j=1}^J \gamma_j M_j^T M_j.$$

Remark: The above theorem also says what $A^T A$ is. This is important, as will be clear from the next theorem.



25 / 30

ESO for $C_1(A)$ Functions and an Arbitrary Sampling

Theorem 10 ([5])

Assume $f \in C^1(A)$ for some real matrix $A \in \mathbb{R}^{m \times n}$, let \hat{S} be an arbitrary sampling and let $P = P(\hat{S})$ be its probability matrix. If $v = (v_1, \dots, v_n)$ satisfies

$$P \bullet A^T A \preceq \text{Diag}(p \bullet v),$$

then

$$(f, \hat{S}) \sim \text{ESO}(v).$$



26 / 30

Sampling Identity for a Quadratic

In order to prove Theorem 10, we will need the following lemma.

Lemma 11 ([5])

Let G be any real $n \times n$ matrix and \hat{S} an arbitrary sampling. Then for any $h \in \mathbb{R}^n$ we have

$$\mathbf{E} [h_{\hat{S}}^T G h_{\hat{S}}] = h^T \left(P(\hat{S}) \bullet G \right) h, \quad (17)$$

where \bullet denotes the Hadamard (elementwise) product of matrices, and $P(\hat{S})$ is the probability matrix of \hat{S} .



27 / 30

Proof of Lemma 11

Proof.

Let 1_{ij} be the indicator random variable of the event $i \in \hat{S} \ \& \ j \in \hat{S}$:

$$1_{ij} = \begin{cases} 1 & \text{if } i \in \hat{S} \ \& \ j \in \hat{S}, \\ 0 & \text{otherwise.} \end{cases}$$

Note that $\mathbf{E}[1_{ij}] = \mathbf{P}(i \in \hat{S} \ \& \ j \in \hat{S}) = p_{ij}$. We now have

$$\begin{aligned} \mathbf{E} [h_{\hat{S}}^T G h_{\hat{S}}] &\stackrel{(2)}{=} \mathbf{E} \left[\sum_{i \in \hat{S}} \sum_{j \in \hat{S}} G_{ij} h_i h_j \right] = \mathbf{E} \left[\sum_{i=1}^n \sum_{j=1}^n 1_{ij} G_{ij} h_i h_j \right] \\ &= \sum_{i=1}^n \sum_{j=1}^n \mathbf{E}[1_{ij}] G_{ij} h_i h_j = \sum_{i=1}^n \sum_{j=1}^n p_{ij} G_{ij} h_i h_j \\ &= h^T \left(P(\hat{S}) \bullet G \right) h. \end{aligned}$$



28 / 30

Proof of Theorem 10

Having established Lemma 11, we are now ready prove Theorem 10.

Proof.

Fixing any $x, h \in \mathbb{R}^n$, we have

$$\begin{aligned} \mathbf{E} [f(x + h_{\hat{S}})] &\stackrel{(15)}{\leq} \mathbf{E} [f(x) + \langle \nabla f(x), h_{\hat{S}} \rangle + \frac{1}{2} h_{\hat{S}}^T A^T A h_{\hat{S}}] \\ &= f(x) + \langle \nabla f(x), h \rangle_p + \frac{1}{2} \mathbf{E} [h_{\hat{S}}^T A^T A h_{\hat{S}}] \\ &\stackrel{(\text{Lemma 11})}{=} f(x) + \langle \nabla f(x), h \rangle_p + \frac{1}{2} h^T (P \bullet A^T A) h \\ &\leq f(x) + \langle \nabla f(x), h \rangle_p + \frac{1}{2} \underbrace{h^T \text{Diag}(p \bullet v) h}_{= \|h\|_{p \bullet v}^2}. \end{aligned}$$

□



29 / 30

References I

- [1] Yurii Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341-362, 2012
- [2] Peter Richtárik and Martin Takáč. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Mathematical Programming*, 144(2):1–38, 2014
- [3] Peter Richtárik and Martin Takáč. Parallel coordinate descent methods for big data optimization. *Mathematical Programming*, 1–52, 2015
- [4] Peter Richtárik and Martin Takáč. On optimal probabilities in stochastic coordinate descent methods. *Optimization Letters*, 2015
- [5] Zheng Qu and Peter Richtárik. Coordinate descent with arbitrary sampling II: expected separable overapproximation, *arXiv:1412.8063*, 2014



30 / 30