

Mathematics of Data: From Theory to Computation

Prof. Volkan Cevher
volkan.cevher@epfl.ch

Lecture 6: Time-data tradeoffs and variance reduction

Laboratory for Information and Inference Systems (LIONS)
École Polytechnique Fédérale de Lausanne (EPFL)

EE-556 (Fall 2021)



License Information for Mathematics of Data Slides

- ▶ This work is released under a [Creative Commons License](#) with the following terms:
- ▶ **Attribution**
 - ▶ The licensor permits others to copy, distribute, display, and perform the work. In return, licensees must give the original authors credit.
- ▶ **Non-Commercial**
 - ▶ The licensor permits others to copy, distribute, display, and perform the work. In return, licensees may not use the work for commercial purposes – unless they get the licensor's permission.
- ▶ **Share Alike**
 - ▶ The licensor permits others to distribute derivative works only under a license identical to the one that governs the licensor's work.
- ▶ [Full Text of the License](#)

Outline

- ▶ This class
 1. Time-data trade-offs
 2. Rate iteration-cost trade-offs
 3. Variance reduction
- ▶ Next class
 1. Deep learning introduction

A simple *regression* model

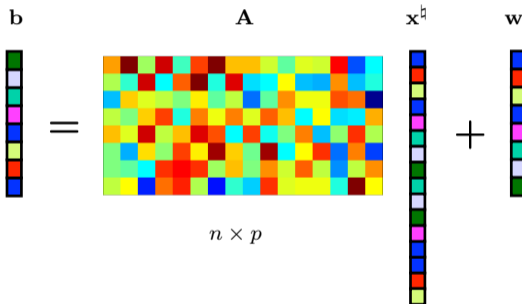
$$b_i = h_{\mathbf{x}^{\natural}}(\mathbf{a}_i)$$

\mathbf{x}^{\natural} : unknown function parameters

\mathbf{a}_i : input

b_i : response / output

Linear model:



$$\mathbf{b}_i = \langle \mathbf{a}_i, \mathbf{x}^{\natural} \rangle + \mathbf{w}_i$$

Applications: **Compressive sensing, machine learning, theoretical computer science...**

A simple *regression* model and many *practical* questions

$$\mathbf{b}_i = \langle \mathbf{a}_i, \mathbf{x}^{\natural} \rangle + \mathbf{w}_i$$

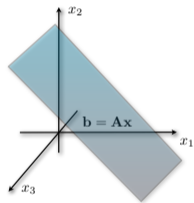
\mathbf{x}^{\natural} : unknown function parameters

\mathbf{a}_i : input

\mathbf{b}_i : response / output

\mathbf{w}_i : perturbations / noise

- Estimation: find \mathbf{x}^* to minimize $\|\mathbf{x}^* - \mathbf{x}^{\natural}\|$
- Prediction: find \mathbf{x}^* to minimize $L(\langle \mathbf{a}_i, \mathbf{x}^* \rangle, \langle \mathbf{a}_i, \mathbf{x}^{\natural} \rangle)$
- Decision: choose \mathbf{a}_i for estimation or prediction



A difficult estimation challenge when $n < p$:

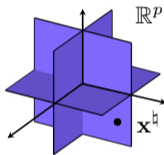
Nullspace (null) of \mathbf{A} : $\mathbf{x}^{\natural} + \mathbf{v} \rightarrow \mathbf{b}, \quad \forall \mathbf{v} \in \text{null}(\mathbf{A})$

- Needle in a haystack: **We need additional information on \mathbf{x}^{\natural} !**

A natural signal model

Definition (s -sparse vector)

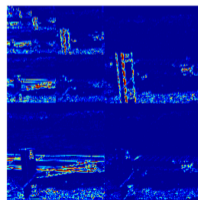
A vector $\mathbf{x} \in \mathbb{R}^p$ is s -sparse if it has at most s non-zero entries.



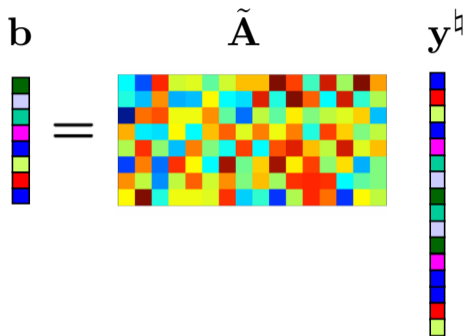
$$\mathbf{y}^h = \Psi \mathbf{x}^h$$

Sparse representations

- \mathbf{x}^h : *sparse* transform coefficients
- Basis representations $\Psi \in \mathbb{R}^{p \times p}$
 - ▶ *Wavelets*, DCT, ...
- Frame representations $\Psi \in \mathbb{R}^{m \times p}$, $m > p$
 - ▶ Gabor, curvelets, shearlets, ...
- Other *dictionary* representations...

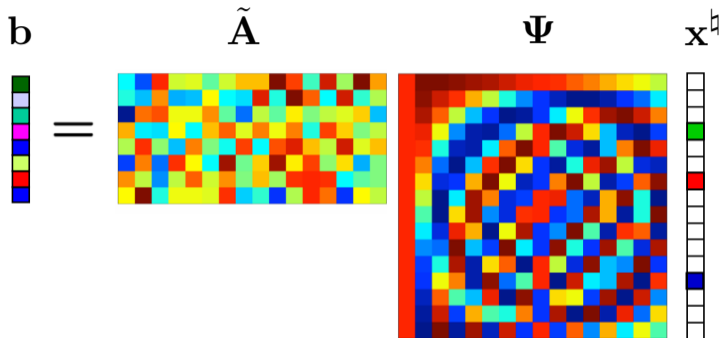


Sparse representations strike back!

$$\mathbf{b} = \tilde{\mathbf{A}} \mathbf{y}^{\natural}$$


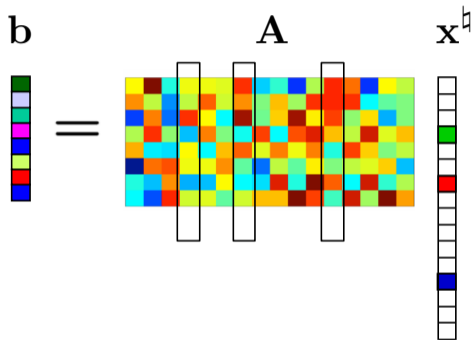
- $\mathbf{b} \in \mathbb{R}^n$, $\tilde{\mathbf{A}} \in \mathbb{R}^{n \times p}$, and $n < p$

Sparse representations strike back!



- $\mathbf{b} \in \mathbb{R}^n$, $\tilde{\mathbf{A}} \in \mathbb{R}^{n \times p}$, and $n < p$
- $\Psi \in \mathbb{R}^{p \times p}$, $\mathbf{x}^{\natural} \in \mathbb{R}^p$, and $\|\mathbf{x}^{\natural}\|_0 \leq s < n$

Sparse representations strike back!



◦ $\mathbf{b} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{n \times p}$, and $\mathbf{x}^h \in \mathbb{R}^p$, and $\|\mathbf{x}^h\|_0 \leq s < n < p$

Sparse representations strike back!

$$\mathbf{b} = \mathbf{A} \mathbf{x}^{\natural}$$

$n \times 1$ $n \times s$ $s \times 1$

- Observations:**
- The matrix \mathbf{A} effectively becomes *overcomplete*.
 - We could solve for \mathbf{x}^{\natural} if we knew *the location of the non-zero entries of \mathbf{x}^{\natural}* .

Enter sparsity

A combinatorial approach for estimating \mathbf{x}^{\natural} from $\mathbf{b} = \mathbf{A}\mathbf{x}^{\natural} + \mathbf{w}$

We may consider the estimator with the least number of non-zero entries. That is,

$$\mathbf{x}^{\star} \in \arg \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ \|\mathbf{x}\|_0 : \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2 \leq \kappa \right\} \quad (\mathcal{P}_0)$$

with some $\kappa \geq 0$. If $\kappa = \|\mathbf{w}\|_2$, then \mathbf{x}^{\natural} is a feasible solution.

Enter sparsity

A combinatorial approach for estimating \mathbf{x}^{\natural} from $\mathbf{b} = \mathbf{A}\mathbf{x}^{\natural} + \mathbf{w}$

We may consider the estimator with the least number of non-zero entries. That is,

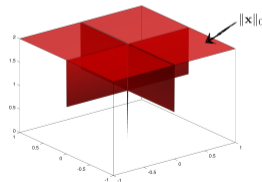
$$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ \|\mathbf{x}\|_0 : \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2 \leq \kappa \right\} \quad (\mathcal{P}_0)$$

with some $\kappa \geq 0$. If $\kappa = \|\mathbf{w}\|_2$, then \mathbf{x}^{\natural} is a feasible solution.

o \mathcal{P}_0 has the following characteristics:

- ▶ sample complexity: $\mathcal{O}(s)$
- ▶ computational effort: NP-Hard
- ▶ stability: No

$\|\mathbf{x}\|_0$ over the unit ℓ_∞ -ball



Enter sparsity

A combinatorial approach for estimating \mathbf{x}^{\natural} from $\mathbf{b} = \mathbf{A}\mathbf{x}^{\natural} + \mathbf{w}$

We may consider the estimator with the least number of non-zero entries. That is,

$$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ \|\mathbf{x}\|_0 : \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2 \leq \kappa \right\} \quad (\mathcal{P}_0)$$

with some $\kappa \geq 0$. If $\kappa = \|\mathbf{w}\|_2$, then \mathbf{x}^{\natural} is a feasible solution.

o \mathcal{P}_0 has the following characteristics:

- ▶ sample complexity: $\mathcal{O}(s)$
- ▶ computational effort: NP-Hard
- ▶ stability: No

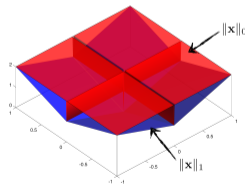
o **Tightest convex relaxation:**

- ▶ $\|\mathbf{x}\|_0^{**}$ is the **biconjugate**
- ▶ i.e., Fenchel conjugate of Fenchel conjugate

o **Fenchel conjugate:**

- ▶ $f^*(\mathbf{y}) := \sup_{\mathbf{x} \in \text{dom}(f)} \mathbf{x}^T \mathbf{y} - f(\mathbf{x})$.

$\|\mathbf{x}\|_1$ is the **convex envelope** of $\|\mathbf{x}\|_0$



A technicality: Restrict $\mathbf{x}^{\natural} \in [-1, 1]^p$.

The role of convexity

A convex candidate solution for $\mathbf{b} = \mathbf{A}\mathbf{x}^\dagger + \mathbf{w}$

$$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ \|\mathbf{x}\|_1 : \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2 \leq \|\mathbf{w}\|_2, \|\mathbf{x}\|_\infty \leq 1 \right\}. \quad (\text{SOCP})$$

Theorem (A model recovery guarantee [17])

Let $\mathbf{A} \in \mathbb{R}^{n \times p}$ be a matrix of i.i.d. Gaussian random variables with zero mean and variances $1/n$. For any $t > 0$ with probability at least $1 - 6 \exp(-t^2/26)$, we have

$$\|\mathbf{x}^* - \mathbf{x}^\dagger\|_2 \leq \left[\frac{2 \sqrt{2s \log(\frac{p}{s}) + \frac{5}{4}s}}{\sqrt{n} - \sqrt{2s \log(\frac{p}{s}) + \frac{5}{4}s} - t} \right] \|\mathbf{w}\|_2 := \varepsilon, \quad \text{when } \|\mathbf{x}^\dagger\|_0 \leq s.$$

- Observations:**
- perfect recovery (i.e., $\varepsilon = 0$) with $n \geq 2s \log(\frac{p}{s}) + \frac{5}{4}s$ whp when $\mathbf{w} = 0$.
 - ε -accurate solution in $k = \mathcal{O}\left(\sqrt{2p+1} \log(\frac{1}{\varepsilon})\right)$ iterations via IPM with a total complexity of $\mathcal{O}(n^2 p^{1.5} \log(\frac{1}{\varepsilon}))$ with each iteration requiring the solution of a structured $n \times 2p$ linear system.
 - robust to noise.

A Time-Data conundrum — I

A computational dogma

Running time of a learning algorithm increases with the size of the data.

A Time-Data conundrum — I

A computational dogma

Running time of a learning algorithm increases with the size of the data.

- Misaligned goals in the statistical and optimization disciplines

Discipline	Goal	Metric
Optimization	reaching numerical ϵ -accuracy	$\ \mathbf{x}^k - \mathbf{x}^*\ \leq \epsilon$
Statistics	learning ϵ -accurate model	$\ \mathbf{x}^* - \mathbf{x}^{\natural}\ \leq \epsilon$

- Main issue: ϵ and ϵ are **NOT** the same but should be treated jointly!

A Time-Data conundrum — II

A stylized formalization of the time-data tradeoff

The goals of optimization and statistical modeling are tightly connected:

$$\underbrace{\|\mathbf{x}^k - \mathbf{x}^\dagger\|}_{\text{learning quality}} \leq \underbrace{\|\mathbf{x}^k - \mathbf{x}^*\|}_{\varepsilon: \text{ needs "time" } t(k)} + \underbrace{\|\mathbf{x}^* - \mathbf{x}^\dagger\|}_{\varepsilon: \text{ needs "data" } n},$$

\mathbf{x}^\dagger : true model in \mathbb{R}^p
 \mathbf{x}^* : statistical model estimate
 \mathbf{x}^k : numerical solution at iteration k

o As the number of data samples n increases with a fixed optimization formulation,

$$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ \|\mathbf{x}\|_1 : \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2 \leq \|\mathbf{w}\|_2, \|\mathbf{x}\|_\infty \leq 1 \right\}$$

▶ numerical methods take longer time t to reach ε -accuracy

▶ e.g., per-iteration time to solve an $n \times 2p$ linear system

▶ statistical model estimates ε become more precise when $\|\mathbf{w}\|_2 = \mathcal{O}(\sqrt{n})$

▶ $\varepsilon = \frac{2 \sqrt{2s \log(\frac{p}{s}) + \frac{5}{4}s}}{\sqrt{n} - \sqrt{2s \log(\frac{p}{s}) + \frac{5}{4}s - t}} \|\mathbf{w}\|_2$, with probability $1 - 6\exp(-t^2/26)$.

A Time-Data conundrum — II

A stylized formalization of the time-data tradeoff

The goals of optimization and statistical modeling are tightly connected:

$$\underbrace{\|\mathbf{x}^k - \mathbf{x}^\dagger\|}_{\leq \bar{\varepsilon}(t(k), n)} \leq \underbrace{\|\mathbf{x}^k - \mathbf{x}^*\|}_{\varepsilon: \text{ needs "time" } t(k)} + \underbrace{\|\mathbf{x}^* - \mathbf{x}^\dagger\|}_{\varepsilon: \text{ needs "data" } n},$$

\mathbf{x}^\dagger :	true model in \mathbb{R}^p
\mathbf{x}^* :	statistical model estimate
\mathbf{x}^k :	numerical solution at iteration k
$\bar{\varepsilon}(t(k), n)$:	actual learning quality at time $t(k)$ with n samples

- o As the number of data samples n increases with a fixed optimization formulation,

$$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ \|\mathbf{x}\|_1 : \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2 \leq \|\mathbf{w}\|_2, \|\mathbf{x}\|_\infty \leq 1 \right\}$$

- ▶ numerical methods take longer time t to reach ε -accuracy
 - ▶ e.g., per-iteration time to solve an $n \times 2p$ linear system
- ▶ statistical model estimates ε become more precise when $\|\mathbf{w}\|_2 = \mathcal{O}(\sqrt{n})$

- ▶ $\varepsilon = \frac{2 \sqrt{2s \log(\frac{p}{s}) + \frac{5}{4}s}}{\sqrt{n} - \sqrt{2s \log(\frac{p}{s}) + \frac{5}{4}s - t}} \|\mathbf{w}\|_2$, with probability $1 - 6\exp(-t^2/26)$.

“Time” effort has significant diminishing returns on ε in the underdetermined case* (cf., [8, 5, 19, 7, 6])

* “Data” effort also exhibits a similar behavior in the overdetermined case when a signal prior is used due to noise!

Data as a computational resource

A stylized formalization of the time-data tradeoff

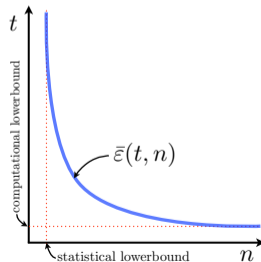
The goals of optimization and statistical modeling are tightly connected:

$$\underbrace{\|\mathbf{x}^{k(t)} - \mathbf{x}^{\natural}\|}_{\leq \bar{\epsilon}(t, n)} \leq \underbrace{\|\mathbf{x}^{k(t)} - \mathbf{x}^{\star}\|}_{\epsilon: \text{ needs "time" } t} + \underbrace{\|\mathbf{x}^{\star} - \mathbf{x}^{\natural}\|}_{\epsilon: \text{ needs "data" } n},$$

\mathbf{x}^{\natural} : true model in \mathbb{R}^p

$\bar{\epsilon}(t, n)$: actual model precision at time t with n samples

- Rest of the lecture:
- o estimator formulation and sample complexity
 - o a “continuous” time-data tradeoff
 - o a different, algorithmic tradeoff with SGD



Sample complexity analysis

Convex optimization formulation for the estimator

$$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^p} \{f(\mathbf{x}) : \mathbf{b} = \mathbf{A}\mathbf{x}\},$$

where $f : \mathbb{R}^p \rightarrow \mathbb{R} \cup \{-\infty, \infty\}$ is a convex function.

Sample complexity

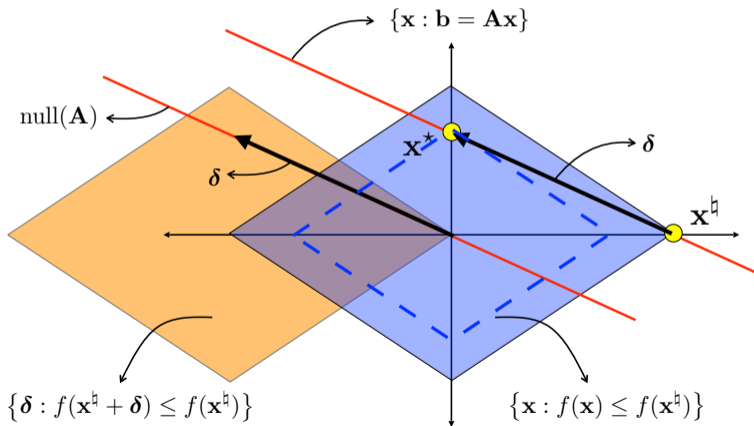
Assume that $A \in \mathbb{R}^{n \times p}$ is a matrix of independent identically distributed (i.i.d.) Gaussian random variables.

What is the minimum number of samples n such that $\mathbf{x}^* = \hat{\mathbf{x}}$ with high probability?

Characterization of the error vector

$$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^p} \{f(\mathbf{x}) : \mathbf{b} = \mathbf{A}\mathbf{x}\}$$

Define the error vector $\delta := \mathbf{x}^* - \mathbf{x}^{\natural}$.

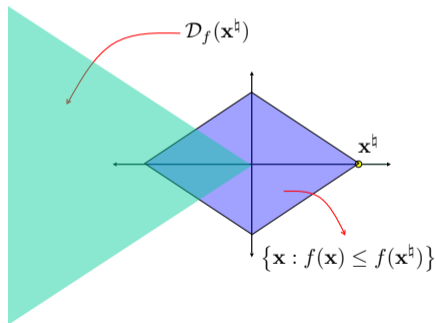


Descent cone

Definition (Descent cone)

Let $f : \mathbb{R}^p \rightarrow \mathbb{R} \cup \{-\infty, \infty\}$ be a proper lower-semicontinuous function. The **descent cone** of f at \mathbf{x}^h is defined as

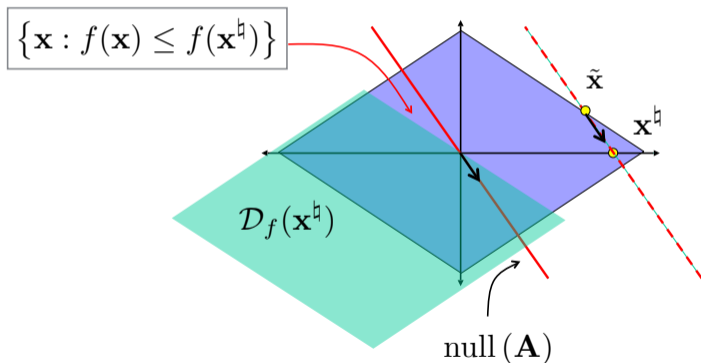
$$\mathcal{D}_f(\mathbf{x}^h) := \text{cone} \left(\left\{ \delta : f(\mathbf{x}^h + \delta) \leq f(\mathbf{x}^h) \right\} \right).$$



Condition for exact recovery in the *noiseless* case

Proposition (Condition for exact recovery)

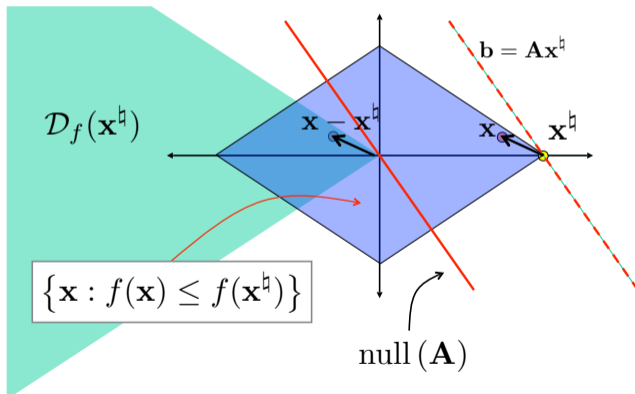
We have successful recovery, i.e., $\delta := \mathbf{x}^* - \mathbf{x}^{\natural} = 0$ with $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^p} \{f(\mathbf{x}) : \mathbf{b} = \mathbf{A}\mathbf{x}\}$, if and only if $\text{null}(\mathbf{A}) \cap \mathcal{D}_f(\mathbf{x}^{\natural}) = \{0\}$.



Condition for exact recovery in the *noiseless* case

Proposition (Condition for exact recovery)

We have successful recovery, i.e., $\delta := \mathbf{x}^* - \mathbf{x}^{\natural} = 0$ with $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^p} \{f(\mathbf{x}) : \mathbf{b} = \mathbf{A}\mathbf{x}\}$, if and only if $\text{null}(\mathbf{A}) \cap \mathcal{D}_f(\mathbf{x}^{\natural}) = \{0\}$.



Statistical dimension and approximate kinematic formula

Now we have

$$\mathbb{P} \{ \mathbf{x}^* = \mathbf{x}^{\natural} \} = \mathbb{P} \{ \text{null}(\mathbf{A}) \cap \mathcal{D}_f(\mathbf{x}^{\natural}) = \{0\} \}.$$

Definition (Statistical dimension [3]¹)

Let $\mathcal{C} \subseteq \mathbb{R}^p$ be a closed convex cone. The *statistical dimension* of \mathcal{C} is defined as

$$d(\mathcal{C}) := \mathbb{E} \left[\|\text{proj}_{\mathcal{C}}(\mathbf{g})\|_2^2 \right].$$

Theorem (Approximate kinematic formula [3])

Let $A \in \mathbb{R}^{n \times p}$, $n < p$, be a matrix of i.i.d. standard Gaussian random variables, and let $\mathcal{C} \subseteq \mathbb{R}^p$ be a closed convex cone. Let $\eta \in (0, 1)$. Then

$$\begin{aligned} n \geq d(\mathcal{C}) + c_{\eta} \sqrt{p} &\Rightarrow \mathbb{P} \{ \text{null}(\mathbf{A}) \cap \mathcal{C} = \{0\} \} \geq 1 - \eta; \\ n \leq d(\mathcal{C}) - c_{\eta} \sqrt{p} &\Rightarrow \mathbb{P} \{ \text{null}(\mathbf{A}) \cap \mathcal{C} = \{0\} \} \leq \eta, \end{aligned}$$

where $c_{\eta} := \sqrt{8 \log(4/\eta)}$.

¹The statistical dimension is closely related to the Gaussian complexity [4], Gaussian width [9], and Gaussian squared complexity [8].

Probability of exact recovery

Corollary

For any $\eta \in (0, 1)$,

$$n \geq d(\mathcal{D}_f(\mathbf{x}^{\natural})) + c_\eta \sqrt{p} \quad \Rightarrow \quad \mathbb{P} \left\{ \mathbf{x}^* = \mathbf{x}^{\natural} \right\} \geq 1 - \eta;$$

$$n \leq d(\mathcal{D}_f(\mathbf{x}^{\natural})) - c_\eta \sqrt{p} \quad \Rightarrow \quad \mathbb{P} \left\{ \mathbf{x}^* = \mathbf{x}^{\natural} \right\} \leq \eta,$$

where $c_\eta := \sqrt{8 \log(4/\eta)}$.

- There is a *phase transition* at $n \approx d(\mathcal{D}_f(\mathbf{x}^{\natural}))$.

Examples ([3])

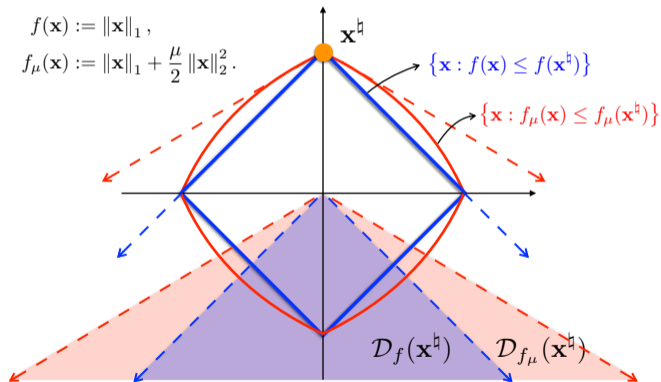
- Let $f(\mathbf{x}) := \|\mathbf{x}\|_1$, and let $\mathbf{x}^{\natural} \in \mathbb{R}^p$ be s -sparse. Then $d(\mathcal{D}_f(\mathbf{x}^{\natural})) \leq 2s \log(p/s) + (5/4)s$.
- Let $f(\mathbf{x}) := \|\mathbf{X}\|_*$, and let $\mathbf{X}^{\natural} \in \mathbb{R}^{p \times p}$ of rank r . Then $d(\mathcal{D}_f(\mathbf{x}^{\natural})) \leq 3r(2p - r)$.

Smoothing increases the statistical dimension

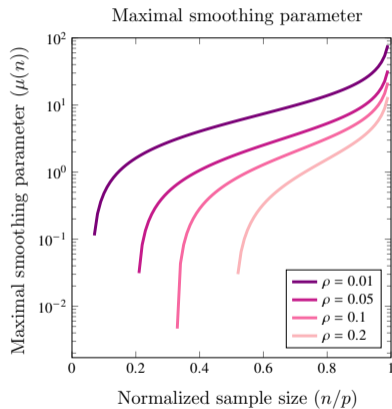
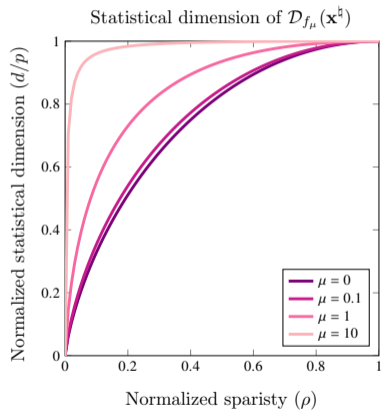
Key properties of the statistical dimension [3]

- The statistical dimension is invariant under unitary transformations (rotations).
- Let \mathcal{C}_1 and \mathcal{C}_2 be closed convex cones. If $\mathcal{C}_1 \subseteq \mathcal{C}_2$, then $d(\mathcal{C}_1) \leq d(\mathcal{C}_2)$.

The larger the statistical dimension is, the more number of observations is required.



Numerical results for the statistical dimension and $\mu(n)$



Smoothing decreases the computational cost

- Consider the estimator,

$$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f_\mu(\mathbf{x}) : \mathbf{b} = \mathbf{A}\mathbf{x}, \|\mathbf{x}\|_\infty \leq \|\mathbf{x}^*\|_\infty \right\}, \quad \mu \in [0, \infty).$$

Proposition

Let $\mu > 0$ and $f(\mathbf{x}) = \|\mathbf{x}\|_1$. Consider solving (1) with a primal-dual method as in [6, 7]. The output after the k -th iteration, \mathbf{x}^k , satisfies

$$\|\mathbf{x}^* - \mathbf{x}^k\|_2 \leq \frac{4p\kappa(\mathbf{A}) \left[\rho(1 + \mu \|\mathbf{x}^*\|_\infty)^2 + (1 - \rho) \right]}{\mu k} \propto \frac{1}{\mu k} \Big|_{\rho \ll 1},$$

where $\rho := s/p$, s being the number of non-zero entries in \mathbf{x}^* , and $\kappa(\mathbf{A})$ denotes the restricted condition number of \mathbf{A} .

- Observations:**
- When $\rho \ll 1$, the number of iterations k to achieve the required precision decreases.
 - In fact, we need $1/(\mu\varepsilon)$ iterations to have an error bound $\|\mathbf{x}^* - \mathbf{x}^k\|_2 \leq \varepsilon$ for a fixed $\varepsilon > 0$.

Time-data tradeoff

- Define the maximal smoothing parameter

$$\mu(n) := \arg \max_{\mu > 0} \left\{ \mu : d \left(\mathcal{D}_{f_\mu}(\mathbf{x}^{\natural}) \right) \leq n \right\}.$$

- Consider the “conservative” estimator in probability,

$$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f_\mu(\mathbf{x}) \Big|_{\mu = \frac{1}{4} \mu(n)} : \mathbf{b} = \mathbf{A}\mathbf{x} \right\}.$$

Corollary

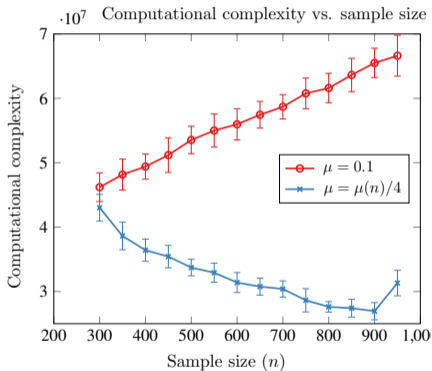
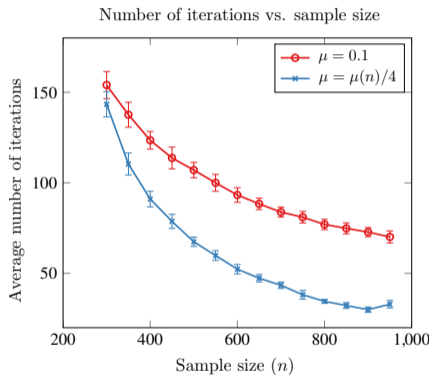
Let $\rho := s/p \ll 1$. Then we have, with high probability, $\mathbf{x}^* = \mathbf{x}^{\natural}$, and

$$\left\| \mathbf{x}^{\natural} - \mathbf{x}^k \right\|_2 \propto \frac{1}{\mu(n)k}.$$

Therefore, to achieve the error bound, $\left\| \mathbf{x}^{\natural} - \mathbf{x}^k \right\|_2 \leq \varepsilon$ for a fixed $\varepsilon > 0$, it suffices to choose

$$k = O \left(\frac{1}{\mu(n)} \right).$$

A numerical result for the time-data tradeoff



Another trade-off in optimization

- Statistics vs Optimization:

Discipline	Goal	Metric
Optimization	reaching numerical ϵ -accuracy	$\ \mathbf{x}^k - \mathbf{x}^*\ \leq \epsilon$
Statistics	learning ϵ -accurate model	$\ \mathbf{x}^* - \hat{\mathbf{x}}\ \leq \epsilon$

Remarks: *As data sample size gets larger we have seen that:*

- Algorithms take **longer to reach ϵ** accuracy.
- However, **statistical error ϵ decreases** as the estimation is more precise.

Another trade-off in optimization

- Statistics vs Optimization:

Discipline	Goal	Metric
Optimization	reaching numerical ϵ -accuracy	$\ \mathbf{x}^k - \mathbf{x}^*\ \leq \epsilon$
Statistics	learning ϵ -accurate model	$\ \mathbf{x}^* - \hat{\mathbf{x}}\ \leq \epsilon$

Remarks: *As data sample size gets larger we have seen that:*

- Algorithms take **longer to reach ϵ** accuracy.
- However, **statistical error ϵ decreases** as the estimation is more precise.

Similar analogy exists between **per-iteration cost** and **convergence rate** for optimization algorithms

Another trade-off in optimization

- Statistics vs Optimization:

Discipline	Goal	Metric
Optimization	reaching numerical ϵ -accuracy	$\ \mathbf{x}^k - \mathbf{x}^*\ \leq \epsilon$
Statistics	learning ϵ -accurate model	$\ \mathbf{x}^* - \hat{\mathbf{x}}\ \leq \epsilon$

Remarks: *As data sample size gets larger we have seen that:*

- Algorithms take **longer to reach ϵ** accuracy.
- However, **statistical error ϵ decreases** as the estimation is more precise.

Similar analogy exists between **per-iteration cost** and **convergence rate** for optimization algorithms

Understanding this trade-off helps us reduce total complexity!

Recall: GD vs. SGD

Problem (Unconstrained convex minimization)

Deterministic setting

$$f^* = \min_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x})$$

- $f(\mathbf{x})$ is a *proper, closed, convex and smooth*.
- The solution set
 $\mathcal{S}^* := \{\mathbf{x}^* \in \text{dom}(f) : f(\mathbf{x}^*) = f^*\} \neq \emptyset$.

Stochastic programming

$$f^* = \min_{\mathbf{x} \in \mathbb{R}^p} \{f(\mathbf{x}) := \mathbb{E}[f(\mathbf{x}, \theta)]\}$$

- $f(\mathbf{x})$ is *proper, closed, convex and smooth*.
- The solution set
 $\mathcal{S}^* := \{\mathbf{x}^* \in \text{dom}(f) : f(\mathbf{x}^*) = f^*\} \neq \emptyset$.
- θ is a random vector, supported on set Θ .

Algorithms

Gradient Descent

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k \nabla f(\mathbf{x}^k)$$

- $\alpha_k < 2/L$.

Stochastic Gradient Descent

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k G(\mathbf{x}^k, \theta_k)$$

- $\alpha_k = \mathcal{O}(1/\sqrt{k})$
- $\mathbb{E}[G(\mathbf{x}^k, \theta_k)] = \nabla f(\mathbf{x}^k)$

Example: Convex optimization with finite sum

- Consider the finite sum (e.g., ERM) setting

$$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) := \frac{1}{n} \sum_{j=1}^n f_j(\mathbf{x}) \right\}.$$

Algorithms in the finite sum setting

Gradient Descent

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k \nabla f(\mathbf{x}^k)$$

- $\nabla f(\mathbf{x}^k) = \frac{1}{n} \sum_{j=1}^n \nabla f_j(\mathbf{x}^k)$

Stochastic Gradient Descent

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k G(\mathbf{x}^k, \theta_k)$$

- $G(\mathbf{x}^k, \theta_k) = \nabla f_j(\mathbf{x}^k), j \sim \text{Uniform}(\{1, \dots, n\})$

Example: Convex optimization with finite sum

- Consider the finite sum (e.g., ERM) setting

$$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) := \frac{1}{n} \sum_{j=1}^n f_j(\mathbf{x}) \right\}.$$

Algorithms in the finite sum setting

Gradient Descent

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k \nabla f(\mathbf{x}^k)$$

- $\nabla f(\mathbf{x}^k) = \frac{1}{n} \sum_{j=1}^n \nabla f_j(\mathbf{x}^k)$

Stochastic Gradient Descent

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k G(\mathbf{x}^k, \theta_k)$$

- $G(\mathbf{x}^k, \theta_k) = \nabla f_j(\mathbf{x}^k), j \sim \text{Uniform}(\{1, \dots, n\})$

- $f(\mathbf{x})$: **convex** and L -Lipschitz gradient

	rate	cost per iteration	iteration complexity	total complexity
GD	$1/k$	n	$1/\epsilon$	n/ϵ
SGD	$1/\sqrt{k}$	1	$1/\epsilon^2$	$1/\epsilon^2$

Example: Convex optimization with finite sum

- Consider the finite sum (e.g., ERM) setting

$$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) := \frac{1}{n} \sum_{j=1}^n f_j(\mathbf{x}) \right\}.$$

Algorithms in the finite sum setting

Gradient Descent

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k \nabla f(\mathbf{x}^k)$$

- $\nabla f(\mathbf{x}^k) = \frac{1}{n} \sum_{j=1}^n \nabla f_j(\mathbf{x}^k)$

Stochastic Gradient Descent

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k G(\mathbf{x}^k, \theta_k)$$

- $G(\mathbf{x}^k, \theta_k) = \nabla f_j(\mathbf{x}^k), j \sim \text{Uniform}(\{1, \dots, n\})$

- $f(\mathbf{x})$: μ -strongly convex and L -Lipschitz gradient

	rate	cost per iteration	iteration complexity	total complexity
GD	ρ^k	n	$\log(1/\epsilon)$	$n \log(1/\epsilon)$
SGD	$1/k$	1	$1/\epsilon$	$1/\epsilon$

When f is μ -strongly convex and L -Lipschitz gradient

Finite sums

$$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) := \frac{1}{n} \sum_{j=1}^n f_j(\mathbf{x}) \right\}.$$

	rate	cost per iteration	iteration complexity	total complexity
GD	ρ^k	n	$\log(1/\epsilon)$	$n \log(1/\epsilon)$
SGD	$1/k$	1	$1/\epsilon$	$1/\epsilon$

- Remarks:**
- o SGD trades off **convergence rate** with **low per-iteration cost**.
 - o When n is large, SGD proves to be effective.
 - o To **control variance** of the stochastic gradient estimate, SGD **decreases step size** at a certain rate.
 - o In turn, convergence deteriorates from **linear** to **sublinear**.

An observation of GD vs. SGD step

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k \nabla f(\mathbf{x}^k) \quad (\text{GD})$$

Lemma

Assume f is Lipschitz smooth with constant L . Then,

$$f(\mathbf{x}^{k+1}) - f(\mathbf{x}^k) \leq \left(\frac{\gamma_k^2 L}{2} - \gamma_k \right) \|\nabla f(\mathbf{x}^k)\|^2.$$

An observation of GD vs. SGD step

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k G(\mathbf{x}^k, \theta_k) \quad (\text{SGD})$$

Lemma

Assume f is Lipschitz smooth with constant L . Then,

$$\mathbb{E}[f(\mathbf{x}^{k+1}) - f(\mathbf{x}^k)] \leq \left(\frac{\gamma_k^2 L}{2} - \gamma_k \right) \mathbb{E}[\|\nabla f(\mathbf{x}^k)\|^2] + \frac{L\gamma_k^2}{2} \mathbb{E}[\|G(\mathbf{x}^k, \theta_k) - \nabla f(\mathbf{x}^k)\|^2]$$

An observation of GD vs. SGD step

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k G(\mathbf{x}^k, \theta_k) \quad (\text{SGD})$$

Lemma

Assume f is Lipschitz smooth with constant L . Then,

$$\mathbb{E}[f(\mathbf{x}^{k+1}) - f(\mathbf{x}^k)] \leq \left(\frac{\gamma_k^2 L}{2} - \gamma_k \right) \mathbb{E}[\|\nabla f(\mathbf{x}^k)\|^2] + \frac{L\gamma_k^2}{2} \mathbb{E}[\|G(\mathbf{x}^k, \theta_k) - \nabla f(\mathbf{x}^k)\|^2]$$

- The variance of gradient estimate dominates as $\nabla f(\mathbf{x}^k) \rightarrow 0$.
- To ensure convergence we need to control variance.

$\gamma_k \rightarrow 0 \implies$ Slow convergence!

Can we decrease the variance while using a constant step-size?

An observation of GD vs. SGD step

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k G(\mathbf{x}^k, \theta_k) \quad (\text{SGD})$$

Lemma

Assume f is Lipschitz smooth with constant L . Then,

$$\mathbb{E}[f(\mathbf{x}^{k+1}) - f(\mathbf{x}^k)] \leq \left(\frac{\gamma_k^2 L}{2} - \gamma_k \right) \mathbb{E}[\|\nabla f(\mathbf{x}^k)\|^2] + \frac{L\gamma_k^2}{2} \mathbb{E}[\|G(\mathbf{x}^k, \theta_k) - \nabla f(\mathbf{x}^k)\|^2]$$

- The variance of gradient estimate dominates as $\nabla f(\mathbf{x}^k) \rightarrow 0$.
- To ensure convergence we need to control variance.

$\gamma_k \rightarrow 0 \implies$ Slow convergence!

Can we decrease the variance while using a constant step-size?

Choose a stochastic gradient, s.t. $\mathbb{E}[\|G(\mathbf{x}^k; \theta_k)\|^2] \rightarrow 0$.

A simple approach: Mini-batch SGD

- More samples \rightarrow better estimate for full gradient.

SGD with mini batches

Let $G(\mathbf{x}, \theta)$ be an unbiased gradient estimate ($\mathbb{E}[G(\mathbf{x}, \theta)] = \nabla f(\mathbf{x})$) and B_k be the batch size. Then,

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k \frac{1}{B_k} \sum_{j=1}^{B_k} G(\mathbf{x}^k, \theta_{k,j})$$

Theorem

Let $B_k > 0$ be the batch size and $G(\mathbf{x}, \theta)$ be an unbiased gradient estimate with bounded variance, i.e., $\mathbb{E}[\|G(\mathbf{x}, \theta) - \nabla f(\mathbf{x})\|^2 | \mathbf{x}] \leq \sigma^2$. Then, the mini-batch estimate has the following properties:

$$\mathbb{E} \left[\frac{1}{B_k} \sum_{j=1}^{B_k} G(\mathbf{x}, \theta_{k,j}) \right] = \nabla f(\mathbf{x}) \quad \text{and} \quad \mathbb{E} \left[\left\| \frac{1}{B_k} \sum_{j=1}^{B_k} G(\mathbf{x}, \theta_{k,j}) - \nabla f(\mathbf{x}) \right\|^2 \mid \mathbf{x} \right] \leq \frac{\sigma^2}{B_k}$$

- Remarks:**
- We might need to increase the batch size over time to take variance to 0.
 - We can come up with a “smarter” estimate for $\nabla f(\mathbf{x})$.

How to construct a new estimate $G(\mathbf{x}^k; \theta_k)$? [10]

Finite sum structure:	SGD update rule:
$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) := \frac{1}{n} \sum_{j=1}^n f_j(\mathbf{x}) \right\}$	$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k \nabla f_j(\mathbf{x}^k)$

- Let $X = \nabla f_j(\mathbf{x}^k)$ be a random variable (due to $j \sim \text{Uniform}(\{1, \dots, n\})$).
- Let $Y = \nabla f_j(\tilde{\mathbf{x}})$ be another random variable, and $\tilde{\mathbf{x}}$ is a particularly selected point.

Remarks:

- We want X and Y to be correlated (we will see why!).
- Given Y , we should be able to estimate $\mathbb{E}[X]$ with more confidence.

Observations:

- Choice of $\tilde{\mathbf{x}}$ affects how correlated X and Y are.
- We can compute $\mathbb{E}[Y] = \frac{1}{n} \sum_{j=1}^n \nabla f_j(\tilde{\mathbf{x}}) = \nabla f(\tilde{\mathbf{x}})$.

Goal:

- Find a **good** estimate of $\mathbb{E}[X] = \frac{1}{n} \sum_{j=1}^n \nabla f_j(\mathbf{x}^k) = \nabla f(\mathbf{x}^k)$.

How to construct a new estimate $G(\mathbf{x}^k; \theta_k)$? [10]

Finite sum structure:	SGD update rule:
$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) := \frac{1}{n} \sum_{j=1}^n f_j(\mathbf{x}) \right\}$	$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k \nabla f_j(\mathbf{x}^k)$

- Let $X = \nabla f_j(\mathbf{x}^k)$ be a random variable (due to $j \sim \text{Uniform}(\{1, \dots, n\})$).
- Let $Y = \nabla f_j(\tilde{\mathbf{x}})$ be another random variable, and $\tilde{\mathbf{x}}$ is a particularly selected point.

A generalized estimator: $R_\alpha = \alpha(X - Y) + \mathbb{E}[Y]$

- $\mathbb{E}[R_\alpha] = \alpha\mathbb{E}[X] + (1 - \alpha)\mathbb{E}[Y]$
- $\text{Var}(R_\alpha) = \alpha^2(\text{Var}(X) + \text{Var}(Y) - 2\text{Cov}(X, Y))$

- Observations:**
- When $\alpha = 1$, R_α becomes unbiased, i.e., $\mathbb{E}[R_\alpha] = \mathbb{E}[X]$.
 - If $\text{Cov}(X, Y)$ is large enough (X and Y are correlated enough), $\text{Var}(R_\alpha) \leq \text{Var}(X)$.

How could we use this information to construct our estimate?

Variance reduction techniques: SVRG

- Select the stochastic gradient ∇f_{i_k} , and compute a gradient estimate

$$\mathbf{r}_k = \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}) + \nabla f(\tilde{\mathbf{x}}).$$

- As $\tilde{\mathbf{x}} \rightarrow \mathbf{x}^*$ and $\mathbf{x}^k \rightarrow \mathbf{x}^*$,

$$\nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}) + \nabla f(\tilde{\mathbf{x}}) \rightarrow 0.$$

- Therefore,

$$\mathbb{E}[\|\nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}) + \nabla f(\tilde{\mathbf{x}})\|^2] \rightarrow 0.$$

Remarks:

- Remember the generalized estimator: $R_\alpha = \alpha(X - Y) + \mathbb{E}[Y]$.
- For SVRG, $\alpha = 1$, $X = \nabla f_{i_k}(\mathbf{x}^k)$ and $Y = \nabla f_{i_k}(\tilde{\mathbf{x}})$.
- We will see how $\tilde{\mathbf{x}}$ is computed!

Stochastic gradient algorithm with variance reduction

Stochastic gradient with variance reduction (SVRG) [12, 21]

1. Choose $\tilde{\mathbf{x}}^0 \in \mathbb{R}^p$ as a starting point and $\gamma > 0$ and $q \in \mathbb{N}_+$.

2. For $s = 0, 1, 2, \dots$, perform:

2a. $\tilde{\mathbf{x}} = \tilde{\mathbf{x}}^s$, $\tilde{\mathbf{v}} = \nabla f(\tilde{\mathbf{x}})$, $\mathbf{x}^0 = \tilde{\mathbf{x}}$.

2b. For $k = 0, 1, \dots, q-1$, perform:

$$\begin{cases} \text{Pick } i_k \in \{1, \dots, n\} \text{ uniformly at random} \\ \mathbf{r}_k = \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}) + \tilde{\mathbf{v}} \\ \mathbf{x}^{k+1} := \mathbf{x}^k - \gamma \mathbf{r}_k, \end{cases} \quad (1)$$

2c. Update $\tilde{\mathbf{x}}^{s+1} = \frac{1}{m} \sum_{j=0}^{q-1} \mathbf{x}^j$.

Features

- The SVRG method uses a multistage scheme to reduce the **variance** of the **stochastic gradient** \mathbf{r}_k .
- **Learning rate** γ does not necessarily tend to 0 while \mathbf{x}^k and $\tilde{\mathbf{x}}^s$ tend to \mathbf{x}_* .
- Each stage, SVRG uses $n + 2q$ component **gradient** evaluations.
- n for the **full gradient** at the beginning of each stage, and $2q$ for each of the q **stochastic gradient steps**.

Convergence analysis

Assumption A5.

- (i) f is μ -strongly convex
- (ii) The learning rate $0 < \gamma < 1/(4L_{\max})$, where $L_{\max} = \max_{1 \leq j \leq n} L_j$.
- (iii) q is large enough such that

$$\kappa = \frac{1}{\mu\gamma(1 - 4\gamma L_{\max})q} + \frac{4\gamma L_{\max}(q + 1)}{(1 - 4\gamma L_{\max})q} < 1.$$

Theorem

Assumptions:

- The sequence $\{\tilde{\mathbf{x}}^s\}_{k \geq 0}$ is generated by SVRG.
- Assumption A5 is satisfied.

Conclusion: Linear convergence is obtained:

$$\mathbb{E}f(\tilde{\mathbf{x}}^s) - f(\mathbf{x}^*) \leq \kappa^s (f(\tilde{\mathbf{x}}^0) - f(\mathbf{x}^*)).$$

Choice of γ and q , and complexity

Chose γ and q such that $\kappa \in (0, 1)$:

For example

$$\gamma = 0.1/L_{\max}, q = 100(L_{\max}/\mu) \implies \kappa \approx 5/6.$$

Complexity

$$\mathbb{E}f(\tilde{\mathbf{x}}^s) - f(\mathbf{x}^*) \leq \epsilon, \quad \text{when } s \geq \log((f(\tilde{\mathbf{x}}^0) - f(\mathbf{x}^*))/\epsilon) / \log(\kappa^{-1})$$

- o Each stage needs $n + 2q$ **component gradient evaluations**
- o With $q = \mathcal{O}(L_{\max}/\mu)$, we obtain an **overall complexity** of

$$\mathcal{O}\left((n + L_{\max}/\mu) \log(1/\epsilon)\right).$$

Comparison: GD vs. SGD vs. SVRG

- GD update:

$$\{ \mathbf{x}^{k+1} := \mathbf{x}^k - \gamma \nabla f(\mathbf{x}^k),$$

- SGD update:

$$\{ \mathbf{x}^{k+1} := \mathbf{x}^k - \gamma \nabla f_{i_k}(\mathbf{x}^k),$$

- SVRG update:

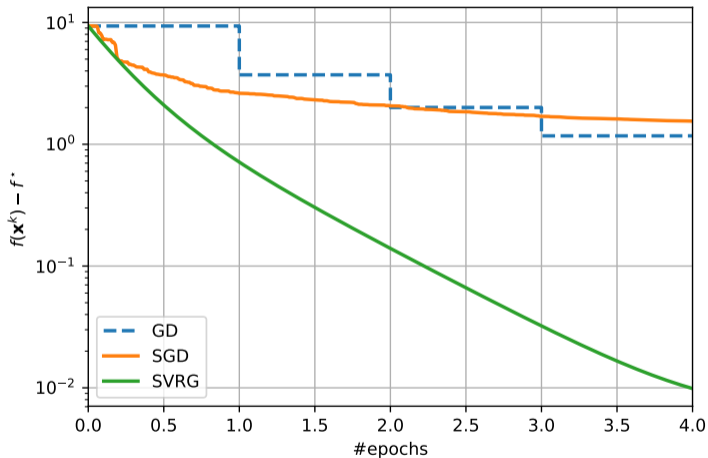
$$\begin{cases} \mathbf{r}_k = \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}) + \nabla f(\tilde{\mathbf{x}}) \\ \mathbf{x}^{k+1} := \mathbf{x}^k - \gamma \mathbf{r}_k, \end{cases}$$

	SGD	SVRG	GD
Requires gradient storage?	no	no	no
Epoch-based	no	yes	no
Parameters	stepsize	stepsize & epoch length	stepsize
Gradient evaluations	1 per iteration	$n + 2q$ per epoch	n per iteration

Table: Comparisons of SGD, SVRG and GD [10]

- Recall that $q = \mathcal{O}(L_{\max}/\mu)$ is the epoch length for SVRG.

Example: ℓ_2 -regularized least squares with synthetic data



Taxonomy of algorithms

$$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) := \frac{1}{n} \sum_{j=1}^n f_j(\mathbf{x}) \right\}.$$

- $f(\mathbf{x}) = \frac{1}{n} \sum_{j=1}^n f_j(\mathbf{x})$: μ -strongly convex with L -Lipschitz continuous gradient.

SVRG	GD	SGD
Linear	Linear	Sublinear

Table: Rate of convergence.

- $\kappa = L/\mu$.

SVRG	AGD	SGD
$\mathcal{O}((n + \kappa) \log(1/\varepsilon))$	$\mathcal{O}(n\kappa \log(1/\varepsilon))$	$1/\varepsilon$

Table: Complexity to obtain ε -solution.

The variance reduction zoo

Setting	Algorithm	Lower bound	Complexity bound
L -smooth f_i 's with bounded variance	Gradient descent	$L\Delta_0 \min\{\sigma/\epsilon^3, \sqrt{n}/\epsilon^2\}$ [11]	$nL\Delta_0/\epsilon^2$
	SVRG ($B_k = 1$) [18]	$L\Delta_0 \min\{\sigma/\epsilon^3, \sqrt{n}/\epsilon^2\}$ [11]	$nL\Delta_0/\epsilon^2$
	SVRG ($B_k = \Omega(n^{2/3})$) [18]	$L\Delta_0 \min\{\sigma/\epsilon^3, \sqrt{n}/\epsilon^2\}$ [11]	$n^{2/3}L\Delta_0/\epsilon^2$
	SAGA ($B_k = 1$) [18]	$L\Delta_0 \min\{\sigma/\epsilon^3, \sqrt{n}/\epsilon^2\}$ [11]	$nL\Delta_0/\epsilon^2$
	SAGA ($B_k = \Omega(n^{2/3})$) [18]	$L\Delta_0 \min\{\sigma/\epsilon^3, \sqrt{n}/\epsilon^2\}$ [11]	$n^{2/3}L\Delta_0/\epsilon^2$
	SpiderBoost [20]	$L\Delta_0 \min\{\sigma/\epsilon^3, \sqrt{n}/\epsilon^2\}$ [11]	$\sqrt{n}L\Delta_0/\epsilon^2$
	SpiderBoost-M [20]	$L\Delta_0 \min\{\sigma/\epsilon^3, \sqrt{n}/\epsilon^2\}$ [11]	$\sqrt{n}L\Delta_0/\epsilon^2$
	Spider [11] PAGE [14]	$L\Delta_0 \min\{\sigma/\epsilon^3, \sqrt{n}/\epsilon^2\}$ [11] $L\Delta_0 \min\{\sigma/\epsilon^3, \sqrt{n}/\epsilon^2\}$ [11]	$L\Delta_0 \min\{\sigma/\epsilon^3, \sqrt{n}/\epsilon^2\}$ $L\Delta_0 \min\{\sigma/\epsilon^3, \sqrt{n}/\epsilon^2\}$
f is μ -SCVX and L -smooth f_i 's are average L -smooth	KatyushaX [2]	$(n + n^{3/4} \sqrt{\frac{L}{\mu}}) \log \frac{\Delta_0}{\epsilon}$ [22]	$(n + n^{3/4} \sqrt{\frac{L}{\mu}}) \log \frac{\Delta_0}{\epsilon}$
f is CVX and L -smooth f_i 's are average L -smooth	KatyushaX [2]	$n + n^{3/4} \sqrt{\frac{LD_0^2}{\epsilon}}$ [23]	$n + n^{3/4} \sqrt{\frac{LD_0^2}{\epsilon}}$
f is α -weakly CVX and L -smooth f_i 's are average L -smooth	Spider [11]	$\frac{\Delta_0}{\epsilon^2} \min\{n^{3/4} \sqrt{\alpha L}, \sqrt{n}L\}$ [23]	$\frac{\Delta_0}{\epsilon^2} \min\{n^{3/4} \sqrt{\alpha L}, \sqrt{n}L\}$
f_i 's are α -weakly CVX and L -smooth	Natasha [1]	$\frac{\Delta_0}{\epsilon^2} \min\{\sqrt{n\alpha L}, L\}$ [23]	$\frac{\Delta_0}{\epsilon^2} \min\{\sqrt{n\alpha L}, \sqrt{n}L\}$

- Remarks:**
- Complexity (nonCVX f): total number of stochastic first-order oracle calls to find $\hat{\mathbf{x}}$ with $\mathbb{E}[\|\nabla f(\hat{\mathbf{x}})\|^2] \leq \epsilon^2$
 - Complexity ((S)CVX f): total number of stochastic first-order oracle calls to find $\hat{\mathbf{x}}$ with $\mathbb{E}[f(\hat{\mathbf{x}}) - f(\mathbf{x}^*)] \leq \epsilon$
 - $\Delta_0 = f(\mathbf{x}^0) - f^*$, $D_0 = \|\mathbf{x}^0 - \mathbf{x}^*\|$
 - Bounded variance: $\mathbb{E}_i[\|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|^2] \leq \sigma^2 \quad \forall \mathbf{x}$
 - Average L -smooth: $\mathbb{E}_i[\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\|^2] \leq L^2 \|\mathbf{x} - \mathbf{y}\|^2 \quad \forall \mathbf{x}, \mathbf{y}$
 - $f(\mathbf{x})$ is α -weakly convex if $f(\mathbf{x}) + \frac{\alpha}{2} \|\mathbf{x}\|^2$ is convex $\forall \mathbf{x}$.

Wrap up!

- Please finalize Homework 1 on Friday!
- Deep learning next week!

* Calculation of $d(\mathcal{D}_f(\mathbf{x}^{\natural}))$ and $d(\mathcal{D}_{f_\mu}(\mathbf{x}^{\natural}))$

Lemma ([3])

Let f be a proper lower-semicontinuous convex function, and let $\mathbf{x} \in \text{dom}(f)$. We have

$$d(\mathcal{D}_f(\mathbf{x})) \leq \inf_{\tau > 0} \mathbb{E} [\text{dist}^2(\mathbf{g}, \tau \partial f(\mathbf{x}))],$$

where \mathbf{g} is a vector of i.i.d. standard Gaussian random variables.

The upper bounds on $d(\mathcal{D}_f(\mathbf{x}^{\natural}))$ and $d(\mathcal{D}_{f_\mu}(\mathbf{x}^{\natural}))$ can be derived based on above.

Proposition

Let \mathbf{x}^{\natural} be an s -sparse vector. We have

$$d(\mathcal{D}_{f_\mu}(\mathbf{x}^{\natural})) \leq \inf_{\tau > 0} \left\{ s(1 + \tau^2) + 2\mu f_\mu(\mathbf{x}^{\natural})\tau^2 + (p - s) \sqrt{\frac{2}{\pi}} \int_{\tau}^{\infty} (u - \tau)^2 e^{-u^2/2} du \right\}.$$

Note that $f = f_\mu|_{\mu=0}$.

*Variance reduction techniques: SAGA

- Select the stochastic gradient \mathbf{r}_k as

$$\mathbf{r}_k = \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}_{i_k}^k) + \frac{1}{n} \sum_{j=1}^n \nabla f_j(\tilde{\mathbf{x}}_j^k),$$

where, at each iteration, $\tilde{\mathbf{x}}$ is updated as $\tilde{\mathbf{x}}_{i_k}^k = \mathbf{x}^k$ and $\tilde{\mathbf{x}}_j^k$ stays the same for $j \neq i_k$.

- As $\tilde{\mathbf{x}}_j^k \rightarrow \mathbf{x}^*$ and $\mathbf{x}^k \rightarrow \mathbf{x}^*$,

$$\nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}_{i_k}^k) + \frac{1}{n} \sum_{j=1}^n \nabla f_j(\tilde{\mathbf{x}}_j^k) \rightarrow 0.$$

- Therefore,

$$\mathbb{E} \left[\left\| \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}_{i_k}^k) + \frac{1}{n} \sum_{j=1}^n \nabla f_j(\tilde{\mathbf{x}}_j^k) \right\|^2 \right] \rightarrow 0.$$

*Variance reduction techniques: SAGA

Stochastic Average Gradient (SAGA) [10]

- 1a.** Choose $\tilde{\mathbf{x}}_i^0 = \mathbf{x}^0 \in \mathbb{R}^p, \forall i, q \in \mathbb{N}_+$ and stepsize $\gamma > 0$.
- 1b.** Store $\nabla f_i(\tilde{\mathbf{x}}_i^0)$ in a table data-structure with length n .
- 2.** For $k = 0, 1 \dots$ perform:
 - 2a.** Pick $i_k \in \{1, \dots, n\}$ uniformly at random
 - 2b.** Take $\tilde{\mathbf{x}}_{i_k}^{k+1} = \mathbf{x}^k$, store $\nabla f_{i_k}(\tilde{\mathbf{x}}_{i_k}^{k+1})$ in the table and leave other entries the same.
 - 2c.** $\mathbf{r}_k = \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}_{i_k}^k) + \frac{1}{n} \sum_{j=1}^n \nabla f_j(\tilde{\mathbf{x}}_j^k)$
- 3.** $\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma \mathbf{r}_k$

Recipe:

In each iteration:

- ▶ Store last gradient evaluated at each datapoint.
- ▶ Previous gradient for datapoint j is $\nabla f_j(\tilde{\mathbf{x}}_j^k)$.
- ▶ Perform SG-iterations with the following stochastic gradient

$$\mathbf{r}_k = \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}_{i_k}^k) + \frac{1}{n} \sum_{j=1}^n \nabla f_j(\tilde{\mathbf{x}}_j^k).$$

*Convergence of SAGA

$$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) := \frac{1}{n} \sum_{j=1}^n f_j(\mathbf{x}) \right\}.$$

Theorem (Convergence of SAGA [10])

Suppose that f is μ -strongly convex and that the stepsize is $\gamma = \frac{1}{2(\mu n + L)}$ with

$$\rho = 1 - \frac{\mu}{2(\mu n + L)} < 1,$$

$$C = \|\mathbf{x}^0 - \mathbf{x}^*\|^2 + \frac{n}{\mu n + L} [f(\mathbf{x}^0) - \langle \nabla f(\mathbf{x}^*), \mathbf{x}^0 - \mathbf{x}^* \rangle - f(\mathbf{x}^*)]$$

Then

$$\mathbb{E}[\|\mathbf{x}^k - \mathbf{x}^*\|^2] \leq \rho^k C.$$

- Allows the constant step-size.
- Obtains linear rate convergence.

*Variance reduction techniques: SARAH

- o Select the stochastic gradient \mathbf{r}_k

$$\mathbf{r}_k = \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\mathbf{x}^{k-1}) + \mathbf{r}_{k-1},$$

- o The variance reduction in SARAH can be characterized as

$$\mathbb{E}[\|\mathbf{r}_k\|^2] \leq \left[1 - \left(\frac{2}{\gamma L} - 1\right)\mu^2\gamma^2\right]^k \mathbb{E}[\|\nabla f(\mathbf{x}^0)\|^2].$$

*Variance reduction techniques: SARAH

Stochastic Recursive Gradient Algorithm (SARAH) [16]

1. Choose $\bar{\mathbf{x}}^0 \in \mathbb{R}^p$, $q \in \mathbb{N}_+$ and stepsize $\gamma > 0$.
2. For $k = 0, 1 \dots$ perform:
 2. $\mathbf{x}^0 = \bar{\mathbf{x}}^k$, $\mathbf{r}_0 = \frac{1}{n} \sum_{j=1}^n f_j(\bar{\mathbf{x}}^0)$
 - 2a. $\mathbf{x}^1 = \mathbf{x}^0 - \gamma \mathbf{r}_0$
 - 2b. For $l = 1 \dots, q - 1$, perform:
$$\begin{cases} \text{pick } i_l \in \{1, \dots, n\} \text{ uniformly at random,} \\ \mathbf{r}_l = \nabla f_{i_l}(\mathbf{x}^l) - \nabla f_{i_l}(\mathbf{x}^{l-1}) + \mathbf{r}_{l-1}, \\ \mathbf{x}^{l+1} = \mathbf{x}^l - \gamma \mathbf{r}_l. \end{cases}$$
- 3 Update $\bar{\mathbf{x}}^{k+1} = \mathbf{x}^l$ where l is chosen uniformly at random from $\{0, \dots, q\}$.

Recipe: *In a cycle of q inner iterations:*

- o Compute stochastic step direction by recursively adding and subtracting component gradients.

$$\mathbf{r}_l = \nabla f_{i_l}(\mathbf{x}^l) - \nabla f_{i_l}(\mathbf{x}^{l-1}) + \mathbf{r}_{l-1}.$$

- o Perform q SG-iterations with \mathbf{r}_l .
- o Update next iteration by picking uniformly at random from q previous iterations.

*Convergence of SARAH

$$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) := \frac{1}{n} \sum_{j=1}^n f_j(\mathbf{x}) \right\}.$$

Theorem (Convergence of SARAH [16])

Suppose that f is μ -strongly convex and that the stepsize γ and number of inner iterations q satisfies

$$\rho_q = \frac{1}{\mu\gamma(1+q)} + \frac{L_{\max}\gamma}{2 - L_{\max}\gamma} < 1.$$

Then

$$\mathbb{E}[\|\nabla f(\bar{\mathbf{x}}^k)\|^2] \leq \rho_q^k \|\nabla f(\bar{\mathbf{x}}^0)\|^2.$$

*An abridged variance reduction results for distributed optimization

Setting	Algorithm	Complexity bound
L -smooth f_i 's	Gradient descent	$mL\Delta_0/\epsilon^2$
	SVRG ($B_k = \Omega(n^{2/3})$) [18]	$n + n^{2/3}L\Delta_0/\epsilon^2$
	SpiderBoost [20]	$n + \sqrt{n}L\Delta_0/\epsilon^2$
	Spider [11]	$n + \sqrt{n}L\Delta_0/\epsilon^2$
	SARAH [16]	$n + \sqrt{n}L\Delta_0/\epsilon^2$
	PAGE [14]	$n + \sqrt{n}L\Delta_0/\epsilon^2$
	ZeroSARAH [15] ($B_0 = n$ and then $B_k = \sqrt{n}$)	$n + \sqrt{n}L\Delta_0/\epsilon^2$
Distributed with L -smooth $f_{i,j}$'s	Gradient descent	$mL\Delta_0/\epsilon^2$
	SCAFFOLD [13]	$m + \frac{m}{n^{1/3}} \frac{L\Delta_0}{\epsilon^2}$
	Spider [11]	$m + \frac{\sqrt{m}}{\sqrt{n}} \frac{L\Delta_0}{\epsilon^2}$
	SARAH [16]	$m + \frac{\sqrt{m}}{\sqrt{n}} \frac{L\Delta_0}{\epsilon^2}$
	ZeroSARAH [15] ($B_0 = m$ and then $B_k = \sqrt{m}$)	$m + \frac{\sqrt{m}}{\sqrt{n}} \frac{L\Delta_0}{\epsilon^2}$

Distributed: $f_i(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m f_{i,j}(\mathbf{x})$ loss on client or device i with m data samples

References I

- [1] Zeyuan Allen-Zhu.
Natasha: Faster non-convex stochastic optimization via strongly non-convex parameter.
In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 89–97. PMLR, 06–11 Aug 2017.
- [2] Zeyuan Allen-Zhu.
Katyusha x: Simple momentum method for stochastic sum-of-nonconvex optimization.
In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 179–185. PMLR, 10–15 Jul 2018.
- [3] Dennis Amelunxen, Martin Lotz, Michael B. McCoy, and Joel A. Tropp.
Living on the edge: Phase transitions in convex programs with random data.
2014.
arXiv:1303.6672v2 [cs.IT].
- [4] Peter L. Bartlett and Shahar Mendelson.
Rademacher and Gaussian complexities: Risk bounds and structural results.
J. Mach. Learn. Res., 3, 2002.
- [5] Léon Bottou and Oliver Bousquet.
The tradeoffs of large scale learning.
In *Advances in Neural Information Processing Systems*, 2007.

References II

- [6] John J Bruer, Joel A Tropp, Volkan Cevher, and Stephen Becker.
Time–data tradeoffs by aggressive smoothing.
In Advances in Neural Information Processing Systems, pages 1664–1672, 2014.
- [7] John J Bruer, Joel A Tropp, Volkan Cevher, and Stephen R Becker.
Designing statistical estimators that balance sample size, risk, and computational cost.
IEEE Journal of Selected Topics in Signal Processing, 9(4):612–624, 2015.
- [8] Venkat Chandrasekaran and Michael I. Jordan.
Computational and statistical tradeoffs via convex relaxation.
Proc. Natl. Acad. Sci., 110(13):E1181–E1190, 2013.
- [9] Venkat Chandrasekaran, Benjamin Recht, Pablo A. Parrilo, and Alan S. Willsky.
The convex geometry of linear inverse problems.
Found. Comput. Math., 12:805–849, 2012.
- [10] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien.
Saga: A fast incremental gradient method with support for non-strongly convex composite objectives.
In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, Advances in Neural Information Processing Systems 27, pages 1646–1654. Curran Associates, Inc., 2014.

References III

- [11] Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang.
SPIDER: near-optimal non-convex optimization via stochastic path-integrated differential estimator.
In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 687–697, 2018.
- [12] Rie Johnson and Tong Zhang.
Accelerating stochastic gradient descent using predictive variance reduction.
In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 315–323. Curran Associates, Inc., 2013.
- [13] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh.
Scaffold: Stochastic controlled averaging for federated learning.
In *International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020.
- [14] Zhize Li, Hongyan Bao, Xiangliang Zhang, and Peter Richtarik.
Page: A simple and optimal probabilistic gradient estimator for nonconvex optimization.
In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 6286–6295. PMLR, 18–24 Jul 2021.

References IV

- [15] Zhize Li and Peter Richtárik.
ZeroSARAH: Efficient nonconvex finite-sum optimization with zero full gradient computation.
arXiv preprint arXiv:2103.01447, 2021.
- [16] Lam M. Nguyen, Jie Liu, Katya Scheinberg, and Martin Takac.
Sarah: A novel method for machine learning problems using stochastic recursive gradient, 2017.
- [17] Samet Oymak, Christos Thrampoulidis, and Babak Hassibi.
Simple bounds for noisy linear inverse problems with exact side information.
2013.
arXiv:1312.0641v2 [cs.IT].
- [18] Sashank J Reddi, Suvrit Sra, Barnabás Póczos, and Alex Smola.
Stochastic frank-wolfe methods for nonconvex optimization.
arXiv preprint arXiv:1607.08254, 2016.
- [19] Shai Shalev-Shwartz and Nathan Srebro.
Svm optimization: inverse dependence on training set size.
In *Proceedings of the 25th international conference on Machine learning*, pages 928–935, 2008.

References V

- [20] Zhe Wang, Kaiyi Ji, Yi Zhou, Yingbin Liang, and Vahid Tarokh.
Spiderboost and momentum: Faster stochastic variance reduction algorithms.
In Advances in Neural Information Processing Systems, 2019.
- [21] Lin Xiao and Tong Zhang.
A proximal stochastic gradient method with progressive variance reduction.
SIAM Journal on Optimization, 24, 03 2014.
- [22] Guangzeng Xie, Luo Luo, and Zhihua Zhang.
A general analysis framework of lower complexity bounds for finite-sum optimization, 2019.
- [23] Dongruo Zhou and Quanquan Gu.
Lower bounds for smooth nonconvex finite-sum optimization.
In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, Proceedings of the 36th International Conference on Machine Learning, volume 97 of *Proceedings of Machine Learning Research*, pages 7574–7583. PMLR, 09–15 Jun 2019.