

Randomized Algorithms

Zhifeng Wang

Department of Statistics
Florida State University

Outline

- ▶ (Block) Coordinate Descent
 - Selecting Coordinates
 - Rates of Convergence
- ▶ Stochastic Gradient Descent
 - Stochastic Optimization
 - Empirical Risk Minimization
 - Momentum and Adaptive Learning Rates
 - Incremental Gradient Methods

Coordinate Descent

- ▶ Our goal is to solve

$$\min_{\beta \in \mathbb{R}^p} f(\beta),$$

where f is convex and smooth (f is continuously differentiable and gradient is Lipschitz continuous).

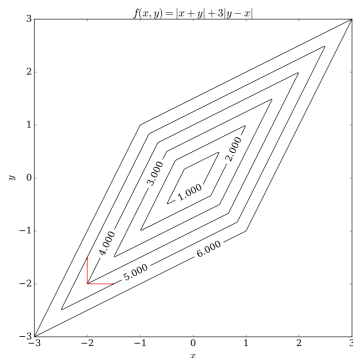
- ▶ The general idea for a coordinate descent algorithm is

Algorithm 1 Coordinate Descent/Minimization

- 1: Given initial $\beta^{(0)} \in \mathbb{R}^p$, $t \leftarrow 0$
 - 2: **repeat**
 - 3: Choose $j \in \{1, 2, \dots, p\}$
 - 4: $\beta_j^{(t+1)} \leftarrow \arg \min_{\beta_j} f(\beta_j, \beta_{-j}^{(t)})$, $\beta_{-j}^{(t+1)} \leftarrow \beta_{-j}^{(t)}$
 - 5: $t \leftarrow t + 1$
 - 6: **until** some convergence criterion is met
-

Remarks

- ▶ The function values are non-increasing.
- ▶ If f is **strictly convex** and **smooth**, CD converges to the global minimum.
- ▶ If f is nonconvex or not even smooth, CD may not converge at all or get stuck at a non-stationary point.



Coordinate Gradient Descent

- ▶ A popular alternative is to do one-step gradient descent for $\min_{\beta_j} f(\beta_j, \boldsymbol{\beta}_{-j}^{(t)})$, i.e., minimizing a quadratic approximation of $f(\beta_j, \boldsymbol{\beta}_{-j}^{(t)})$.

Algorithm 2 Coordinate Gradient Descent

- 1: Given initial $\boldsymbol{\beta}^{(0)} \in \mathbb{R}^p$, $t \leftarrow 0$
 - 2: **repeat**
 - 3: Choose $j \in \{1, 2, \dots, p\}$
 - 4: $\boldsymbol{\beta}^{(t+1)} \leftarrow \boldsymbol{\beta}^{(t)} - \alpha_t \nabla_j f(\boldsymbol{\beta}^{(t)}) \mathbf{e}_j$ for some $\alpha_t > 0$
 - 5: $t \leftarrow t + 1$
 - 6: **until** some convergence criterion is met
-

BCD

- ▶ The framework can be generalized for block updates, which is usually known as **block coordinate descent**.
- ▶ In the case when we have two blocks, it simply reduces to the **alternating minimization**.

BCD Convergence

- ▶ Consider the following unconstrained optimization problems with block structure

$$\min_{\boldsymbol{\beta}=[\boldsymbol{\beta}_1^\top, \dots, \boldsymbol{\beta}_k^\top]^\top \in \mathbb{R}^{p_1} \times \dots \times \mathbb{R}^{p_k}} f(\boldsymbol{\beta}), \quad (1)$$

where $\sum_j p_j = p$.

- ▶ Suppose that f is continuously differentiable over the convex closed set \mathcal{X} , and for each j ,

$$g(\boldsymbol{\beta}') := f(\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_{j-1}, \boldsymbol{\beta}', \boldsymbol{\beta}_{j+1}, \dots, \boldsymbol{\beta}_k)$$

is **strictly convex**. Then every limit point of BCD iterates is a **stationary point** of (1).

Rules for selecting coordinates

- ▶ **Cyclic Order:** run all coordinates in cyclic order, i.e. $1 \rightarrow 2 \rightarrow \dots \rightarrow p$ at each iteration
- ▶ **Random Sampling:** Randomly select some coordinate to update (uniformly or according to some distribution)
- ▶ **Gauss-Southwell:** At each iteration, pick coordinate $j = \arg \max_{1 \leq j \leq p} |\nabla_j f(\boldsymbol{\beta}^{(t)})|$.
- ▶ **Random Permutation:** run cyclic order on a permuted index (sample without replacement)

Convergence Rates

- ▶ Most convergence-rate analysis in the literature focus on block coordinate gradient descent, so we introduce some rate results for BCGD.
- ▶ Decompose the identity matrix $I_{p \times p} = [\mathbf{U}_1 \ \mathbf{U}_2 \ \dots \ \mathbf{U}_k]$, where \mathbf{U}_j is a matrix of size of $p \times p_j$. Then

$$\boldsymbol{\beta} = \sum_j \mathbf{U}_j \boldsymbol{\beta}_j, \quad \boldsymbol{\beta}_j = \mathbf{U}_j^\top \boldsymbol{\beta}.$$

Assumptions

- ▶ f is convex and smooth, ∇f is L -Lipschitz continuous for some $L \geq 0$.
- ▶ $\nabla_j f(\boldsymbol{\beta})$ is L_j -Lipschitz continuous.
- ▶ Denote $L_{\max} = \max_j L_j$, $L_{\min} = \min_j L_j$. Then

$$L_j \leq L \leq \sum_j L_j \leq kL_{\max}.$$

GS-BCGD

- ▶ **Gauss-Southwell** Block Coordinate Gradient Descent
- ▶ At t -th iteration, pick up $j = \arg \max_j |\nabla f(\boldsymbol{\beta}^{(t)})|$,

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - \frac{1}{L} \mathbf{U}_j \nabla_j f(\boldsymbol{\beta}^{(t)})$$

- ▶ If f is convex, then

$$f(\boldsymbol{\beta}^{(T)}) - f(\boldsymbol{\beta}^o) \leq \frac{2Lk \|\boldsymbol{\beta}^{(0)} - \boldsymbol{\beta}^o\|_2^2}{T}, \quad T > 0$$

- ▶ If f is μ -strongly convex, then

$$f(\boldsymbol{\beta}^{(T)}) - f(\boldsymbol{\beta}^o) \leq \left(1 - \frac{\mu}{Lk}\right)^T (f(\boldsymbol{\beta}^{(0)}) - f(\boldsymbol{\beta}^o)), \quad T > 0$$

GS-BCGD

- ▶ The previous GS analysis under strong convexity loss too much when we use $\|\nabla_j f(\beta^{(t)})\|_2^2 \geq \|\nabla f(\beta^{(t)})\|/k$.
- ▶ A refined result can be derived if we measure the strong convexity in the ℓ_1 -norm, i.e.,

$$f(\gamma) \geq f(\beta) + \langle \nabla f(\beta), \gamma - \beta \rangle + \frac{\mu_1}{2} \|\gamma - \beta\|_1^2$$

- ▶ Minimizing both sides w.r.t. γ , we obtain $f(\beta^o) \geq f(\beta) - \frac{1}{2\mu_1} \|\nabla f(\beta)\|_\infty^2 \geq f(\beta) - \frac{1}{2\mu_1} \max_j \|\nabla_j f(\beta)\|_2^2$.
- ▶ Then

$$f(\beta^{(T)}) - f(\beta^o) \leq \left(1 - \frac{\mu_1}{L}\right)^T (f(\beta^{(0)}) - f(\beta^o))$$

- ▶ Faster than previous result if $\mu_1 > \mu/k$.

Randomized-BCGD

- ▶ **Randomized** Block Coordinate Gradient Descent
- ▶ At t -th iteration, pick up j *randomly* from $\{1, 2, \dots, k\}$,

$$\beta^{(t+1)} = \beta^{(t)} - \frac{1}{L} \mathbf{U}_j \nabla_j f(\beta^{(t)})$$

- ▶ If f is convex, then

$$\mathbb{E}[f(\beta^{(T)}) - f(\beta^o)] \leq \frac{2LkR^2(\beta^{(0)})}{T}, \quad T > 0,$$

where $R(\beta^{(0)}) = \max\{\|\beta - \beta^o\|_2 : f(\beta) \leq f(\beta^{(0)})\}$.

- ▶ If f is μ -strongly convex, then

$$\mathbb{E}[f(\beta^{(T)}) - f(\beta^o)] \leq \left(1 - \frac{\mu}{Lk}\right)^T (f(\beta^{(0)}) - f(\beta^o)), \quad T > 0$$

Randomized-BCGD

- ▶ How to improve the rate?
- ▶ Use larger step sizes: $\alpha_t = 1/L_{j_t}$
- ▶ Use non-uniform sampling distribution: $p_j = L_j / \sum_j L_j$
- ▶ These two ways give an improved convergence bound for convex f :

$$\mathbb{E}[f(\boldsymbol{\beta}^{(t)}) - f(\boldsymbol{\beta}^o)] \leq \frac{2 \sum_j L_j R^2(\boldsymbol{\beta}^{(0)})}{T}$$

Cyclic-BCGD

- ▶ **Cyclic** Block Coordinate Gradient Descent
- ▶ At t -th iteration, let $\beta_0^{(t)} = \beta^{(t)}$, for $j = 1, 2, \dots, k$, do

$$\beta_j^{(t)} = \beta_{j-1}^{(t)} - \frac{1}{L} \mathbf{U}_j \nabla_j f(\beta_{j-1}^{(t)}),$$

Set $\beta^{(t+1)} = \beta_k^{(t)}$.

- ▶ Note that the Cyclic-BCGD is a *deterministic* algorithm.
- ▶ The cost for Cyclic-BCGD is $\mathcal{O}(k)$ times larger than R-BCGD, since it needs to go through entire blocks.

Cyclic-BCGD

- ▶ If f is convex, then

$$f(\boldsymbol{\beta}^{(T)}) - f(\boldsymbol{\beta}^o) \leq \frac{4L(k+1)R^2(\boldsymbol{\beta}^{(0)})}{T}, \quad T > 0$$

- ▶ If a larger step sizes $\alpha_t = 1/L_{j_t}$ is adopted, the rate can be improved to

$$f(\boldsymbol{\beta}^{(T)}) - f(\boldsymbol{\beta}^o) \leq \frac{4L_{\max}(kL^2/L_{\min}^2 + 1)R^2(\boldsymbol{\beta}^{(0)})}{T}$$

- ▶ If f is μ -strongly convex, then

$$f(\boldsymbol{\beta}^{(T)}) - f(\boldsymbol{\beta}^o) \leq \left(1 - \frac{\mu}{2(k+1)L}\right)^T (f(\boldsymbol{\beta}^{(0)}) - f(\boldsymbol{\beta}^o)), \quad T > 0$$

CD v.s. GD

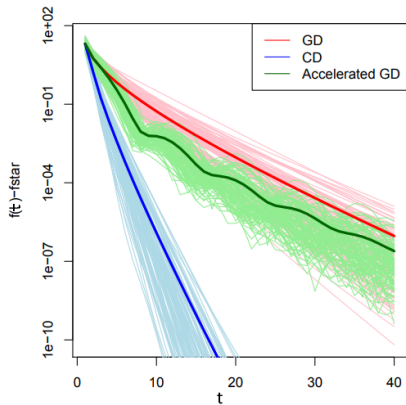


Figure: Linear regression with $n = 100, p = 20$

Summary

Algorithm	$\alpha_t = 1/L$	$\alpha_t = 1/L_{j_t}$
Gauss-Southwell	$\mathcal{O}(\frac{kL}{T})$	$\mathcal{O}(\frac{kL_{\max}}{T})$
Cyclic	$\mathcal{O}(\frac{kL}{T})$	$\mathcal{O}(\frac{L_{\max}(kL^2/L_{\min}^2+1)}{T})$
Randomized	$\mathcal{O}(\frac{kL}{T})$	$\mathcal{O}(\frac{\sum_j L_j}{T})$

Table: Convergence rates for three variants of BCGD

- ▶ Open questions: Are these bounds optimal? Under what condition can we drop the dependence of k ? When cyclic version is better than randomized version?

Deterministic v.s. Stochastic Optimization

- ▶ Deterministic optimization problem

$$\min_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n f_i(\beta). \quad (2)$$

- ▶ Stochastic optimization problem

$$\min_{\beta \in \mathbb{R}^p} \mathbb{E}_{\xi} [F(\beta, \xi)]. \quad (3)$$

where ξ is a random variable.

- ▶ (2) can be written as (3) if ξ is a uniform discrete random variable with $\mathbb{P}(\xi = i) = 1/n, i = 1, \dots, n$ and $F(\beta, \xi) = f_{\xi}(\beta)$.

Expected Risk

- ▶ In many supervised learning problems, input-output pairs $x, y \in \mathcal{X} \times \mathcal{Y}$ follow an unknown probability distribution $\mathcal{P}(x, y)$.
- ▶ The typical learning setting relies on
 - A loss function: $l : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$,
 $l(\hat{y}, y)$ measures the discrepancy between predicted \hat{y} and true y
 - A predictor function $\phi : \mathcal{X} \rightarrow \mathcal{Y}$.
 $\phi(x)$ maps an input x to the predicted output \hat{y} .

Expected Risk

- ▶ Given a loss function l , expected risk of a predictor function ϕ is defined by

$$\mathbf{E}(\phi) := \int l(\phi(x), y) d\mathcal{P}(x, y).$$

- ▶ The **ideal** goal:

$$\phi^*(x) := \arg \min_{\hat{y} \in \mathcal{Y}} \int l(\hat{y}, y) d\mathcal{P}(y|x).$$

- ▶ This is a typical *Stochastic Optimization* problem.
- ▶ We do not know the source distribution $\mathcal{P}(x, y)$. Instead, we have access to **samples** from the distribution. (*Sample Average Approximation*)

Empirical Risk

- ▶ Given a training dataset $\{(x_1, y_i)\}_{i=1}^n$ (which is assumed to be drawn iid from $\mathcal{P}(x, y)$)
- ▶ The empirical risk $\mathbf{E}_n(\phi)$ as a proxy to $\mathbf{E}(\phi)$:

$$\mathbf{E}_n(\phi) := \frac{1}{n} \sum_{i=1}^n l(\phi(x_i), y_i)$$

- ▶ Further restriction: ϕ belongs to a specific class \mathcal{F} , e.g., linear functions $\phi_w(x) = w^\top x$.
- ▶ \mathcal{F} together with the choice of l amounts to the choice of a specific machine learning technique.

Empirical Risk Minimization

- ▶ How to assess the quality of an empirical predictor $\hat{\phi}^n$?
- ▶ Empirical minimizer $\phi_{\mathcal{F}}^n := \arg \min_{\phi \in \mathcal{F}} \mathbf{E}_n(\phi)$.
- ▶ Expected minimizer $\phi_{\mathcal{F}}^* := \arg \min_{\phi \in \mathcal{F}} \mathbf{E}(\phi)$.
- ▶ The goal of a machine learner is to minimize the **excess error** $\mathbb{E}[\mathbf{E}(\hat{\phi}^n) - \mathbf{E}(\phi^*)]$.

- ▶
$$\mathbb{E}[\mathbf{E}(\hat{\phi}^n) - \mathbf{E}(\phi^*)] = \underbrace{\mathbb{E}[(\mathbf{E}(\phi_{\mathcal{F}}^*) - \mathbf{E}(\phi^*))]}_{\text{approximation error}} + \underbrace{\mathbb{E}[\mathbf{E}(\phi_{\mathcal{F}}^n) - \mathbf{E}(\phi_{\mathcal{F}}^*)]}_{\text{estimation error}} + \underbrace{\mathbb{E}[\mathbf{E}(\hat{\phi}^n) - \mathbf{E}(\phi_{\mathcal{F}}^n)]}_{\text{optimization error}}$$
$$\underbrace{\hspace{15em}}_{\mathbb{E}[\mathbf{E}(\hat{\phi}^n) - \mathbf{E}(\phi_{\mathcal{F}}^*)]}$$

- ▶ $\phi_{\mathcal{F}}^*$ reduces to $\phi_{\mathcal{F}}^n$ if $\mathcal{P}(x, y)$ is a uniform discrete distribution.

Stochastic Gradient Descent

- ▶ Originally designed for stochastic optimization, but still applicable for deterministic optimization (finite n).
- ▶ Recall the **finite sum problem**

$$\min_{\beta \in \mathbb{R}^p} f(\beta) := \frac{1}{n} \sum_{i=1}^n f_i(\beta)$$

- ▶ SGD update: $\beta^{(t+1)} = \beta^{(t)} - \alpha_t \underbrace{\nabla f_{i_t}(\beta^{(t)})}_{\text{stoch. grad.}}$
- ▶ Advantage: can deal with **huge** data and **streaming** data (online scenario)

Stochastic Gradient Descent

- ▶ **Unbiasedness:** the stochastic gradient is unbiased, i.e., $\mathbb{E}[\nabla f_{i_t}(\beta^{(t)})] = \nabla f(\beta^{(t)})$.
- ▶ **Decreasing step sizes:** since the stochastic gradient is random, we cannot guarantee $\nabla f_i(\beta^*) = 0, \forall i$. Hence, we need $\alpha_t \rightarrow 0$ as $t \rightarrow \infty$.
- ▶ **Convergence measure:** Due to the randomness, we cannot use previous error functions to measure the optimality. Instead, consider the **expectation**.

R-CD: A Special Case of SGD

- ▶ Randomized Coordinate Descent can be view as a special case of SGD, where $g^{(t)} = p \nabla_{j_t} f(\beta^{(t)}) e_{j_t}$.
- ▶ $\mathbb{E}[g^{(t)}] = \frac{1}{p} \sum_{j=1}^p p \nabla_j f(\beta^{(t)}) e_j = \nabla f(\beta^{(t)})$ certifies **unbiasedness**.
- ▶ The relationship can be seen from **dual** problems.

Example

- ▶ Linear system $Aw = b$, where $A \in \mathbb{R}^{m \times n}$ and $\|A_i\|_2 = 1$.
- ▶ The least-norm solution is found by solving $\min_{w \in \mathbb{R}^n} \frac{1}{2} \|w\|_2^2$ s.t. $Aw = b$, and its **dual problem** is

$$\min_{z \in \mathbb{R}^m} \frac{1}{2} \|A^\top z\|_2^2 - b^\top z.$$

- ▶ Applying **CD** to the dual problem with $\alpha_t = 1$ gives $z^{t+1} = z^t - (A_i \cdot A^\top z^t - b_i) e_i$, which corresponds

$$w^{t+1} = w^t - (A_i \cdot w^t - b_i) A_i^\top. \quad (4)$$

- ▶ (4) can be seen as **SGD** applied to the *least square problem* $\min_w \frac{1}{2} \|Aw - b\|_2^2$.

Convergence Rate of SGD

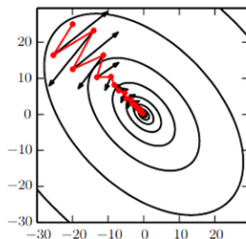
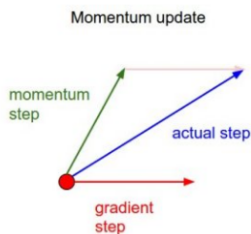
- ▶ Assume that $\mathbb{E}[\|\nabla f_i(\beta)\|_2^2] \leq G^2$. When $f(\beta)$ is μ -**strongly convex**, $\alpha_t = \mathcal{O}(\frac{1}{\mu t})$, we have

$$\mathbb{E}[f(\beta^{(T)}) - f(\beta^*)] \leq \frac{\max(\|\beta^{(1)} - \beta^*\|_2^2, G^2/\mu^2)}{T}$$

- ▶ The $\mathcal{O}(1/T)$ rate is **tight!** $\alpha_t = \mathcal{O}(\frac{1}{t})$ is theoretically best.
- ▶ For a general convex f , the rate of convergence is $\mathcal{O}(\frac{1}{\sqrt{T}})$.
- ▶ SGD is much slower than (accelerated) GD in terms of iteration number, but the per iteration cost is only $\mathcal{O}(1)$.

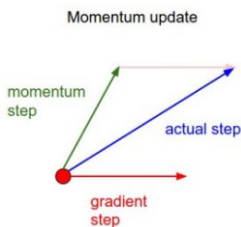
Momentum

- ▶ Momentum introduces a variable v that plays the role of **velocity**, which is set to an exponentially decaying **average** of the negative gradient.
- ▶ $v = \mu v - \text{learning_rate} * dx$
 $x += v$



Nesterov Momentum

- ▶ $x_ahead = x + \mu * v$
 $v = \mu * v - learning_rate * dx_ahead$
 $x += v$
- ▶ The difference between Nesterov momentum and standard momentum is where the gradient is evaluated.



- ▶ Very popular in practice, but **NO improvement** for the rate of convergence.

Adaptive Learning Rates

- ▶ **AdaGrad**

```
cache += dx**2
```

```
x += - learning_rate * dx / (np.sqrt(cache)+eps)
```

- ▶ **RMSprop**

```
cache = decay_rate*cache + (1-decay_rate)*dx**2
```

```
x += - learning_rate * dx / (np.sqrt(cache)+eps)
```

- ▶ **Adam** (simplified version)

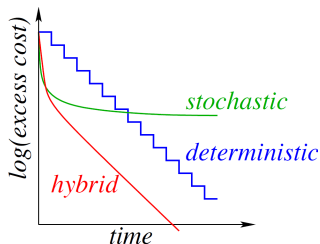
```
m = beta1*m + (1-beta1)*dx
```

```
v = beta2*v + (1-beta2)*(dx**2)
```

```
x += - learning_rate * m / (np.sqrt(v) + eps)
```

Incremental Gradient Methods

- ▶ GD enjoys a linear rate but $\mathcal{O}(n)$ iteration cost.
- ▶ SGD enjoys a sublinear rate but with $\mathcal{O}(1)$ iteration cost.



- ▶ For **finite** n , can we design a better method?
Incremental Gradient Methods (**SAG**, **SAGA**, **SVRG**, **S2GD**, ...)

Variance Reduction

- ▶ In **expectation** both GD and SGD are identical, i.e.,

$$\mathbb{E}[\beta^{(t)}|\beta^{(t-1)}] = \beta^{(t-1)} - \alpha_t \left[\frac{1}{n} \sum_{i=1}^n \nabla f_i(\beta^{(t-1)}) \right].$$

- ▶ But the SG equals $\nabla f(\beta)$ only in expectation and may have **large deviation** from the mean.
- ▶ We observe $\mathcal{O}(1/T)$ rate due to the inherent large variance.

Batching

- ▶ The simplest variant is to use **large** sample.
- ▶ For a fixed batch size B , the rate is still **sublinear**.
Another way is to grow B gradually.
- ▶ Growing B eventually requires $\mathcal{O}(n)$ iteration cost.
- ▶ Can we have 1 gradient per iteration but **linear** convergence rate?
Yes! Stochastic Average Gradient

Variance Reduction Technique

- ▶ Suppose we want to estimate $\Theta = \mathbb{E}[X]$, and we have access to a random variable Y which is **highly correlated** with X , and we can compute $\mathbb{E}[Y]$ easily.
- ▶ Consider the point estimator $\hat{\Theta}_a$ with $a \in [0, 1]$

$$\hat{\Theta}_a = a(X - Y) + \mathbb{E}(Y).$$

- ▶ $\mathbb{E}[\hat{\Theta}_a] = a\mathbb{E}[X] + (1 - a)\mathbb{E}[Y]$
 $\text{Var}[\hat{\Theta}_a] = a^2(\text{Var}[X] + \text{Var}[Y] - 2\text{Cov}[X, Y])$

Variance Reduction Techniques

$$\hat{\Theta}_a = a(X - Y) + \mathbb{E}(Y).$$

$$\mathbb{E}[\hat{\Theta}_a] = a\mathbb{E}[X] + (1 - a)\mathbb{E}[Y]$$

$$\text{Var}[\hat{\Theta}_a] = a^2(\text{Var}[X] + \text{Var}[Y] - 2\text{Cov}[X, Y])$$

- ▶ If $\text{Cov}[X, Y]$ is sufficiently large, then $\text{Var}(\hat{\Theta}_a) < \text{Var}[X]$. $\hat{\Theta}_a$ has **smaller variance** than the direct estimator X
- ▶ $a = 1$, $\hat{\Theta}_a$ is an **unbiased** estimator of X .
- ▶ $a = 0$, $\hat{\Theta}_a = \mathbb{E}[Y]$, **zero variance** but could be biased.
- ▶ As a increase from 0 to 1, the bias *decreases* and the variance *increases*.

Stochastic Average Gradient

- ▶ $\beta^{(t+1)} = \beta^{(t)} - \frac{\alpha_t}{n} \sum_i^n g_i^t$
- ▶ $g_i^t = \nabla f_i(\beta^t)$ from the **last** t where i was selected.
- ▶ **Memory:**

$$\begin{bmatrix} - & - & - & - & - & g_1 & - & - & - & - & - \\ - & - & - & - & - & g_2 & - & - & - & - & - \\ & & & & & \vdots & & & & & \\ - & - & - & - & - & g_n & - & - & - & - & - \end{bmatrix},$$

- ▶ Each g_i keeps track of the last time we randomly picked sample i . **Update** g_i if i is selected in the next iteration.
- ▶ Then take a step in the direction of the **average** of g_i .

Stochastic Average Gradient

- ▶ SAG can be rewritten as

$$\beta^{(t+1)} = \beta^{(t)} - \alpha \underbrace{\left[\underbrace{\frac{1}{n}}_a \left(\underbrace{\nabla f_i(\beta^{(t)})}_X - \underbrace{g_i^{t-1}}_Y \right) + \underbrace{\frac{1}{n} \sum_{i=1}^n g_i^{t-1}}_{\mathbb{E}[Y]} \right]}_{\hat{\Theta}_a}$$

- ▶ $\Theta = \nabla f(\beta^{(t)})$, $X = \nabla f_i(\beta^{(t)})$, $Y = g_i^{t-1}$, $a = 1/n$.
- ▶ SAGA: similar to SGA, only change $a = 1$.

Convergence Rate of SAG

- ▶ If f_i is L -smooth and f is μ -strongly convex with $\alpha = 1/16L$, SAG has

$$\mathbb{E}[f(\beta^{(T)}) - f(\beta^*)] \leq \mathcal{O}\left(\left[1 - \min\left(\frac{\mu}{16L}, \frac{1}{8n}\right)\right]^T\right)$$

- ▶ High memory cost: $\mathcal{O}(n)$!

Stochastic Variance Reduced Gradient

- ▶ SVRG: gets rid of memory by occasionally computing exact gradient.

Algorithm 3 SVRG

- 1: Initialize $\tilde{\beta}^{(0)}$
 - 2: **for** $s = 1, 2, \dots$ **do**
 - 3: $\tilde{\theta} = \sum_{i=1}^n \nabla f_i(\tilde{\beta}^{(s-1)})$
 - 4: $\beta^{(0)} = \tilde{\beta}^{(s-1)}$
 - 5: **for** $t = 1, 2, \dots, m$ **do**
 - 6: Randomly pick $i \in \{1, 2, \dots, n\}$
 - 7: $\beta^{(t)} = \beta^{(t-1)} - \alpha[\nabla f_i(\beta^{(t-1)}) - \nabla f_i(\tilde{\beta}^{(s-1)}) + \tilde{\theta}]$
 - 8: **end for**
 - 9: $\tilde{\beta}^{(s)} = \beta^{(m)}$
 - 10: **end for**
-

Stochastic Variance Reduced Gradient

- ▶ At the s -th epoch, SVRG updates

$$\beta^{(t+1)} = \beta^{(t)} - \alpha \underbrace{\left[\underbrace{\nabla f_i(\beta^{(t)})}_X - \underbrace{\nabla f_i(\tilde{\beta}^{(s-1)})}_Y + \underbrace{\sum_{i=1}^n \nabla f_i(\tilde{\beta}^{(s-1)})}_{\mathbb{E}[Y]} \right]}_{\hat{\Theta}_a}$$

- ▶ $\Theta = \nabla f(\beta^{(t)})$, $X = \nabla f_i(\beta^{(t)})$, $Y = \nabla f_i(\tilde{\beta}^{(s-1)})$, $a = 1$.

Convergence Rate of SVRG

- ▶ Assume f_i is L -smooth and f is μ -strongly convex. If m is sufficiently large and $\alpha < 1/2L$, SVRG enjoys

$$\mathbb{E}[f(\tilde{\beta}^{(S)}) - f(\beta^*)] \leq \mathcal{O}\left(\left[\frac{1}{\mu\alpha(1-2L\alpha)m} + \frac{2L\alpha}{1-2L\alpha}\right]^S\right)$$

- ▶ Low memory cost, but require **two** loops.
- ▶ S2GD is same as SVRG, but with *random* number of inner steps.

Algorithm Comparison

- ▶ Logistic regression:

$$\min_{\beta} f(\beta) := \frac{1}{n} \sum_{i=1}^n \log(1 + \exp -y_i \beta^\top x_i) + \frac{\lambda}{2} \|\beta\|_2^2$$

- ▶ $X \in \mathbb{R}^{n \times p} \sim N(0, I)$, $n = 10000$, $p = 50$, $\lambda = 10^{-4}$. The columns of X are normalized. y_i is generated according to the logistic model.
- ▶ $L = \lambda + \|x_i\|_2^2/4 = \lambda + 1/4$, $\mu = \lambda$.
- ▶ GD: $\alpha = 2/(L + \mu)$
- ▶ SGD: $\alpha_t = 0.1/(\mu(t + 100))$, $B = 50$
- ▶ SVRG: $\alpha = 1/L$, $B = 50$, $m = 3n$

Algorithm Comparison

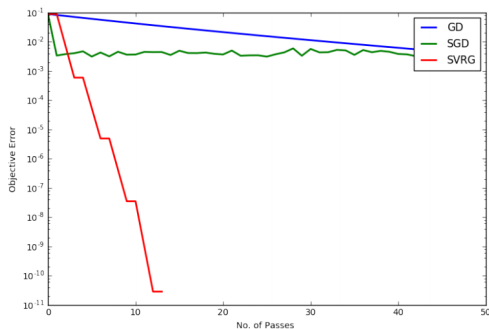


Figure: Comparison of GD, SGD, and SVRG on a synthetic dataset