Takeaways:
1. Template of <u>Variance Reduction</u> analysis:
2. "Warm start": a good starting point matters ✓ (engineering trick)
3. Trade-offs (amount of computation VS performance) and how to benefit from them
   A. switching
   B. mixing (if time allows)

minimize $F(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x)$
$x \in \mathbb{R}^d$

Lecture notation $F(w,b) = \frac{1}{n} \sum_{i=1}^{n} f_i(w,b) + \Psi(w)$

depend on $i^{th}$ sample $(x_i, y_i)$

Vector X: $x^t$ for iterates, $x_i$ for component
Scalar X: $\alpha_t$ for iterates

GD: $x^t = x^{t-1} - \alpha_t \nabla F(x^{t-1})$

SGD: for $t = 1, 2, \cdots$     $t > n$
   $i_t \leftarrow \text{Unif} \{1, 2, \cdots, n\}$
   $x^t = x^{t-1} - \alpha_t \nabla f_{i_t}(x^{t-1})$

$\mathbb{E}[\nabla f_{i_t}(x^{t-1})] = \nabla F(x^{t-1})$    ?

$\mathbb{E}[\nabla f_{i_t}(x)] = \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(x) = \nabla F(x^{t-1})$

$1^{st}$ order   $\nabla F$

$F = (f_1 + f_2) \cdot \frac{1}{2}$
$x^*$   $\nabla f_1(x^*) = 1$
    $\nabla f_2(x^*) = -1$

$\alpha_t \simeq$ constant

$\frac{1}{t}$

$\frac{1}{\sqrt{t}}$

VR template:
   $\nabla f_{i_t}(x)$       $\nabla F(x)$

   Consider an estimator $X$ for parameter $\theta \in \mathbb{R}^n$

   $X$ is Unbiased if $\mathbb{E}X = \theta$

   $\text{Var}(X) := \mathbb{E}[(X - \theta)^2]$

$Z := X - Y$. $X$ unbiased, $\mathbb{E}Y = 0 \Rightarrow$ $\color{blue}{\mathbb{E}Z = \mathbb{E}X = \theta}$ $Z$ is unbiased

"Recall" $\text{Cov}(X, Y) = \mathbb{E}\left[(X - \mathbb{E}X)(Y - \mathbb{E}Y)\right]$

**Fact:** $\text{Var}(Z) = \text{Var}(X) + \text{Var}(Y) - 2\text{Cov}(X, Y)$ ☆

$\text{Var}(Z) \leq \text{Var}(X) \iff \text{Cov}(X, Y) \geq \frac{1}{2}\text{Var}(Y) \color{red}{\geq 0}$

$\color{red}{X, Y \text{ "positively correlated"}}$

$\color{blue}{\text{Example (midterm)}}$

estimator $X$ for parameter $\theta$. $\mathbb{E}X = \theta = 1$

$Z := X - Y$. $\text{Var}(Z) \leq \text{Var}(X)$

$\text{Var}(X) = 0.5$, $\text{Var}(Y) = 0.1$, $\mathbb{E}[Y] = 0$

$\text{Cov}(X, Y) \geq \frac{1}{2}\text{Var}(Y) = 0.05$

---

# SAG (Stochastic averaged gradient, 2013)

Maintain a table of $g_i$ $\quad i = 1, 2, \ldots, n$ ($g$ for <u>g</u>radient)

Initialize $X^0$, $g_i^0 = \nabla f_i(X^0)$ $\forall i$

for $t = 1, 2, \ldots$ do

$\quad i_t \leftarrow \text{Unif}\{1, \ldots, n\}$

$$g_{i_t}^t = \nabla f_i(x^{t-1}), \quad g_i^t = g_i^{t-1} \text{ for } i \neq i_t$$

$$x^t = x^{t-1} - \alpha_t \frac{1}{n} \sum_{i=1}^n g_i^t$$

$$
\begin{array}{c|c|c|c|c|c}
g & 1 & 2 & 3 & \cdots & n \\
\hline
\text{Initialize} & \nabla f_1(x^0) & \nabla f_2(x^0) & \nabla f_3(x^0) & \cdots & \nabla f_n(x^0) \\
 & -//- & \nabla f_2(x^{t-1}) & -//- & & -//- \\
\end{array}
$$

$$g_i^0 := \nabla f_i(x^0)$$

Remark:

① Each $f_i$ contribute a part to gradient estimate

② $x^0$ — How to choose

③ $\frac{1}{n} \sum_{i=1}^n g_i^t$ — too much to average?

$$x^t = x^{t-1} - \alpha_t \left( \frac{1}{n} \sum_{i \neq i_t} g_i^t + \frac{1}{n} g_{i_t}^t \right)$$

$$= x^{t-1} - \alpha_t \left( \frac{1}{n} \sum_{i \neq i_t} g_i^{t-1} + \frac{1}{n} g_{i_t}^{t-1} - \frac{1}{n} g_{i_t}^{t-1} + \frac{1}{n} g_{i_t}^t \right)$$

$$= x^{t-1} - \alpha_t \left( \frac{1}{n} g_{it}^t - \frac{1}{n} g_{it}^{t-1} + \frac{1}{n} \sum_{i=1}^{n} g_i^{t-1} \right)$$

old average

$$\underbrace{g_{it}^t}_{\nabla f_{it}(x^t)} - \left( g_{it}^{t-1} - \sum_{i=1}^{n} g_i^{t-1} \right)$$

$$z := X - \underbrace{Y}$$

$$\mathbb{E} X = \nabla F(x^t) \quad - \quad \text{unbiased}$$

$$\mathbb{E} Y \neq 0 \quad \Rightarrow \quad z \text{ is biased}$$

$X, Y$ correlated

As $t \to \infty$, $X - Y \to 0$

$$g_{it}^t \approx g_{it}^{t-1}$$

$$\sum_{i=1}^{n} g_i \approx 0$$

SAGA (2014)

If $\mathbb{E}Y = 0$, $Z$ is unbiased

$$\underbrace{g_{i_t}^t}_{X} - \underbrace{\left(g_{i_t}^{t-1} - \frac{1}{n}\sum_{i=1}^{n} g_i^{t-1}\right)}_{Y}$$

$i_t \leftarrow \text{Unif}\{1, \cdots, n\}$

$\mathbb{E}Y = 0$ because $i_t$ chosen u.a.r from $\{1, 2, \cdots, n\}$.

Exercise: Var SAGA $n^2$ time larger than SAG

Note: SAGA still needs a careful choice of $x^0$

Choose $x^{00}$, run 1 epoch of SGD to give $x^0$

(2013)

## Stochastic variance reduction gradient (SVRG)

Template:

$$\underbrace{g_{i_t}^t}_{X} - \underbrace{\nabla f_{i_t}(x^{old}) + \nabla F(x^{old})}_{Y}$$

Exercise: ① $\mathbb{E} Y = 0$

② $X, Y$ positively correlated

$X - Y \to 0$, because . . .

③ Need to compute $\nabla F(x^{old})$ ocassionally

Algorithm: Initialize $X$

Until convergence do

$x^{old} \leftarrow$ newest iterate $X$

$x^o \leftarrow x^{old}$, Compute $\nabla F(x^{old})$

$$F = \frac{1}{n} \sum_{i=1}^{n} f_i$$

$$\text{for } t = 0, 1, \cdots, m-1 \text{ do}$$

$$i_t \leftarrow \text{Unif} \{1, \cdots, n\}$$

$$\hat{\nabla}^t = \nabla f_{i_t}(x^t) - \nabla f_{i_t}(x^{old}) + \nabla F(x^{old})$$

$$x^{t+1} = x^t - \eta \hat{\nabla}^t$$

Remark: ① Constant $\eta$

② Small $m \rightarrow$ Compute $\nabla F$ more often

③ $m \gg \eta$

# Computation    SVRG $\approx$ SGD

# SARAH:

Initialize $X$

Until convergence do

$x^{old} \leftarrow \text{newest iterate } x$          $\hat{\nabla}^0 = \nabla F(\text{newest } x)$

$x^0 \leftarrow x^{old}$, Compute $\nabla F(x^{old})$

$$\text{for } t = 0, 1, \cdots, m-1 \text{ do}$$

$$i_t \leftarrow \text{Unif} \{1, \cdots, n\}$$

$$\hat{\nabla}^t = \nabla f_{i_t}(x^t) - \nabla f_{i_t}(x^{t-1}) + \hat{\nabla}^{t-1}$$

$$x^{t+1} = x^t - \eta \hat{\nabla}^t$$

Note that I misspoke in the lecture - SARAH still computes full gradient, the same as SVRG. It just doesn't use it again and again in the inner loop update, but use it as an occasional correction
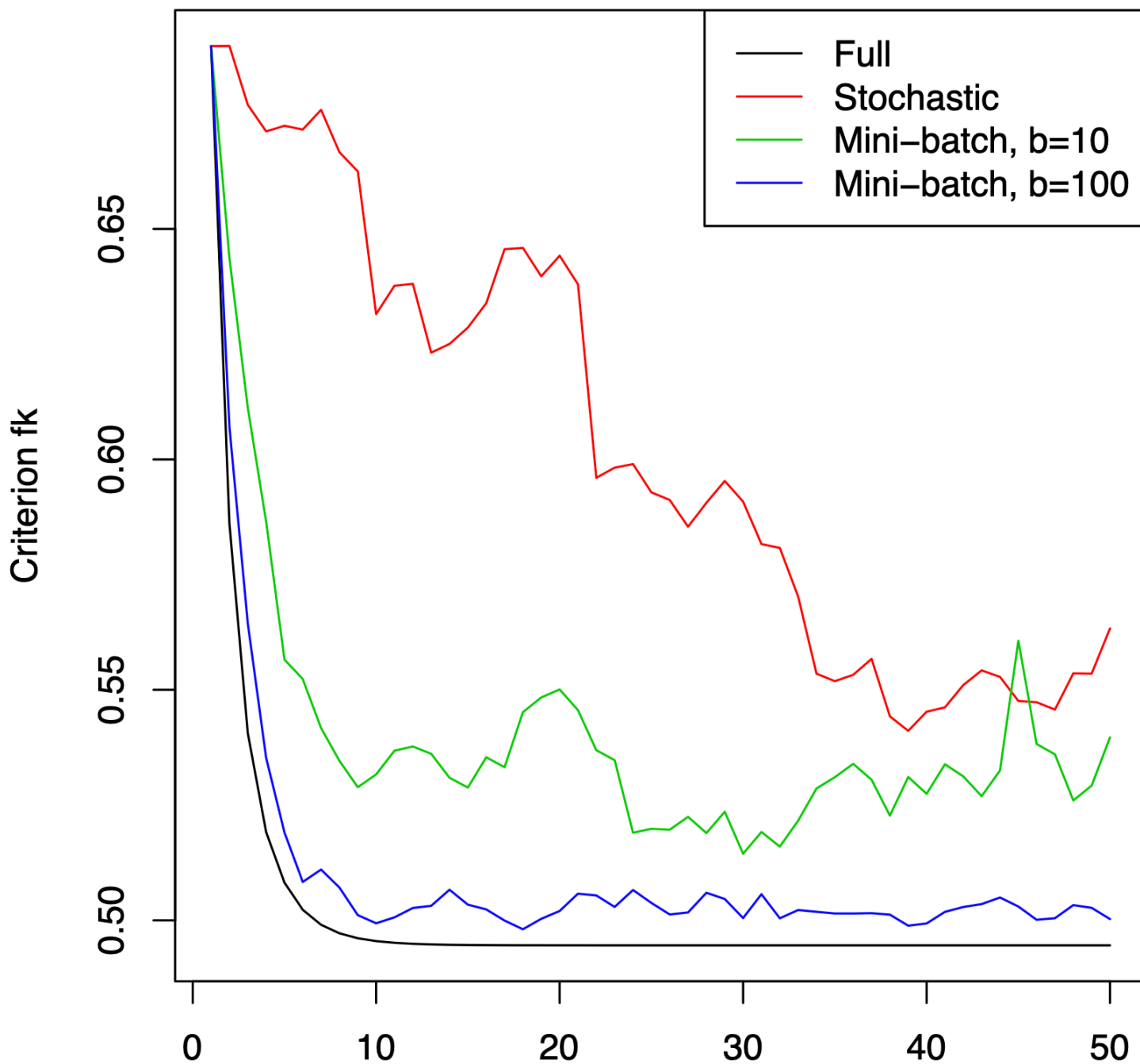
PAGE: Each step do                    (2020)

$b'$ - mini batch SARAH       w.p. $P_t$

$b$ - mini batch SGD         w.p. $1 - P_t$

$$P_t \equiv \frac{b}{b+b'}, \quad b' \ll b$$

The following pictures are from http://www.stat.cmu.edu/~ryantibs/convexopt/lectures/modern-sgd.pdf.

This note is prepared bas http://www.stat.cmu.edu/~ryantibs/convexopt/lectures/modern-sgd.pdf and http://www.princeton.edu/~yc5/ele522_optimization/lectures/variance_reduction.pdf