# Adaptive Sampling Strategies for Stochastic Optimization

Raghu (Vijaya Raghavendra) Bollapragada<sup>1</sup> Richard Byrd<sup>2</sup> Jorge Nocedal<sup>1</sup>

> <sup>1</sup>Northwestern University <sup>2</sup>University of Colorado Boulder

> > January 8, 2018

 $11^{th}~{\rm US}$  - Mexico Workshop on Optimization and its Applications Huatulco, Mexico

Northwestern ENGINEERING

## **Optimization Problem**

$$\min_{x\in\mathbb{R}^d}F(x)=\mathbb{E}_{\zeta}[f(x;\zeta)]$$

## **Optimization Problem**

$$\min_{x\in\mathbb{R}^d}F(x)=\mathbb{E}_{\zeta}[f(x;\zeta)]$$

Structural Risk Minimization

$$\min_{x\in\mathbb{R}^d}F(x)=\int f(x;z,y)dP(z,y)$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三日= のへで

$$\min_{x\in\mathbb{R}^d}F(x)=\mathbb{E}_{\zeta}[f(x;\zeta)]$$

$$\min_{x\in\mathbb{R}^d}F(x)=\int f(x;z,y)dP(z,y)$$

Empirical Risk Minimization

$$\min_{x\in\mathbb{R}^d}R(x)=\frac{1}{n}\sum_{i=1}^nF_i(x)$$

$$\min_{x\in\mathbb{R}^d}F(x)=\mathbb{E}_{\zeta}[f(x;\zeta)]$$

$$\min_{x\in\mathbb{R}^d}F(x)=\int f(x;z,y)dP(z,y)$$

Empirical Risk Minimization

$$\min_{x\in\mathbb{R}^d}R(x)=\frac{1}{n}\sum_{i=1}^nF_i(x)$$

Stochastic Gradient is a popular first order method for solving these problems

$$\min_{x\in\mathbb{R}^d}F(x)=\mathbb{E}_{\zeta}[f(x;\zeta)]$$

$$\min_{x\in\mathbb{R}^d}F(x)=\int f(x;z,y)dP(z,y)$$

Empirical Risk Minimization

$$\min_{x \in \mathbb{R}^d} R(x) = \frac{1}{n} \sum_{i=1}^n F_i(x)$$

- Stochastic Gradient is a popular first order method for solving these problems
- Many stochastic first order variance reduced methods have been proposed for finite sum problem
   SAG[Schmidt et al. 2016], SAGA[Defazio et al. 2014], SVRG[Johnson and Zhang 2013]

$$\min_{x\in\mathbb{R}^d}F(x)=\mathbb{E}_{\zeta}[f(x;\zeta)]$$

$$\min_{x\in\mathbb{R}^d}F(x)=\int f(x;z,y)dP(z,y)$$

Empirical Risk Minimization

$$\min_{\mathbf{x}\in\mathbb{R}^d}R(\mathbf{x})=\frac{1}{n}\sum_{i=1}^nF_i(\mathbf{x})$$

- Stochastic Gradient is a popular first order method for solving these problems
- Many stochastic first order variance reduced methods have been proposed for finite sum problem
   SAG[Schmidt et al. 2016], SAGA[Defazio et al. 2014], SVRG[Johnson and Zhang 2013]
- Require either storage or computation of full gradient

$$\min_{x\in\mathbb{R}^d}F(x)=\mathbb{E}_{\zeta}[f(x;\zeta)]$$

$$\min_{x\in\mathbb{R}^d}F(x)=\int f(x;z,y)dP(z,y)$$

Empirical Risk Minimization

$$\min_{x \in \mathbb{R}^d} R(x) = \frac{1}{n} \sum_{i=1}^n F_i(x)$$

- Stochastic Gradient is a popular first order method for solving these problems
- Many stochastic first order variance reduced methods have been proposed for finite sum problem
   SAG[Schmidt et al. 2016], SAGA[Defazio et al. 2014], SVRG[Johnson and Zhang 2013]
- Require either storage or computation of full gradient
- Achieve linear convergence for strongly convex functions

Raghu Bollapragada (NU)

三日 のへで

$$x_{k+1} = x_k - lpha_k 
abla F_{\mathcal{S}_k}(x_k), \qquad 
abla F_{\mathcal{S}_k}(x_k) = rac{1}{|\mathcal{S}_k|} \sum_{i \in \mathcal{S}_k} 
abla F_i(x_k),$$

where the set  $S_k \subset \{1, 2, \ldots\}$  indexes data points  $(y^i, z^i)$  drawn at random from the distribution P

$$x_{k+1} = x_k - \alpha_k \nabla F_{S_k}(x_k), \qquad \nabla F_{S_k}(x_k) = \frac{1}{|S_k|} \sum_{i \in S_k} \nabla F_i(x_k),$$

where the set  $S_k \subset \{1, 2, \ldots\}$  indexes data points  $(y^i, z^i)$  drawn at random from the distribution P

• Noise in the steps is controlled by sample sizes

$$x_{k+1} = x_k - \alpha_k \nabla F_{S_k}(x_k), \qquad \nabla F_{S_k}(x_k) = \frac{1}{|S_k|} \sum_{i \in S_k} \nabla F_i(x_k),$$

where the set  $S_k \subset \{1, 2, ...\}$  indexes data points  $(y^i, z^i)$  drawn at random from the distribution P

- Noise in the steps is controlled by sample sizes
- These methods can take advantage of parallel frameworks

$$x_{k+1} = x_k - \alpha_k \nabla F_{S_k}(x_k), \qquad \nabla F_{S_k}(x_k) = \frac{1}{|S_k|} \sum_{i \in S_k} \nabla F_i(x_k),$$

where the set  $S_k \subset \{1, 2, ...\}$  indexes data points  $(y^i, z^i)$  drawn at random from the distribution P

- Noise in the steps is controlled by sample sizes
- These methods can take advantage of parallel frameworks
- If sample sizes are increased at geometric rate, R-Linear convergence for strongly convex functions

[Byrd et al. 2012] [Friedlander and Schmidt 2012] [Pasupathy et al. 2015]

- Adaptive Sampling Tests
- 2 Convergence Analysis
- Practical Implementation
- 4 Numerical Experiments

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三日= のへで

5 Summary

## Overview

#### 1 Adaptive Sampling Tests

- Norm test
- Inner Product Test

#### Convergence Analysis

- Orthogonal test
- Linear Convergence

#### 3 Practical Implementation

• Step-Length Strategy

Parameter Selection

#### 4 Numerical Experiments

#### 5 Summary

### Norm Test

Raghu Bollapragada (NU)

#### $\|\nabla F_{\mathcal{S}_k}(x_k) - \nabla F(x_k)\| \le \theta_n \|\nabla F(x_k)\|, \quad \text{for some} \quad \theta_n \in [0, 1).$

$$\|\nabla F_{\mathcal{S}_k}(x_k) - \nabla F(x_k)\| \le \theta_n \|\nabla F(x_k)\|, \quad \text{for some} \quad \theta_n \in [0, 1).$$

$$\frac{\mathbb{E}[\|\nabla F_i(x_k) - \nabla F(x_k)\|^2]}{|S_k|} \leq \theta_n^2 \|\nabla F(x_k)\|^2$$

$$\|
abla F_{\mathcal{S}_k}(x_k) - 
abla F(x_k)\| \le heta_n \|
abla F(x_k)\|, \quad \text{for some} \quad heta_n \in [0, 1).$$

$$\frac{\mathbb{E}[\|\nabla F_i(x_k) - \nabla F(x_k)\|^2]}{|S_k|} \leq \theta_n^2 \|\nabla F(x_k)\|^2$$

• Byrd et al. [2012] proposed this test to control the sample sizes

$$\|
abla F_{\mathcal{S}_k}(x_k) - 
abla F(x_k)\| \le heta_n \|
abla F(x_k)\|, \quad \text{for some} \quad heta_n \in [0, 1).$$

$$\frac{\mathbb{E}[\|\nabla F_i(x_k) - \nabla F(x_k)\|^2]}{|S_k|} \leq \theta_n^2 \|\nabla F(x_k)\|^2$$

- Byrd et al. [2012] proposed this test to control the sample sizes
- Cartis and Scheinberg [2016] ensured this condition is satisfied in probability and analyzed global convergence properties

$$\|
abla F_{\mathcal{S}_k}(x_k) - 
abla F(x_k)\| \le heta_n \|
abla F(x_k)\|, \quad \text{for some} \quad heta_n \in [0, 1).$$

$$\frac{\mathbb{E}[\|\nabla F_i(x_k) - \nabla F(x_k)\|^2]}{|S_k|} \leq \theta_n^2 \|\nabla F(x_k)\|^2$$

- Byrd et al. [2012] proposed this test to control the sample sizes
- Cartis and Scheinberg [2016] ensured this condition is satisfied in probability and analyzed global convergence properties
- Hashemi et al. [2014] similar test in simulation optimization settings

$$\|\nabla F_{\mathcal{S}_k}(x_k) - \nabla F(x_k)\| \le \theta_n \|\nabla F(x_k)\|, \quad \text{for some} \quad \theta_n \in [0, 1).$$

$$\frac{\mathbb{E}[\|\nabla F_i(x_k) - \nabla F(x_k)\|^2]}{|S_k|} \leq \theta_n^2 \|\nabla F(x_k)\|^2$$

- Byrd et al. [2012] proposed this test to control the sample sizes
- Cartis and Scheinberg [2016] ensured this condition is satisfied in probability and analyzed global convergence properties
- Hashemi et al. [2014] similar test in simulation optimization settings
- This test is designed to get more than just descent directions

$$\|
abla F_{\mathcal{S}_k}(x_k) - 
abla F(x_k)\| \le heta_n \|
abla F(x_k)\|, \quad \text{for some} \quad heta_n \in [0, 1).$$

$$\frac{\mathbb{E}[\|\nabla F_i(x_k) - \nabla F(x_k)\|^2]}{|S_k|} \le \theta_n^2 \|\nabla F(x_k)\|^2$$

- Byrd et al. [2012] proposed this test to control the sample sizes
- Cartis and Scheinberg [2016] ensured this condition is satisfied in probability and analyzed global convergence properties
- Hashemi et al. [2014] similar test in simulation optimization settings
- This test is designed to get more than just descent directions
- Sample gradients are unnecessarily close to the true gradients

$$\|
abla F_{\mathcal{S}_k}(x_k) - 
abla F(x_k)\| \le heta_n \|
abla F(x_k)\|, \quad \text{for some} \quad heta_n \in [0, 1).$$

$$\frac{\mathbb{E}[\|\nabla F_i(x_k) - \nabla F(x_k)\|^2]}{|S_k|} \le \theta_n^2 \|\nabla F(x_k)\|^2$$

- Byrd et al. [2012] proposed this test to control the sample sizes
- Cartis and Scheinberg [2016] ensured this condition is satisfied in probability and analyzed global convergence properties
- Hashemi et al. [2014] similar test in simulation optimization settings
- This test is designed to get more than just descent directions
- Sample gradients are unnecessarily close to the true gradients
- Sample sizes are increased at much faster rates than the desired rates

Raghu Bollapragada (NU)

6 / 27

Raghu Bollapragada (NU)

→ 3 → 4 3

Image: A matrix

三日 のへで

First-Order Descent Condition

 $\nabla F_{S_k}(x_k)^T \nabla F(x_k) > 0$ 

EL OQO

First-Order Descent Condition

$$\nabla F_{S_k}(x_k)^T \nabla F(x_k) > 0$$

Holds in Expectation

$$\mathbb{E}\left[\nabla F_{S_k}(x_k)^T \nabla F(x_k)\right] = \|\nabla F(x_k)\|^2 > 0$$

ELE DQC

First-Order Descent Condition

$$\nabla F_{S_k}(x_k)^T \nabla F(x_k) > 0$$

Holds in Expectation

$$\mathbb{E}\left[\nabla F_{S_k}(x_k)^T \nabla F(x_k)\right] = \|\nabla F(x_k)\|^2 > 0$$

For descent condition to hold at most iterations, we impose bounds on the variance

ELE NOR

First-Order Descent Condition

$$\nabla F_{S_k}(x_k)^T \nabla F(x_k) > 0$$

Holds in Expectation

$$\mathbb{E}\left[\nabla F_{S_k}(x_k)^T \nabla F(x_k)\right] = \|\nabla F(x_k)\|^2 > 0$$

For descent condition to hold at most iterations, we impose bounds on the variance

$$\frac{\mathbb{E}\left[\left(\nabla F_i(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2\right)^2\right]}{|S_k|} \leq \theta_{ip}^2 \|\nabla F(x_k)\|^4, \, \theta_{ip} \in [0,1)$$

First-Order Descent Condition

$$\nabla F_{S_k}(x_k)^T \nabla F(x_k) > 0$$

Holds in Expectation

$$\mathbb{E}\left[\nabla F_{S_k}(x_k)^T \nabla F(x_k)\right] = \|\nabla F(x_k)\|^2 > 0$$

For descent condition to hold at most iterations, we impose bounds on the variance

$$\frac{\mathbb{E}\left[\left(\nabla F_i(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2\right)^2\right]}{|S_k|} \leq \theta_{ip}^2 \|\nabla F(x_k)\|^4, \, \theta_{ip} \in [0, 1)$$

• Test is designed to achieve descent directions sufficiently often

First-Order Descent Condition

$$\nabla F_{S_k}(x_k)^T \nabla F(x_k) > 0$$

Holds in Expectation

$$\mathbb{E}\left[\nabla F_{S_k}(x_k)^T \nabla F(x_k)\right] = \|\nabla F(x_k)\|^2 > 0$$

For descent condition to hold at most iterations, we impose bounds on the variance

$$\frac{\mathbb{E}\left[\left(\nabla F_i(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2\right)^2\right]}{|S_k|} \leq \theta_{ip}^2 \|\nabla F(x_k)\|^4, \, \theta_{ip} \in [0, 1)$$

- Test is designed to achieve descent directions sufficiently often
- Samples sizes required to satisfy this condition are smaller than those required for norm condition

Raghu Bollapragada (NU)

Adaptive Sampling Methods



Figure: Given a gradient  $\nabla F$  the shaded areas denote the set of vectors satisfying the deterministic (a) Norm test (b) Inner Product test

ELE NOR

#### Lemma

Let  $|S_{ip}|$ ,  $|S_n|$  represent the minimum number of samples required to satisfy the inner product test and norm test at any given iterate x and any given  $\theta_{ip} = \theta_n < 1$ . Then we have

$$\frac{|S_{ip}|}{|S_n|} = \beta(x) \le 1,$$

where

$$\beta(x) = \frac{\mathbb{E}[\|\nabla F_i(x)\|^2 \cos^2(\gamma_i)] - \|\nabla F(x)\|^2}{\mathbb{E}[\|\nabla F_i(x)\|^2] - \|\nabla F(x)\|^2},$$

and  $\gamma_i$  is the angle made by  $\nabla F_i(x)$  with  $\nabla F(x)$ .

5 1 SQC

### Test Approximation

$$\frac{\mathbb{E}\left[\left(\nabla F_i(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2\right)^2\right]}{|S_k|} \leq \theta_{ip}^2 \|\nabla F(x_k)\|^4, \ \theta_{ip} \in [0,1)$$

< 17 ▶

三日 のへの

## Test Approximation

$$\frac{\mathbb{E}\left[\left(\nabla F_i(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2\right)^2\right]}{|S_k|} \leq \theta_{ip}^2 \|\nabla F(x_k)\|^4, \ \theta_{ip} \in [0, 1)$$

• Computing true gradient is expensive

ELE NOR

#### Test Approximation

$$\frac{\mathbb{E}\left[\left(\nabla F_i(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2\right)^2\right]}{|S_k|} \leq \theta_{ip}^2 \|\nabla F(x_k)\|^4, \ \theta_{ip} \in [0, 1)$$

- Computing true gradient is expensive
- Approximate population variance with sample variance and true gradient with sampled gradient

5 1 SQA
$$\frac{\mathbb{E}\left[\left(\nabla F_i(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2\right)^2\right]}{|S_k|} \leq \theta_{ip}^2 \|\nabla F(x_k)\|^4, \ \theta_{ip} \in [0,1)$$

- Computing true gradient is expensive
- Approximate population variance with sample variance and true gradient with sampled gradient

$$\frac{\operatorname{Var}_{i\in S_k}(\nabla F_i(x_k)^T \nabla F_{S_k}(x_k))}{|S_k|} \leq \theta_{ip}^2 \|\nabla F_{S_k}(x_k)\|^4,$$

$$\frac{\mathbb{E}\left[\left(\nabla F_i(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2\right)^2\right]}{|S_k|} \leq \theta_{ip}^2 \|\nabla F(x_k)\|^4, \ \theta_{ip} \in [0, 1)$$

- Computing true gradient is expensive
- Approximate population variance with sample variance and true gradient with sampled gradient

$$\frac{\operatorname{Var}_{i\in S_k}(\nabla F_i(x_k)^T \nabla F_{S_k}(x_k))}{|S_k|} \leq \theta_{ip}^2 \|\nabla F_{S_k}(x_k)\|^4,$$

$$\operatorname{Var}_{i \in S_k} \left( \nabla F_i(x_k)^T \nabla F_{S_k}(x_k) \right) = \frac{1}{|S_k| - 1} \sum_{i \in S_k} \left( \nabla F_i(x_k)^T \nabla F_{S_k}(x_k) - \| \nabla F_{S_k}(x_k) \|^2 \right)^2$$

$$\frac{\mathbb{E}\left[\left(\nabla F_i(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2\right)^2\right]}{|S_k|} \leq \theta_{ip}^2 \|\nabla F(x_k)\|^4, \ \theta_{ip} \in [0,1)$$

- Computing true gradient is expensive
- Approximate population variance with sample variance and true gradient with sampled gradient

$$\frac{\operatorname{Var}_{i\in S_k}(\nabla F_i(x_k)^T \nabla F_{S_k}(x_k))}{|S_k|} \leq \theta_{ip}^2 \|\nabla F_{S_k}(x_k)\|^4,$$

 $\operatorname{Var}_{i\in S_k}\left(\nabla F_i(x_k)^T \nabla F_{S_k}(x_k)\right) = \frac{1}{|S_k| - 1} \sum_{i\in S_k} \left(\nabla F_i(x_k)^T \nabla F_{S_k}(x_k) - \|\nabla F_{S_k}(x_k)\|^2\right)^2$ 

• Whenever condition is not satisfied, increase sample size to satisfy the condition

# Overview

### Adaptive Sampling Tests

- Norm test
- Inner Product Test

### 2 Convergence Analysis

- Orthogonal test
- Linear Convergence

### 3 Practical Implementation

• Step-Length Strategy

Parameter Selection

# 4 Numerical Experiments

# 5 Summary

Raghu Bollapragada (NU)

三日 のへで

(4回) (4回) (4回)

 Although the Inner Product test is practical, convergence cannot be established because it allows sample gradients that are arbitrarily long relative to ||∇F(x<sub>k</sub>)||

1 = nar

- Although the Inner Product test is practical, convergence cannot be established because it allows sample gradients that are arbitrarily long relative to ||∇F(x<sub>k</sub>)||
- No restriction on near orthogonality of sample gradient and gradient

- Although the Inner Product test is practical, convergence cannot be established because it allows sample gradients that are arbitrarily long relative to ||∇F(x<sub>k</sub>)||
- No restriction on near orthogonality of sample gradient and gradient
- Component of sample gradient orthogonal to gradient is 0 in expectation

- Although the Inner Product test is practical, convergence cannot be established because it allows sample gradients that are arbitrarily long relative to ||∇F(x<sub>k</sub>)||
- No restriction on near orthogonality of sample gradient and gradient
- Component of sample gradient orthogonal to gradient is 0 in expectation
- We control the variance in the orthogonal components of sampled gradients

$$\frac{\mathbb{E}\left[\left\|\nabla F_i(x_k) - \frac{\nabla F_i(x_k)^T \nabla F(x_k)}{\|\nabla F(x_k)\|^2} \nabla F(x_k)\right\|^2\right]}{|S_k|} \leq \nu^2 \|\nabla F(x_k)\|^2, \nu > 0$$

Raghu Bollapragada (NU)

< 🗇 🕨

三日 のへで

#### Theorem

Suppose that F is twice continuously differentiable and that there exist constants  $0 < \mu \leq L$ such that

$$\mu I \preceq \nabla^2 F(x) \preceq LI, \quad \forall x \in \mathbb{R}^d.$$

< 67 ▶

EL OQO

#### Theorem

Suppose that F is twice continuously differentiable and that there exist constants  $0 < \mu \leq L$  such that

$$\mu I \preceq \nabla^2 F(x) \preceq LI, \quad \forall x \in \mathbb{R}^d.$$

Let  $\{x_k\}$  be the iterates generated by subsampled gradient method with any  $x_0$ , where  $|S_k|$  is chosen such that inner product test and orthogonal test are satisfied at each iteration for any given  $\theta_{ip} > 0$  and  $\nu > 0$ .

A ∃ ► A ∃ ► ∃ | = \0 Q Q

#### Theorem

Suppose that F is twice continuously differentiable and that there exist constants  $0 < \mu \leq L$  such that

$$\mu I \preceq \nabla^2 F(x) \preceq LI, \quad \forall x \in \mathbb{R}^d.$$

Let  $\{x_k\}$  be the iterates generated by subsampled gradient method with any  $x_0$ , where  $|S_k|$  is chosen such that inner product test and orthogonal test are satisfied at each iteration for any given  $\theta_{ip} > 0$  and  $\nu > 0$ . Then, if the steplength satisfies

$$\alpha_k = \alpha = \frac{1}{(1 + \theta_{ip}^2 + \nu^2)L},$$

A ∃ ► A ∃ ► ∃ | = \0 Q Q

#### Theorem

Suppose that F is twice continuously differentiable and that there exist constants  $0 < \mu \leq L$  such that

$$\mu I \preceq \nabla^2 F(x) \preceq LI, \quad \forall x \in \mathbb{R}^d.$$

Let  $\{x_k\}$  be the iterates generated by subsampled gradient method with any  $x_0$ , where  $|S_k|$  is chosen such that inner product test and orthogonal test are satisfied at each iteration for any given  $\theta_{ip} > 0$  and  $\nu > 0$ . Then, if the steplength satisfies

$$\alpha_k = \alpha = \frac{1}{(1 + \theta_{ip}^2 + \nu^2)L},$$

we have that

$$\mathbb{E}[F(w_k)-F(w^*)] \leq \rho^k(F(x_0)-F(x^*)),$$

where

$$\rho = 1 - \frac{\mu}{L} \frac{1}{(1 + \theta_{ip}^2 + \nu^2)}$$

Convex Non-Convex

▲ Ξ ▶ ▲ Ξ ▶ Ξ Ξ

# Overview

### Adaptive Sampling Tests

- Norm test
- Inner Product Test

### Convergence Analysis

- Orthogonal test
- Linear Convergence

### ③ Practical Implementation

Step-Length Strategy

Parameter Selection

## 4 Numerical Experiments

### 5 Summary

# Step-Length Selection

$$x_{k+1} = x_k - \alpha_k \nabla F_{\mathcal{S}_k}(x_k)$$

三日 のへの

$$x_{k+1} = x_k - \alpha_k \nabla F_{S_k}(x_k)$$

• Stochastic gradient is employed with diminishing stepsizes

$$x_{k+1} = x_k - \alpha_k \nabla F_{\mathcal{S}_k}(x_k)$$

- Stochastic gradient is employed with diminishing stepsizes
- Exact line search can be performed to determine the stepsize but it is too expensive

$$x_{k+1} = x_k - \alpha_k \nabla F_{\mathcal{S}_k}(x_k)$$

- Stochastic gradient is employed with diminishing stepsizes
- Exact line search can be performed to determine the stepsize but it is too expensive
- Constant stepsize can be employed but one needs to know the Lipschitz constant of the problem

$$\alpha_k = 1/L$$

$$x_{k+1} = x_k - \alpha_k \nabla F_{\mathcal{S}_k}(x_k)$$

- Stochastic gradient is employed with diminishing stepsizes
- Exact line search can be performed to determine the stepsize but it is too expensive
- Constant stepsize can be employed but one needs to know the Lipschitz constant of the problem

$$\alpha_k = 1/L$$

• We propose to estimate the Lipschitz constant as we proceed, resulting in adaptive stepsizes

Input:  $L_{k-1} > 0$ , some  $\eta > 1$ 1: Compute parameter  $\zeta_k > 1$ 2: Set  $L_k = L_{k-1}/\zeta_k$ ;  $\triangleright$  Decrease the Lipschitz constant 3: Compute  $F_{new} = F_{s_k} \left( x_k - \frac{1}{L_k} \nabla F_{s_k}(x_k) \right)$ 4: while  $F_{new} > F_{s_k}(x_k) - \frac{1}{2L_k} \| \nabla F_{s_k}(x_k) \|^2$  do  $\triangleright$  sufficient decrease 5: Set  $L_k = \eta L_{k-1}$   $\triangleright$  Increase the Lipschitz constant 6: Compute  $F_{new} = F_{s_k} \left( x_k - \frac{1}{L_k} \nabla F_{s_k}(x_k) \right)$ 7: end while

ヨト イヨト ヨヨ わすや

Input:  $L_{k-1} > 0$ , some  $\eta > 1$ 1: Compute parameter  $\zeta_k > 1$ 2: Set  $L_k = L_{k-1}/\zeta_k$ ;  $\triangleright$  Decrease the Lipschitz constant 3: Compute  $F_{new} = F_{s_k} \left( x_k - \frac{1}{L_k} \nabla F_{s_k}(x_k) \right)$ 4: while  $F_{new} > F_{s_k}(x_k) - \frac{1}{2L_k} \| \nabla F_{s_k}(x_k) \|^2$  do  $\triangleright$  sufficient decrease 5: Set  $L_k = \eta L_{k-1}$   $\triangleright$  Increase the Lipschitz constant 6: Compute  $F_{new} = F_{s_k} \left( x_k - \frac{1}{L_k} \nabla F_{s_k}(x_k) \right)$ 7: end while

• Beck and Teboulle [2009] proposed this algorithm to estimate Lipschitz constant in deterministic settings

A = A = A = A = A = A

Input:  $L_{k-1} > 0$ , some  $\eta > 1$ 1: Compute parameter  $\zeta_k > 1$ 2: Set  $L_k = L_{k-1}/\zeta_k$ ;  $\triangleright$  Decrease the Lipschitz constant 3: Compute  $F_{new} = F_{s_k} \left( x_k - \frac{1}{L_k} \nabla F_{s_k}(x_k) \right)$ 4: while  $F_{new} > F_{s_k}(x_k) - \frac{1}{2L_k} \| \nabla F_{s_k}(x_k) \|^2$  do  $\triangleright$  sufficient decrease 5: Set  $L_k = \eta L_{k-1}$   $\triangleright$  Increase the Lipschitz constant 6: Compute  $F_{new} = F_{s_k} \left( x_k - \frac{1}{L_k} \nabla F_{s_k}(x_k) \right)$ 7: end while

- Beck and Teboulle [2009] proposed this algorithm to estimate Lipschitz constant in deterministic settings
- Scmidt et al. [2016] adapted this algorithm to stochastic algorithms such as SAG

Input:  $L_{k-1} > 0$ , some  $\eta > 1$ 1: Compute parameter  $\zeta_k > 1$ 2: Set  $L_k = L_{k-1}/\zeta_k$ ;  $\triangleright$  Decrease the Lipschitz constant 3: Compute  $F_{new} = F_{s_k} \left( x_k - \frac{1}{L_k} \nabla F_{s_k}(x_k) \right)$ 4: while  $F_{new} > F_{s_k}(x_k) - \frac{1}{2L_k} \| \nabla F_{s_k}(x_k) \|^2$  do  $\triangleright$  sufficient decrease 5: Set  $L_k = \eta L_{k-1}$   $\triangleright$  Increase the Lipschitz constant 6: Compute  $F_{new} = F_{s_k} \left( x_k - \frac{1}{L_k} \nabla F_{s_k}(x_k) \right)$ 7: end while

• Similar to Line-Search on sampled functions with memory

A = A = A = A = A = A

Input:  $L_{k-1} > 0$ , some  $\eta > 1$ 1: Compute parameter  $\zeta_k > 1$ 2: Set  $L_k = L_{k-1}/\zeta_k$ ;  $\triangleright$  Decrease the Lipschitz constant 3: Compute  $F_{new} = F_{s_k} \left( x_k - \frac{1}{L_k} \nabla F_{s_k}(x_k) \right)$ 4: while  $F_{new} > F_{s_k}(x_k) - \frac{1}{2L_k} \| \nabla F_{s_k}(x_k) \|^2$  do  $\triangleright$  sufficient decrease 5: Set  $L_k = \eta L_{k-1}$   $\triangleright$  Increase the Lipschitz constant 6: Compute  $F_{new} = F_{s_k} \left( x_k - \frac{1}{L_k} \nabla F_{s_k}(x_k) \right)$ 7: end while

- Similar to Line-Search on sampled functions with memory
- Only needs access to sampled function values

▲ Ξ ▶ ▲ Ξ ▶ Ξ Ξ

Raghu Bollapragada (NU)

三日 のへで

### It is well known and easy to show that

$$\mathbb{E}[F(x_{k+1})] - F(x_k) \leq -\alpha_k \|F(x_k)\|^2 + \frac{\alpha_k^2 L}{2} \mathbb{E}[\|\nabla F_{\mathcal{S}_k}(x_k)\|^2]$$

I= nan

It is well known and easy to show that

$$\mathbb{E}[F(x_{k+1})] - F(x_k) \leq -\alpha_k \|F(x_k)\|^2 + \frac{\alpha_k^2 L}{2} \mathbb{E}[\|\nabla F_{\mathcal{S}_k}(x_k)\|^2]$$

Thus, we can obtain a decrease in the true objective, in expectation, if

$$\frac{L\alpha_k^2}{2} \left( \mathbb{E}[\|\nabla F_{\mathcal{S}_k}(x_k) - \nabla F(x_k)\|^2] + \|\nabla F(x_k)\|^2 \right) \le \alpha_k \|\nabla F(x_k)\|^2$$

# $\frac{L\alpha_k^2}{2} \left( \mathbb{E}[\|\nabla F_{\mathcal{S}_k}(x_k) - \nabla F(x_k)\|^2] + \|\nabla F(x_k)\|^2 \right) \le \alpha_k \|\nabla F(x_k)\|^2$

3 × 4 3 × 3 1 × 0 0 0

$$\frac{L\alpha_k^2}{2} \left( \mathbb{E}[\|\nabla F_{\mathcal{S}_k}(x_k) - \nabla F(x_k)\|^2] + \|\nabla F(x_k)\|^2 \right) \le \alpha_k \|\nabla F(x_k)\|^2$$

Using,  $\alpha_k = 1/L_k$ , assuming  $L_{k-1} \ge L$ , and sample approximations

$$L_k \geq rac{L_{k-1}}{2} \left( rac{\operatorname{Var}\left( 
abla F_i(x_k) 
ight)}{|S_k| \| 
abla F_{S_k}(x_k) \|^2} + 1 
ight),$$

where  $\operatorname{Var}(\nabla F_i(x_k)) = \frac{1}{|S_k|-1} \sum_{i \in S_k} \|\nabla F_i(x_k) - \nabla F_{S_k}(x_k)\|^2$ .

▲ Ξ ▶ ▲ Ξ ▶ Ξ Ξ

$$\frac{L\alpha_k^2}{2} \left( \mathbb{E}[\|\nabla F_{\mathcal{S}_k}(x_k) - \nabla F(x_k)\|^2] + \|\nabla F(x_k)\|^2 \right) \le \alpha_k \|\nabla F(x_k)\|^2$$

Using,  $\alpha_k = 1/L_k$ , assuming  $L_{k-1} \ge L$ , and sample approximations

$$L_k \geq rac{L_{k-1}}{2} \left( rac{\operatorname{Var}\left( 
abla F_i(x_k) 
ight)}{|S_k| \| 
abla F_{S_k}(x_k) \|^2} + 1 
ight),$$

where  $\operatorname{Var}(\nabla F_i(x_k)) = \frac{1}{|S_k|-1} \sum_{i \in S_k} \|\nabla F_i(x_k) - \nabla F_{S_k}(x_k)\|^2$ . Therefore,

$$\zeta_k = \max\left(1, \frac{2}{\frac{\mathsf{Var}(\nabla F_i(x_k))}{|S_k| \|\nabla F_{S_k}(x_k)\|^2} + 1}\right)$$

A ∃ ► A ∃ ► ∃ | = \0 Q Q

# Parameter Selection

Raghu Bollapragada (NU)

三日 のへの

$$\frac{\nabla F_{S_k}(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2}{\left(\frac{\sigma}{\sqrt{|S_k|}}\right)} \sim \mathcal{N}(0, 1),$$
  
where  $\sigma^2 = \mathbb{E}\left[\left(\nabla F_i(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2\right)^2\right]$  is the true variance.

三日 のへの

$$\frac{\nabla F_{\mathcal{S}_k}(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2}{\left(\frac{\sigma}{\sqrt{|\mathcal{S}_k|}}\right)} \sim \mathcal{N}(0, 1),$$

where  $\sigma^2 = \mathbb{E}\left[\left(\nabla F_i(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2\right)^2\right]$  is the true variance.

• Parameter  $\theta_{ip}$  is directly proportional to probability of getting a descent direction

$$\frac{\nabla F_{S_k}(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2}{\left(\frac{\sigma}{\sqrt{|S_k|}}\right)} \sim \mathcal{N}(0, 1),$$

where  $\sigma^2 = \mathbb{E}\left[\left(\nabla F_i(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2\right)^2\right]$  is the true variance.

- Parameter  $\theta_{ip}$  is directly proportional to probability of getting a descent direction
- $\theta_{ip} = 0.7$  corresponds to 0.9 probability

$$\frac{\nabla F_{S_k}(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2}{\left(\frac{\sigma}{\sqrt{|S_k|}}\right)} \sim \mathcal{N}(0, 1),$$

where  $\sigma^2 = \mathbb{E}\left[\left(\nabla F_i(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2\right)^2\right]$  is the true variance.

- Parameter  $\theta_{ip}$  is directly proportional to probability of getting a descent direction
- $\theta_{ip} = 0.7$  corresponds to 0.9 probability
- $\theta_{ip} = 0.9$  works well in practice
$$\frac{\nabla F_{\mathcal{S}_k}(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2}{\left(\frac{\sigma}{\sqrt{|\mathcal{S}_k|}}\right)} \sim \mathcal{N}(0, 1),$$

where  $\sigma^2 = \mathbb{E}\left[\left(\nabla F_i(x_k)^T \nabla F(x_k) - \|\nabla F(x_k)\|^2\right)^2\right]$  is the true variance.

- Parameter  $\theta_{ip}$  is directly proportional to probability of getting a descent direction
- $\theta_{ip} = 0.7$  corresponds to 0.9 probability
- $\theta_{ip} = 0.9$  works well in practice
- Orthogonal test is seldom active in practice and we choose  $\nu = \tan(80^{\circ}) = 5.84$  for all problems

# Overview

#### Adaptive Sampling Tests

- Norm test
- Inner Product Test

#### Convergence Analysis

- Orthogonal test
- Linear Convergence

#### 3 Practical Implementation

• Step-Length Strategy

Parameter Selection

#### 4 Numerical Experiments

#### 5 Summary

$$R(x) = \frac{1}{n} \sum_{i=1}^{n} \log(1 + \exp(-z^{i} x^{T} y^{i})) + \frac{\lambda}{2} ||x||^{2}$$

-

ELE NOR



Figure: Norm Test vs. Inner Product Test. Synthetic Dataset (n = 7000)

ELE NOR



Other Datasets

### Results: Adaptive Step-Length Strategy



Raghu Bollapragada (NU)

ъ.

### Results: Adaptive Step-Length Strategy



# Overview

#### Adaptive Sampling Tests

- Norm test
- Inner Product Test

#### Convergence Analysis

- Orthogonal test
- Linear Convergence

#### 3 Practical Implementation

• Step-Length Strategy

Parameter Selection

#### 4 Numerical Experiments



# Summary

Raghu Bollapragada (NU)

#### • Adaptive sampling methods are alternate methods for noise reduction

э

1= 9QC

- Adaptive sampling methods are alternate methods for noise reduction
- These methods can lead to speed ups when implemented in parallel environments

= nac

- Adaptive sampling methods are alternate methods for noise reduction
- These methods can lead to speed ups when implemented in parallel environments
- We propose a practical inner product test which is better at controlling the sample sizes than the existing norm test

- Adaptive sampling methods are alternate methods for noise reduction
- These methods can lead to speed ups when implemented in parallel environments
- We propose a practical inner product test which is better at controlling the sample sizes than the existing norm test
- These methods can use adaptive stepsizes and second-order information can be incorporated

- Adaptive sampling methods are alternate methods for noise reduction
- These methods can lead to speed ups when implemented in parallel environments
- We propose a practical inner product test which is better at controlling the sample sizes than the existing norm test
- These methods can use adaptive stepsizes and second-order information can be incorporated
- Currently working on practical sampling tests to control sample sizes in stochastic quasi-Newton methods

# Questions???

# Thank You

Raghu Bollapragada (NU)

Adaptive Sampling Methods

US - Mexico Workshop - 2018 1 / 13

3

# Appendix

・ 同 ト ・ ヨ ト ・ ヨ ト

Raghu Bollapragada (NU)

(3)

• Sample approximations are not accurate when samples are very small (say 1, 5, 10)

-

- Sample approximations are not accurate when samples are very small (say 1, 5, 10)
- Our tests may not be accurate in controlling the sample sizes in such situations

- Sample approximations are not accurate when samples are very small (say 1, 5, 10)
- Our tests may not be accurate in controlling the sample sizes in such situations
- Need more accurate approximations in such scenarios

- Sample approximations are not accurate when samples are very small (say 1, 5, 10)
- Our tests may not be accurate in controlling the sample sizes in such situations
- Need more accurate approximations in such scenarios

$$g_{avg} \stackrel{\text{def}}{=} \frac{1}{r} \sum_{j=k-r+1}^{k} \nabla F_{S_j}(x_j)$$

- Sample approximations are not accurate when samples are very small (say 1, 5, 10)
- Our tests may not be accurate in controlling the sample sizes in such situations
- Need more accurate approximations in such scenarios

$$g_{avg} \stackrel{\text{def}}{=} \frac{1}{r} \sum_{j=k-r+1}^{k} \nabla F_{\mathcal{S}_j}(x_j)$$

- r should be chosen such that the iterates in the summation are close enough and there are enough samples for  $g_{avg}$  to be accurate (r = 10).
- If ||g<sub>avg</sub>|| < γ ||∇F<sub>Sk</sub>(x<sub>k</sub>)||, for some γ ∈ (0, 1) then we use g<sub>avg</sub> instead of ∇F<sub>Sk</sub>(x<sub>k</sub>) in the tests.

Raghu Bollapragada (NU)

(日) (周) (三) (三)

#### Theorem

(General Convex Objective.) Suppose that F is twice continuously differentiable and convex, and that there exists a constant L > 0 such that

 $\nabla^2 F(x) \preceq LI, \quad \forall x \in \mathbb{R}^d.$ 

#### Theorem

(General Convex Objective.) Suppose that F is twice continuously differentiable and convex, and that there exists a constant L > 0 such that

$$abla^2 F(x) \preceq LI, \quad \forall x \in \mathbb{R}^d.$$

Let  $\{x_k\}$  be the iterates generated by subsampled gradient method with any  $x_0$ , where  $|S_k|$  is chosen such that inner product test and orthogonal test are satisfied at each iteration for any given  $\theta_{ip} > 0$  and  $\nu > 0$ .

A ∃ ► A ∃ ► ∃ | = \0 Q Q

#### Theorem

(General Convex Objective.) Suppose that F is twice continuously differentiable and convex, and that there exists a constant L > 0 such that

$$abla^2 F(x) \preceq LI, \quad \forall x \in \mathbb{R}^d.$$

Let  $\{x_k\}$  be the iterates generated by subsampled gradient method with any  $x_0$ , where  $|S_k|$  is chosen such that inner product test and orthogonal test are satisfied at each iteration for any given  $\theta_{ip} > 0$  and  $\nu > 0$ . Then, if the steplength satisfies

$$\alpha_k = \alpha < \frac{1}{(1 + \theta_{ip}^2 + \nu^2)L},$$

A ∃ ► A ∃ ► ∃ | = \0 Q Q

#### Theorem

(General Convex Objective.) Suppose that F is twice continuously differentiable and convex, and that there exists a constant L > 0 such that

$$abla^2 F(x) \preceq LI, \quad \forall x \in \mathbb{R}^d.$$

Let  $\{x_k\}$  be the iterates generated by subsampled gradient method with any  $x_0$ , where  $|S_k|$  is chosen such that inner product test and orthogonal test are satisfied at each iteration for any given  $\theta_{ip} > 0$  and  $\nu > 0$ . Then, if the steplength satisfies

$$\alpha_k = \alpha < \frac{1}{(1 + \theta_{ip}^2 + \nu^2)L},$$

we have for any positive integer T,

$$\min_{0 \le k \le T-1} \mathbb{E} \left[ F(x_k) \right] - F^* \le \frac{1}{2\alpha cT} \|x_0 - x^*\|^2,$$

where the constant c > 0 is given by  $c = 1 - L\alpha(1 + \theta_{ip}^2 + \nu^2)$ .

Raghu Bollapragada (NU)

3 1 4

#### Theorem

(Nonconvex Objective.) Suppose that F is twice continuously differentiable and bounded below, and that there exist a constant L > 0 such that

 $abla^2 F(x) \preceq LI, \quad \forall x \in \mathbb{R}^d.$ 

Image: A matrix

- ▲ 표 ▶ ▲ 표 ▶ . 토 | 표 . • • • • • •

#### Theorem

(Nonconvex Objective.) Suppose that F is twice continuously differentiable and bounded below, and that there exist a constant L > 0 such that

$$abla^2 F(x) \preceq LI, \quad \forall x \in \mathbb{R}^d.$$

Let  $\{x_k\}$  be the iterates generated by subsampled gradient method with any  $x_0$ , where  $|S_k|$  is chosen such that inner product test and orthogonal test are satisfied at each iteration for any given  $\theta_{ip} > 0$  and  $\nu > 0$ .

ELE NOR

#### Theorem

(Nonconvex Objective.) Suppose that F is twice continuously differentiable and bounded below, and that there exist a constant L > 0 such that

$$abla^2 F(x) \preceq LI, \quad \forall x \in \mathbb{R}^d.$$

Let  $\{x_k\}$  be the iterates generated by subsampled gradient method with any  $x_0$ , where  $|S_k|$  is chosen such that inner product test and orthogonal test are satisfied at each iteration for any given  $\theta_{ip} > 0$  and  $\nu > 0$ . Then, if the steplength satisfies

$$\alpha_k = \alpha \leq \frac{1}{(1 + \theta_{ip}^2 + \nu^2)L},$$

ELE NOR

#### Theorem

(Nonconvex Objective.) Suppose that F is twice continuously differentiable and bounded below, and that there exist a constant L > 0 such that

$$abla^2 F(x) \preceq LI, \quad \forall x \in \mathbb{R}^d.$$

Let  $\{x_k\}$  be the iterates generated by subsampled gradient method with any  $x_0$ , where  $|S_k|$  is chosen such that inner product test and orthogonal test are satisfied at each iteration for any given  $\theta_{ip} > 0$  and  $\nu > 0$ . Then, if the steplength satisfies

$$\alpha_k = \alpha \leq \frac{1}{(1 + \theta_{ip}^2 + \nu^2)L},$$

we have that

$$\lim_{k\to\infty}\mathbb{E}[\|\nabla F(x_k)\|^2]\to 0.$$

Moreover, for any positive integer T we have that

$$\min_{0\leq k\leq T-1}\mathbb{E}[\|\nabla F(x_k)\|^2]\leq \frac{2}{\alpha T}(F(x_0)-F_{\min}).$$

Main-Presentation

▲ Ξ ▶ ▲ Ξ ▶ Ξ Ξ






# Results: Constant Step-Length Strategy



Main-Presentation



Main-Presentation

ELE NOR

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >



Main-Presentation

EL OQO

< ロ > < 同 > < 三 > < 三 > :



< 一型

ELE NOR



Main-Presentation

ELE NOR