

# LECTURE NOTES ON STOCHASTIC GRADIENT DESCENT \*

A. A. ERGÜR, V. T. NGUYEN

**Abstract.** This scribe contains lecture notes on Stochastic Gradient Descent

**Key words.** Stochastic Gradient Descent

**AMS subject classifications.** safely ignore

**1. Stochastic Gradient Descent (SGD).** Given a dataset  $\mathcal{S} = \{\mathbf{x}_i, y_i\}_{i=1}^N$ , where each feature vector  $\mathbf{x}_i \in \mathbb{R}^D$ . We define a loss function  $f(\cdot)$  which measures the error between the prediction and the true label for every single sample  $i$ -th as

$$f_{\boldsymbol{\theta}} : \mathbb{R}^D \times \mathbb{R} \rightarrow \mathbb{R}.$$

The average error which calculated on the entire dataset  $\mathcal{S}$  is called cost function, i.e.,

$$C(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y}) = \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}_i, y_i; \boldsymbol{\theta}),$$

where  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top$ ,  $\mathbf{y} = [y_1, \dots, y_N]^\top$ . For training machine learning, we want to find the parameters  $\boldsymbol{\theta} \in \Omega$  that minimize the cost function, i.e.,

$$\hat{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\theta} \in \Omega} C(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y})$$

**1.1. Recall Gradient descent.** Supposedly,  $f(\cdot)$  is differentiable, then the gradient of the cost function wrt  $\boldsymbol{\theta}$  is given as

$$\nabla_{\boldsymbol{\theta}} C(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y}) = \frac{1}{N} \sum_{i=1}^N \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_i, y_i; \boldsymbol{\theta})$$

Then the update parameters follow as

$$\boldsymbol{\theta}^{(k)} = \boldsymbol{\theta}^{(k)} - \alpha^{(k)} \nabla_{\boldsymbol{\theta}} C(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y}).$$

Since computing the gradient on the entire dataset is costly, people refer stochastic methods that approximate the  $\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y})$ , with fewer computational effort at every step. That is Stochastic Gradient Descent.

**1.2. SGD algorithm.** The idea of the stochastic gradient descent algorithm is to build a function  $g(\cdot)$  which approximates  $C(\cdot)$  and in the mean,  $g(\cdot)$  behaves like  $C(\cdot)$ , i.e.,

$$\mathbb{E}\{g(\mathbf{x}, y; \boldsymbol{\theta})\} = C(\mathbf{x}, y; \boldsymbol{\theta}),$$

$$\mathbb{E}\{\nabla_{\boldsymbol{\theta}} g(\mathbf{x}, y; \boldsymbol{\theta})\} = \nabla_{\boldsymbol{\theta}} C(\mathbf{x}, y; \boldsymbol{\theta}).$$

Then, we can update the model weights as

$$(1.1) \quad \boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \alpha^{(k)} \nabla_{\boldsymbol{\theta}} g(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)}).$$

---

\*We thank Robbins family for supporting the Algorithmic Foundations of Data Science Course

**Algorithm 1.1** SGD: Stochastic Gradient Descent**Require:**  $\alpha^{(k)}$ : Learning rate at iteration  $k$ **Require:**  $g(\mathbf{x}, y; \boldsymbol{\theta})$ : Stochastic approximation of the cost function**Require:**  $\boldsymbol{\theta}^{(0)}$ : Initial model parameters

- 1: **for**  $k = 0, 1, 2, \dots$  until convergence **do**
- 2:   Sample data point  $(\mathbf{x}^{(k)}, y^{(k)})$
- 3:   Compute stochastic gradient  $\nabla_{\boldsymbol{\theta}} g(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)})$
- 4:   Update parameters:

$$\boldsymbol{\theta}^{(k+1)} \leftarrow \boldsymbol{\theta}^{(k)} - \alpha^{(k)} \nabla_{\boldsymbol{\theta}} g(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)})$$

5: **end for**6: **return**  $\boldsymbol{\theta}^{(k+1)}$ 

32       *Stochastic gradient descent for linear regression.* Recall that the cost function of  
 33 linear regression is given as

$$34 \quad C(\mathbf{X}, \mathbf{y}; \boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i \boldsymbol{\theta} - y_i)^2 = \frac{1}{N} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2$$

35 By calculation, we have

$$36 \quad \nabla_{\boldsymbol{\theta}} C(\mathbf{X}, \mathbf{y}; \boldsymbol{\theta}) = \frac{2}{N} (\mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} - \mathbf{X}^\top \mathbf{y})$$

37 Instead of using  $C(\cdot)$  and its gradient, we can use the alternative cost function  $g(\cdot)$  as

$$38 \quad g(\mathbf{x}, y; \boldsymbol{\theta}; i) = (\mathbf{x}_i^\top \boldsymbol{\theta} - y_i)^2,$$

39 where  $i$  is a random variable that follows a uniform distribution, i.e.,  $\Pr(i) = \frac{1}{N}$ . We  
 40 can prove that in the mean, the expectation of  $g(\cdot)$  equals to  $C(\cdot)$  as

$$\begin{aligned}
 41 \quad \mathbb{E}_i \{g(\mathbf{x}, y; \boldsymbol{\theta})\} &= \sum_{i=1}^N \Pr(i) (\mathbf{x}_i^\top \boldsymbol{\theta} - y_i)^2 \\
 42 &= \sum_{i=1}^N \frac{1}{N} (\mathbf{x}_i^\top \boldsymbol{\theta} - y_i)^2 \\
 43 &= \sum_{i=1}^N \frac{1}{N} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 \\
 44 &\triangleq C(\mathbf{X}, \mathbf{y}; \boldsymbol{\theta}).
 \end{aligned}$$

45 Similarly, we have  $\nabla_{\boldsymbol{\theta}} g(\mathbf{x}, y; \boldsymbol{\theta}; i) = 2 (\mathbf{x}_i) (\mathbf{x}_i^\top \boldsymbol{\theta} - y_i)$ . Then consider the expec-

tation of  $\nabla_{\theta} g(\mathbf{x}, y; \theta; i)$  over  $i$ , we have:

$$\begin{aligned}
 \mathbb{E}_i\{\nabla_{\theta} g(\mathbf{x}, y; \theta; i)\} &= \mathbb{E}_i\{2(\mathbf{x}_i)(\mathbf{x}_i^{\top} \theta - y_i)\} \\
 &= \sum_{i=1}^N 2 \frac{1}{N} \mathbf{x}_i (\mathbf{x}_i^{\top} \theta - y_i) \\
 &= \frac{2}{N} \sum_{i=1}^N \mathbf{x}_i (\mathbf{x}_i^{\top} \theta - y_i) \\
 &= \frac{2}{N} (\mathbf{X}^{\top} \mathbf{X} \theta - \mathbf{X} \mathbf{y}) \triangleq \nabla_{\theta} C(\mathbf{X}, \mathbf{y}; \theta)
 \end{aligned}$$

*Hinge loss function.* Given dataset  $\mathcal{S} = \{\mathbf{x}_i, y_i\}_{i=1}^N$ , where  $\mathbf{x}_i \in \mathbb{R}^D, y_i \in \{-1, 1\}$ . We want to find the hyperplane parameterized by  $\theta \in \mathbb{R}^D$ , such that

$$C(\mathbf{X}, \mathbf{y}; \theta) = \frac{1}{N} \sum_{i=1}^N \alpha \max(1 - y_i \mathbf{x}_i^{\top} \theta, 0) + \frac{\lambda}{2} \|\theta\|_2^2$$

is minimized, where  $\alpha, \lambda$  are hyper-parameters. Instead of using the above cost function  $C(\cdot)$  to do optimization, we can use the following approximation, i.e.,

$$g(\mathbf{x}, y; \theta; i) = \alpha \max(1 - y_i \mathbf{x}_i^{\top} \theta, 0) + \frac{\lambda}{2} \|\theta\|_2^2.$$

Again, the index  $i$  is chosen uniformly in  $\{1, \dots, N\}$ . Then we have

$$\mathbb{E}_i\{g(\mathbf{x}, y; \theta)\} = \sum_{i=1}^N \frac{1}{N} \alpha \max(1 - y_i \mathbf{x}_i^{\top} \theta, 0) + \frac{\lambda}{2} \|\theta\|_2^2 \triangleq C(\mathbf{X}, \mathbf{y}; \theta).$$

The gradient of  $g(\cdot)$  is given as

$$\nabla_{\theta} g(\mathbf{x}, y; \theta; i) = \alpha \max(-y_i \mathbf{x}_i, 0) + \lambda \theta.$$

The gradient of  $C(\cdot)$  is given as

$$\nabla_{\theta} C(\mathbf{X}, \mathbf{y}; \theta) = \frac{1}{N} \sum_{i=1}^N \alpha \max(-y_i \mathbf{x}_i, 0) + \lambda \theta.$$

Obviously, we have  $\mathbb{E}_i\{\nabla_{\theta} g(\mathbf{x}, y; \theta; i)\} \triangleq \nabla_{\theta} C(\mathbf{X}, \mathbf{y}; \theta)$ .

**1.3. Loss function is convex.** Given a convex function  $f(\cdot)$ , and we apply SGD as presented in (1.1), with the approximate function  $g(\mathbf{x}, y; \theta; i) = f(\mathbf{x}_i, y_i; \theta)$ , where the index  $i \in [1, N]$  chosen uniformly, and supposedly  $\hat{\theta}$  is the optimal point, subtracting both sides for the optimal point and take their squared magnitudes, we have

$$\|\theta^{(k+1)} - \hat{\theta}\|_2^2 = \|\left(\theta^{(k)} - \hat{\theta}\right) - \alpha^{(k)} \nabla_{\theta} f(\mathbf{x}_i^{(k)}, y_i^{(k)}; \theta^{(k)})\|_2^2.$$

The above equality is expressed as

$$\begin{aligned}
 \|\theta^{(k+1)} - \hat{\theta}\|_2^2 - \|\theta^{(k)} - \hat{\theta}\|_2^2 &= \left(\alpha^{(k)}\right)^2 \|\nabla_{\theta} f(\mathbf{x}_i^{(k)}, y_i^{(k)}; \theta^{(k)})\|_2^2 \\
 &\quad - 2\alpha^{(k)} \nabla_{\theta}^{\top} f(\mathbf{x}^{(k)}, y^k; \theta^{(k)}) (\theta^{(k)} - \hat{\theta}).
 \end{aligned}$$

73 Since

$$\begin{aligned} 74 \quad \mathbb{E}_i\{\nabla_{\boldsymbol{\theta}}^\top f(\mathbf{x}_i^{(k)}, y_i^{(k)}; \boldsymbol{\theta}^{(k)})(\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}})\} &= \mathbb{E}_i\{\nabla_{\boldsymbol{\theta}}^\top f(\mathbf{x}_i^{(k)}, y_i^{(k)}; \boldsymbol{\theta}^{(k)})\}(\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}) \\ 75 \quad &= \mathbb{E}\left\{\nabla^\top f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)})\left(\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}\right)\right\}, \end{aligned}$$

76 and

$$77 \quad \mathbb{E}_i\{\|\nabla_{\boldsymbol{\theta}} g(\mathbf{x}, y; \boldsymbol{\theta}^{(k)}); i\|^2\} = \mathbb{E}\{\|\nabla_{\boldsymbol{\theta}} f(\mathbf{x}, y; \boldsymbol{\theta}^{(k)})\|^2\}$$

78 , we have

$$79 \quad \dots$$

80 **1.4. Randomized Kaczmarz method.** This problem is set up that we are  
81 given a training dataset  $\{\mathbf{A}, \mathbf{b}\} \in \mathbb{R}^{N \times D} \times \mathbb{R}^N$  which establish the linear system:

$$82 \quad \mathbf{A}\mathbf{x} = \mathbf{b}.$$

83 Our goal is to find the best  $\hat{\mathbf{x}}$ , which is the root of the system. If  $\mathbf{b} \in \text{Col}(\mathbf{A})$ , then  
84 the linear system above has one or more solutions; otherwise, there is no solution  
85 for the system above. Instead of dealing with the two cases above, we can solve an  
86 optimization problem instead, i.e.,

$$87 \quad \hat{\mathbf{x}} = \text{argmin}_{\mathbf{x} \in \mathbb{R}^D} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2.$$

88 The Randomized Kaczmarz [6] is an iterative method for solving the complex-  
89 value linear system. The Randomized Kaczmarz starts with initializing  $\mathbf{x}_0$  to be an  
90 arbitrary complex-valued initial approximation, then the update formula is given as

$$91 \quad \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \frac{b_i - \langle \mathbf{a}_i^{(k)}, \mathbf{x}^{(k)} \rangle}{\|\mathbf{a}_i^{(k)}\|_2^2} \mathbf{a}_i^{(k)},$$

$$92 \quad \text{with } \Pr\{i = k\} = \frac{\|\mathbf{a}_k\|_2^2}{\sum_{i=1}^N \|\mathbf{a}_i\|_2^2}.$$

93 Let  $\kappa(\mathbf{A}) = \|\mathbf{A}^{-1}\| \|\mathbf{A}\|_F = \frac{1}{\sigma_{\min}(\mathbf{A})} \|\mathbf{A}\|_F$ , then at the  $k$ -th iteration, it is guar-  
94 anteed that

$$95 \quad \mathbb{E}\{\|\mathbf{x}_k - \hat{\mathbf{x}}\|_2^2\} \leq (1 - \kappa(\mathbf{A})^{-2})^k \|\mathbf{x}_0 - \hat{\mathbf{x}}\|_2^2$$

96

97 *Proof.* Since  $\sum_{i=1}^m |\langle \mathbf{z}, \mathbf{a}_i \rangle|^2 \leq \sigma_{\min}(\mathbf{A})^2 \|\mathbf{z}\|^2$ , and  $\|\mathbf{A}\|_F^2 = \sum_{i=1}^m \|\mathbf{a}_i\|_2^2$  then

$$\begin{aligned} 98 \quad &\sum_{i=1}^m \frac{1}{\|\mathbf{A}\|_F^2} |\langle \mathbf{z}, \mathbf{a}_i \rangle|^2 \geq \kappa(\mathbf{A})^{-1} \|\mathbf{z}\|^2 \\ 99 \quad &\Leftrightarrow \sum_{i=1}^m \frac{\|\mathbf{a}_i\|_2^2}{\|\mathbf{A}_F^2\|} |\langle \mathbf{z}, \frac{\mathbf{a}_i}{\|\mathbf{a}_i\|_2} \rangle|^2 \geq \kappa(\mathbf{A})^{-2} \|\mathbf{z}\|^2. \end{aligned}$$

100 Let  $\Pr(e = i) = \frac{\|\mathbf{a}_i\|_2^2}{\sum_{i=1}^m \|\mathbf{a}_i\|_2^2}$ , and  $\mathbf{u}_e = \frac{\mathbf{a}_e}{\|\mathbf{a}_e\|}$ , then

$$101 \quad \mathbb{E}\{|\langle \mathbf{z}, \mathbf{u}_e \rangle|^2\} \geq \kappa(\mathbf{A})^{-2} \|\mathbf{z}\|^2.$$

102 Note that  $\mathbf{x}^{(k)}$  is the projection of  $\mathbf{x}$  on the span of  $\mathbf{a}_{k-1}$ , thus

$$103 \quad \|\mathbf{x}^{(k)} - \hat{\mathbf{x}}\|^2 = \|\mathbf{x}^{(k)} - \hat{\mathbf{x}}\|^2 - \|\mathbf{x}_{k-1} - \mathbf{x}_k\|^2$$

104 , and since  $\hat{\mathbf{x}}$  satisfies  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , then

$$105 \quad \|\mathbf{x}^{(k-1)} - \mathbf{x}^{(k)}\|_2^2 = \left\langle \mathbf{x}^{(k-1)} - \hat{\mathbf{x}}, \frac{\mathbf{a}^{(k)}}{\|\mathbf{a}^{(k)}\|} \right\rangle.$$

106 Then we have

$$107 \quad \mathbb{E}_{\{\mathbf{u}_1, \dots, \mathbf{u}_k\}} \|\mathbf{x}^{(k)} - \hat{\mathbf{x}}\|^2 \leq \left(1 - \mathbb{E}\{|\langle \mathbf{x}^{(k-1)} - \hat{\mathbf{x}}, \mathbf{u}^{(k)} \rangle|^2\}\right) \|\mathbf{x}^{(k-1)} - \hat{\mathbf{x}}\|^2$$

108 And since  $\{\mathbf{u}_i\}$  are independent, then

$$109 \quad \mathbb{E}\{\|\mathbf{x}^{(k)} - \hat{\mathbf{x}}\|_2^2\} \leq (1 - \kappa(\mathbf{A})^{-2}) \|\mathbf{x}^{(k-1)} - \hat{\mathbf{x}}\|_2^2$$

110 holds. □

## 111 2. Quasi-Newton methods.

112 **2.1. Newton method.** The Newton method is a second-order method that uses  
 113 the Hessian of the objective to compute the update in SGD. At every iteration, the  
 114 new vector  $\boldsymbol{\theta}$  is updated as

$$115 \quad \boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \alpha_k \left( \nabla_{\boldsymbol{\theta}}^2 f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)}) \right)^{-1} \nabla_{\boldsymbol{\theta}} f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)})$$

---

### Algorithm 2.1 Newton Method for Optimization

---

**Require:**  $\alpha_k$ : Stepsize at iteration  $k$

**Require:**  $f(\mathbf{x}, y; \boldsymbol{\theta})$ : Objective function

**Require:**  $\boldsymbol{\theta}^{(0)}$ : Initial parameter vector

- 1: **for**  $k = 0, 1, 2, \dots$  until convergence **do**
- 2:   Sample data point  $(\mathbf{x}^{(k)}, y^{(k)})$
- 3:   Compute gradient  $\nabla_{\boldsymbol{\theta}} f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)})$
- 4:   Compute Hessian  $\nabla_{\boldsymbol{\theta}}^2 f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)})$
- 5:   Update parameters:

$$\boldsymbol{\theta}^{(k+1)} \leftarrow \boldsymbol{\theta}^{(k)} - \alpha_k \left( \nabla_{\boldsymbol{\theta}}^2 f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)}) \right)^{-1} \nabla_{\boldsymbol{\theta}} f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)})$$

6: **end for**

7: **return**  $\boldsymbol{\theta}^{(k+1)}$

---

116

117 The Newton method is good but not an ideal one in terms of efficiency since it  
 118 requires the computation of the Hessian matrix and its inversion. Instead, people  
 119 invented intermediate methods in the middle of standard SGD and Newton, called  
 120 Quasi-Newton. Specifically, they used an approximate version of the Hessian matrix  
 121 whose inverse is easily computed, i.e.,

$$122 \quad \boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \alpha^{(k)} \left( \mathbf{B}^{(k)} \right)^{-1} \nabla_{\boldsymbol{\theta}} f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)}).$$

123 , where the  $\mathbf{B}^{(k)}$  an *easier* version of the Hessian matrix  $\nabla_{\boldsymbol{\theta}}^2 f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)})$ , where  
 124  $(\mathbf{B})^{-1}$  is computed quickly but still approximate well the inverse of the Hessian.

125 Next, we discuss several versions of this Quasi-Newton method, including Adap-  
 126 tive Preconditioned Gradient Descent (AdaGrad), AdaDelta, and Root Mean Square  
 127 Propagation.

**2.2. Adaptive Preconditioned Gradient Descent.** The update of the Adaptive Preconditioned Gradient Descent (aka, AdaGrad) [1] is given as

$$(2.1) \quad \boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \alpha^{(k)} \mathbf{G}^{-1} \nabla_{\boldsymbol{\theta}} f(\mathbf{x}^k, y^k; \boldsymbol{\theta}^{(k)}),$$

where

$$\mathbf{G}^{(t)} = \left( \sum_{i=0}^t \nabla_{\boldsymbol{\theta}} f(\mathbf{x}^{(i)}, y^{(i)}; \boldsymbol{\theta}^{(i)}) \left( f(\mathbf{x}^{(i)}, y^{(i)}; \boldsymbol{\theta}^{(i)}) \right)^\top \right)^{1/2}.$$

A practical version of AdaGrad employs the diagonal matrix  $\mathbf{D}^{(t)}$ , with entries  $\sum_{i=1}^n (\|\nabla_{\boldsymbol{\theta}} f(\mathbf{x}^{(i)}, y^{(i)}; \boldsymbol{\theta}^{(i)})\| + \delta)$ .

**2.2.1. Analysis.** Obviously  $\mathbf{G}^{(k)}$  is a positive semi-definite matrix. By subtracting both sides of (2.1) for  $\hat{\boldsymbol{\theta}}$  and taking the squared  $\mathbf{A}$ -norm<sup>1</sup>, we have

$$\begin{aligned} \|\boldsymbol{\theta}^{(k+1)} - \hat{\boldsymbol{\theta}}\|_{\mathbf{G}^{(k)}}^2 &= \left\| \boldsymbol{\theta}^{(k)} - \alpha^{(k)} \left( \mathbf{G}^{(k)} \right)^{-1} \nabla_{\boldsymbol{\theta}} f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)}) - \hat{\boldsymbol{\theta}} \right\|_{\mathbf{G}^{(k)}}^2 \\ &= \left\| \boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}} \right\|_{\mathbf{G}^{(k)}}^2 - 2\alpha^{(k)} \left\langle \nabla_{\boldsymbol{\theta}} f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)}), \boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}} \right\rangle \\ &\quad + \left( \alpha^{(k)} \right)^2 \left\| \left( \mathbf{G}^{(k)} \right)^{-1} \nabla_{\boldsymbol{\theta}} f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)}) \right\|_{\mathbf{G}^{(k)}}^2. \end{aligned}$$

Therefore,

$$\begin{aligned} \nabla_{\boldsymbol{\theta}}^\top f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)}) (\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}) &= \frac{1}{2\alpha^{(k)}} \left( \|\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}\|_{\mathbf{G}^{(k)}}^2 - \|\boldsymbol{\theta}^{(k+1)} - \hat{\boldsymbol{\theta}}\|_{\mathbf{G}^{(k)}}^2 \right) \\ &\quad + \frac{\alpha}{2} \left\| \left( \mathbf{G}^{(k)} \right)^{-1} \nabla_{\boldsymbol{\theta}} f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)}) \right\|_{\mathbf{G}^{(k)}}^2 \end{aligned}$$

We sum up from  $k = 0$ , to the current step  $K$ , and taking its average, and since

$$f\left(\frac{1}{K+1} \sum_{k=0}^K \mathbf{x}^{(k)}, \frac{1}{K+1} \sum_{k=0}^K y^{(k)}; \frac{1}{K+1} \sum_{k=0}^K \boldsymbol{\theta}^{(k)}\right) \leq \frac{1}{K+1} \sum_{k=0}^K f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)}),$$

and supposedly fixed  $\alpha$ , i.e.,  $\alpha^{(k)} = \bar{\alpha} \forall k \in [0, K]$ , then

$$\begin{aligned} &f\left(\frac{1}{K+1} \sum_{k=0}^K \mathbf{x}^{(k)}, \frac{1}{K+1} \sum_{k=0}^K y^{(k)}; \frac{1}{K+1} \sum_{k=0}^K \boldsymbol{\theta}^{(k)}\right) \\ &\leq \frac{1}{t} \sum_{k=0}^K \left[ \frac{1}{2\alpha^{(k)}} \left( \|\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}\|_{\mathbf{G}^{(k)}}^2 - \|\boldsymbol{\theta}^{(k+1)} - \hat{\boldsymbol{\theta}}\|_{\mathbf{G}^{(k)}}^2 \right) \right. \\ &\quad \left. + \frac{\alpha^{(k)}}{2} \left\| \left( \mathbf{G}^{(k)} \right)^{-1} \nabla_{\boldsymbol{\theta}} f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)}) \right\|_{\mathbf{G}^{(k)}}^2 \right] \\ &= \frac{1}{2\bar{\alpha}t} \sum_{k=0}^K \left( \|\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}\|_{\mathbf{G}^{(k)}}^2 - \|\boldsymbol{\theta}^{(k+1)} - \hat{\boldsymbol{\theta}}\|_{\mathbf{G}^{(k)}}^2 \right) \\ &\quad + \frac{\alpha}{2K} \sum_{k=0}^K \left\| \left( \mathbf{G}^{(k)} \right)^{-1} \nabla_{\boldsymbol{\theta}} f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)}) \right\|_{\mathbf{G}^{(k)}}^2. \end{aligned}$$

---

<sup>1</sup>  $\|\mathbf{v}\|_{\mathbf{A}}^2 = \mathbf{v}^\top \mathbf{A} \mathbf{v} = \|\mathbf{A}^{1/2} \mathbf{v}\|_2^2$

151 By Elad Hazan's lemma [2], i.e.,

$$152 \quad \sum_{i=1}^b \left\| \left( \mathbf{G}^{(i)} \right)^{-1} \nabla_{\boldsymbol{\theta}} f(\mathbf{x}^{(i)}, y^{(i)}; \boldsymbol{\theta}^{(k)}) \right\|_{\mathbf{G}^{(k)}}^2 \leq 2\text{Tr}(\mathbf{G}^{(b)}),$$

153 with  $\mathbf{G}^{(b)}$  is defined as (2.2), and terms in the summation

$$154 \quad \sum_{k=0}^K \left( \|\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}\|_{\mathbf{G}^{(k)}}^2 - \|\boldsymbol{\theta}^{(k+1)} - \hat{\boldsymbol{\theta}}\|_{\mathbf{G}^{(k)}}^2 \right)$$

155 cancel out each other, i.e.,

$$156 \quad \sum_{k=0}^K \left( \|\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}\|_{\mathbf{G}^{(k)}}^2 - \|\boldsymbol{\theta}^{(k+1)} - \hat{\boldsymbol{\theta}}\|_{\mathbf{G}^{(k)}}^2 \right) = \|\boldsymbol{\theta}^{(0)} - \hat{\boldsymbol{\theta}}\|_{\mathbf{G}^{(0)}}^2 - \|\boldsymbol{\theta}^{(K+1)} - \hat{\boldsymbol{\theta}}\|_{\mathbf{G}^{(K)}}^2$$

157 Therefore, we have

$$158 \quad f\left(\frac{1}{K+1} \sum_{k=0}^K \mathbf{x}^{(k)}, \frac{1}{K+1} \sum_{k=0}^K y^{(k)}, \frac{1}{K+1} \sum_{k=0}^K \boldsymbol{\theta}^{(k)}\right) \\ 159 \quad \leq \|\boldsymbol{\theta}^{(0)} - \hat{\boldsymbol{\theta}}\|_{\mathbf{G}^{(0)}}^2 - \|\boldsymbol{\theta}^{(K+1)} - \hat{\boldsymbol{\theta}}\|_{\mathbf{G}^{(K)}}^2 + 2\text{Tr}(\mathbf{G}^{(b)})$$

160 **2.3. AdaDelta.** Idea: Accumulate over a restricted time window.

161 Let

$$162 \quad \mathbb{E}_t\{\mathbf{g}^\top \mathbf{g}\}_t = p\mathbb{E}_{t-1}\{\mathbf{g}^\top \mathbf{g}\} + (1-p)\mathbf{g}_t^\top \mathbf{g}_t,$$

163 then the updated parameters is followed as

$$164 \quad \boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \alpha^{(k)} \left( \mathbb{E}_t\{\mathbf{g}^\top \mathbf{g}\} + \epsilon \mathbf{I} \right)^{-1/2} \mathbf{g}^{(k)},$$

165 with  $\mathbf{g}^{(k)} \triangleq \nabla_{\boldsymbol{\theta}} f(\mathbf{x}^{(k)}, y^{(k)}; \boldsymbol{\theta}^{(k)})$ .

166 Another idea is picking  $\alpha^{(k)}$  to imitate Newton method, i.e., from the Newton  
167 method, we have

$$168 \quad \boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \alpha^{(k)} \left( \mathbf{H}^{(k)} \right)^{-1} \mathbf{g}^{(k)}$$

169 , where  $\mathbf{H}^{(k)} \triangleq \nabla_{\boldsymbol{\theta}}^2 g^{(k)}$ . Notice that the vector  $(\mathbf{H}^{(k)})^{-1} \mathbf{g}^{(k)}$  can be estimated as

170  $\frac{\Delta \boldsymbol{\theta}^{(k)}}{\Delta \mathbf{g}^{(k)}}$ , so  $\alpha^{(k)}$  can be suggested as  $\|\Delta \boldsymbol{\theta}^{(k)}\|$ . From this idea, we pick

$$171 \quad \alpha^{(k)} = \sqrt{\mathbb{E}_{k-1}\{\|\Delta \boldsymbol{\theta}\|^2\}},$$

172 where  $\mathbb{E}_k\{\|\Delta \boldsymbol{\theta}\|^2\} = p\mathbb{E}_{k-1}\{\|\Delta \boldsymbol{\theta}\|^2\} + (1-p)\mathbb{E}\{\|\Delta \boldsymbol{\theta}^{(k)}\|^2\}$ . Putting together, we  
173 have the updated formulaas follows

$$174 \quad \boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \sqrt{\mathbb{E}_{k-1}\{\|\Delta \boldsymbol{\theta}\|^2\}} \left( \mathbb{E}_t\{\mathbf{g}^\top \mathbf{g} + \epsilon \mathbf{I}\} \right)^{-1/2} \mathbf{g}^{(k)}$$

175 **2.4. Root Mean Square Propagation.** RMSProp [3] was provided by Geof-  
176 frey Hinton, where at every iteration, the new parameter is updated as:

$$177 \quad \boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \gamma \left( \mathbb{E}_t\{\mathbf{g}^\top \mathbf{g}\} + \epsilon \mathbf{I} \right)^{-1} \mathbf{g}^{(k)}$$

### 3. Accelerated descent methods.

**3.1. Issue with (Stochastic) Gradient Descent.** Get stuck at the local minimum.

**3.2. Gradient descent with Momentum.** It is also called the Polyak Heavy Ball method.

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \alpha^{(k)} \mathbf{g}^{(k)} + \mu^{(k)} (\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^{(k-1)})$$

*Issue.* Let  $f(x, y; \theta)$  is defined by its gradient as

$$\nabla f_{\theta}(x, y, \theta) = \begin{cases} 25\theta & \text{if } \theta < 1, \\ \theta + 24 & \text{if } 1 \leq \theta < 2, \\ 25\theta - 24 & \text{if } \theta \geq 2. \end{cases}$$

Then the below graph displays the solution of the Polyak Heavy Ball method over 50 iterations, with  $\alpha = \frac{1}{9}$ ,  $\mu = \frac{4}{9}$ , and initial  $x_0 = 3.3$ .

### 3.3. Nestorov momentum.

$$\begin{aligned} \boldsymbol{\theta}^{(k+1)} &= \boldsymbol{\theta}^{(k)} - \mu (\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^{(k-1)}) \\ &\quad - \gamma \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^{(k)} + \mu (\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^{(k-1)})) \end{aligned}$$

**3.4. Optimality of Nestorov momentum.** Theorem 2.1.7 of [5] said that for any  $\boldsymbol{\theta}_0 \in \mathbb{R}^D$  and  $k \leq \frac{d-1}{2}$ , there exists an  $L$ -smooth convex function  $f(\cdot)$  such that for any first-order algorithm the output sequence satisfies

$$f(\boldsymbol{\theta}^{(k)}) - f(\hat{\boldsymbol{\theta}}) \geq \frac{3L\|\boldsymbol{\theta}_0 - \hat{\boldsymbol{\theta}}\|_2^2}{32(k+1)^2}.$$

The goal is now to have an algorithm with

$$\begin{aligned} |f(\boldsymbol{\theta}^{(k)}) - f(\hat{\boldsymbol{\theta}})| &= \mathcal{O}\left(\frac{L}{k^2}\right) \\ \Leftrightarrow |f(\boldsymbol{\theta}^{(k)}) - f(\hat{\boldsymbol{\theta}})| &\leq \epsilon \quad \text{for } k \geq \sqrt{\frac{L}{\epsilon}}. \end{aligned}$$

Nestorov proved if  $f(\boldsymbol{\theta}^{(k)}) \leq f(\boldsymbol{\zeta}^{(k)}) - \frac{L}{2} \|\nabla f(\boldsymbol{\zeta}^{(k)})\|_2^2$ , we have the optimal convergence rate  $\mathcal{O}\left(\frac{L}{k^2}\right)$

**3.5. Accelerated Gradient Descent.** The update formula of accelerated Gradient Descent is given as

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\zeta}^{(k)} - \epsilon \nabla f(\boldsymbol{\zeta}^{(k)}),$$

where  $\boldsymbol{\zeta}^{(k)} = \boldsymbol{\theta}^{(k)} + \frac{k+1}{k+2}(\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^{(k-1)})$ . The illustration of updating process is demonstrated in Figure 1 (based on the blog of Wibisono<sup>2</sup>).

When  $\|\nabla^2 f\| \leq \frac{1}{\epsilon}$ , the optimal rate of accelerated gradient descent is  $\mathcal{O}\left(\frac{1}{\epsilon k^2}\right)$  while gradient descent only achieve the optimal rate of  $\mathcal{O}\left(\frac{1}{\epsilon k}\right)$ .

<sup>2</sup><https://awibisono.github.io/2016/06/20/accelerated-gradient-descent.html>



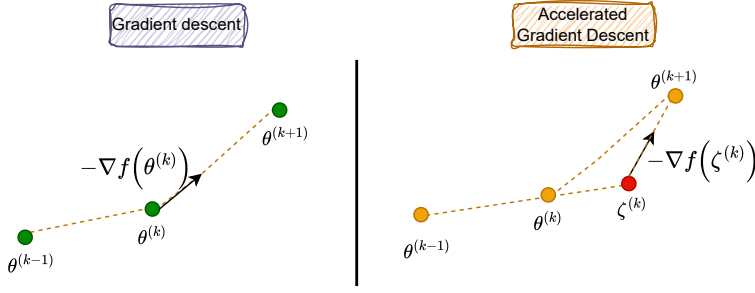


FIG. 1. Illustration of parameter update of Gradient Descent versus Accelerated GD.

206 *Proof of convergence rate.* Consider a Lyapunov function<sup>3</sup>  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . Let  
 207  $\alpha \triangleq \epsilon, \beta \triangleq \frac{k+1}{k+2}$ , from accelerated GD method we have

$$\begin{aligned} 208 \quad \theta^{(k+1)} &= \zeta^{(k)} - \alpha \nabla f(\zeta^{(k)}) \\ 209 \quad \theta^{(k)} &= \zeta^{(k-1)} - \alpha \nabla f(\zeta^{(k-1)}) \end{aligned}$$

210 Subtracting both sides of the two equalities above for the optimal point  $\hat{\theta}$ , we have

$$211 \quad \begin{bmatrix} \theta^{(k+1)} - \hat{\theta} \\ \theta^{(k)} - \hat{\theta} \end{bmatrix} = \mathbf{T} \begin{bmatrix} \theta^{(k)} - \hat{\theta} \\ \theta^{(k-1)} - \hat{\theta} \end{bmatrix}.$$

212 **Lemma:**  $p(\mathbf{T}) := \max_{\mathbf{T}\theta = \lambda\theta} |\lambda|$ .

213 Then  $p(\mathbf{T}) < p \Leftrightarrow \exists$  a positive semidefinite matrix  $\mathbf{K}$  such that

$$214 \quad p^2 \mathbf{K} - \mathbf{T}^\top \mathbf{K} \mathbf{T} \geq 0$$

215 For  $\mathbf{T}$  and corresponding  $\mathbf{K}$ , define

$$216 \quad V(\theta) = \theta^\top \mathbf{Q} \theta.$$

217 We analyze the decrease of  $V(\theta)$  on  $\{\theta^{(t)}\}_{t=0}^k$ . By using the lemma above, we can  
 218 show that

$$219 \quad \|\theta^{(k)} - \hat{\theta}\| \leq \sqrt{2\text{cond}(\mathbf{Q})} p(\mathbf{T})^k \|\theta^{(0)} - \hat{\theta}\|.$$

220 The optimal values of  $\alpha, \beta$  are given as

$$221 \quad \alpha = \left( \frac{2}{\sqrt{\sigma_1} + \sqrt{\sigma_n}} \right)^2, \quad \beta = \left( \frac{\sqrt{\sigma_1} - \sqrt{\sigma_n}}{\sqrt{\sigma_1} + \sqrt{\sigma_n}} \right)^2,$$

222 where  $\sigma_1, \sigma_n$  is the largest and smallest singular values. In practice, these values are  
 223 chosen as fixed  $\beta = 0.99$  and tune value of  $\alpha$ .

224 **4. Adaptive Moment Estimation (ADAM).** The Adaptive Moment Esti-  
 225 mation (ADAM) [4] combines all advantages of the momentum method, and the  
 226 AdaGrad, RMSProp.

227 **5. Exercises.**

<sup>3</sup>[https://en.wikipedia.org/wiki/Lyapunov\\_function](https://en.wikipedia.org/wiki/Lyapunov_function)

**Algorithm 4.1** Adam: Adaptive Moment Estimation**Require:**  $\alpha$ : Stepsize**Require:**  $\beta_1, \beta_2 \in [0, 1)$ : Exponential decay rates for moment estimates**Require:**  $f(\theta)$ : Stochastic objective function with parameters  $\theta$ **Require:** Initial parameters  $\theta_0$ 


---

```

1:  $m_0 \leftarrow 0$  {Initialize 1st moment vector}
2:  $v_0 \leftarrow 0$  {Initialize 2nd moment vector}
3:  $t \leftarrow 0$  {Initialize timestep}
4: while  $\theta_t$  not converged do
5:    $t \leftarrow t + 1$ 
6:    $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$  {Compute gradients}
7:    $m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$  {Update biased 1st moment}
8:    $v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$  {Update biased 2nd moment}
9:    $\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$  {Bias-corrected 1st moment}
10:   $\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$  {Bias-corrected 2nd moment}
11:   $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$  {Update parameters}
12: end while
13: return  $\theta_t$  {Resulting parameters}

```

---

*Exercise 1.* Suppose we have  $n$  data points with each having 100 features and unit norm:  $\mathbf{x}_i \in \mathbb{R}^{100}, \|\mathbf{x}_i\|_2 = 1, \forall i \in [n]$ . We want to classify these  $n$  data points into 5 categories. We'll use logistic regression for this purpose. Analyze the number of expected steps for stochastic gradient descent to make the loss go below 0.01.

**Solution:**

Let the softmax logistic regression loss be defined as:

$$L(w) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, w^T x_i),$$

where  $\ell$  is the cross-entropy loss between the predicted probabilities and the true labels  $y_i \in \{1, 2, 3, 4, 5\}$ . Let assume: The gradient  $\nabla L(w)$  is  $L$ -Lipschitz continuous. We use a constant learning rate  $\eta$ . The stochastic gradients have bounded variance:  $\mathbb{E}[\|\nabla L(w) - \nabla \tilde{L}(w)\|^2] \leq \sigma^2$ . Desired squared gradient norm  $\mathbb{E}[\|\nabla L(w)\|^2] \leq \varepsilon = 0.01$ .

From the convergence theory of SGD for non-convex objectives, the number of iterations  $T$  required to achieve this level of accuracy is:

$$T = \mathcal{O}\left(\frac{1}{\varepsilon^2}\right)$$

Substituting the target accuracy:

$$T = \mathcal{O}\left(\frac{1}{0.01^2}\right) = \mathcal{O}(10^4)$$

Therefore, SGD is expected to require approximately 10,000 iterations to reduce the loss gradient norm below 0.01.

*Exercise 2.* Suppose we have  $n$  data points  $\mathbf{x}_i \in \mathbb{R}^{100}, \|\mathbf{x}_i\| = 1, \forall i \in [n]$  and an objective function  $h(\cdot)$  is written as

$$h(\mathbf{x}) = \sum_{i=1}^n f(\mathbf{x}_i).$$

That is  $f(\cdot)$ , a specific function is applied to the data point  $\mathbf{x}_i$  separately, one by one.  
We assume Gaussian noise in our computation, that is

$$(\nabla f(\mathbf{x}_i))_j = (\nabla f(\mathbf{x}_i))_j + (\epsilon_i)_j, \quad \forall i \in [n], j \in [100],$$

where  $(\epsilon_j)_j \sim \mathcal{N}(0, \sigma^2)$ . We sample  $m$  data points from our database uniformly randomly without replacement, suppose this set of data points is called  $\mathcal{S}$ . Let

$$g(\mathbf{x}) = \frac{1}{m} \sum_{i \in \mathcal{S}} \nabla f(\mathbf{x}_i).$$

Please derive a bound on the error term

$$\|\nabla f(\mathbf{x})\| - \mathbb{E}_{\mathcal{S}}\{g(\mathbf{x})\}.$$

**Solution:** By linearity of expectation, we have:

$$\mathbb{E}_{\mathcal{S}}[g(x)] = \frac{1}{m} \sum_{i \in \mathcal{S}} \mathbb{E}[\nabla f(x_i)] = \nabla f(x)$$

So the expected bias is zero, and we compute the variance of the estimator:

$$\mathbb{E}_{\mathcal{S}} \left[ \|\nabla f(x) - g(x)\|^2 \right] = \mathbb{E}_{\mathcal{S}} \left[ \left\| \frac{1}{m} \sum_{i \in \mathcal{S}} \xi_i \right\|^2 \right]$$

Each noise vector  $\xi_i \in \mathbb{R}^{100}$  has components  $\xi_{ij} \sim \mathcal{N}(0, \sigma^2)$ , so:

$$\mathbb{E} \left[ \|\xi_i\|^2 \right] = \sum_{j=1}^{100} \mathbb{E}[\xi_{ij}^2] = 100 \cdot \sigma^2$$

Since we are averaging  $m$  i.i.d. noise vectors:

$$\text{Var}(g(x)) = \frac{1}{m^2} \cdot m \cdot \mathbb{E}[\|\xi_i\|^2] = \frac{100\sigma^2}{m}$$

Taking the square root to get the standard deviation (expected deviation):

$$\mathbb{E} [\|\nabla f(x) - g(x)\|] \leq \sqrt{\frac{100 \cdot \sigma^2}{m}} = \frac{10\sigma}{\sqrt{m}}$$

Thus, the expected error norm is bounded by:

$$\mathbb{E}_{\mathcal{S}} [\|\nabla f(x) - g(x)\|] \leq \frac{10\sigma}{\sqrt{m}}.$$

## REFERENCES

- [1] J. DUCHI, E. HAZAN, AND Y. SINGER, Adaptive subgradient methods for online learning and stochastic optimization., Journal of machine learning research, 12 (2011).
- [2] E. HAZAN ET AL., Introduction to online convex optimization, Foundations and Trends® in Optimization, 2 (2016), pp. 157–325.
- [3] G. HINTON, N. SRIVASTAVA, AND K. SWERSKY, Neural networks for machine learning lecture 6a overview of mini-batch gradient descent, Cited on, 14 (2012), p. 2.
- [4] D. P. KINGMA AND J. BA, Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980, (2014).
- [5] Y. NESTEROV ET AL., Lectures on convex optimization, vol. 137, Springer, 2018.
- [6] T. STROHMER AND R. VERSHYNIN, A randomized kaczmarz algorithm with exponential convergence, Journal of Fourier Analysis and Applications, 15 (2009), pp. 262–278.