

SGD: Stochastic Gradient Methods

Jingchang Liu

University of Science and Technology of China

April 6, 2017

Table of Contents

Stochastic Gradient Methods

Lipschitz-Continuous Gradient

Analyses of SGD

Table of Contents

Stochastic Gradient Methods

Lipschitz-Continuous Gradient

Analyses of SGD

Motivation

Empirical Risk minimize

- ▶ Regression: $F(w) = \frac{1}{n} \sum_{i=1}^n (f(x_i; w) - y_i)^2 \triangleq \frac{1}{n} \sum_{i=1}^n f_i(w)$
- ▶ Classification: $F(w) = \frac{1}{n} \sum_{k=1}^n I(f(x_i; w) = y_i) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(w)$
- ▶ Neural Networks: $E = \frac{1}{n} \sum_{k=1}^n E_k = \frac{1}{n} \sum_{k=1}^n \left[\frac{1}{2} \sum_{j=1}^l (\hat{y}_j^k - y_j^k)^2 \right]$

Regularization

- ▶ $F(w) = \frac{1}{n} \sum_{i=1}^n (f(x_i; w) - y_i)^2 + \frac{\lambda}{2} \|w\|^2 \triangleq \frac{1}{n} \sum_{i=1}^n f_i(w)$

Algorithm 1 Stochastic Gradient(SG) Method

- 1: Choose an initial iterate w_1 .
 - 2: **for** $k = 1, 2, \dots$ **do**
 - 3: Generate a realization of the random variable i_k .
 - 4: Compute a stochastic vector $\nabla_{i_k} f(w_k)$
 - 5: Choose a stepsize $\alpha_k > 0$
 - 6: Set the new iterate as $w_{k+1} \leftarrow w_k - \alpha_k \nabla_{i_k} f(w_k)$
 - 7: **end for**
-

Gradient Methods

Full Gradient

$$w_{k+1} \leftarrow w_k + \alpha_k \nabla F(w) = w_k - \frac{\alpha_k}{n} \sum_{i=1}^n \nabla f_i(w_k)$$

Stochastic Gradient

$$w_{k+1} \leftarrow w_k + \alpha_k \nabla f_{l_k}(w)$$

Batch Gradient

$$w_{k+1} \leftarrow w_k - \frac{\alpha_k}{|S_k|} \sum_{i \in S_k} \nabla f_i(w_k)$$

Gradient Methods

Formulation

$$\min_w F(w) = \min_w \frac{1}{n} \sum_{i=1}^n f_i(w)$$

Gradient

$$g(w_k, \xi_k) = \begin{cases} \nabla f(w_k; \xi_k) & \text{Stochastic Gradient} \\ \frac{1}{|S_k|} \sum_{i=1}^{|S_k|} \nabla f(w_k; \xi_{k,i}) & \text{Batch Gradient} \\ \nabla F(w_k) & \text{Full Gradient} \end{cases}$$

Table of Contents

Stochastic Gradient Methods

Lipschitz-Continuous Gradient

Analyses of SGD

Assumption

Lipschitz-continuous gradient

The objective function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is continuously differentiable and the gradient function of f , namely, $\nabla f : \mathbb{R}^d \rightarrow \mathbb{R}^d$, is Lipschitz continuous with Lipschitz constant $L > 0$, i.e.,

$$\|\nabla f(w) - \nabla f(\bar{w})\|_2 \leq L \|w - \bar{w}\|_2 \quad \text{for all } \{w, \bar{w}\} \subset \mathbb{R}^d \quad (1)$$

Strongly convex

f is strongly convex with parameter $c > 0$ if

$$g(x) = f(x) - \frac{c}{2} x^T x \quad \text{is convex}$$

Equal Properties

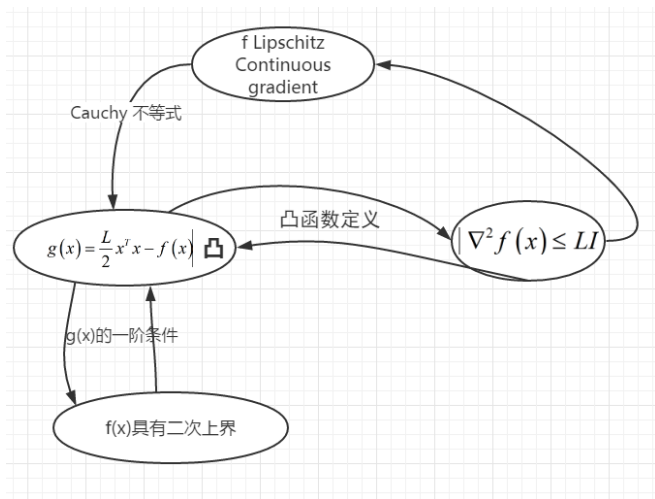


Figure: Equal Properties introduced from Lipschitz

Tips 1

1.

$$\|\nabla f(x) - \nabla f(y)\|_2 \leq L \|x - y\|_2$$

2.

$$\begin{aligned} (\nabla f(x) - \nabla f(y))^T (x - y) &\leq \|\nabla f(x) - \nabla f(y)\|_2 \|x - y\|_2 \\ &\leq L \|x - y\|_2^2 \end{aligned}$$

3.

$$\text{def. } g(x) = \frac{L}{2} x^T x - f(x) \quad \text{Then } \nabla g(x) = Lx - \nabla f(x)$$

Monotonicity of gradient

A differentiable function g is convex if and only if $\text{dom } g$ is convex and

$$(\nabla g(x) - \nabla g(y))^T (x - y) \geq 0 \quad \text{for all } x, y \in \text{dom } g$$

Equal Properties

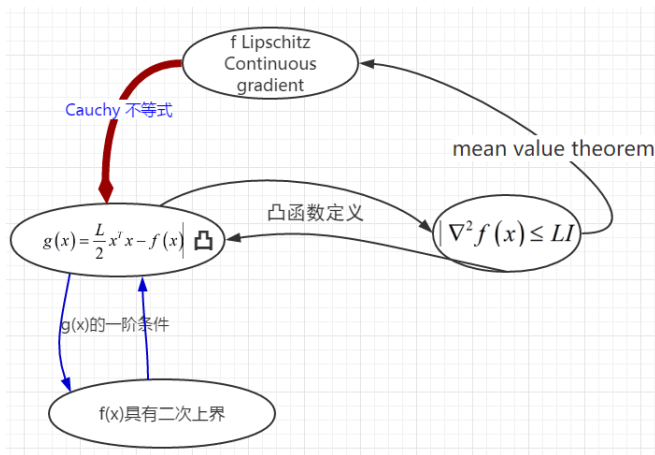


Figure: Equal Properties introduced from Lipschitz

Tips 2

First-order Condition of Convex Function g

$$g(y) \geq g(x) + \nabla g(x)^T (y - x) \quad \text{for all } x, y \in \text{dom } g$$

Quadratic upper bound on f

$$g(x) = \frac{L}{2} x^T x - f(x) \quad \text{Then} \quad \nabla g(x) = Lx - \nabla f(x)$$

$$f(y) \leq f(x) + \nabla f(x)^T (y - x) + \frac{L}{2} \|y - x\|_2^2$$

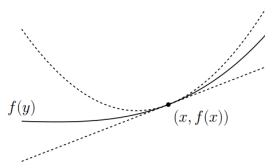


Figure: Quadratic upper bound

Equal Properties

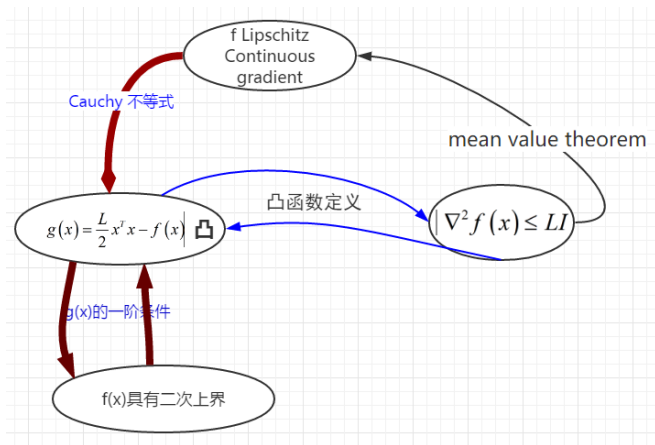


Figure: Equal Properties introduced from Lipschitz

Tip3

Second-order Condition of Convex Function g

$$\nabla^2 g(x) \succeq 0$$

Second-order Gradient of f

$$\nabla g(x) = Lx - \nabla f(x)$$

$$\nabla^2 g(x) = L - \nabla^2 f(x)$$

$$\nabla^2 f(x) \preceq L$$

Equal Properties

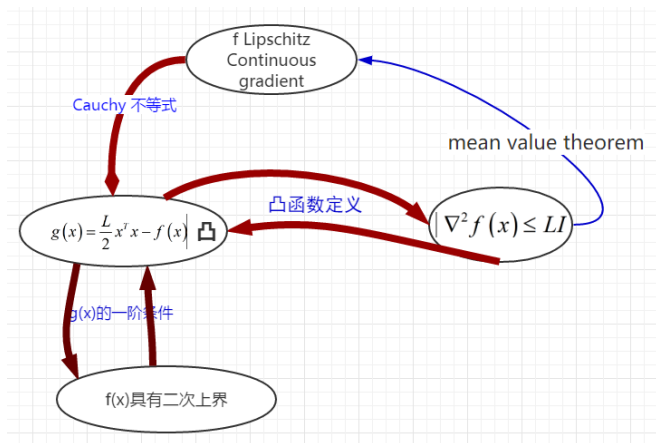


Figure: Equal Properties introduced from Lipschitz

Tips 4

Mean value theorem of f'

A function f is continuous on the closed interval $[a, b]$, and differentiable on the open interval (a, b) , then there exists a point c in (a, b) such that:

$$f'(c) = \frac{f(b) - f(a)}{b - a}$$

Mean value theorem of $\nabla^2 f$

$$\frac{\nabla f(x) - \nabla f(y)}{x - y} = \nabla^2 f(\zeta) \leq L$$

Equal Properties

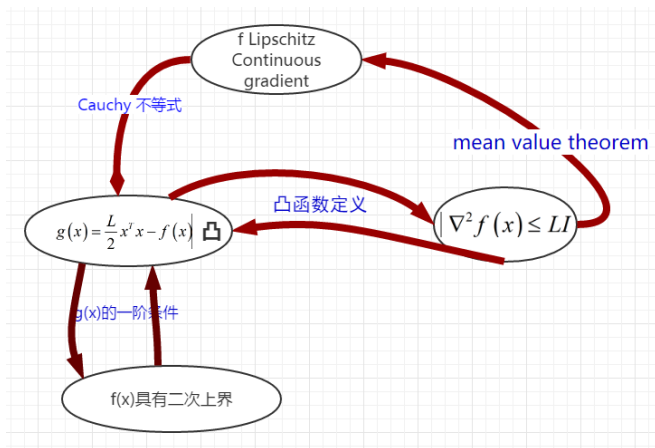


Figure: Equal Properties introduced from Lipschitz

Table of Contents

Stochastic Gradient Methods

Lipschitz-Continuous Gradient

Analyses of SGD

Theorem

Th

Under *Lipschitz-continuous* Assumption, the iterations of SG satisfy the following inequality for all $k \in N$:

$$\begin{aligned} E_{\xi_k} [F(w_{k+1})] - F(w_k) &\leq -\alpha_k \nabla F(w_k)^T E_{\xi_k} [g(w_k, \xi_k)] \\ &\quad + \frac{1}{2} \alpha_k^2 L E_{\xi_k} [\|g(w_k, \xi_k)\|_2^2] \end{aligned}$$

Proof

By Lipschitz-continuous Assumption, the iterates generated by SG satisfy

$$\begin{aligned} F(w_{k+1}) - F(w_k) &\leq \nabla F(w_k)^T (w_{k+1} - w_k) + \frac{1}{2} L \|w_{k+1} - w_k\|_2^2 \\ &\leq -\alpha_k \nabla F(w_k)^T g(w_k, \xi_k) + \frac{1}{2} \alpha_k^2 L \|g(w_k, \xi_k)\|_2^2 \end{aligned}$$

Noting that w_{k+1} but not w_k depends on ξ_k . Take expectation with respect to ξ_k , we'll get

$$\begin{aligned} E_{\xi_k} [F(w_{k+1})] - F(w_k) &\leq -\alpha_k \nabla F(w_k)^T E_{\xi_k} [g(w_k, \xi_k)] \\ &\quad + \frac{1}{2} \alpha_k^2 L E_{\xi_k} [\|g(w_k, \xi_k)\|_2^2] \end{aligned}$$

Variance Reduce

Trade-Offs of Mini-Batch

$$g(w_k, \xi_k) = \frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(w_k)$$

$$\mathbb{V}_{\xi_k} [g(w_k, \xi_k)] = \frac{\mathbb{V}_{\xi_k} [\nabla_{i_k} f(w_k)]}{|S_k|}$$

Popular Variance Reduce Methods

SAG, SAGA, SVRG

Convergence Rate of Full Gradient

Th

Assume F is convex and L -Lipschitz continuous gradient. With step size $\alpha = \frac{1}{L}$. Then

$$F(x_{k+1}) - F^* \leq \frac{2L \|x_1 - x^*\|_2^2}{k}$$

Th

Assume F is c -strongly convex and L -Lipschitz continuous gradient. With step size $\alpha = \frac{2}{M+L}$. Then

$$F(x_{k+1}) - F^* \leq \frac{L}{2} \exp\left(-\frac{4k}{Q+1}\right) \|x_1 - x^*\|_2^2$$

where $Q = \frac{L}{c}$

Convergence Rate of SG, Fixed Stepsize

Th

Assume F is *c-strongly convex* and *L-Lipschitz continuous gradient*. Under some mild assumption. Suppose that the SG method is run with a *fixed stepsize*, $\alpha_k = \bar{\alpha}$ for all $k \in \mathbb{N}$, satisfying

$$0 < \bar{\alpha} \leq \frac{\mu}{LM_G}$$

Then, the expected optimality gap satisfies the following inequality for all $k \in \mathbb{N}$:

$$\mathbb{E}[F(w_k) - F_*] \leq \frac{\alpha \bar{L} M}{2c\mu} + (1 - \bar{\alpha} c \mu)^{k-1} \left(F(w_1) - F_* - \frac{\bar{\alpha} L M}{2c\mu} \right)$$
$$\xrightarrow{k \rightarrow \infty} \frac{\alpha \bar{L} M}{2c\mu}$$

Convergence Rate of SG, Diminishing Stepsize

Assume F is c -strongly convex and L -Lipschitz continuous gradient. Under some mild assumption. Suppose that the SG method is run with a stepsize sequence such that, for all $k \in \mathbb{N}$,

$$\alpha_k = \frac{\beta}{\gamma + k} \quad \text{for some } \beta > \frac{1}{c\mu} \text{ and } \gamma > 0 \quad \text{such that } \alpha_1 \leq \frac{\mu}{LM_G}$$

Then, for all $k \in \mathbb{N}$, the expected optimality gap satisfies

$$\mathbb{E}[F(w_k) - F_*] \leq \frac{v}{\gamma + k}$$

where

$$v := \max \left\{ \frac{\beta^2 LM}{2(\beta c\mu - 1)}, (\gamma + 1)(F(w_1) - F_*) \right\}$$

Q&A