

Statistical Machine Learning

Instructor: Jie Shen

Dept. of Computer Science

February 14, 2020

Gradient Descent

Consider minimization without constraints:

$$\min_{\mathbf{w}} F(\mathbf{w}), \mathbf{w} \in \mathbb{R}^d$$

Gradient Descent:

1. Initialize \mathbf{w}^0 arbitrarily, e.g. $\mathbf{w}^0 = \mathbf{0}$
2. For $t = 1, 2, \dots$

$$\mathbf{w}^t = \mathbf{w}^{t-1} - \eta_t \nabla F(\mathbf{w}^{t-1}) \quad (1)$$

Goal:

- $\mathbf{w}^t \rightarrow \mathbf{w}^*$, where $\mathbf{w}^* = \arg \min F(\mathbf{w})$
- in few iterations (cheap computation)

Why GD “decreases” objective value (under proper conditions)?

Smooth: $F(\mathbf{w})$ is smooth if for any $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^d$

$$\|\nabla F(\mathbf{w}_2) - \nabla F(\mathbf{w}_1)\|_2 \leq L \|\mathbf{w}_2 - \mathbf{w}_1\|_2$$

Examples:

When GD fails to find global optimum

- when it terminates
- when it gets stuck at a non-optimal point

The Big Picture

Convex set: $\mathcal{C} \subset \mathbb{R}^d$ is said to be a convex set if for any $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ and any $0 \leq \lambda \leq 1$, $\lambda \mathbf{u} + (1 - \lambda) \mathbf{v} \in \mathcal{C}$

- illustration, examples

Convex Function

Convex function: $F(\mathbf{w})$ is said to be a convex function if the set $\mathcal{E} = \{(\mathbf{w}, y) : y \geq F(\mathbf{w})\}$ is convex

- \mathcal{E} is called the **epigraph** of $F(\mathbf{w})$
- illustration

Characterization of Convex Functions

Theorem 1

Suppose that $F : \mathbb{R}^d \rightarrow \mathbb{R}$ is twice differentiable. The following are equivalent:

- ① F is convex;
- ② $F(\mathbf{w}_2) \geq F(\mathbf{w}_1) + \langle \nabla F(\mathbf{w}_1), \mathbf{w}_2 - \mathbf{w}_1 \rangle$;
- ③ $\nabla^2 F(\mathbf{w})$ is positive semi-definite.

Typically use 3 to check the convexity.

Let f and h be convex functions.

- $a \cdot f + b \cdot h$ is convex when $a \geq 0$ and $b \geq 0$
- $f(h)$ may NOT be convex

$$\min_{\mathbf{w}} F(\mathbf{w}), \quad \mathbf{w} \in \mathcal{C}.$$

Convex Program: both $F(\mathbf{w})$ and \mathcal{C} are convex

- optimality: local optimum \iff global optimum
- works well
- easy to solve

Theorem 2

Suppose $F(\mathbf{w})$ is convex and L -smooth. Pick $0 < \eta \leq 1/L$. Then for all $t \geq 1$,

$$F(\mathbf{w}^t) - F(\mathbf{w}^*) \leq \frac{\|\mathbf{w}^0 - \mathbf{w}^*\|_2^2}{2\eta} \cdot \frac{1}{t}$$

In particular, picking $\eta = 1/L$ gives

$$F(\mathbf{w}^t) - F(\mathbf{w}^*) \leq \frac{2L \|\mathbf{w}^0 - \mathbf{w}^*\|_2^2}{t}.$$

Implications

1. $F(\mathbf{w}^t) - F(\mathbf{w}^*)$ v.s. $\|\mathbf{w}^t - \mathbf{w}^*\|_2$

2. Iteration complexity

3. Estimate L

Faster Rate of Convergence

Strongly Convex: for any $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^d$

$$\|\nabla F(\mathbf{w}_2) - \nabla F(\mathbf{w}_1)\|_2 \geq \alpha \|\mathbf{w}_2 - \mathbf{w}_1\|_2$$

- functions satisfying SC
- not satisfying

Theorem 3

Suppose $F(\mathbf{w})$ is α -strongly convex and L -smooth. Let $\{\mathbf{w}^t\}_{t \geq 1}$ be the iterates generated by GD where $0 < \eta \leq 2/(\alpha + L)$. Then for all $t \geq 1$,

$$\|\mathbf{w}^t - \mathbf{w}^*\|_2 \leq \sqrt{1 - \frac{2\eta\alpha L}{\alpha + L}} \|\mathbf{w}^{t-1} - \mathbf{w}^*\|_2.$$

In particular, picking $\eta = 2/(\alpha + L)$ gives

$$\|\mathbf{w}^t - \mathbf{w}^*\|_2 \leq \left(1 - \frac{2}{c + 1}\right) \|\mathbf{w}^{t-1} - \mathbf{w}^*\|_2$$

where $c \stackrel{\text{def}}{=} L/\alpha$ is the *condition number*.

- converges linearly / geometric rate of convergence
- typically the best one can hope

Implications

For all $t \geq 1$,

$$\begin{aligned}\|\mathbf{w}^t - \mathbf{w}^*\|_2 &\leq \left(1 - \frac{2}{c+1}\right)^t \|\mathbf{w}^0 - \mathbf{w}^*\|_2 \\ &\leq e^{-\frac{2t}{c+1}} \|\mathbf{w}^0 - \mathbf{w}^*\|_2 \quad (\text{by } 1 + x \leq e^x)\end{aligned}$$

For any pre-defined error $0 < \epsilon < 1$,

$$\underbrace{t = c \log \frac{\|\mathbf{w}^0 - \mathbf{w}^*\|_2}{\epsilon}}_{\text{iteration complexity}} \implies \|\mathbf{w}^t - \mathbf{w}^*\|_2 \leq \epsilon$$

3. Estimate α

Overall Computational Complexity

Below #Iter hides the dependence on $\|\mathbf{w}^0 - \mathbf{w}^*\|_2$, $c = L/\alpha$

Condition	Guarantee	#Iter
α -SC, L -smooth	$\ \mathbf{w}^t - \mathbf{w}^*\ _2 \leq \epsilon$	$c \log(1/\epsilon)$
L -smooth	$F(\mathbf{w}^t) - F(\mathbf{w}^*) \leq \epsilon$	L/ϵ

- illustration
- GD solves linear regression efficiently

$$d^2(n + d) \quad \text{v.s.} \quad nd \cdot c \log(1/\epsilon)$$

- note on c

Improve Gradient Descent

Program

$$\min_{\mathbf{w}} F(\mathbf{w}), \quad \text{s.t. } \mathbf{w} \in \mathbb{R}^d.$$

- $F(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$ for $y \in \mathbb{R}$
- $F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \max\{1 - y_i \mathbf{x}_i \cdot \mathbf{w}, 0\} + 0.5\lambda \|\mathbf{w}\|_2^2$ for $y \in \{+1, -1\}$

GD:

- $O(nd)$ to evaluate $\nabla F(\mathbf{w})$
- always converges to opt

If computational cost is major concern,

can we boost the efficiency?

Explore the problem structure

Investigate Problem Structure

Suppose

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

- linear regression

$$F(\mathbf{w}) = \frac{1}{n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2, \quad f_i(\mathbf{w}) = (y_i - \mathbf{x}_i \cdot \mathbf{w})^2$$

- SVM

$$F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \max\{1 - y_i \mathbf{x}_i \cdot \mathbf{w}, 0\} + 0.5\lambda \|\mathbf{w}\|_2^2$$
$$f_i(\mathbf{w}) = \max\{1 - y_i \mathbf{x}_i \cdot \mathbf{w}, 0\} + 0.5\lambda \|\mathbf{w}\|_2^2$$

- any sample-wise loss

Stochastic Gradient Descent

- 1 Initialize \mathbf{w}^0 , say $\mathbf{w}^0 = \mathbf{0}$
- 2 For $t = 1, 2, \dots$

Uniformly draw i_t from $\{1, 2, \dots, n\}$, and update

$$\mathbf{w}^t = \mathbf{w}^{t-1} - \eta_t \nabla f_{i_t}(\mathbf{w}^{t-1}) \quad (2)$$

- example, intuition
- time cost per iteration is $O(d)$ (GD needs nd)
- total time = cost/iter \cdot #iter

Convergence Rate for SGD

α -SC, Lipschitz

$$\eta_t \leq \frac{1}{\alpha t}, \quad \mathbb{E}[\|\mathbf{w}^t - \mathbf{w}^*\|_2] \leq \frac{\log t}{t}$$

convex, Lipschitz

$$\eta_t \leq \frac{1}{\sqrt{t}}, \quad \mathbb{E}[F(\mathbf{w}^t) - F(\mathbf{w}^*)] \leq \frac{\log t}{\sqrt{t}}$$

We can modify SGD for faster rate (a rich literature).

Table 1: Overall computational cost to obtain ϵ opt. error

Condition	GD	SGD
SC	$n \log(1/\epsilon)$	$1/\epsilon$
Convex	n/ϵ	$(1/\epsilon)^2$

SGD wins if

- large-scale data

GD wins if

- small data set
- need high accuracy (i.e. ϵ is small)

In Practice...

SGD is used in



and more...

Other Practical Concerns

- Storage
- real-time decision-making

Initialize the model.

for $t = 1, 2, \dots$

- receive \mathbf{x}_t
- make prediction $\hat{\mathbf{y}}_t$
- receive \mathbf{y}_t
- evaluate loss $\ell(\hat{\mathbf{y}}_t, \mathbf{y}_t)$
- update model

Compare to SGD

March 12, 18:30 - 21:00 EST

- linear algebra, calculus
- advanced probability
- linear regression
- gradient descent
- stochastic gradient
- online learning
- **statistical machine learning** (Mar. 5)