

# Statistical Machine Learning

Instructor: Jie Shen

Dept. of Computer Science

February 14, 2020

# Gradient Descent

$$\epsilon \in (0, 1)$$

Consider minimization without constraints: Goal:  $F(w^t) - F(w^*) \leq \epsilon$

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

$$t = ?$$

$$t = f(\epsilon)$$

## Gradient Descent:

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

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

## Goal:

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

# Informal Analysis

$$\forall t, \quad F(w^{t+1}) - F(w^t) < 0 \Leftrightarrow F(w^{t+1}) < F(w^t)$$

Why GD "decreases" objective value (under proper conditions)?

$$w^{t+1} = w^t - \eta \cdot \nabla F(w^t)$$

$$F(w^{t+1}) - F(w^t) = \langle \nabla F(w), \underbrace{w^{t+1} - w^t} \rangle$$

$$\begin{aligned} \eta \cdot \nabla F(w) &\stackrel{\eta=0}{=} \langle \nabla F(w), -\eta \cdot \nabla F(w^t) \rangle \\ &= -\eta \cdot \underbrace{\langle \nabla F(w), \nabla F(w^t) \rangle}_{\approx \nabla F(w^t)} \end{aligned}$$

$$\begin{array}{c} \xrightarrow{w^t} \xrightarrow{w} \xrightarrow{w^{t+1}} \\ \text{---} \end{array}$$

$$\approx 0 \quad \eta \rightarrow 0. \quad w \rightarrow w^t \Rightarrow \nabla F(w) \approx \nabla F(w^t)$$

Mean-value Thm,

$$F: \mathbb{R} \rightarrow \mathbb{R}$$

$$\forall x, y \in \mathbb{R},$$

$$F(x) - F(y) = \underbrace{F'(z)}_{\mathbb{R}^d} \cdot \underbrace{(x-y)}_{\mathbb{R}^d}.$$

$z \in [x, y]$



$$F: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\forall x, y \in \mathbb{R}^d$$

$$F(x) - F(y) = \langle \nabla F(z), x-y \rangle$$

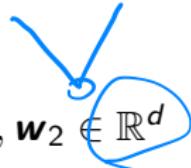
$$\mathbb{R} \quad \mathbb{R} \quad \mathbb{R}^d \cdot \mathbb{R}^d$$

$$z = \lambda x + (1-\lambda)y, \quad \lambda \in [0, 1]$$

# Smoothness



(0, +∞)



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:

$$F(w) = w^2$$

$$F(w) = |w|.$$

$$\nabla F(w_2) = 2w_2$$

$$w_1 = -\epsilon, \quad w_2 = \epsilon$$

$$\nabla F(w_1) = 2w_1$$

$$\nabla F_1 = -1, \quad \nabla F_2 = 1$$

$$\text{LHS} = 2 \cdot \|w_2 - w_1\|$$

$$\text{LHS} = 2 \quad \text{RHS} = L \cdot 2\epsilon$$

$$L = 2, \quad 2\text{-smooth.}$$

can we find L.

$$\begin{aligned} & 2 < 2L\epsilon ? \\ \Leftrightarrow & L \cdot \epsilon > 1 \end{aligned}$$

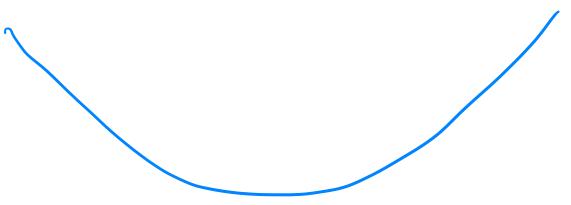
$$L = 1 \text{ M. } 10^9$$

$$\epsilon = \frac{1}{10^9} \quad \frac{1}{10^{12}}$$



$$w_1 \quad w_2 \quad \underline{\epsilon^{-12}}$$

$$F(w) = w^4$$



$$\nabla = 4w^3.$$

$$\text{LHS} = 4 \cdot |w_1^3 - w_2^3|$$

$$\text{RHS} = L \cdot |w_1 - w_2|$$

can find  $\lambda > 0$ . s.t.  $\forall w_1, w_2 \in \mathbb{R}$

$$4 \cdot |w_1^3 - w_2^3| < L \cdot |w_1 - w_2|$$

$$\Leftrightarrow \lambda > 4 \cdot \underbrace{|w_1^2 + w_1 w_2 + w_2^2|}_{\forall w_1, w_2 \in \mathbb{R}}$$

non-smooth.

$$\text{domain} = [-1, 1] \quad \checkmark$$

## When GD fails to find global optimum

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

# The Big Picture

# Convex Optimization

**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

# Convex Program

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

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

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

# Convergence Analysis

$$t \in (0, 1) , \quad t = ? \quad F(w^t) - F(w^*) \leq \epsilon$$

## Theorem 2

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

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

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

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

Input:  $\mathcal{G} \subset (0, 1)$

$$\overbrace{\mathcal{F}(w)}$$

Output:  $w^t, \mathcal{F}(w^t) - \mathcal{F}(w^*) \leq \epsilon$

Alg:  $\eta = \frac{1}{L}$

$$I = \frac{1}{c}$$

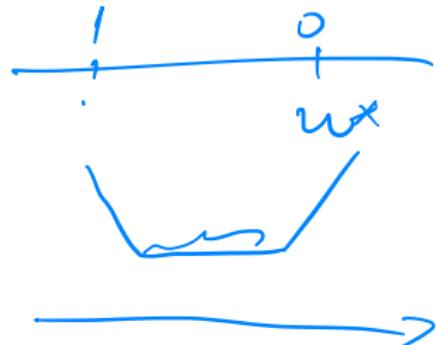
$$t = 1, \dots, I$$

$$w^{t+1} \leftarrow w^t - \eta \cdot \nabla \mathcal{F}(w^t)$$

# Implications

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

$$F(\mathbf{w}) = 1$$



# Implications

## 2. Iteration complexity

$$t = \frac{1}{\epsilon}$$

$$t \geq \frac{F(w^t) - F(w^*)}{\epsilon} \leq \frac{\epsilon}{\epsilon} = 1$$

$t$

# Implications

$$\sup_{w_1, w_2} \frac{\|\nabla^2 F(w) \cdot (w_2 - w_1)\|}{\|w_2 - w_1\|} = \sup_v \frac{\|\nabla^2 F(w) \cdot v\|}{\|v\|}$$

3. Estimate  $L$

$L = \max$ -eigen value of Hessian

$\lambda \in \mathbb{R}, v \in \mathbb{R}^d$   
 $(\lambda, v), \underline{M}$ .

$$\frac{\|\nabla F(w_2) - \nabla F(w_1)\|}{\|w_2 - w_1\|} \leq L \quad \text{if } M \cdot v = \lambda \cdot v$$

$$\|M \cdot v\| = |\lambda| \cdot \|v\|$$

$$\|w_2 - w_1\|$$

$$\Leftrightarrow L = \sup_{w_1, w_2} \frac{\|\nabla F(w_2) - \nabla F(w_1)\|}{\|w_2 - w_1\|}$$

$$\frac{\|M \cdot v\|}{\|v\|} = |\lambda|$$

$$\frac{\|\nabla F(w) \cdot v\|}{\|\nabla F(w) \cdot v\|} = \lambda$$

$$F: \mathbb{R} \rightarrow \mathbb{R} \quad \checkmark$$

$$F: \mathbb{R}^d \rightarrow \mathbb{R} \quad \checkmark$$

$$\underline{g}: \mathbb{R}^d \rightarrow \mathbb{R}^d$$

$$\cdot w_1, w_2$$

matrix-ve-  
multi.  
 $\downarrow$

$$g(w_2) - g(w_1) = \underline{\nabla g(w) \cdot (w_2 - w_1)}$$

$$\mathbb{R}^d$$

$$\mathbb{R}^d$$

$$\mathbb{R}^{d \times d}$$

$$\mathbb{R}^d$$

$$g \stackrel{\Delta}{=} \nabla F$$

$$\underline{\nabla F(w_2) - \nabla F(w_1)} = \underline{\nabla^2 F(w) \cdot (w_2 - w_1)}$$

Given  $\bar{F}(w)$

$\nabla^2 \bar{F}(w)$  by yourself.



python / Matlab

# Faster Rate of Convergence

$w^*$  is unique.

$$F(w) = w^T X$$

Not unique.  $w_1 \neq w^*$ .

$$\frac{\|\nabla F(w_1) - \nabla F(w^*)\|_2}{\|w_1 - w^*\|} = \frac{4(w_1^2 + w_1 w^* + w^*^2)}{\alpha}$$

Strongly Convex: for any  $w_1, w_2 \in \mathbb{R}^d$

$$\nabla^2 F \succeq \lambda \cdot I \quad \|\nabla F(w_2) - \nabla F(w_1)\|_2 \geq \alpha \|w_2 - w_1\|_2$$

- functions satisfying SC

$$0 \geq \alpha \|w_2 - w^*\|_2 = 1/2$$

- not satisfying

$$\nabla^2 \leq 0 \Rightarrow w_2 = w^*$$

7:40

$$F(w) = w^T \lambda \cdot I \quad \lambda = \text{min-eigenvalue}$$

2-SC, convex ↑ of Hessian.

$$F(w) = g(w) + \lambda \cdot \|w\|^2$$

# Theoretical Guarantee

## Theorem 3

Suppose  $F(\mathbf{w})$  is  $\alpha$ -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  $t \geq 1$ ,

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

$\alpha \leq L$ .  
 $c > 1$

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

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

$c \in (0, 1)$

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

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

# Implications

$$\epsilon = 10^{-6} \quad \gamma \epsilon = 1M$$

$\log 1/\epsilon = b$  App. of Thm 3 + telescoping

For all  $t \geq 1$ ,

$i = 1, \dots, t$

$$\|\mathbf{w}^t - \mathbf{w}^*\|_2 \leq \left(1 - \frac{2}{c+1}\right)^t \|\mathbf{w}^0 - \mathbf{w}^*\|_2$$

$$\leq e^{-\frac{2t}{c+1}} \|\mathbf{w}^0 - \mathbf{w}^*\|_2$$

(by  $1+x \leq e^x$ )

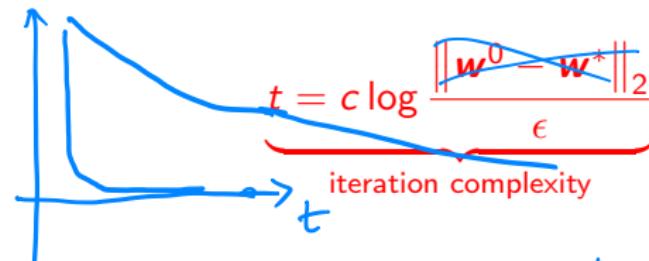
$$\|\mathbf{w}^t - \mathbf{w}^*\|$$

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

$$= \epsilon$$

$$x = -\frac{2}{c+1}$$

$$t = c \log \frac{\|\mathbf{w}^0 - \mathbf{w}^*\|_2}{\epsilon} \Rightarrow \|\mathbf{w}^t - \mathbf{w}^*\|_2 \leq \epsilon$$



$$t = c \cdot \underline{\log \frac{1}{\epsilon}}$$

$$t = \underline{\underline{\frac{1}{\epsilon}}}$$

# Implications

3. Estimate  $\alpha$

$$\overline{D^2}.$$

$$L \geq \alpha$$

$$\eta \leq \frac{2}{\alpha + qL}.$$

$$\begin{aligned}\eta &= \left( \frac{2}{L+L} \right) \leq \frac{2}{\alpha+L} \\ &= \frac{1}{C}.\end{aligned}$$

# Overall Computational Complexity

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

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

$$\{(x_i, y_i)\} \subseteq \mathbb{R}^d \times \mathbb{R}$$

- illustration

$$F(w) = \frac{1}{n} \sum_{i=1}^n (y_i - w \cdot x_i)^2 + \lambda \cdot \|w\|^2$$

- GD solves linear regression efficiently

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

$d^3 \quad \curvearrowright \quad d.$

- 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

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

new assumption

- 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

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

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

$$\rightarrow \nabla f_3(w)$$

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

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

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

$$\nabla f_1(w) \cdot \frac{1}{n} + \nabla f_2(w) \cdot \frac{1}{n} + \nabla f_3(w) \cdot \frac{1}{n}$$

$$+ \dots + \nabla f_n(w) \cdot \frac{1}{n} = \nabla F(w)$$

Necessary condition:

$$w^{t+1} = w^t - \eta_t \cdot \nabla f_{i_t}(w^t)$$



$$w^*$$

$$w^*$$

$$\Rightarrow \eta_t \cdot \nabla f_{i_t}(w^*) = 0$$

$$\Rightarrow \eta_t = 0 \text{ as } t \rightarrow \infty$$

$$\nabla f_{i_t}(w^*) \neq 0$$

$$F(w) = \underbrace{\frac{1}{2} \cdot w^2}_{f_1(w)} + \underbrace{\frac{1}{2} \cdot (w-1)^2}_{f_2(w)}$$

$$w^* = \frac{1}{2}$$

$$\nabla f_1(w^*) = \frac{1}{2} \quad \nabla f_2(w^*) = -\frac{1}{2}$$

# Convergence Rate for SGD

①  $\nexists M > 0. \quad \forall |F(w_2) - F(w_1)| \leq M \cdot \|w_2 - w_1\|$

②  $\exists M > 0. \quad \|\nabla F(w)\| \leq M$   
 $\alpha$ -SC, Lipschitz

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

SGD total cost =  $\boxed{d \cdot \frac{1}{\epsilon}} \cdot C.$

convex, Lipschitz

G.D.

$$n^2 d \quad \eta_t \leq \frac{1}{\sqrt{t}}, \quad \mathbb{E}[F(w^t) - F(w^*)] \leq \frac{\log t}{\sqrt{t}}$$

G.D.  $\boxed{\frac{C \cdot n d \cdot (\log \frac{1}{\epsilon})}{O(n)}} \Rightarrow \frac{(C + n d) \log \frac{1}{\epsilon}}{O(n)} \quad \boxed{\sqrt{t}}$

We can modify SGD for faster rate (a rich literature).  $t = \frac{1}{\epsilon^2}$

Rie Johnson & SVRG

T-ng Zhang '14 stochastic variance reduced

# GD v.s. SGD

Oct 19.  
TA

Oct. 26  
6:30 - 9:00 pm

Table 1: Overall computational cost to obtain  $\epsilon$  opt. error

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

Week 1 - Oct 12

SGD wins if

- large-scale data

GD wins if

- small data set
- need high accuracy (i.e.  $\epsilon$  is small)

Oct 21

Open book.

No calculator

# In Practice...

SGD is used in



and more...

## Other Practical Concerns

- Storage
- real-time decision-making

# Online Learning

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

# Mid-term

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)