

Optimization for Machine Learning

September 27, 2023

Today's roadmap:

- Gradient descent and convex optimization
- Accelerated methods

⇒ Notes + illustration notebook available on the course webpage after the class

GRADIENT METHODS AND CONVEX OPTIMIZATION

Problem of interest:

minimize $f(w)$
 $w \in \mathbb{R}^d$

$f: \mathbb{R}^d \rightarrow \mathbb{R}$
convex

Assumption: f belongs to the class of $C_L^{1,1}$ functions, i.e.

• f is C^1 (continuously differentiable)

• At every $w \in \mathbb{R}^d$, $\exists \nabla f(w) \in \mathbb{R}^d$ (gradient of f at w) that represents how the function varies locally

• The gradient mapping $\nabla f: \mathbb{R}^d \rightarrow \mathbb{R}^d$
 $w \mapsto \nabla f(w)$ is L -Lipschitz continuous where $L > 0$.

$$\forall (w, v) \in (\mathbb{R}^d)^2, \quad \|\nabla f(w) - \nabla f(v)\| \leq L \|w - v\|$$

Examples of $C_L^{1,1}$ functions:

- linear least squares
- logistic regression objective "logistic loss"
- quadratic functions

Remark: $C_L^{1,1}$ assumption is a simplifying assumption

- Functions are not always C^1 (see next week!)
- Functions / gradients are not always Lipschitz continuous on the entire space \mathbb{R}^d
- ⇒ Possible to use local Lipschitz constants

↳ Sometimes $C_L^{1,1}$ functions are called L -smooth

Two key properties

If f is $C_L^{1,1}$ and convex, then

$\forall w \in \mathbb{R}^d, \forall v \in \mathbb{R}^d,$

(1)
$$f(v) \leq f(w) + \underbrace{\nabla f(w)^T (v-w)}_{\text{linear function in } v} + \underbrace{\frac{L}{2} \|v-w\|^2}_{\text{quadratic function of } v}$$

$\|v-w\|^2 = \sum_{i=1}^d (v_i - w_i)^2$

and
(2)
$$f(v) \geq f(w) + \underbrace{\nabla f(w)^T (v-w)}_{\text{linear function of } v}$$

(1) \Rightarrow Upper bound on $f(v)$ with a quadratic function of v

(2) \Rightarrow Lower bound on $f(v)$ with a linear function of v

(2) is actually a characterization of convexity for C^1 functions

Towards gradient descent

- Suppose that you are at $w \in \mathbb{R}^d$ and you know $f(w), \nabla f(w)$
- If $\|\nabla f(w)\| = 0$, then w is a global minimum (f convex)
- When $\|\nabla f(w)\| \neq 0$, we can find v such that

$$f(w) + \nabla f(w)^T (v-w) + \frac{L}{2} \|v-w\|^2 < f(w)$$

Need this < 0

implying that $f(v) \stackrel{(1)}{\leq} f(w) + \nabla f(w)^T (v-w) + \frac{L}{2} \|v-w\|^2 < f(w)$

- We can then replace w by v and repeat the process

Main idea behind the gradient descent method

① Gradient descent ($f \in C_2^{2,1}$, convex)

Gradient descent (GD) algorithm

$\{w_k\}_k$: iterates

Initialization: Choose $w_0 \in \mathbb{R}^d$.

For $k=0, 1, \dots$

- (i) Evaluate $\nabla f(w_k)$. If $\|\nabla f(w_k)\| = 0$, terminate and output w_k .
- (ii) Compute a stepsize $\alpha_k > 0$.
- (iii) Define $w_{k+1} = w_k - \alpha_k \nabla f(w_k)$

"Gradient descent iteration"

End For

- In an implementation, the for loop is replaced by a while loop involving a convergence criterion and a budget criterion

Ex) while $[\|\nabla f(w_k)\| \neq 0]$ and $[k \leq K]$

End while

↑
convergence
criterion

↑
budget criterion.
→ Number of iterations
→ Number of function/
gradient evaluations
→ CPU time

Other convergence criteria

- $\|\nabla f(w_k)\| < \epsilon$ where ϵ is a small tolerance/precision (e.g. 10^{-4} , 10^{-5} , 10^{-16})

- $f(w_k) - \min_{w \in \mathbb{R}^d} f(w) < \epsilon$

"close to the best possible objective"

⚠ Most of the time the optimal value is unknown!

But for convex functions, can provide theoretical bounds on the cost of satisfying this condition

• $\|w_k - w^*\| < \epsilon$

where $w^* \in \underset{w \in \mathbb{R}^d}{\text{argmin}} f(w)$ (w^* solution of the problem)

"close to a solution"

⚠ In general optimal solutions are not known!

⇒ But for certain convex functions can provide bounds on the cost of satisfying $\|w_k - w^*\| < \epsilon$

● → How do we select α_k ?

• One choice that works for $C_L^{1,1}$ functions: $\alpha_k = \frac{1}{L}$

Proposition: Suppose that we are at iteration k of GD and $\|\nabla f(w_k)\| \neq 0$. Then, if $\alpha_k = \frac{1}{L}$, we have

$$f(w_k - \alpha_k \nabla f(w_k)) \leq f(w_k) - \frac{1}{2L} \|\nabla f(w_k)\|^2 \leq f(w_k)$$

Proof: Apply (1) with $v = w_k - \alpha_k \nabla f(w_k)$ and $w = w_k$

$$f(w_k - \alpha_k \nabla f(w_k)) \leq f(w_k) + \nabla f(w_k)^T (w_k - \alpha_k \nabla f(w_k) - w_k) + \frac{L}{2} \|w_k - \alpha_k \nabla f(w_k) - w_k\|^2$$

$\forall w, v \in \mathbb{R}^d$
 $\|v\|^2 = \sum_{i=1}^d v_i^2$
 $u^T v = \sum_{i=1}^d u_i v_i$
 $v^T v = \|v\|^2$

$$\begin{aligned} &= f(w_k) + \nabla f(w_k)^T (-\alpha_k \nabla f(w_k)) + \frac{L}{2} \|\alpha_k \nabla f(w_k)\|^2 \\ &= f(w_k) - \underbrace{\alpha_k \nabla f(w_k)^T \nabla f(w_k)}_{\|\nabla f(w_k)\|^2} + \frac{L}{2} \alpha_k^2 \|\nabla f(w_k)\|^2 \\ &= f(w_k) - \alpha_k \|\nabla f(w_k)\|^2 + \frac{L}{2} \alpha_k^2 \|\nabla f(w_k)\|^2 \\ &= f(w_k) - \frac{1}{L} \|\nabla f(w_k)\|^2 + \frac{L}{2} \times \frac{1}{L^2} \|\nabla f(w_k)\|^2 \\ &= f(w_k) + \left[-\frac{1}{L} + \frac{1}{2L}\right] \|\nabla f(w_k)\|^2 = f(w_k) - \frac{1}{2L} \|\nabla f(w_k)\|^2 \end{aligned}$$

More general result: For any $\alpha_k \in (0, \frac{2}{L})$, one has

$$f(w_k - \alpha_k \nabla f(w_k)) < f(w_k)$$

Remarks:
→ Suggests a problem-dependent choice for α_k
→ Requires knowledge of the Lipschitz constant (assumed to be known in this lecture)

Q Under our assumptions, what can we prove about the algorithm?

Theorem: Suppose that f is convex and $C_L^{1,1}$.
Suppose that we run GD with $\alpha_k = \frac{1}{L} \forall k$.

a) Convergence rate

For any $K \geq 1$, after K iterations of GD, we have

$$f(w_K) - \min_{w \in \mathbb{R}^d} f(w) \leq O\left(\frac{1}{K}\right)$$

"worst case behavior"

"big-O" notation

A constant times $\frac{1}{K}$ where the constant does not depend on K

b) Complexity bound

For any $\epsilon > 0$, GD reaches an iterate w_K such that

$$f(w_K) - \min_{w \in \mathbb{R}^d} f(w) \leq \epsilon$$

after at most $O\left(\frac{1}{\epsilon}\right)$ iterations.

$$\epsilon \Rightarrow \frac{\epsilon}{10}$$

$$O(\epsilon^{-1}) \Rightarrow O(10\epsilon^{-1})$$

↑
constant that depends on w_0, L , distance between w_0 and argmin f

↳ CV rates and complexity bounds are used as a guidance/as an indicator of algorithmic performance.
 They can be close to practical performance on simple examples and these rates/bounds are attained on some (pathological) examples.

↳ Two ways to get better rates/bounds

- Add assumptions on the function \Rightarrow strong convexity
- Change the algorithm \Rightarrow acceleration

② Strongly convex functions

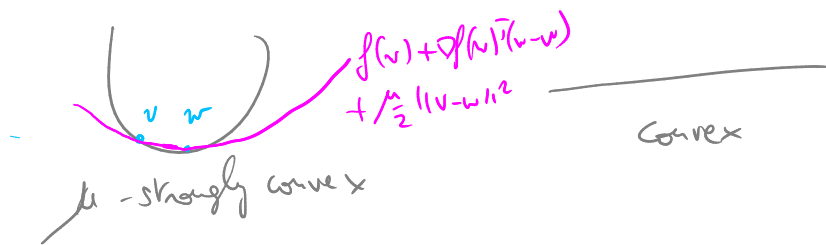
Def: $f: \mathbb{R}^d \rightarrow \mathbb{R}$ C^1

f is called μ -strongly convex (where $\mu > 0$) if

$$\forall (v, w) \in (\mathbb{R}^d)^2, \quad \underline{f(v) \geq f(w) + \nabla f(w)^T (v-w) + \frac{\mu}{2} \|v-w\|^2}$$

Corresponds to (2)

↑
Additional quadratic term compared to (2)



NB: μ -strongly convex \Rightarrow convex

Property: If f is C^1 and μ -strongly convex, then it has a unique global minimum, which is the solution of $\|\nabla f(w)\| = 0$

unique

Consequences on GD

↳ Suppose that $f \in \mathcal{F}_{L, \mu}^{1,1}$

$\mathcal{F}_{L, \mu}^{1,1} = \{ C_{L, \mu}^{1,1} \text{ and } \mu\text{-strongly convex} \}$

$\Rightarrow L \geq \mu$

$\{w^*\} = \underset{w}{\operatorname{argmin}} f(w)$

The running GD with $\alpha = \frac{1}{L}$ gives the following guarantees:

a) CV rate After $K \geq 1$ iterations

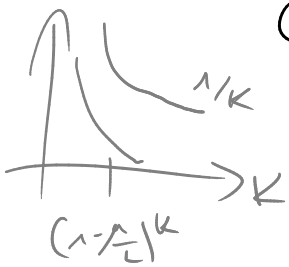
$$\underbrace{\|w_K - w^*\|}_{\text{distance between the last iterate and the solution}} \leq \underbrace{\left(1 - \frac{\mu}{L}\right)^K}_{\in [0, 1]} \underbrace{\|w_0 - w^*\|}_{\text{distance between the initial point and the solution}}$$

$\left(1 - \frac{\mu}{L}\right)^K \xrightarrow{K \rightarrow \infty} 0$

This is called a linear rate of CV or an exponential rate of CV

$$\left(1 - \frac{\mu}{L}\right)^K = e^{K \log\left(1 - \frac{\mu}{L}\right)}$$

"linear" $\Rightarrow \log(\|w_K - w^*\|)$ decreases as a linear function of K



⊕ Better than the result for convex functions

• Stronger criterion: $\|w_K - w^*\|$ vs $f(w_K) - \min_{w \in \mathcal{R}^d} f(w)$

• Better rate

$\left(1 - \frac{\mu}{L}\right)^K$ goes to zero faster than $1/K$

⊕ In practice, GD is faster on strongly convex functions than on convex ones

For certain strongly convex functions, one can define a better stepsize than $\frac{1}{L}$ (ex: $\frac{2}{L + \mu}$ for quadratic functions)

b) Complexity bound

$$\|w_K - w^*\| \leq \epsilon \text{ holds}$$

after at most $O\left(\frac{L}{\mu} \log(1/\epsilon)\right)$ iterations

"Condition number of the function"

$$\log(1/\epsilon) < 1/\epsilon \text{ for } \epsilon \text{ sufficiently small}$$

③ Acceleration

Question: Is there an algorithm that has better guarantees than GD while doing the same amount of work per iteration (i.e. evaluating 1 gradient)?

Recall: GD iteration

$$\forall k \geq 0, \quad w_{k+1} = w_k - \alpha_k \nabla f(w_k) \quad (\text{think } \alpha_k = 1/L)$$

To define the step $w_{k+1} - w_k = -\alpha_k \nabla f(w_k)$, we only rely on information related to w_k

and we ignore what we did in the previous steps

$$w_k - w_{k-1}, \dots$$



Intuition: Maybe $w_k - w_{k-1}$ would also be a good step from w_k

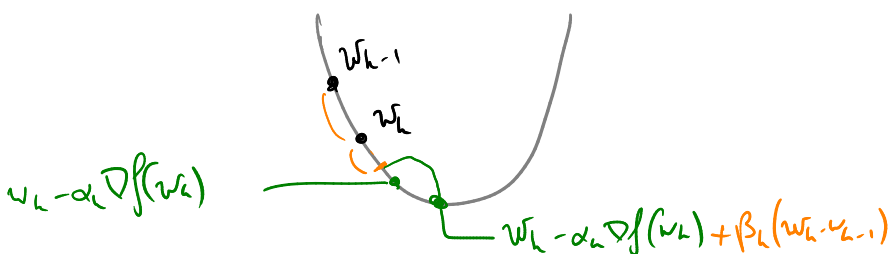
Polyak (1964): Heavy-ball method / Gradient descent with momentum

$$\forall k \geq 0, \quad w_{k+1} = w_k - \alpha_k \nabla f(w_k) + \beta_k (w_k - w_{k-1})$$

with $\beta_0 = 0 \quad w_{-1} = w_0$

↑ Previous step
"Momentum term"

$$\beta_k \geq 0$$



↳ This is optimal on strongly convex quadratic functions (better than GD and cannot do better with only 1 gradient evaluation per iteration)

$$\text{CV rate } (1 - \sqrt{\frac{\mu}{L}})^K \leq \underbrace{(1 - \frac{\mu}{L})^K}_{\text{GD}}$$

↳ But Heavy-ball can fail on strongly convex functions that are not quadratic

Nesterov (1983): Accelerated gradient / Nesterov's method

$$\forall k \geq 0, \quad w_{k+1} = w_k - \alpha_k \nabla f(w_k + \beta_k (w_k - w_{k-1})) + \beta_k (w_k - w_{k-1})$$

≠ with Heavy ball:

- First compute $w_k + \beta_k (w_k - w_{k-1})$
- Then do a gradient step

First do $w_k - \alpha_k \nabla f(w_k)$

then add $\beta_k (w_k - w_{k-1})$

from $w_k + \beta_k (w_k - w_{k-1})$

Guarantees for Nesterov's method (optimal: cannot do better)

• If $f \in \mathcal{F}_{L,\mu}^{1,1}$, set $\alpha_k = \frac{1}{L}$ and $\beta_k = \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}$

Then

CV rate $\|w_k - w^*\| \leq (1 - \sqrt{\frac{\mu}{L}})^k \|w_0 - w^*\|$
 $\forall k \geq 1$

↓ problem dependent

↑ improves over $\frac{\mu}{L}$ for GD

complexity

$\|w_k - w^*\| \leq \epsilon$ after $O\left(\sqrt{\frac{L}{\mu}} \log\left(\frac{1}{\epsilon}\right)\right)$ iterations
 ↑ improves over $\frac{L}{\mu}$ for GD

• If $f \in C_L^{1,1}$ and convex

Set $\alpha_k = \frac{1}{L}$ and $\beta_k = \frac{t_k - 1}{t_{k+1}}$ where $\begin{cases} t_0 = 0 \\ t_{k+1} = \frac{1}{2}(1 + \sqrt{1 + 4t_k^2}) \end{cases}$

↓
independent of the problem (or even w_0)

CV rate $f(w_{1/k}) - \min_{w \in \mathbb{R}^d} f(w) \leq O\left(\frac{1}{k^2}\right)$

Complexity $f(w_{1/k}) - \min_{w \in \mathbb{R}^d} f(w) \leq \varepsilon$ after at most
 $O(\varepsilon^{-1/2})$ iterations

$\frac{1}{k^2} \ll \frac{1}{k}$ for large k and $\varepsilon^{-1/2} \ll \varepsilon^{-1}$ for small ε