

OPTIMIZATION FOR MACHINE LEARNING

September 23, 2024

Today: Gradient descent (theory)

Next session: Lab on gradient descent (bring laptops!)



Tomorrow September 24, 8.30 - 11.45
(Lab room TBA)

Bo48

GRADIENT DESCENT

Q: How do you solve an optimization problem in practice?
↑
exactly?
approximately?
↑
→ on a computer
→ using an algorithm

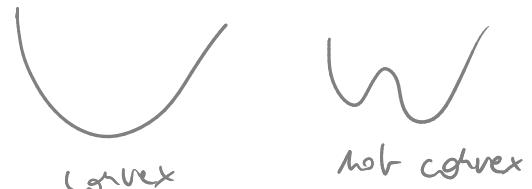
Problem of interest

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

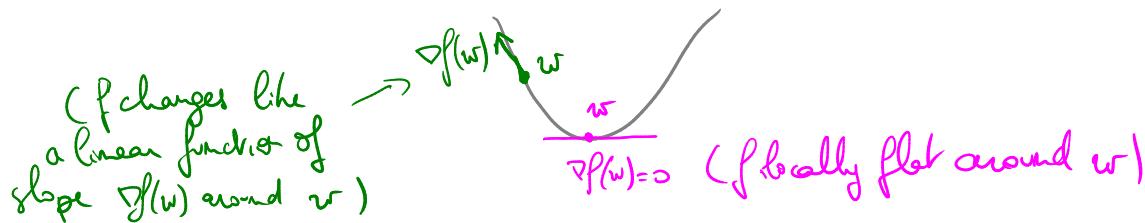
$f: \mathbb{R}^d \rightarrow \mathbb{R}$ objective function
 $w \in \mathbb{R}^d$ variables

Assumptions:

- f is convex
- f is C^1



$\forall w \in \mathbb{R}^d, \exists \nabla f(w) \in \mathbb{R}^d$ that indicates how the function changes around w



$\text{argmin } f$: set of global minima of f

Recall: Optimality conditions for $f C^1$ convex

$$[\bar{w} \in \underset{w \in \mathbb{R}^d}{\text{argmin}} f(w)] \iff [\nabla f(\bar{w}) = 0_{\mathbb{R}^d}]$$

$$\begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_d \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

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

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

$$f \in C^1$$

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

$$w \mapsto \begin{bmatrix} \frac{\partial f}{\partial w_1}(w) \\ \vdots \\ \frac{\partial f}{\partial w_d}(w) \end{bmatrix} = \nabla f(w)$$

↳ In general, cannot directly compute a vector $\bar{w} \in \mathbb{R}^d$ such that $\nabla f(\bar{w}) = 0_{\mathbb{R}^d}$

- Either there is no mathematical expression for the solution of the equation $\nabla f(\bar{w}) = 0_{\mathbb{R}^d}$ (e.g. deep neural networks)
- Or the solution is too expensive to compute in practice

Ex: Linear regression

$$\text{minimize}_{w \in \mathbb{R}^d} \frac{1}{2m} \|Xw - y\|^2 = \frac{1}{2m} \sum_{i=1}^m (x_i^T w - y_i)^2$$

$$\text{Data: } \{(x_i, y_i)\}_{i=1..m} \quad x_i \in \mathbb{R}^d$$

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_m^T \end{bmatrix} \in \mathbb{R}^{m \times d}$$

$$y \in \mathbb{R}$$

m : Number of data points

↳ On paper, we can give an explicit formula for the set of $\bar{w} \in \mathbb{R}^d$ such that $\nabla f(\bar{w}) = 0_{\mathbb{R}^d}$

d : Dimension of the parameter space

↳ However, if $m \gg 1$, computing these expressions is not possible!

⇒ The large number of data points makes the calculations too expensive

(inverse of a matrix in $\mathbb{R}^{m \times m}$, storage of such a matrix, ...)

↳ Instead of solving $\nabla f(\bar{w}) = 0_{\mathbb{R}^d}$ directly, we will build an algorithm that converges to a point for which $\nabla f(\bar{w}) = 0_{\mathbb{R}^d}$

→ Gradient descent

Idea: Given any $w \in \mathbb{R}^d$, compute $\nabla f(w) \in \mathbb{R}^d$

① If $\nabla f(w) = 0_{\mathbb{R}^d}$, then $w \in \arg\min_{w \in \mathbb{R}^d} f(w)$ (f convex)
(problem solved!)

② If $\nabla f(w) \neq 0_{\mathbb{R}^d}$, use $\nabla f(w)$ to find a better point $w^+ \in \mathbb{R}^d$ such that $f(w^+) < f(w)$

\Rightarrow Such a point necessarily exists!

\Rightarrow Mathematically, one can show that there exists $\alpha > 0$ such that

$$f(w - \alpha \nabla f(w)) < f(w)$$

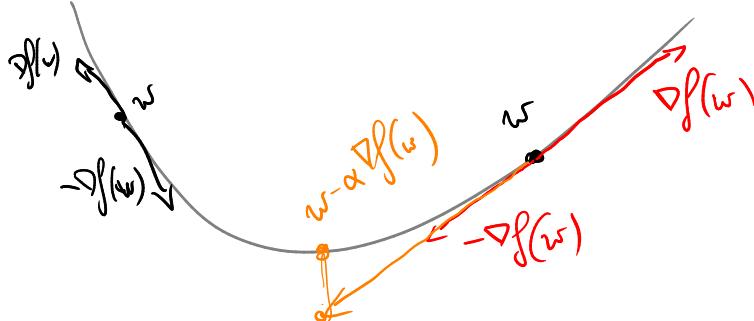
Not true
for any $\alpha > 0$

Opposite/negative
gradient direction

When $\nabla f(w) \neq 0_{\mathbb{R}^d}$

$\nabla f(w)$ always defines a direction in which f increases

$-\nabla f(w)$ always defines a direction in which f decreases



Generic Gradient Descent algorithm

Pseudo-code:

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

$f \in C^2$
convex

Initialization: Choose $w_0 \in \mathbb{R}^d$, set $k = 0$. \leftarrow iteration index
initial point

Loop: For $k = 0, 1, 2, \dots$) infinite loop

① Compute $\nabla f(w_k) \in \mathbb{R}^d$

② If $\nabla f(w_k) = 0_{\mathbb{R}^d}$, stop and return w_k .

③ Otherwise, compute $\alpha_k > 0$ $\xrightarrow{\text{stepsize}}$ and

set $w_{k+1} = w_k - \alpha_k \nabla f(w_k)$

) Gradient
descent
update

↳ The algorithm above is an infinite loop and a conceptual method
⇒ To implement it, we change the stopping criterion and we
specify how to compute the stepsize.

Stopping the algorithm (or making sure that it stops)

↳ Two ways to stop the algorithm

. Convergence criterion: Checks that the algorithm has
computed a solution or a good enough approximation

Ex) $\nabla f(w_k) = 0_{\mathbb{R}^d}?$ ($\Leftrightarrow \|\nabla f(w_k)\| = 0?$)

\uparrow
exact
CV criterion

$$\|v\| = \sqrt{\sum_{j=1}^d v_j^2}$$

$\|\nabla f(w_k)\| \leq \varepsilon$ where $\varepsilon > 0$ is a small real
number ($10^{-3}, 10^{-4}, \dots$)

- . Budget criterion: Guarantees that the algorithm will stop in a finite amount of time

Ex) Do the for loop over $K \geq 1$ iterations.

\Rightarrow Do at most K iterations / Compute at most K gradients

$$\begin{aligned} & \|w_{k+1} - w_k\| \\ &= \alpha_k \|\nabla f(w_k)\| \\ &\text{can decrease} \\ &\text{faster than} \\ &\|\nabla f(w_k)\| \end{aligned}$$

- Stop when $\|w_{k+1} - w_k\|$ is smaller than some tolerance (e.g. $10^{-8}, 10^{-12}, \dots$)
 \Rightarrow Algorithm is not changing the weight point significantly

(pick $\alpha_0 = 10^{-16}$, $\|\nabla f(w_0)\| = 1$, $\varepsilon = 10^{-6}$

$\|\nabla f(w_0)\| < 10^{-6}$ not satisfied

$\|w_1 - w_0\| = 10^{-16}$ should stop
 because this norm is of order of machine precision

- Stop after a CPU/GPU time limit is exceeded

\Rightarrow In practice, the algorithm includes a convergence criterion and one or several budget criteria.

COMPUTING THE STEPSIZE α_k (aka the learning rate)

↳ Theory: If $\nabla f(w) \neq 0_{\mathbb{R}^d}$, $\exists \bar{\alpha} > 0$ such that

$$f(w - \alpha \nabla f(w)) < f(w) \quad \forall \alpha \in (0, \bar{\alpha}]$$

\Rightarrow There exist infinitely many values that produce a decrease

↳ Empirically, there are 3 main strategies to choose the stepsize in the algorithm

a) Constant stepsize: $\alpha_k = \alpha > 0 \quad \forall k = 0, 1, \dots$



- ⊕ Easy to compute a priori
- ⊕ Actually works on some standard problems when α is well chosen

⊖ It is difficult to set a good value for α in general (α too large \Rightarrow increase the function value)
(α too small \Rightarrow numerical issues, little or no progress)

\Rightarrow The method is very sensitive to the choice of α .

b) Decreasing stepsize

Idea: Define a priori a sequence $\{\alpha_k\}_k$ such that

$$\alpha_k > 0 \quad \text{and} \quad \alpha_k \xrightarrow[k \rightarrow \infty]{} 0 \quad \text{with} \quad \alpha_0 \geq \alpha_1 \geq \alpha_2 \dots$$

Classical choices: $\alpha_k = \frac{1}{k+1}$, $\alpha_k = \frac{1}{\sqrt{k+1}}$, $\alpha_k = \frac{\alpha_0}{k+1}$, $\alpha_0 > 0$

⊕ the theory guarantees that α_k is a good choice for k sufficiently large \Rightarrow Eventually, the function values only decrease

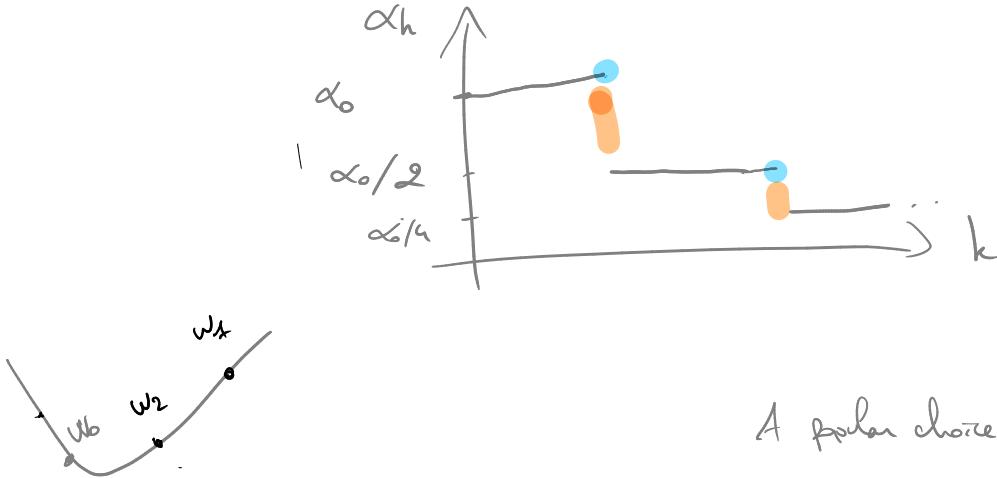
⊖ the stepsize can become unnecessarily small very quickly
 \Rightarrow No further progress

c) Hybrid stepsize approach

↳ Combination of constant/decreasing

↳ Use a constant value for α_k during a given number

of iterations, then decrease the value and repeat



Implementation

- Define a rule or when to decrease α_h (often heuristic)
- Define how to decrease it (e.g. divide by 2 or 10)

A popular choice in ML: 0.1, 0.01, 0.001, ...

d) Adaptive step size

↳ Choose α_h according to $w_h, f(w_h), \nabla f(w_h)$

⇒ Goal: Find a good value adapted to the current point

↳ One example: Backtracking line search

$$(w_h, \nabla f(w_h), f(w_h))$$

NB: 1 and 1/2 are arbitrary values, they can be replaced by a positive value

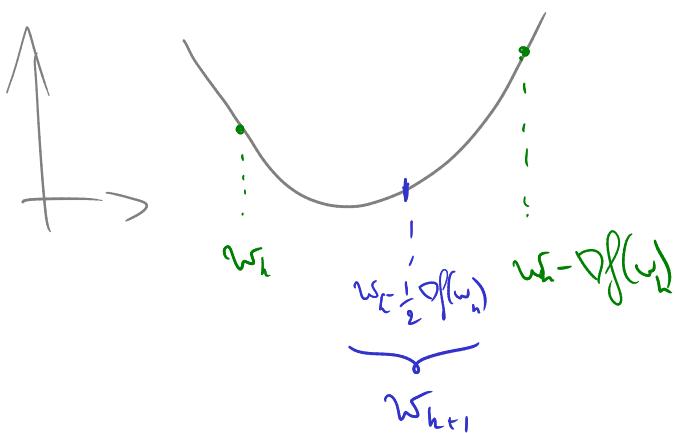
and a number in (0,1)

① Set $\alpha = 1$.

② If $f(w_h - \alpha \nabla f(w_h)) < f(w_h)$

Set $\alpha_h = \alpha$ and stop

Otherwise set $\alpha = \frac{\alpha}{2}$ and repeat ②.



⊕ Theory guarantees that the backtracking line-search process will stop (with a good stepsize)

⊖ Extra cost in function values

(in the worst case: proportional to $-\log(\|\nabla f(w_h)\|)$)

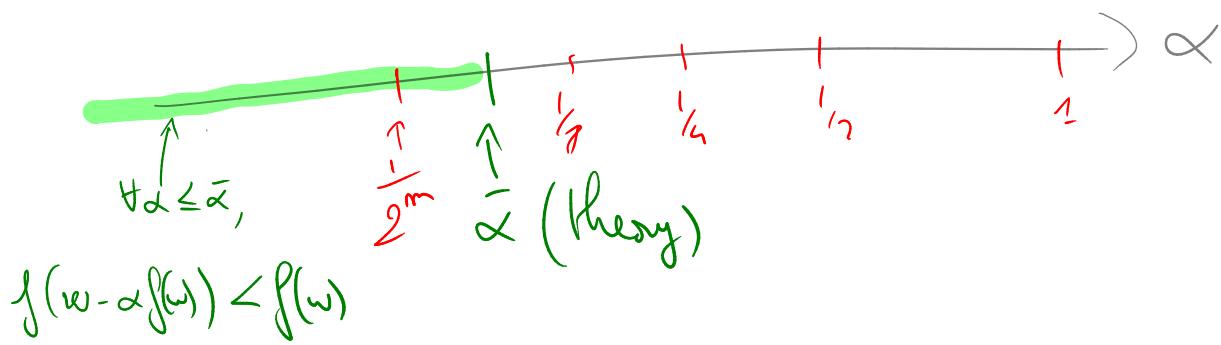
What is the best choice of α_n ?

It depends!

- If function values are expensive (like in deep learning), then predefined step sizes (constant, decreasing, some hybrid) are preferable
- But if function values can be afforded, adaptive step sizes give better results at the price of extra function value calculations

⚠ The gradient descent algorithm only computes gradient values and not function values in general

$$\text{Iteration: } w_{n+1} = w_n - \alpha_n \nabla f(w_n)$$



↳ For particular classes of problems (such as linear regression, logistic regression, linear SVM, ...), we can compute a constant value for α_n that guarantees decrease at every iteration and guarantees convergence to an approximate solution in finite time

$C_L^{1,1}$ functions and gradient descent

Def. f: $\mathbb{R}^d \rightarrow \mathbb{R}$ is $C_L^{1,1}$ if it is C^1 and $\nabla f: \mathbb{R}^d \rightarrow \mathbb{R}^d$ is L-Lipschitz continuous for some $L > 0$, that is

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

"The change in the gradients can be bounded by the change in the points."

Theorem Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be $C_L^{1,1}$ and let $w \in \mathbb{R}^d$ such that $\nabla f(w) \neq 0_{\mathbb{R}^d}$. Then,

$$f\left(w - \frac{1}{L} \nabla f(w)\right) \leq f(w) - \frac{1}{2L} \|\nabla f(w)\|^2 < f(w)$$

Corollary: If f is $C_L^{1,1}$, we can implement gradient descent on f with $\alpha_h = \frac{1}{L} \|\nabla f(w_h)\|$, and this will guarantee (without checking the function value!) that $f(w_{h+1}) < f(w_h)$ as long as $\nabla f(w_h) \neq 0_{\mathbb{R}^d}$

Corollary (Rate of convergence for gradient descent) — gradient descent
 Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be $C_L^{1,1}$ convex. Suppose that we run GD with $\alpha_h = \frac{1}{L} \|\nabla f(w_h)\|$ starting with $w_0 \in \mathbb{R}^d$. Then, for any $K \geq 1$,

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

↑ "Big-O of $(\frac{1}{K})$ "

↑ "Best possible value for f "

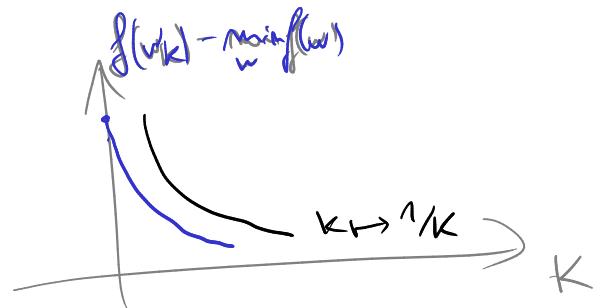
↑ "Iterate after K iterations"

where $\mathcal{O}\left(\frac{1}{K}\right) = \text{constant} \times \frac{1}{K}$, where the constant is positive and does not depend on K

(Here the constant depends on w_0 , L , $\arg\min_{w \in \mathbb{R}^d} f(w)$)

$$\frac{L}{2} \max_{w^* \in \arg\min f} \|w_0 - w^*\|$$

\Rightarrow We say that GD converges at the rate of $\frac{1}{K}$ for $C_L^{1,1}$ and convex functions



Equivalent result : Complexity bound

Under the same assumptions as the corollary, let $\varepsilon > 0$.

Then GD computes w_k such that $f(w_k) - \min_{w \in \mathbb{R}^d} f(w) \leq \varepsilon$

in at most $\mathcal{O}\left(\frac{1}{\varepsilon}\right)$ iterations.