

# Gradient Descent Algorithm

Andersen Ang

Mathématique et recherche opérationnelle  
UMONS, Belgium

[manshun.ang@umons.ac.be](mailto:manshun.ang@umons.ac.be)    Homepage: [angms.science](http://angms.science)

First draft : August 2, 2017

Last update : November 17, 2018

- 1 Optimization Problem
- 2 Descent Algorithm
- 3 Gradient Descent Algorithm

# Optimization Problem

Minimization (maximization) problem aims to find the optimal value of a given function  $f(\mathbf{x})$ .

Mathematically, given a function  $f : \text{dom } f = \mathcal{Q} \rightarrow \mathbb{R}$ , a minimization problem is to :

- find the optimal value  $f^*$

$$f^* = \min_{\mathbf{x} \in \mathcal{Q} \subset \mathbb{R}^n} f(\mathbf{x}).$$

- find the optimal decision variable  $\mathbf{x}$  that minimizes the function  $f$

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathcal{Q} \subset \mathbb{R}^n} f(\mathbf{x}).$$

- $\mathbf{x}$  : the decision variable / optimization variable
- $\text{dom } f = \mathcal{Q}$  : the domain of the function  $f$
- $\mathcal{Q}$  is also called the constraint set
- A solution is feasible if  $\mathbf{x} \in \mathcal{Q}$
- A solution is infeasible if  $\mathbf{x} \notin \mathcal{Q}$

# Unconstrained Optimization Problem

An optimization problem is unconstrained if  $Q = \mathbb{R}^n$ . That is

$$\min_{\mathbf{x}} f(\mathbf{x}).$$

That is, all  $\mathbf{x}$  are feasible for the problem.

Remarks :

- All  $\mathbf{x}$  are feasible for the problem
- More accurately, "min" should be replaced by "inf"

$$f^* = \inf_{\mathbf{x} \in Q} f(\mathbf{x})$$

- The optimal value  $f^*$  can be  $-\infty$ . (e.g. minimizes  $-\|\mathbf{x}\|_2^2$ )
- Some problems have bounded optimal value  $f^* > -\infty$

# Solving unconstrained optimization problem

One way to solve an unconstrained optimization problem is to come up with an algorithm that produces a sequence

$$\{\mathbf{x}_k\}_{k \in \mathbb{N}} : \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k, \dots$$

where  $k > 0$  is the iteration number, such that

$$\lim_{k \rightarrow \infty} f(\mathbf{x}_k) = f^*.$$

**$\epsilon$ -accuracy:** we stop the algorithm when the function value is  $\epsilon$ -close to the optimal value  $f^*$  :

$$0 \leq |f(\mathbf{x}_k) - f^*| \leq \epsilon.$$

- $\epsilon$  is user-defined.
- This stopping condition requires the knowledge of  $f^*$ , which is usually not known in advance.

# Iterative descent algorithm

One iterative algorithm is called *Descent Algorithm*, which iterates

$$\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \Delta_k$$

where

- $k$  : iteration counter
- $\mathbf{x}_k$  : the current decision variable
- $\mathbf{x}_{k+1}$  : the decision variable of the next iteration
- $t_k \geq 0$  : the step size
- $\Delta_k \in \mathbb{R}^n$  : the descent direction
  - ▶  $\|\Delta_k\| \neq 1$ , which is not "direction" in the common sense.

"Descent" is defined by the monotonic non-increasing of function value per iteration. i.e.

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) \quad \forall k$$

# Iterative descent algorithm

In conceptual level, iterative algorithm can be generalized as

---

## Algorithm 1: General framework of descent algorithm

---

**Result:** A solution  $\mathbf{x}$  that approximately solve  $\min_{\mathbf{x} \in \mathcal{Q}} f(\mathbf{x})$

**Initialization** pick initial point  $\mathbf{x}_0 \in \mathcal{Q}$  and initial parameters  $\theta$  (if any)

**while** *stopping condition is not met* **do**

    Find a descent direction  $\Delta_k$  in the form  $\Delta_k = p(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_k; \theta)$

    Pick a step size  $t_k$  in the form  $t_k = q(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_k; \Delta_k, \theta)$

    Update  $\mathbf{x}_{k+1}$  in the form  $r(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_k; t_k, \Delta_k, \theta)$

    Update  $\theta$  in the form  $s(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_k, \mathbf{x}_{k+1}; t_k, \Delta_k, \theta)$

**end**

---

General questions :

- How to initialize? (How to pick  $\mathbf{x}_0$ ?)
- How to define stopping condition? (What does "converge" means?)
- How to pick descent direction? (How to pick  $\Delta x_k$ ?)
- How to select the step size? (How to pick  $t_k$ ?)
- How to update? (How to get  $\mathbf{x}_{k+1}$  from the available information?)

## Theorem on the descent direction

**Theorem.** If  $f$  is convex and differentiable<sup>1</sup>, then the descent algorithm requires the descent direction  $\Delta x_k$  satisfies the following inequality

$$\nabla f(\mathbf{x}_k)^\top \Delta_k \leq 0$$

Proof:  $f$  is convex,  $\forall a, b \in \text{dom } f$  we have  $f(a) \geq f(b) + \nabla f(b)^\top (a - b)$ .  
Put  $a = \mathbf{x}_{k+1}$ ,  $b = \mathbf{x}_k$

$$f(\mathbf{x}_{k+1}) \geq f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top (\mathbf{x}_{k+1} - \mathbf{x}_k)$$

By definition of update  $\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \Delta_k$  thus  $\mathbf{x}_{k+1} - \mathbf{x}_k = t_k \Delta_k$  and

$$f(\mathbf{x}_{k+1}) \geq f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top t_k \Delta_k = f(\mathbf{x}_k) + t_k \nabla f(\mathbf{x}_k)^\top \Delta_k$$

Re-arrange, we have  $f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \geq t_k \nabla f(\mathbf{x}_k)^\top \Delta_k$ .

Using the definition of step size  $t_k \geq 0$  and the definition of descent  $f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k)$ , the proof is completed. □

<sup>1</sup>If  $f$  is not differentiable, replaces  $\nabla f$  by sub-gradient, the theorem still works.



## Gradient descent direction

Now we have  $\nabla f(\mathbf{x}_k)^\top \Delta_k \leq 0$ , then how to pick  $\Delta_k$  ?

One way is  $\Delta_k = -\nabla f(\mathbf{x}_k)$ , since  $\nabla f(\mathbf{x}_k)^\top \Delta_k = -\|\nabla f(\mathbf{x}_k)\|_2^2 \leq 0$ .

This is the descent direction used in **Gradient Descent** (GD).

It can be shown (next page), GD with a **small enough** step size  $t_k$  will converge to a point.

Remark : Note that the theorem on descent direction does not require

- $f$  to be convex
- stationary point to be local minimum of the  $f$

# GD converges to a point with step size $t \leq \frac{2}{\beta}$

**Theorem** For a  $\beta$ -smooth function  $f$  with optimal value  $f^* > -\infty$ , GD with step size  $t_k \leq \frac{2}{\beta}$  always converges to a stationary point.

Proof.  $f$  is  $\beta$ -smooth implies  $\forall a, b \in \text{dom} f$  we have

$$f(b) \leq f(a) + \nabla f(a)^T(b - a) + \frac{\beta}{2} \|a - b\|_2^2.$$

Put  $a = \mathbf{x}_k$ ,  $b = \mathbf{x}_{k+1}$  :

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T(\mathbf{x}_{k+1} - \mathbf{x}_k) + \frac{\beta}{2} \|\mathbf{x}_k - \mathbf{x}_{k+1}\|_2^2.$$

By definition of GD update  $\mathbf{x}_{k+1} - \mathbf{x}_k = -t \nabla f(\mathbf{x}_k)$

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leq -t \|\nabla f(\mathbf{x}_k)\|_2^2 + \frac{\beta}{2} t^2 \|\nabla f(\mathbf{x}_k)\|_2^2 = -t \left(1 - \frac{\beta t}{2}\right) \|\nabla f(\mathbf{x}_k)\|_2^2$$

Note  $0 \leq t \leq \frac{2}{\beta} \implies t \left(1 - \frac{\beta t}{2}\right) \geq 0$ . Rearrange and let  $p = \frac{1}{t(1 - \frac{\beta t}{2})}$ , we get

$$\|\nabla f(x_k)\|_2^2 \leq p (f(x_k) - f(x_{k+1})),$$

which forms a telescoping series



## Summary :

- Minimization problem
- Descent algorithm iterates  $\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \Delta_k$
- Descent algorithm requires  $\nabla f(\mathbf{x}_k)^T \Delta_k \leq 0$
- Gradient Descent algorithm picks  $\Delta_k = -\nabla f(\mathbf{x}_k)$
- For  $\beta$ -smooth function, the sequence  $\{\mathbf{x}_k\}_{k \in \mathbb{N}}$  produced by Gradient Descent algorithm with sufficient small step size ( $t_k \leq \frac{2}{\beta}$ ) converges to a point

## Not discussed :

- The convergence of  $\{\mathbf{x}_k\}_{k \in \mathbb{N}}$  to a (first-order) stationary point of  $f$ .
- Convergence rate of gradient descent on different  $f$ .

End of document