

Introductory Course on Non-smooth Optimisation

Lecture 01 - Gradient methods

Jingwei Liang

Department of Applied Mathematics and Theoretical Physics

- 1 Unconstrained smooth optimisation
- 2 Descent methods
- 3 Gradient of convex functions
- 4 Gradient descent
- 5 Heavy-ball method
- 6 Nesterov's optimal schemes
- 7 Dynamical system

Convex set

A set $S \subset \mathbb{R}^n$ is convex if for any $\theta \in [0, 1]$ and two points $x, y \in S$,

$$\theta x + (1 - \theta)y \in S.$$

Convex function

Function $F : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if $\text{dom}(F)$ is convex and for all $x, y \in \text{dom}(F)$ and $\theta \in [0, 1]$,

$$F(\theta x + (1 - \theta)y) \leq \theta F(x) + (1 - \theta)F(y).$$

- Proper convex: $F(x) < +\infty$ at least for one x and $F(x) > -\infty$ for all x .
- 1st-order condition: F is continuous differentiable

$$F(y) \geq F(x) + \langle \nabla F(x), y - x \rangle, \forall x, y \in \text{dom}(F).$$

- 2nd-order condition: if F is twice differentiable

$$\nabla^2 F(x) \succeq 0, \forall x \in \text{dom}(F).$$

Problem

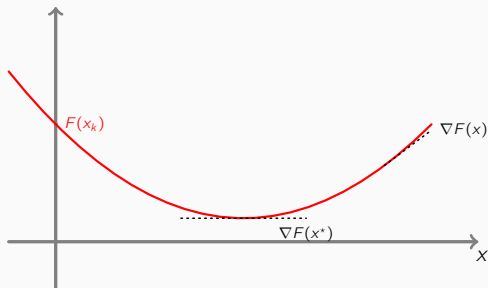
Unconstrained smooth optimisation

$$\min_{x \in \mathbb{R}^n} F(x),$$

where $F : \mathbb{R}^n \rightarrow \mathbb{R}$ is proper convex and smooth differentiable.

- Optimality condition: let x^* be a minimiser of $F(x)$, then

$$0 = \nabla F(x^*).$$



Quadratic programming

General quadratic programming problem

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} x^T A x + b^T x + c,$$

where $A \in \mathbb{R}^{n \times n}$ is symmetric positive definite, $b \in \mathbb{R}^n$ and $c \in \mathbb{R}$.

- Optimality condition:

$$0 = A x^* + b.$$

Special case least square

$$\|Ax - b\|^2 = x^T (A^T A) x - 2(A^T b)^T x + b^T b.$$

Optimality condition

$$A^T A x^* = A^T b.$$

Geometric programming

$$\min_{x \in \mathbb{R}^n} \log \left(\sum_{i=1}^m \exp(a_i^T x + b_i) \right).$$

Optimality condition:

$$0 = \frac{1}{\sum_{i=1}^m \exp(a_i^T x^* + b_i)} \sum_{i=1}^m \exp(a_i^T x^* + b_i) a_i.$$

1 Unconstrained smooth optimisation

2 Descent methods

3 Gradient of convex functions

4 Gradient descent

5 Heavy-ball method

6 Nesterov's optimal schemes

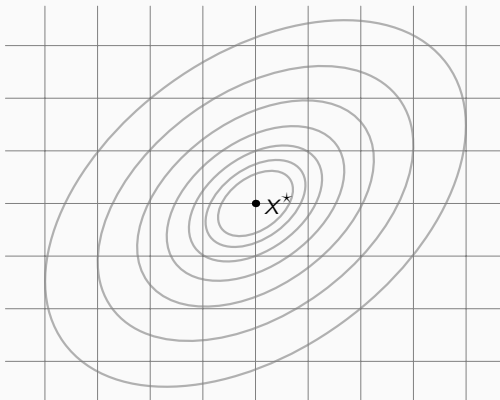
7 Dynamical system

Unconstrained smooth optimisation

Consider minimising

$$\min_{x \in \mathbb{R}^n} F(x),$$

where $F : \mathbb{R}^n \rightarrow \mathbb{R}$ is proper convex and smooth differentiable.



Unconstrained smooth optimisation

Consider minimising

$$\min_{x \in \mathbb{R}^n} F(x),$$

where $F : \mathbb{R}^n \rightarrow \mathbb{R}$ is proper convex and smooth differentiable.

- The set of minimisers, *i.e.*

$$\text{Argmin}(F) = \{x \in \mathbb{R}^n : F(x) = \min_{x \in \mathbb{R}^n} F(x)\}$$

is non-empty.

- However, given $x^* \in \text{Argmin}(F)$, no closed form expression.
- Iterative strategy to find one $x^* \in \text{Argmin}(F)$: start from x_0 and generate a train of sequence $\{x_k\}_{k \in \mathbb{N}}$ such that

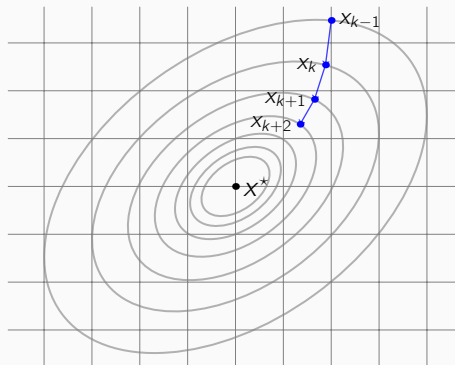
$$\lim_{k \rightarrow \infty} x_k = x^* \in \text{Argmin}(F).$$

Unconstrained smooth optimisation

Consider minimising

$$\min_{x \in \mathbb{R}^n} F(x),$$

where $F : \mathbb{R}^n \rightarrow \mathbb{R}$ is proper convex and smooth differentiable.



Iterative scheme

For each $k = 1, 2, \dots$, find $\gamma_k > 0$ and $d_k \in \mathbb{R}^n$ and then

$$x_{k+1} = x_k + \gamma_k d_k,$$

where

- d_k is called search/descent direction.
- γ_k is called step-size.

Descent methods

An algorithm is called descent method, if there holds

$$F(x_{k+1}) < F(x_k).$$

NB: if $x_k \in \text{Argmin}(F)$, then $x_{k+1} = x_k \dots$

From convexity of F , we have

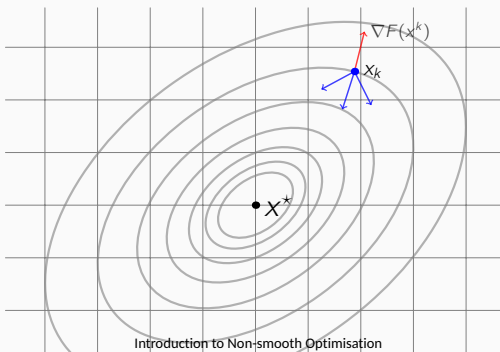
$$F(x_{k+1}) \geq F(x_k) + \langle \nabla F(x_k), x_{k+1} - x_k \rangle,$$

which gives

$$\langle \nabla F(x_k), x_{k+1} - x_k \rangle \geq 0 \implies F(x_{k+1}) \geq F(x_k).$$

Since $x_{k+1} - x_k = \gamma_k d_k$, the direction d_k should be such that

$$\langle \nabla F(x_k), d_k \rangle < 0.$$



General descent method

initial : $x_0 \in \text{dom}(F)$;

repeat :

1. Find a descent direction d_k .
2. Choose a step-size γ_k : line search.
3. Update $x_{k+1} = x_k + \gamma_k d_k$.

until : stopping criterion is satisfied.

Stopping criterion: $\epsilon > 0$ is the tolerance,

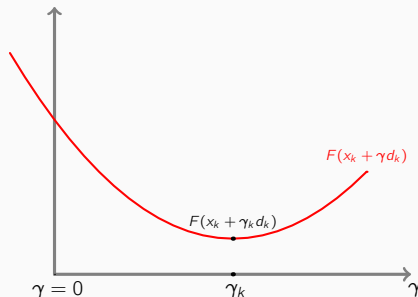
- Function value: $F(x_{k+1}) - F(x_k) \leq \epsilon$ (can be time consuming).
- Sequence: $\|x_{k+1} - x_k\| \leq \epsilon$.
- Optimality condition: $\|\nabla F(x_k)\| \leq \epsilon$.

Exact line search

Suppose that the direction d_k is given. Choose γ_k such that $F(x)$ is minimised along the ray $x_k + \gamma d_k, \gamma > 0$:

$$\gamma_k = \operatorname{argmin}_{\gamma > 0} F(x_k + \gamma d_k).$$

- Useful when the minimisation problem for γ_k is simple.
- γ_k can be found analytically for special cases.



Backtracking line search

Suppose that the direction d_k is given. Choose $\delta \in]0, 0.5[$ and $\beta \in]0, 1[$, let $\gamma = 1$

while $F(x_k + \gamma d_k) > F(x_k) + \delta \gamma \langle \nabla F(x_k), d_k \rangle : \gamma = \beta \gamma$.

- Reduce F enough along the direction d_k .

- Since d_k is a descent direction

$$\langle \nabla F(x_k), d_k \rangle < 0.$$

- Stopping criterion for backtracking:

$$F(x_k + \gamma d_k) \leq F(x_k) + \delta \gamma \langle \nabla F(x_k), d_k \rangle.$$

- When γ is small enough

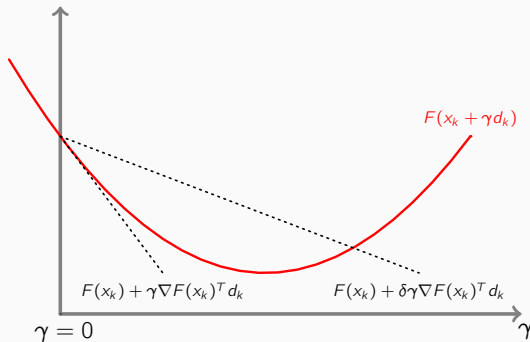
$$F(x_k + \gamma d_k) \approx F(x_k) + \gamma \langle \nabla F(x_k), d_k \rangle < F(x_k) + \delta \gamma \langle \nabla F(x_k), d_k \rangle,$$

whcih means backtracking eventually will stop.

Backtracking line search

Suppose that the direction d_k is given. Choose $\delta \in]0, 0.5[$ and $\beta \in]0, 1[$, let $\gamma = 1$

while $F(x_k + \gamma d_k) > F(x_k) + \delta \gamma \langle \nabla F(x_k), d_k \rangle$: $\gamma = \beta \gamma$.



1 Unconstrained smooth optimisation

2 Descent methods

3 Gradient of convex functions

4 Gradient descent

5 Heavy-ball method

6 Nesterov's optimal schemes

7 Dynamical system

Monotonicity of gradient

Let $F : \mathbb{R}^n \rightarrow \mathbb{R}$ be proper convex and smooth differentiable, then

$$\langle \nabla F(x) - \nabla F(y), x - y \rangle \geq 0, \forall x, y \in \text{dom}(F).$$

- C^1 : proper convex and smooth differentiable functions on \mathbb{R}^n .

Proof Owing to convexity, given $x, y \in \text{dom}(F)$, we have

$$F(y) \geq F(x) + \langle \nabla F(x), y - x \rangle$$

and

$$F(x) \geq F(y) + \langle \nabla F(y), x - y \rangle.$$

Summing them up yields

$$\langle \nabla F(x) - \nabla F(y), x - y \rangle \geq 0.$$

NB: Let $F \in C^1$, F is convex if and only if $\nabla F(x)$ is monotone.

Lipschitz continuity

The gradient of F is L -Lipschitz continuous if there exists $L > 0$ such that

$$\|\nabla F(x) - \nabla F(y)\| \leq L\|x - y\|, \forall x, y \in \text{dom}(F).$$

- C_L^1 : proper convex functions with L -Lipschitz continuous gradient on \mathbb{R}^n .

If $F \in C_L^1$, then

$$H(x) \stackrel{\text{def}}{=} \frac{L}{2}\|x\|^2 - F(x)$$

is convex.

Hint: monotonicity of $\nabla H(x)$, i.e.

$$\begin{aligned}\langle \nabla H(x) - \nabla H(y), x - y \rangle &= L\|x - y\|^2 - \langle \nabla F(x) - \nabla F(y), x - y \rangle \\ &\geq L\|x - y\|^2 - L\|x - y\|^2 \\ &= 0.\end{aligned}$$

Descent lemma, quadratic upper bound

Let $F \in C_L^1$, then there holds

$$F(y) \leq F(x) + \langle \nabla F(x), y - x \rangle + \frac{L}{2} \|y - x\|^2, \forall x, y \in \text{dom}(F).$$

Proof Define $H(t) = F(x + t(y - x))$, then

$$\begin{aligned} F(y) - F(x) &= H(1) - H(0) = \int_0^1 \nabla H(t) dt = \int_0^1 (y - x)^T \nabla F(x + t(y - x)) dt \\ &\leq \int_0^1 (y - x)^T \nabla F(x) dt + \int_0^1 |(y - x)^T (\nabla F(x + t(y - x)) - \nabla F(x))| dt \\ &\leq (y - x)^T \nabla F(x) + \int_0^1 \|y - x\| \|\nabla F(x + t(y - x)) - \nabla F(x)\| dt \\ &\leq (y - x)^T \nabla F(x) + \|y - x\| \int_0^1 tL \|y - x\| dt \\ &= (y - x)^T \nabla F(x) + \frac{L}{2} \|y - x\|^2. \end{aligned}$$

NB: first-order condition of convexity for $H(x) \stackrel{\text{def}}{=} \frac{L}{2} \|x\|^2 - F(x)$.

Corollary

Let $F \in C_L^1$ and $x^* \in \text{Argmin}(F)$, then

$$\frac{1}{2L} \|\nabla F(x)\|^2 \leq F(x) - F(x^*) \leq \frac{L}{2} \|x - x^*\|^2, \forall x \in \text{dom}(F).$$

Proof Right-hand inequality: $\nabla F(x^*) = 0$,

$$F(x) \leq F(x^*) + \langle \nabla F(x^*), x - x^* \rangle + \frac{L}{2} \|x - x^*\|^2, \forall x \in \text{dom}(F).$$

Left-hand inequality:

$$\begin{aligned} F(x^*) &\leq \min_{y \in \text{dom}(F)} \left\{ F(x) + \langle \nabla F(x), y - x \rangle + \frac{L}{2} \|y - x\|^2 \right\} \\ &= F(x) - \frac{1}{2L} \|\nabla F(x)\|^2. \end{aligned}$$

The corresponding y is $y = x - \frac{1}{L} \nabla F(x)$.

Co-coercivity

Let $F \in C_L^1$, then

$$\langle x - y, \nabla F(x) - \nabla F(y) \rangle \geq \frac{1}{L} \|\nabla F(x) - \nabla F(y)\|^2.$$

- Co-coercivity implies Lipschitz continuity

- For $F \in C_L^1$, $H(x) \stackrel{\text{def}}{=} \frac{L}{2} \|x\|^2 - F(x)$

Lipschitz continuity of $\nabla F \implies$ Convexity of $H(x)$

\implies Co-coercivity of $\nabla F(x)$

\implies Lipschitz continuity of ∇F

Co-coercivity

Let $F \in C_L^1$, then

$$\langle x - y, \nabla F(x) - \nabla F(y) \rangle \geq \frac{1}{L} \|\nabla F(x) - \nabla F(y)\|^2.$$

Proof Define $R(z) = F(z) - \langle \nabla F(x), z \rangle$, then $\nabla R(x) = 0$.

Recall the lemma

$$F \in C_L^1 \text{ and } x^* \in \text{Argmin}(F): \frac{1}{2L} \|\nabla F(x)\|^2 \leq F(x) - F(x^*) \leq \frac{L}{2} \|x - x^*\|^2.$$

Then we have

$$\begin{aligned} F(y) - F(x) - \langle \nabla F(x), y - x \rangle &= R(y) - R(x) \geq \frac{1}{2L} \|\nabla R(y)\|^2 \\ &= \frac{1}{2L} \|\nabla F(y) - \nabla F(x)\|^2. \end{aligned}$$

Similarly, define $S(z) = F(z) - \langle \nabla F(y), z \rangle$, then

$$F(x) - F(y) - \langle \nabla F(y), x - y \rangle = S(x) - S(y) \geq \frac{1}{2L} \|\nabla F(x) - \nabla F(y)\|^2.$$

Strong convexity

Function $F : \mathbb{R}^n \rightarrow \mathbb{R}$ is strongly convex if $\text{dom}(F)$ is convex and for all $x, y \in \text{dom}(F)$ and $\theta \in [0, 1]$, there exists $\alpha > 0$ such that

$$F(\theta x + (1 - \theta)y) \leq \theta F(x) + (1 - \theta)F(y) - \frac{\alpha}{2}\theta(1 - \theta)\|x - y\|^2.$$

- F is strongly convex with parameter $\alpha > 0$ if

$$G(x) \stackrel{\text{def}}{=} F(x) - \frac{\alpha}{2}\|x\|^2$$

is convex.

- Monotonicity:

$$\langle \nabla F(x) - \nabla F(y), x - y \rangle \geq \alpha\|x - y\|^2, \forall x, y \in \text{dom}(F).$$

- Second-order condition for strong convexity: if $F \in C^2$,

$$\nabla^2 F(x) \succeq \alpha \text{Id}, \forall x \in \text{dom}(F).$$

Quadratic lower bound

Let $F \in C^1$ and strongly convex, then

$$F(y) \geq F(x) + \langle \nabla F(x), y - x \rangle + \frac{\alpha}{2} \|y - x\|^2, \forall x, y \in \text{dom}(F).$$

Proof First-order condition of convexity for $G(x) \stackrel{\text{def}}{=} F(x) - \frac{\alpha}{2} \|x\|^2$.

Corollary

Let $F \in C^1$ be α -strongly convex and $x^* \in \text{Argmin}(F)$, then

$$\frac{\alpha}{2} \|x - x^*\|^2 \leq F(x) - F(x^*) \leq \frac{1}{2\alpha} \|\nabla F(x)\|^2, \forall x \in \text{dom}(F).$$

Proof Left-hand inequality: quadratic lower bound.

Right-hand inequality:

$$\begin{aligned} F(x^*) &\geq \min_{y \in \text{dom}(F)} \left\{ F(y) + \langle \nabla F(x), y - x \rangle + \frac{\alpha}{2} \|y - x\|^2 \right\} \\ &= F(x) - \frac{1}{2\alpha} \|\nabla F(x)\|^2. \end{aligned}$$

If $F \in C_L^1$ and α -strongly convex, then

$$G(x) \stackrel{\text{def}}{=} F(x) - \frac{\alpha}{2} \|x\|^2$$

is convex, and ∇G is $L - \alpha$ -Lipschitz continuous.

The co-coercivity of ∇G yields

$$\langle \nabla F(x) - \nabla F(y), x - y \rangle \geq \frac{\alpha L}{\alpha + L} \|x - y\|^2 + \frac{1}{\alpha + L} \|\nabla F(x) - \nabla F(y)\|^2$$

for all $x, y \in \text{dom}(F)$.

$S_{\alpha, L}^1$: functions in C_L^1 that are α -strongly convex.

- Sequence x_k converges linearly to x^* if

$$\lim_{k \rightarrow +\infty} \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|} = \rho$$

holds for $\rho \in]0, 1[$, and ρ is called the rate of convergence.

- If x_k converges, let $\rho_k = \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|}$,
 - if $\lim_{k \rightarrow +\infty} \rho_k = 0$: super-linear convergence.
 - if $\lim_{k \rightarrow +\infty} \rho_k = 1$: sub-linear convergence.

- Superlinear convergence: $q > 1$

$$\lim_{k \rightarrow +\infty} \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|^q} < \eta$$

for some $\eta \in]0, 1[$.

- $q = 2$: quadratic convergence.
- $q = 3$: cubic convergence.

- 1 Unconstrained smooth optimisation
- 2 Descent methods
- 3 Gradient of convex functions
- 4 Gradient descent**
- 5 Heavy-ball method
- 6 Nesterov's optimal schemes
- 7 Dynamical system

Unconstrained smooth optimisation

Consider minimising

$$\min_{x \in \mathbb{R}^n} F(x),$$

where $F : \mathbb{R}^n \rightarrow \mathbb{R}$ is proper convex and smooth differentiable.

Assumptions:

- $F \in C^1$ is convex.
- $\nabla F(x)$ is L -Lipschitz continuous for some $L > 0$.
- Set of minimisers is non-empty, i.e. $\text{Argmin}(F) \neq \emptyset$.

Descent direction: let $d = -\nabla F(x)$, then

$$\langle \nabla F(x), d \rangle = -\|\nabla F(x)\|^2 \leq 0.$$

Gradient descent

initial : $x_0 \in \text{dom}(F)$;

repeat :

1. Choose step-size $\gamma_k > 0$
2. Update $x_{k+1} = x_k - \gamma_k \nabla F(x_k)$

until : stopping criterion is satisfied.

- Owing to the quadratic upper bound

$$\begin{aligned}F(x_{k+1}) &\leq F(x_k) + \langle \nabla F(x_k), x_{k+1} - x_k \rangle + \frac{L}{2} \|x_{k+1} - x_k\|^2 \\&= F(x_k) - \gamma \|\nabla F(x_k)\|^2 + \frac{\gamma^2 L}{2} \|\nabla F(x_k)\|^2 \\&= F(x_k) - \gamma \left(1 - \frac{\gamma L}{2}\right) \|\nabla F(x_k)\|^2.\end{aligned}$$

Hence

$$F(x_k) - F(x_{k+1}) \geq \gamma \left(1 - \frac{\gamma L}{2}\right) \|\nabla F(x_k)\|^2.$$

- Let $\gamma \in]0, 2/L[$,

$$\gamma \left(1 - \frac{\gamma L}{2}\right) \sum_{i=0}^k \|\nabla F(x_i)\|^2 \leq F(x_0) - F(x_{k+1}) \leq F(x_0) - F(x^*).$$

- $F(x^*) > -\infty$, rhs is a positive constant.
- for lhs, let $k \rightarrow +\infty$,

$$\lim_{k \rightarrow +\infty} \|\nabla F(x_k)\|^2 = 0.$$

NB: convexity is not required here.

- Let $\gamma \in]0, 1/L]$, then $\gamma(1 - \frac{\gamma L}{2}) \geq \frac{\gamma}{2}$, and

$$\begin{aligned} F(x_{k+1}) &\leq F(x_k) - \frac{\gamma}{2} \|\nabla F(x_k)\|^2 \\ (\text{cvx of } F \text{ at } x_k) &\leq F(x^*) + \langle \nabla F(x_k), x_k - x^* \rangle - \frac{\gamma}{2} \|\nabla F(x_k)\|^2 \\ &= F(x^*) + \frac{1}{2\gamma} (\|x_k - x^*\|^2 - \|x_k - x^* - \gamma \nabla F(x_k)\|^2) \\ &= F(x^*) + \frac{1}{2\gamma} (\|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2). \end{aligned}$$

- Summability of $F(x_k) - F(x^*)$,

$$\begin{aligned} \sum_{i=1}^k (F(x_i) - F(x^*)) &\leq \frac{1}{2\gamma} \sum_{i=1}^k (\|x_{i-1} - x^*\|^2 - \|x_i - x^*\|^2) \\ &= \frac{1}{2\gamma} (\|x_0 - x^*\|^2 - \|x_{k+1} - x^*\|^2) \\ &\leq \frac{1}{2\gamma} \|x_0 - x^*\|^2. \end{aligned}$$

- Since $F(x_k) - F(x^*)$ is decreasing

$$F(x_k) - F(x^*) \leq \frac{1}{k} \left(\sum_{i=1}^k (F(x_i) - F(x^*)) \right) \leq \frac{1}{2\gamma k} \|x_0 - x^*\|^2.$$

- Besides the basic assumptions, let's further assume $F \in S_{\alpha, L}^1$.

- Recall that, for all $x, y \in \text{dom}(F)$

$$\langle \nabla F(x) - \nabla F(y), x - y \rangle \geq \frac{\alpha L}{\alpha + L} \|x - y\|^2 + \frac{1}{\alpha + L} \|\nabla F(x) - \nabla F(y)\|^2.$$

- Analysis for constant step-size: let $\gamma \in]0, 2/(\alpha + L)[$

$$\begin{aligned} \|x_{k+1} - x^*\|^2 &= \|x_k - \gamma \nabla F(x_k) - x^*\|^2 \\ &= \|x_k - x^*\|^2 - 2\gamma \langle \nabla F(x_k), x_k - x^* \rangle + \gamma^2 \|\nabla F(x_k)\|^2 \\ (\nabla F(x^*) = 0) &\leq \left(1 - \frac{2\gamma\alpha L}{\alpha + L}\right) \|x_k - x^*\|^2 + \gamma \left(\gamma - \frac{2}{\alpha + L}\right) \|\nabla F(x_k)\|^2 \\ &\leq \left(1 - \frac{2\gamma\alpha L}{\alpha + L}\right) \|x_k - x^*\|^2. \end{aligned}$$

Distance to minimiser: $\rho = 1 - \frac{2\gamma\alpha L}{\alpha+L}$

$$\|x_k - x^*\|^2 \leq \rho^k \|x_0 - x^*\|^2.$$

- linear convergence

- for $\gamma = \frac{2}{\alpha+L}$,

$$\rho = \left(\frac{L-\alpha}{L+\alpha}\right)^2.$$

Convergence rate of objective function value:

$$F(x_k) - F(x^*) \leq \frac{L}{2} \|x_k - x^*\|^2 \leq \frac{\rho^k L}{2} \|x_0 - x^*\|^2.$$

Numer of iterations k needed for $F(x_k) - F(x^*) \leq \epsilon$

- $F \in C_L^1$: $O(1/\epsilon)$.

- $F \in S_{\alpha,L}^1$: $O(\log(1/\epsilon))$.

First-order method: x_k is an element from the set

$$x_0 + \text{span}\{\nabla F(x_0), \dots, \nabla F(x_i), \dots, \nabla F(x_{k-1})\}. \quad 4.1$$

Problem class: C_L^1

Nesterov's lower bound

For every integer $k \leq (n-1)/2$ and every x_0 , there exist functions in the problem class such that for any first-order method satisfies (4.1),

$$F(x_k) - F(x^*) \geq \frac{3}{32} \frac{L \|x_0 - x^*\|^2}{(k+1)^2},$$
$$\|x_k - x^*\|^2 \geq \frac{1}{8} \|x_0 - x^*\|^2.$$

- Suggests $O(1/k)$ is not the optimal rate.
- Accelerated gradient methods can achieve $O(1/k^2)$ rate.

1 Unconstrained smooth optimisation

2 Descent methods

3 Gradient of convex functions

4 Gradient descent

5 Heavy-ball method

6 Nesterov's optimal schemes

7 Dynamical system

Gradient descent:

$$-\gamma \nabla F(x_k) = x_{k+1} - x_k.$$

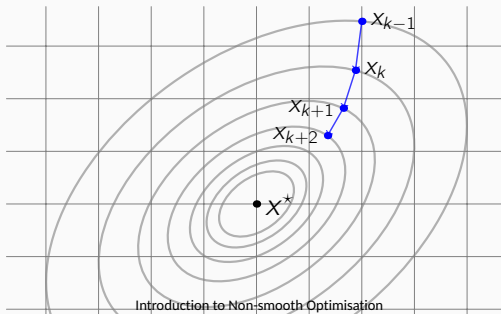
Consider the angle: $\theta_k \stackrel{\text{def}}{=} \text{angle}(\nabla F(x_{k+1}), \nabla F(x_k))$,

$$\lim_{k \rightarrow +\infty} \theta_k = 0.$$

Exercise: prove this claim for least square.

Let $a > 0$ be some constant,

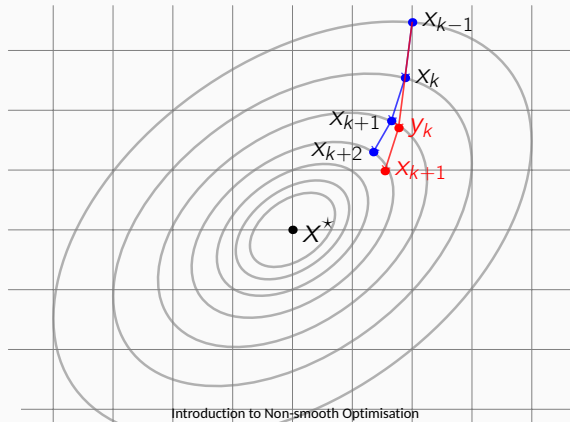
$$-\nabla F(x_{k+1}) \approx a(x_{k+1} - x_k).$$



Heavy-ball method (Polyak)

Initial : $x_0 \in \text{dom}(F)$ and $\gamma \in]0, 2/L[$;

$$y_k = x_k + a_k(x_k - x_{k-1}), \quad a_k \in [0, 1],$$
$$x_{k+1} = y_k - \gamma \nabla F(x_k).$$



Heavy-ball method (Polyak)

Initial : $x_0 \in \text{dom}(F)$ and $\gamma \in]0, 2/L[$;

$$\begin{aligned}y_k &= x_k + a_k(x_k - x_{k-1}), \quad a_k \in [0, 1], \\x_{k+1} &= y_k - \gamma \nabla F(x_k).\end{aligned}$$

- $x_k - x_{k-1}$ is called the inertial term or momentum term.
- a_k is called the inertial parameter.
- Convergence can be proved by studying the Lyapunov function

$$\mathcal{E}(x_k) \stackrel{\text{def}}{=} F(x_k) + \frac{a_k}{2\gamma} \|x_k - x_{k-1}\|^2.$$

- In general, no convergence rate for $F \in C_L^1$. Local rate for $F \in S_{\alpha,L}^2$.

Theorem

Let x^* be a (local) minimiser of F such that $\alpha \text{Id} \preceq \nabla^2 F(x^*) \preceq L \text{Id}$ and choose a, γ with $a \in [0, 1[, \gamma \in]0, 2(1+a)/L[$. There exists $\underline{\rho} < 1$ such that if $\underline{\rho} < \rho < 1$ and if x_0, x_1 are close enough to x^* , one has

$$\|x_k - x^*\| \leq C\rho^k.$$

Moreover, if

$$a = \left(\frac{\sqrt{L} - \sqrt{\alpha}}{\sqrt{L} + \sqrt{\alpha}} \right)^2, \quad \gamma = \frac{4}{(\sqrt{L} + \sqrt{\alpha})^2} \quad \text{then} \quad \underline{\rho} = \frac{\sqrt{L} - \sqrt{\alpha}}{\sqrt{L} + \sqrt{\alpha}}.$$

- Starting points need to be close enough to x^*
- Almost the optimal rate can be achieved by gradient method (or first-order method)
- Gradient descent

$$\underline{\rho} = \frac{L - \alpha}{L + \alpha}.$$

- Taylor expansion

$$x_{k+1} = x_k + a(x_k - x_{k-1}) - \gamma \nabla^2 F(x^*)(x_k - x^*) + o(\|x_k - x^*\|).$$

- Let $z_k = (x_k - x^*, x_{k-1} - x^*)^T$ and $H = \nabla^2 F$, then

$$z_{k+1} = \underbrace{\begin{bmatrix} (1+a)\text{Id} - aH & -a\text{Id} \\ \text{Id} & 0 \end{bmatrix}}_M z_k + o(\|z_k\|).$$

- Spectral radius $\rho(M)$, $\eta = 1 - \gamma\alpha$

$$0 = \rho^2 - (a + \eta)\rho + a\eta.$$

- $\rho(M)$ is a function of a and η (essentially γ).

- 1 Unconstrained smooth optimisation
- 2 Descent methods
- 3 Gradient of convex functions
- 4 Gradient descent
- 5 Heavy-ball method
- 6 Nesterov's optimal schemes**
- 7 Dynamical system

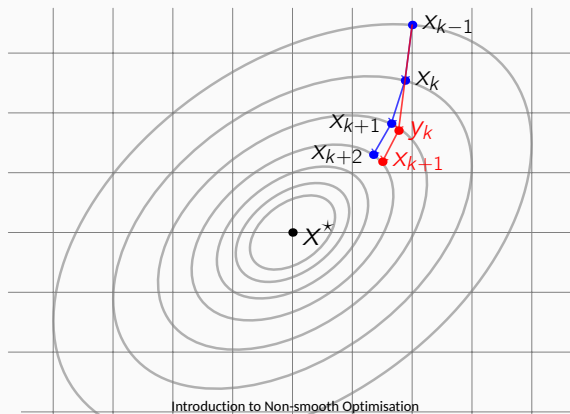
Gradient descent with constant step-size:

- $F \in \mathcal{C}_L^1$

$$F(x_k) - F(x^*) \leq \frac{L\|x_0 - x^*\|^2}{k+4}.$$

- $F \in \mathcal{S}_{\alpha,L}^1$

$$F(x_k) - F(x^*) \leq \frac{L}{2} \left(\frac{L-\alpha}{L+\alpha} \right)^2 \|x_0 - x^*\|^2.$$



Optimal scheme with constant step-size

initial : Choose $x_0 \in \mathbb{R}^n$, $\phi_0 \in]0, 1[$; Let $y_0 = x_0$ and $q = \alpha/L$.

repeat :

1. Compute $\phi_{k+1} \in]0, 1[$ from equation

$$\phi_{k+1}^2 = (1 - \phi_{k+1})\phi_k^2 + q\phi_{k+1}.$$

Let $a_k = \frac{\phi_k(1-\phi_k)}{\phi_k^2 + \phi_{k+1}}$ and

$$y_k = x_k + a_k(x_k - x_{k-1}).$$

2. Update x_{k+1} by

$$x_{k+1} = y_k - \frac{1}{L}\nabla F(y_k).$$

until : stopping criterion is satisfied.

Convergence rate

Let $\phi_0 \geq \sqrt{\alpha/L}$, then

$$F(x_k) - F(x^*) \leq \min \left\{ \left(1 - \sqrt{\frac{\alpha}{L}}\right)^k, \frac{4L}{(2\sqrt{L} + k\sqrt{\nu})^2} \right\} \\ \times \left(F(x_0) - F(x^*) + \frac{\nu}{2} \|x_0 - x^*\|^2 \right),$$

where $\nu = \frac{\phi_0(\phi_0 L - \alpha)}{1 - \phi_0}$.

Parameter choices:

- $F \in \mathcal{C}_L^1$: $\phi_0 = 1$,

$$q = 0, \quad \phi_k \approx \frac{2}{k+1} \rightarrow 0 \quad \text{and} \quad a_k \approx \frac{1 - \phi_k}{1 + \phi_k} \rightarrow 1.$$

- $F \in \mathcal{S}_{\alpha, L}^1$: $\phi_0 = \sqrt{\alpha/L}$

$$q = \sqrt{\frac{\alpha}{L}}, \quad \phi_k \equiv \sqrt{\frac{\alpha}{L}} \quad \text{and} \quad a_k \equiv \frac{\sqrt{L} - \sqrt{\alpha}}{\sqrt{L} + \sqrt{\alpha}}.$$

- 1 Unconstrained smooth optimisation
- 2 Descent methods
- 3 Gradient of convex functions
- 4 Gradient descent
- 5 Heavy-ball method
- 6 Nesterov's optimal schemes
- 7 Dynamical system**

From gradient descent

$$\frac{x_{k+1} - x_k}{\gamma} = -\nabla F(x_k).$$

Let γ be small enough

$$\dot{X}(t) + \nabla F(X(t)) = 0.$$

Discretisation

- Explicit Euler method

$$\dot{X}(t) = \frac{X(t+h) - X(t)}{h}.$$

- Implicit Euler method

$$\dot{X}(t) = \frac{X(t) - X(t-h)}{h}.$$

Given a 2nd order dynamical system

$$\ddot{X}(t) + \lambda(t)\dot{X}(t) + \nabla F(X(t)) = 0.$$

Discretisation:

- 2nd order term

$$\ddot{X}(t) = \frac{X(t+h) - 2X(t) + X(t-h)}{h^2}.$$

- Implicit Euler method

$$\dot{X}(t) = \frac{X(t) - X(t-h)}{h}.$$

Combine together:

$$X(t+h) - X(t) - (1 - h\lambda(t))(X(t) - X(t-h)) + h^2\nabla F(X(t)) = 0.$$

Choices:

- Heavy-ball: $h\lambda(t) \in]0, 1[$.
- Nesterov: $\lambda(t) = \frac{d}{t}$, $d > 3$.

- S. Boyd and L. Vandenberghe. "Convex optimization". Cambridge university press, 2004.
- B. Polyak. "Introduction to optimization". Optimization Software, 1987.
- Y. Nesterov. "Introductory lectures on convex optimization: A basic course". Vol. 87. Springer Science & Business Media, 2013.
- W. Su, S. Boyd, and E. Candès. "A differential equation for modeling Nesterov's accelerated gradient method: Theory and insights". Advances in Neural Information Processing Systems. 2014.