

① Continuous optimization : an Introduction

References: { Convex optimization (Stephen Boyd and Lieven Vandenberghe)
Numerical optimization (Nocedal and Wright)

* Mathematical problem

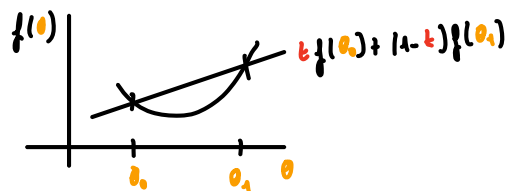
$$\left\{ \begin{array}{l} \min_{\theta} f_0(\theta) \\ \text{st.} \quad f_i(\theta) \leq 0 \\ i = 1, \dots, m \end{array} \right. \quad \text{Similarly} \quad \Rightarrow \quad \left\{ \begin{array}{l} \min_{\theta \in C} f_0(\theta) \\ \text{where } C = \{ \theta \mid f_i(\theta) \leq 0 \} \end{array} \right.$$

{ Names
+ Optimization
+ Mathematical programming

* Convex optimization

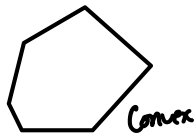
* We ask { $f_0(\theta)$ to be convex
 C to be a convex set

* Convex function: $f(t\theta_0 + (1-t)\theta_1) \leq t f(\theta_0) + (1-t)f(\theta_1) \quad \forall \theta_0, \theta_1, \forall t \in [0,1]$



* Convex set: $t\theta_0 + (1-t)\theta_1 \in C$

$\forall \theta_0, \theta_1, \forall t \in [0,1]$ {
+ Examples {
Linear constraints: $A\theta \leq b$
Euclidian norms: $\|\theta\|^2 \leq r$
Bounds: $\theta_{\min} \leq \theta \leq \theta_{\max}$



* Not trivial to always identify convex sets

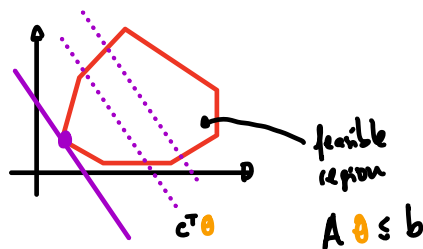
- * Some operators preserve convexity
- * There exists tricks to easily identify convex sets

* Why convexity important?

- If the function is convex, then a local minimum is a global minimum
- Sometimes we only have local convexity \Rightarrow global optimum not guaranteed

* First example: linear programming (convex problem)

$$\begin{cases} \min_{\theta} & c^T \theta \\ \text{s.t.} & A \theta \leq b \end{cases}$$

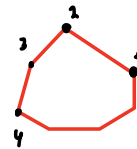


* lots of applications

* Most of the times θ is really large

* Algorithms

- * Simplex: optimal solution traversing the edges
- * Interior point: optimal solution



Unconstrained optimization

$\min_{\theta} f(\theta)$

- f might be convex
- $f \in C^1, \exists \nabla f(\theta) \forall \theta \rightarrow \|\theta\|$ is not differentiable
- In some applications, constraints are added as a penalization in the cost function
- lots of examples in ML

* First order optimality condition

If f is convex and C^1 , then

$$\forall \alpha, \theta \quad f(\theta) + \nabla f^T(\theta)(\alpha - \theta) \leq f(\alpha)$$



* Proof in 1D : $f((1-t)\theta + t\alpha) \leq (1-t)f(\theta) + tf(\alpha)$

$$\frac{f((1-t)\theta + t\alpha) - f(\theta)}{t} + f(\theta) \leq f(\alpha) \xrightarrow{\text{take } t \rightarrow 0} f(\theta) + \nabla f^T(\theta)(\alpha - \theta) \leq f(\alpha)$$

* Taylor approximation of $f(\alpha)$ at θ

+ If the function is convex,

local information like function and gradient at a point θ provides global information (global underestimator)

* If $\exists \theta$ s.t. $\nabla f(\theta) = 0 \Rightarrow f(\theta) \leq f(\alpha) \quad \forall \alpha \Rightarrow \theta$ global minimizer of f

* If f is not convex, but locally convex and C^2 , a Taylor expansion gives

$$f(\alpha) = f(\theta) + \nabla f^T(\theta)(\alpha - \theta) + \frac{1}{2}(\alpha - \theta)^T \nabla^2 f(\theta)(\alpha - \theta) + \dots$$

If $\exists \theta^*$ s.t. $\nabla f(\theta^*) = 0$, then for $\forall \alpha$ s.t. $(\alpha - \theta^*)^T \nabla^2 f(\theta^*)(\alpha - \theta^*) \geq 0$

$$\text{we have } f(\theta^*) = f(\alpha) - (\alpha - \theta^*)^T \nabla^2 f(\theta^*)(\alpha - \theta^*) \leq f(\alpha)$$

* Second example: least squares problem

$$\min_{\theta} f_0(\theta) = \|A\theta - b\|_2^2 = (A\theta - b)^T (A\theta - b) = (A_{ij}\theta_j - b_i)(A_{ik}\theta_k - b_i)$$

$$\begin{aligned} \text{Gateaux derivative : } \frac{\partial f_0}{\partial \theta_m}(\tilde{\theta}) &= A_{ij} \delta_{jm} \tilde{\theta} (A_{ik} \theta_k - b_i) + (A_{ij} \theta_j - b_i) (A_{ik} \delta_{km} \tilde{\theta}) \\ &= [A_{im} (A_{ik} \theta_k - b_i)] + A_{im} [A_{ik} \theta_k - b_i] = 2[A^T (A\theta - b)] \tilde{\theta} = 0 \end{aligned}$$

We have to solve $[A^T A] \theta = A^T b$ (Normal equations)

* $A^T A$ is always square and symmetric.

* Case A square \Rightarrow same as $A \theta = b$

* Case A has vertical shape (overdetermined) $A = \begin{bmatrix} \end{bmatrix}$

↳ Minimize the error in the columns of A

* Case A has horizontal shape (underdetermined) $A = \begin{bmatrix} \end{bmatrix}$

↳ $A^T A$ has not full rank \Rightarrow Ill-posed problem

↳ Plenty of solution

↳ In ML terms is an example of overfitting

* Algorithms

* Gradient method

* $\theta_{k+1} = \theta_k - \eta \nabla f(\theta_k)$ where $\eta > 0$ is small enough

* η is called line-search (or learning rate in ML)

* Descent direction: $f(\theta_{k+1}) = f(\theta_k) + \nabla f(\theta_k)^T (\theta_{k+1} - \theta_k) + \dots$
 $= f(\theta_k) - \eta \|\nabla f(\theta_k)\|^2 \leq f(\theta_k)$

* Steepest descent: find $\Delta \theta$ s.t.
(Euclidean norm)

$$\begin{cases} \min \nabla f(\theta_k)^T \Delta \theta \\ \text{s.t. } \|\Delta \theta\| \leq 1 \end{cases} \Rightarrow \|\Delta \theta\| = \frac{-\nabla f(\theta_k)}{\|\nabla f(\theta_k)\|}$$

* Stopping criteria: $\|\nabla f(\theta_k)\| \leq \eta$

* Line-search (learning-rate)

* Exact line-search

Find z along the ray $\{\theta + z\Delta\theta; z \geq 0\}$:

$$z = \arg \min_{z \geq 0} f(\theta + z\Delta\theta)$$

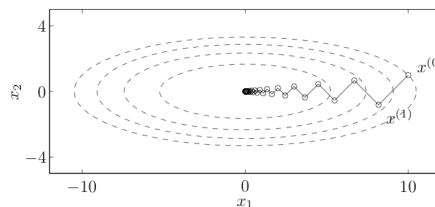
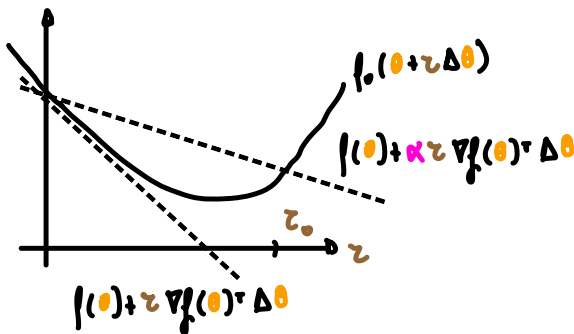
- Usually used when computing the cost is much cheaper than computing the gradient or direction
- Sometimes can be computed analytically or very efficiently

* Backtracking line-search (Inexact)

* Approximately minimize f_0 or reduce f_0 "enough"

* Given $\Delta\theta$, take $z = z^*$; $\alpha = (0, 0.5)$, $\beta \in (0, 1)$

$$\left[\begin{array}{l} \text{while } f(\theta_k + z\Delta\theta_k) > f_0(\theta_k) + \alpha z \nabla f(\theta_k)^T \Delta\theta_k; \\ \quad z = \beta z \end{array} \right.$$



For small z :

$$f_0(\theta_k + z\Delta\theta_k) \approx f_0(\theta_k) + z \overbrace{\nabla f_0(\theta_k)^T \Delta\theta_k}^{< 0} < f_0(\theta_k) + z \alpha \nabla f_0(\theta_k)^T \Delta\theta_k$$

α = $\left\{ \begin{array}{l} \text{Fraction of the decrease in } f \\ \text{predicted by linear extrapolation} \\ \text{that will be accept} \\ \text{Typically between } \alpha \in [0.01, 0.2] \end{array} \right.$

β = $\left\{ \begin{array}{l} \text{Between } \beta \in [0.1, 0.9] \\ \text{Typical number } \beta = 0.5 \end{array} \right.$

* Steepest descent for quadratic norms

Quadratic norm: $\|\theta\|_P = (\theta^T P \theta)^{1/2} = \|P^{1/2} \theta\|_2$

• $\begin{cases} \min \nabla f(\theta_k)^T \Delta \theta \\ \text{s.t. } \|\Delta \theta\|_P \leq 1 \end{cases} \Rightarrow \Delta \theta = -P^{-1} \nabla f(\theta_k)$

$\theta_{k+1} = \theta_k - \tau P^{-1} \nabla f(\theta_k)$ $\left\{ \begin{array}{l} \cdot \text{ If } P \text{ is best squares of linear system of equation: } P \text{ is called preconditioner} \\ \cdot \text{ Interesting when } P^{-1} \text{ is cheap} \\ \cdot \text{ change of coordinates to better condition the problem} \end{array} \right.$



* Newton's method

• $\theta_{k+1} = \theta_k - [\nabla^2 f(\theta_k)]^{-1} \nabla f(\theta_k)$

* Descent direction:

$f(\theta_{k+1}) = f(\theta_k) + \nabla f^T(\theta_k) \Delta \theta_k + \dots$ if $\nabla^2 f(\theta_k)$ s.p.d.
 $= f(\theta_k) - \nabla f^T(\theta_k) [\nabla^2 f(\theta_k)]^{-1} \nabla f(\theta_k) \leq f(\theta_k)$

* Minimization of second-order approximation

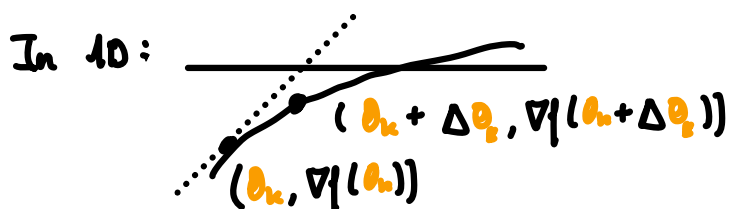
$\left\{ \begin{array}{l} f(\theta_{k+1}) \approx f(\theta_k) + \nabla f^T(\theta_k) \Delta \theta_k + \Delta \theta_k^T [\nabla^2 f(\theta_k)] \Delta \theta_k \\ \min_{\Delta \theta_k} f(\theta_{k+1}) \approx \min_{\Delta \theta_k} f(\theta_k) + \nabla f^T(\theta_k) \Delta \theta_k + \Delta \theta_k^T [\nabla^2 f(\theta_k)] \Delta \theta_k \\ \hookrightarrow \Delta \theta_k = -[\nabla^2 f(\theta_k)]^{-1} \nabla f(\theta_k) \end{array} \right.$

* Steepest descent in Hessian norm

$$\left\{ \begin{array}{l} \min \nabla f(\theta_k)^T \Delta \theta \\ \text{s.t. } \|\Delta \theta\|_{\nabla^2 f} \leq 1 \end{array} \right\} \quad \Delta \theta = - \underbrace{[\nabla^2 f(\theta_k)]^{-1}}_{\text{it can be seen as pseudoinverse}} \nabla f(\theta_k)$$

* Solution of linearized optimality condition

$$\nabla f(\theta_k + \Delta \theta_k) \approx \nabla f(\theta_k) + \nabla^2 f(\theta_k) \Delta \theta_k = 0 \Rightarrow \Delta \theta = -[\nabla^2 f(\theta_k)]^{-1} \nabla f(\theta_k)$$



* { Highly used in computational engineering
Rarely used in ML

* { Convergence when close to the solution
Quadratic convergence (at most 6/7 iterates)
Stopping criteria: $\lambda(\theta_k) = \nabla f^T(\theta_k) [\nabla^2 f(\theta_k)]^{-1} \nabla f(\theta_k) \leq \epsilon$

* Implicit method

{ * Descent method has some similarities with ODEs
* Gradient is explicit
* Newton and quadratic are implicit
* { Explicit { Each iter fast
many iterations
Implicit { Each iter expensive (in general)
Few iterations

* Quasi-newton
{ Approximation of the hessian with products