

§ Multidimensional optimisation (without constraints)

I. Reminders... $\forall \mathbb{R}$

1D

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

f' intuition \rightarrow slope of the tangent

f''

\vdots

Taylor expansion formula (order m)

$$f(x_0+h) = \sum_{k=0}^m \frac{f^{(k)}(x_0)}{k!} \times h^k + o(h^k)$$

polynomial (h)

\downarrow
order 1

$$f(x_0+h) = f(x_0) + \underbrace{f'(x_0) \cdot h}_{(*)} + o(h)$$

order 1

linear fct

tangent

tangent

1D \sim mD

mD — optimization

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

/ general

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

Order 1 derivatives

$$f' = \lim_{h \rightarrow 0} \frac{f(\vec{x}_0 + \vec{h}) - f(\vec{x}_0)}{h}$$

Partial derivatives

$$\frac{\partial f}{\partial x_i}(\vec{x}^0)$$

\downarrow

$$f(x_1, \dots, x_n)$$

all variables are fixed except ith one

$$x_i \mapsto f(x_1^0, \dots, x_{i-1}^0, \underbrace{x_i}_{(*)}, x_{i+1}^0, \dots, x_n^0)$$

$$\mathbb{R} \rightarrow \mathbb{R}$$

Differential

(*) f differentiable:

$$f(\vec{x}_0 + \vec{h}) = f(\vec{x}_0) +$$

$$Df(\vec{x}_0) \cdot \vec{h}$$

linear fct

differential fct

\downarrow
linear fct...

$Df(\vec{x}_0)$ linear ...

\downarrow
matrix

$$Df(\vec{x}_0): \mathbb{R}^m \rightarrow \mathbb{R}$$

(basis \rightarrow m vekt
 $f_{e_1} \dots f_{e_m}$
size of its matrix?

$$m \times m \begin{pmatrix} \cdot & \cdot & \cdot & \cdot \end{pmatrix} \begin{matrix} \text{m cols} \\ \downarrow \end{matrix}$$

$$Df(e_1) \downarrow \dots \downarrow Df(e_m) \in \mathbb{R}$$

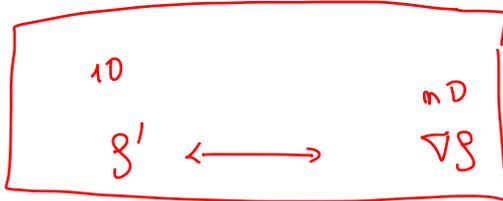
$Df(\vec{x}_0) \longleftrightarrow 1 \times m$ matrix

canonical basis:

$$\left(\frac{\partial f}{\partial x_1} \quad \dots \quad \frac{\partial f}{\partial x_m} \right)$$

$$\nabla f(\vec{x}_0) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_m} \end{pmatrix}$$

matrix of Df
 \downarrow
 ∇f^t



Taylor formula (order 1) linear

$$f(\vec{x}_0 + \vec{h}) = f(\vec{x}_0) + \overbrace{Df(\vec{x}_0)}^{\text{linear}} \cdot \vec{h} + o(\|\vec{h}\|)$$

$$\parallel$$

$$\nabla f(\vec{x}_0)^t \times \vec{h}$$

$$\parallel$$

$$\langle \nabla f(\vec{x}_0), \vec{h} \rangle$$

Remainders

\vec{u}, \vec{v} vectors

$\langle \vec{u}, \vec{v} \rangle$ 2 notations $\vec{u} \cdot \vec{v}$

$$\parallel$$

$$\vec{u}^t \times \vec{v}$$

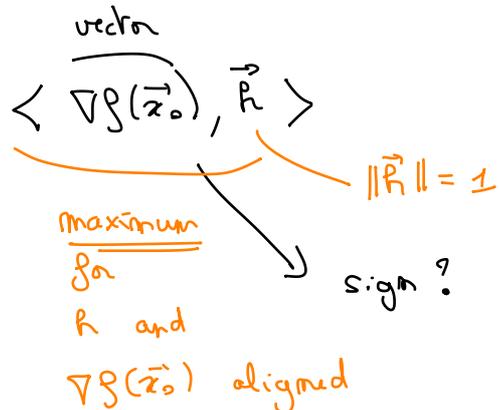
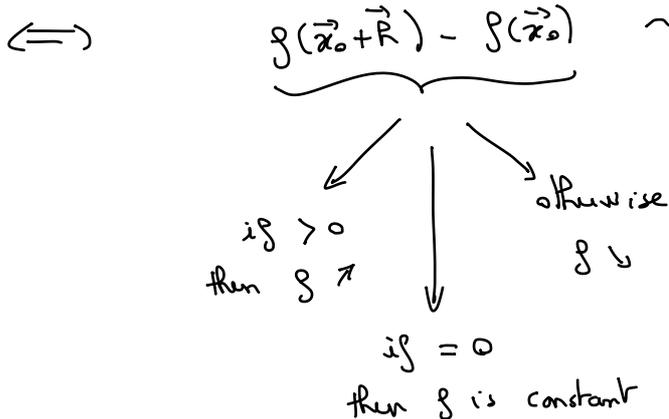
$\left. \begin{matrix} \} \\ (A \times B) \end{matrix} \right\}$

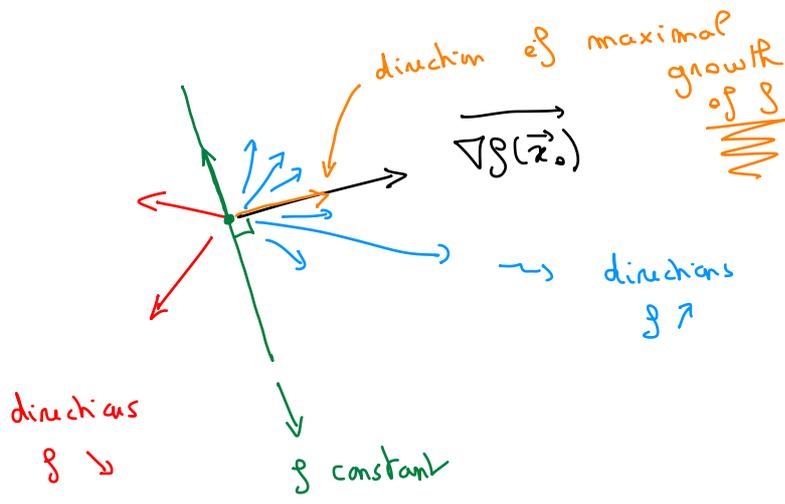
$$\begin{matrix} u^t \\ (\dots) \end{matrix} \times \begin{pmatrix} v \\ \vdots \\ \vdots \end{pmatrix} = \begin{matrix} 1 \times 1 \leftrightarrow \mathbb{R} \\ (\cdot) \end{matrix}$$

Intuition behind the gradient

Taylor formula

$$f(\vec{x}_0 + \vec{h}) \sim f(\vec{x}_0) + \langle \nabla f(\vec{x}_0), \vec{h} \rangle$$

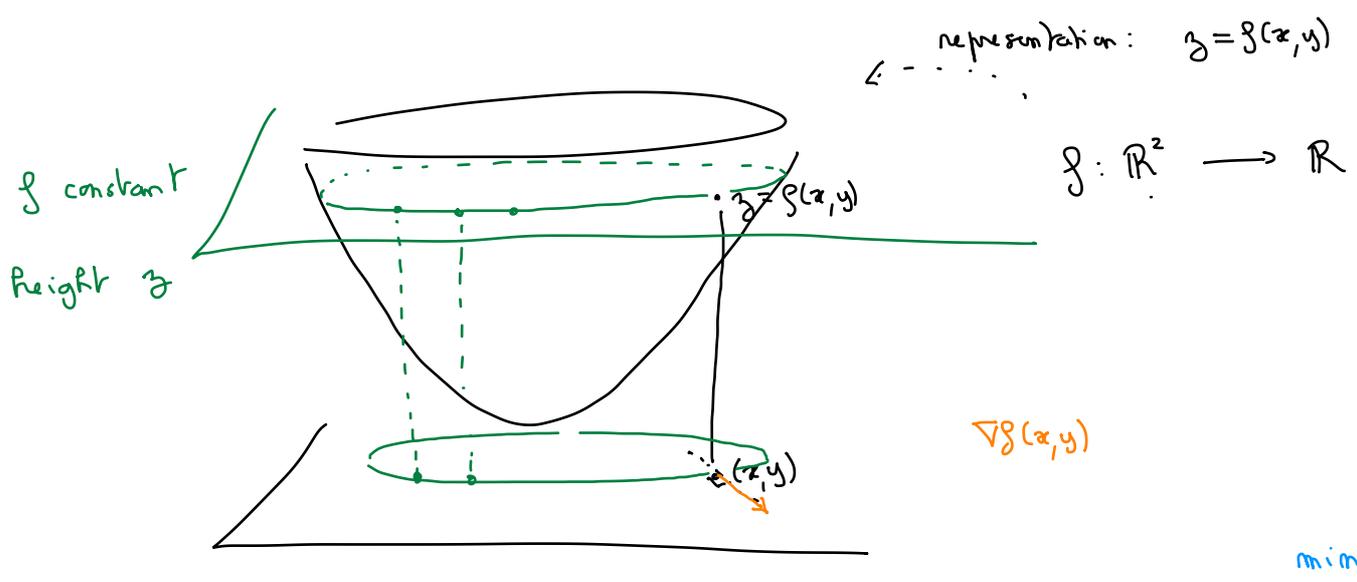




\vec{h} st $\langle \nabla g(\vec{x}_0), \vec{h} \rangle > 0$

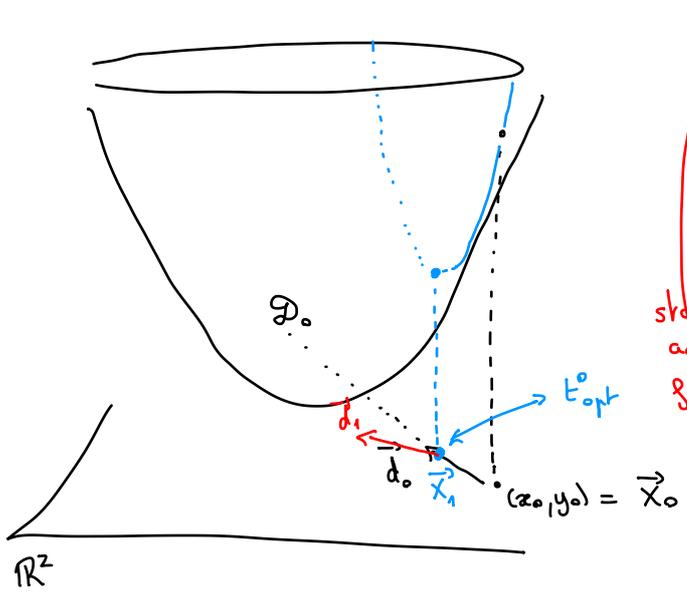
\vec{h} st $\dots < 0$

\vec{h} st $\dots = 0$



II. Descent algorithms

Example $g: \mathbb{R}^2 \rightarrow \mathbb{R}$ represented by a surface



g is decreasing along \vec{d}_0

\mathcal{D}_0 : half-line (starting from \vec{x}_0 with direction \vec{d}_0)

$x_0 + t \cdot \vec{d}_0 \quad t \in \mathbb{R}^+$

find the min of g along this line

x_1

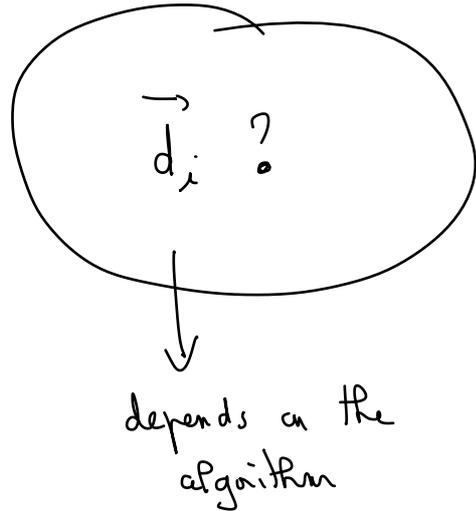
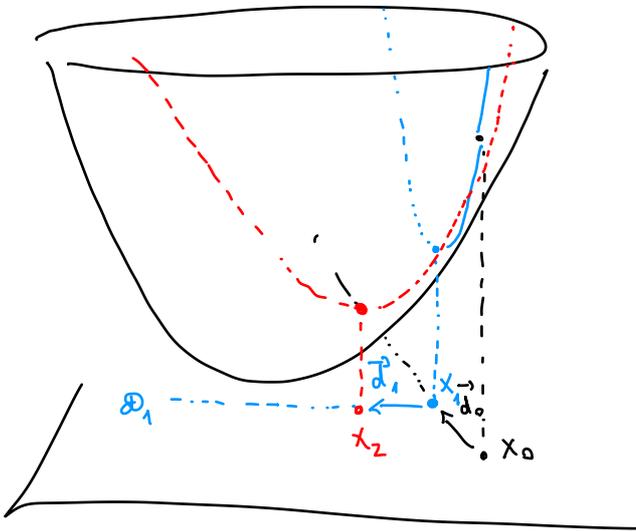
(*) 1D PE

(*) find t s.t

$g(\vec{x}_0 + t \cdot \vec{d}_0)$ minimum

$$\begin{cases} t \mapsto g(\vec{x}_0 + t \cdot \vec{d}_0) & \text{--- min?} \\ \mathbb{R} \longrightarrow \mathbb{R} \end{cases}$$

1D optimization
§ 1



Most basic

gradient descent algorithm

$$\vec{d}_i = - \nabla g(\vec{x}_i)$$

→ greedy algorithm

- ⊖ not the fastest — many iterations
- ⊕ numerically very stable

order 1 methods



flat functions — challenging for opt. algorithms ...

conjugate gradient descent algorithm

⊕ faster

⊖ may fail with non convex functions ...

Order 2 methods

based on ^(mD) Newton-Raphson

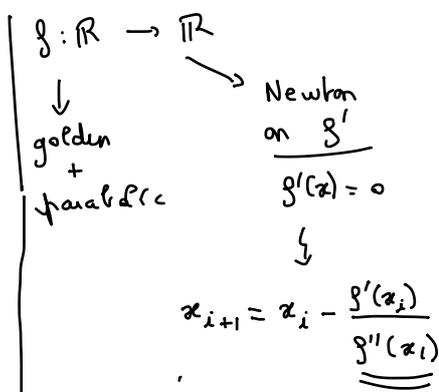
Newton-Raphson method

Quasi-Newton methods

approximations of Hessian

$$x_{i+1} = x_i - H(g)(x_i)^{-1} \times \nabla g(x_i)$$

mD Newton



intuition

" \mathbb{R}^m "

$$x_{i+1} = x_i - \frac{\nabla g(x_i)}{\text{second order derivative of } g}$$

matrix - Hessian matrix
 $H(g)(x_i)$

$$\left(\frac{\partial^2 g}{\partial x_i \partial x_j} \right)$$

① Gradient descent alg.

gradient-descent ($g, \nabla g, x_0, \epsilon, \delta, N$)

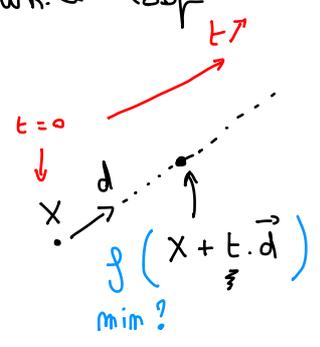
$x \leftarrow x_0$ // x : current point

iter $\leftarrow 0$

$x_{prev} \leftarrow x + \text{ones}(\text{size}(x))$ // just to enter the while loop

while ($(\|x - x_{prev}\| > \epsilon) \ \&\& \ (\text{iter} < N)$)

{
 $d \leftarrow -\nabla g(x)$
 // Optimize along the Ray (x, d)
 $\varphi \leftarrow @(\text{t}) \ g(x + \text{t} \cdot \vec{d})$



// 1D optim (\rightarrow Brent (parabolic + golden) \rightarrow $g_{min} \text{ bnd}$)
 here: 1Doptim($\varphi, \text{interval}$ s.t. the fun is unimodal)

// Sweeping method to init the interval

$t_0 \leftarrow 0$
 $t_1 \leftarrow \delta$
 $y_0 \leftarrow \varphi(t_0)$
 $y_1 \leftarrow \varphi(t_1)$
 while ($y_0 \geq y_1$) // move right

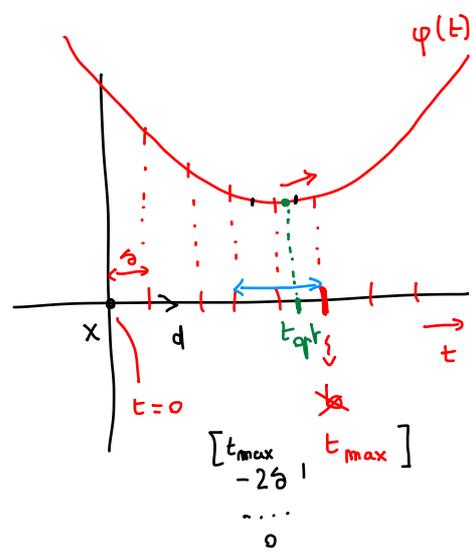
$t_0 \leftarrow t_1$
 $y_0 \leftarrow y_1$
 $t_1 \leftarrow t_0 + \delta$
 $y_1 \leftarrow \varphi(t_1)$

end

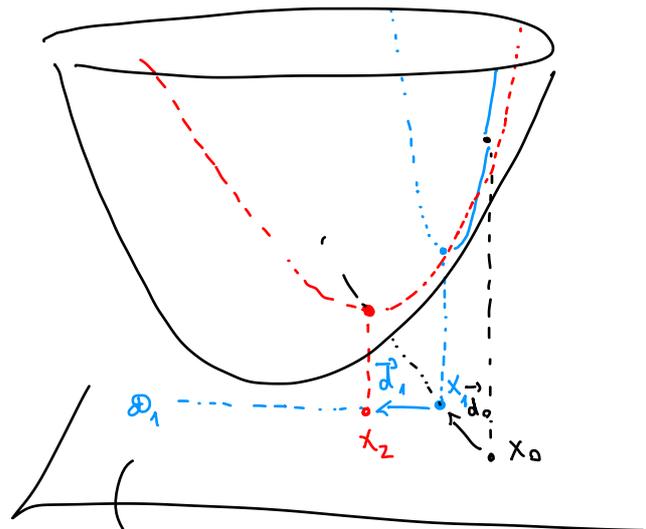
// $t_1 \leftarrow t_{\max} - \frac{t_1 - 2\delta}{2}$

$t_{opt} \leftarrow \text{1Doptim}(\varphi, 0, t_1, \epsilon, (N))$

$X_{prev} \leftarrow X$
 $X \leftarrow X + t_{opt} \cdot d$
 end



end



stop: $\|X_i - X_{i+1}\| < \epsilon$

$\nabla \varphi$
 $\nabla \varphi$
 X_i Cauchy sequence

stop $\leftrightarrow \|X_i - X^*\| < \epsilon$

t_{opt} — optimal for φ

optimal step gradient descent

$d_i \perp d_{i+1}$

Property
 $\forall i$

Prop:

$X_i \rightsquigarrow$ direction: $d_i = -\nabla \varphi(X_i)$

$$\varphi(t) = f(x_i + t \cdot d_i)$$



$$x_{i+1} = x_i + t_{opt} \cdot d_i \leftarrow$$

min of φ

$$\varphi'(t_{opt}) = 0$$

$$\varphi'(t) = \overbrace{f'(x_i + t \cdot d_i)}^{\text{vech}} \cdot \underbrace{d_i}_{\text{vech}}$$

\downarrow

\cdot

intuitively $\sim 1D$

$$(f \circ g)' = f' \circ g \times g'$$

$$\varphi'(t) = \nabla f(x_i + t \cdot d_i) \cdot d_i$$

At $t_{opt} \implies \varphi'(t_{opt}) = \nabla f(x_i + t_{opt} \cdot d_i) \cdot \vec{d}_i$

\downarrow

x_{i+1}

$$\nabla f(x_{i+1}) = -\vec{d}_{i+1}$$

$$\Leftrightarrow \boxed{0 = \nabla f(x_{i+1}) \cdot \vec{d}_i}$$

\downarrow

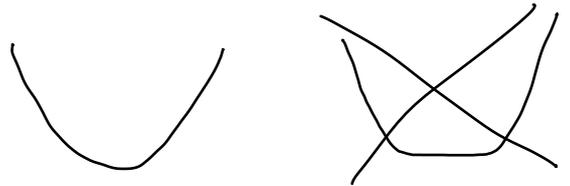
Convergence

$\int \nabla^2 f$ f is \mathcal{C}^2 and elliptic then the algorithm converges

$$\left(\begin{array}{l} \exists \alpha \text{ st. } \forall x, y \\ \|\nabla f(x) - \nabla f(y)\| \geq \alpha \|x - y\| \end{array} \right)$$

\nwarrow

"curvature"



Reminder - quadratic forms / polynomials of d^2 exactly

ex: \mathbb{R}^3 $f(x,y,z) = x^2 - y^2 + 2z^2 + xz + 2yz$

can be represented by a matrix

$X^t \cdot A \cdot X$

$A = \begin{pmatrix} x & y & z \\ y & & \\ z & & \end{pmatrix}$
 symmetric

squared terms

$\rightarrow xz \rightarrow$ split symmetrically

$\begin{pmatrix} x & y & z \\ 1 & 0 & \frac{1}{2} \\ 0 & -1 & 1 \\ \frac{1}{2} & 1 & 2 \end{pmatrix}$

$f(x,y,z) = f(X) = X^t \cdot A \cdot X$

$f(X) = A \cdot X$

$\nabla f(X) = A^t$

Gradient of such a function:

$\nabla f = 2AX$

! is A symmetric

If A not symmetric

$\nabla f = (A + A^t) \cdot X$

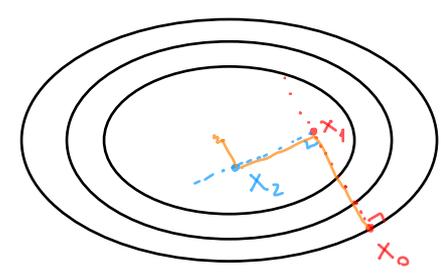
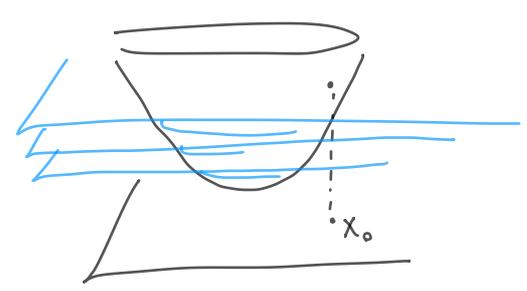
Problem: practical

$f_1(x,y) = \alpha x^2 + \beta y^2$ $\alpha, \beta > 0$

steepest gradient descent converges to min

for this BASIC function

Expected: STOPs at the min



steps

The problem comes from:

$\rightarrow \vec{d}_i \perp \vec{d}_{i+1}$ \rightsquigarrow the choice is "limited"
 \rightsquigarrow always the same directions
 \rightsquigarrow if \vec{d}_0 does not cross the min
 no \vec{d}_i will cross it ...

\rightarrow the alg. is greedy ...

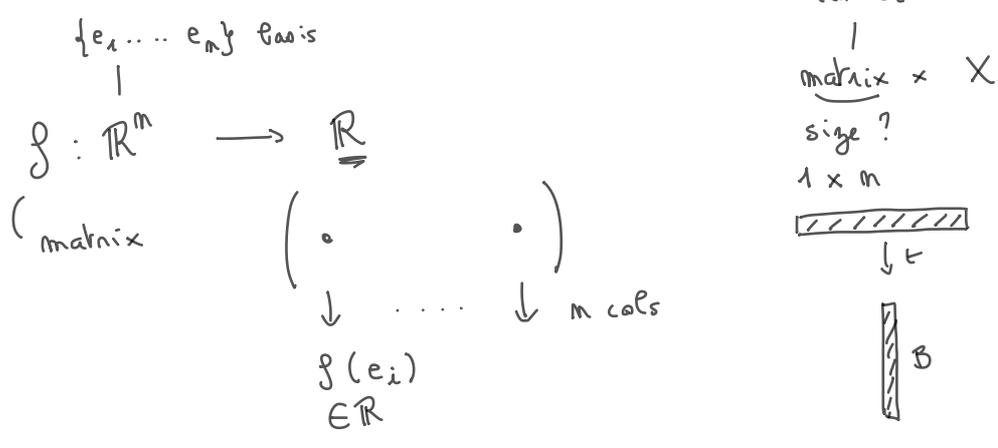
Conjugate gradient descent improves the convergence

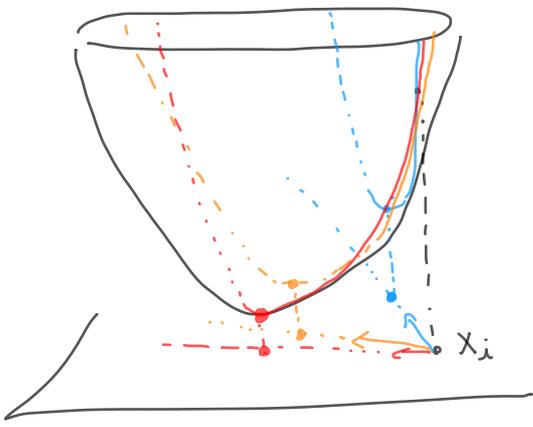
\hookrightarrow designed for quadratic functions
 formulae come from such functions — polynomial $d^{\circ} 2$
 \hookrightarrow can be extended to \mathbb{C}^2 elliptic functions

② Conjugate gradient descent

Idea: f is a quadratic function

$$f(x) = \underbrace{x^t \cdot A \cdot x}_{\text{order 2}} + \underbrace{t B \cdot x}_{\substack{\text{order 1} \\ \text{linear}}} + \underbrace{\frac{R}{c}}_{\substack{\text{order 0} \\ \text{constant}}}$$





each \vec{d} \rightarrow a min
 (the min of g
 along line (x_i, \vec{d}))
 \downarrow
 which \vec{d} gives the best minimum?
 ∇
 $-\nabla g(x_i)$

For quadratic function we can compute an exact formula for this "optimal" \vec{d}

This "optimal" \vec{d} is a linear combination $\begin{cases} \nabla g(x_i) \\ \text{previous direction } \vec{d}_{i-1} \end{cases}$

Optimal direction:

$$\vec{d}_{i+1} = \nabla g(x_{i+1}) + \frac{\|\nabla g(x_{i+1})\|^2}{\|\nabla g(x_i)\|^2} \vec{d}_i$$

g quadratic $g(x) = \frac{1}{2} x^T \cdot A \cdot x + B^T \cdot x + c$
 x_0 : close to the minimum ...

Alg:

$$\nabla g(x) = Ax + B$$

$$\varphi(t) = g(x_i + t \cdot \vec{d}_i)$$

Gradient Steepest descent

$$\vec{d}_i = -\nabla g(x_i)$$

$$x \leftarrow x_0$$

$$x_{\text{prev}} \leftarrow x + \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

$$\text{iter} \leftarrow 0$$

$$d \leftarrow \nabla g(x) \quad \triangle ! \quad d : - \text{descent direction}$$

while $((\|x - x_{\text{prev}}\| > \epsilon) \ \&\& \ (\text{iter} < N))$

$$\left\{ \begin{array}{l} \varphi(t) = g(x - t \cdot d) \\ \text{// min } \varphi? \end{array} \right.$$

// if g is quadratic

$$t_{\text{opt}} = \frac{\|\nabla g(x_i)\|^2}{x_i^T \cdot A \cdot x_i}$$

// if $g \in \mathbb{R}^2$ elliptic

// 1D optim \rightarrow find min φ

$$t_0 = 0$$

$$t_1 = \delta$$

$$y_0 = \varphi(t_0)$$

$$y_1 = \varphi(t_1)$$

while ($y_0 > y_1$)

$$t_0 = t_1$$

$$y_0 = y_1$$

$$t_1 = t_0 + \delta$$

$$y_1 = f(t_1)$$

end

$$t_{opt} = \text{IDoptim}(\varphi, \theta, t_1)$$

$$X_{prev} \leftarrow X$$

$$X \leftarrow X - t_{opt} \cdot d$$

$$d \leftarrow \nabla f(X) + \frac{\|\nabla f(X)\|^2}{\|\nabla f(X_{prev})\|^2} \cdot d$$

end

Convergence ...
 → case f quadratic
 → case f \mathbb{E}^2 elliptic → the alg. converges

When f quadratic and A definite positive :

all eigenvalues are > 0

\Leftrightarrow

$$\forall X \cdot A \cdot X > 0 \quad \forall X \neq 0$$

Property

$\forall i \quad \forall j < i$

\vec{d}_i, \vec{d}_j are conjugate with respect to A



$$\langle A d_j, d_j \rangle = 0$$

conjugate \approx orthogonal for the A -dot product

\vec{d}_i conjugate with $\vec{d}_0 \dots \vec{d}_{i-1}$

Question

is $f: \mathbb{R}^n \rightarrow \mathbb{R}$

quadratic elliptic / A sym. definite positive

the algorithm converges in at most n steps.

steps

For C^2 functions \Rightarrow faster than gradient descent ...

Help... $f(x) = \frac{1}{2} \cdot x^T \cdot A \cdot x + B^T \cdot x + c$

$g(x) = x^T \cdot A \cdot x + B^T x + c$
 $\nabla g(x) = 2AX + BX$

$\nabla f(x) = A \cdot x + B$

In the definition of \vec{d} :

$\| \nabla f(x_i) \|^2 = \| Ax_i + B \|^2 = \langle Ax_i + B, Ax_i + B \rangle$

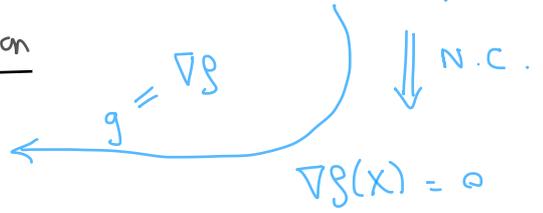
Libraries \rightsquigarrow many variants ...

pre-conditioned conjugate gradients
 bi-conjugate gradients

III . Order 2 methods \equiv Newton

① Newton in \mathbb{R}^m ... \rightarrow Newton-Raphson

$g: \mathbb{R}^m \rightarrow \mathbb{R}$
 Optimize $g \mid X: \min$



Goal \rightarrow compute the zeros of g
 roots

$g: \mathbb{R}^m \rightarrow \mathbb{R}^m$

$\nabla g: \mathbb{R}^m \rightarrow \mathbb{R}^m$
 $x \mapsto \nabla g(x)$

ex: $\mathbb{R}^3 \rightarrow \mathbb{R}^3$
 $g \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} x^2 + \cos(yz) \\ 2x + \ln(z) \\ xyz \end{pmatrix}$

Newton \leftrightarrow solve non linear systems

$g \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \vec{0} \iff \begin{cases} x^2 + \cos(yz) = 0 \\ 2x + \ln(z) = 0 \\ xyz = 0 \end{cases}$

Algorithm ...

Intuitively 1D \rightsquigarrow mD

1D

$$g(x) = 0$$

Newton iteration:

$$x_{i+1} = x_i - \frac{g(x_i)}{g'(x_i)}$$

mD

$$\vec{x}_{i+1} = \vec{x}_i - \frac{g(\vec{x}_i)}{g'(\vec{x}_i)}$$

vector

$$g: \mathbb{R}^m \rightarrow \mathbb{R}^m$$

~~$\frac{\partial g}{\partial x_i} \in \mathbb{R}^m$~~
Jacobian: $J(g)$

$$g: \mathbb{R}^m \rightarrow \mathbb{R}$$
$$\nabla g = \begin{pmatrix} \frac{\partial g}{\partial x_1} \\ \vdots \\ \frac{\partial g}{\partial x_m} \end{pmatrix}$$

Jacobian — "set" of partial derivatives
↓
matrix

$$g: \mathbb{R}^m \rightarrow \mathbb{R}^m$$

$g(\vec{x})$ vector dim m

$$\begin{pmatrix} g_1(\vec{x}) \\ \vdots \\ g_m(\vec{x}) \end{pmatrix}$$

$$g_i: \mathbb{R}^m \rightarrow \mathbb{R}$$

$$J(g) = \left(\frac{\partial g_i}{\partial x_j} \right)_{i,j \in \{1, \dots, m\}} = \begin{pmatrix} \nabla g_1^t & \dots \\ \vdots & \vdots \\ \nabla g_m^t & \dots \end{pmatrix}$$

special case
 $g: \mathbb{R}^n \rightarrow \mathbb{R}$
 ∇g linear $\mathbb{R}^n \rightarrow \mathbb{R}^m$
↕
matrix
Jacobian
↓
Gradient ∇g

ex:

$$\mathbb{R}^3 \rightarrow \mathbb{R}^3$$
$$g \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} x^2 + \cos(yz) \\ 2x + \ln(z) \\ xyz \end{pmatrix} \begin{matrix} \leftarrow g_1 \\ \leftarrow g_2 \\ \leftarrow g_3 \end{matrix}$$

$$g_i: \mathbb{R}^3 \rightarrow \mathbb{R}$$

↙ ∇g_i

Back to mD-Newton:

mD

$$\vec{x}_{i+1} = \vec{x}_i - \frac{g(\vec{x}_i)}{J(g)(\vec{x}_i)}$$

vector
matrix

$$J(g)(\vec{x}_i)^{-1} \times \dots$$

Newton - Raphson — $g: \mathbb{R}^n \rightarrow \mathbb{R}^n$ / $\boxed{g(x) = \vec{0}}$

x_0 : close to the solution

while $(\|x_i - x_{i-1}\| > \epsilon)$

$$x_{i+1} = x_i - \underbrace{J(g)(x_i)^{-1}} \times g(x_i)$$

end

$$\sim \frac{g(x_i)}{g'(x_i)} \text{ } 10 \dots$$

Newton for optimization

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

\rightsquigarrow min f ?

\downarrow necessary condition

$$\nabla f(x) = \vec{0}$$

f - Newton

$J(\nabla f)$?

$$\nabla f(x) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix} \begin{matrix} \leftarrow g_1 \\ \vdots \\ \leftarrow g_n \end{matrix}$$

$$\rightsquigarrow g_i = \frac{\partial f}{\partial x_i}$$

$$g = \nabla f$$

$$J(\nabla f) = \begin{pmatrix} \frac{\partial g_i}{\partial x_j} \end{pmatrix} = \begin{pmatrix} \frac{\partial^2 f}{\partial x_i \partial x_j} \end{pmatrix}_{i,j \in \{1..n\}} = H(f)$$

derivate f'

Hessian matrix of f

Newton for optimization

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

\rightsquigarrow min f ?

\downarrow necessary condition

$$\nabla f(x) = \vec{0}$$

f - Newton

x_0 — close from the minimum

while $(\|x_i - x_{i-1}\| > \epsilon)$

$$x_{i+1} = x_i - \underbrace{H(f)(x_i)^{-1}} \times \nabla f(x_i)$$

end

invert a $n \times n$ matrix

Quasi-Newton methods ...
→ BFGS
→
heuristics

approximate $H(\mathbf{g})^{-1}$
from $\nabla \mathcal{L}$