$$f(x) = 4x_1^2 + 2x_1 x_2 + 5x_2^2 + x_1 + 3x_2$$

1) Write the taylor expansion at $\binom{0}{0}$ of order 2

$$f(x) = f\binom{0}{0} + \langle \nabla f\binom{0}{0}, x \rangle + \frac{1}{2}\langle H f\binom{0}{0}x, x \rangle + O(\|x\|^2)$$

$$\nabla f(x) = \begin{pmatrix} 8x_1 + 2x_2 + 1 \\ 2x_1 + 10x_2 + 3 \end{pmatrix} \qquad Hf(x) = \begin{pmatrix} 8 & 2 \\ 2 & 10 \end{pmatrix}$$

$$\det(Hf(x) - \lambda I) = (8-\lambda)(10-\lambda) - 4 = \lambda^2 - 18\lambda - 4$$

$$\Delta = 18^2 - 16 > 0$$

$$\lambda_i = \frac{18 \pm \sqrt{\Delta}}{2} > 0$$

$$\nabla f\binom{0}{0} = \binom{1}{3} \qquad Hf\binom{0}{0} = \begin{pmatrix} 8 & 2 \\ 2 & 10 \end{pmatrix}$$

$$f(x) = 0 + \langle \binom{1}{3}, x \rangle + \frac{1}{2}\langle Hf\binom{0}{0}x, x \rangle + \;?\; \boxed{0}$$

because $Hf(x) = cte$

$$= \frac{1}{2}\langle Ax, x \rangle + \langle b, x \rangle + c \qquad A \in S_{++}^n$$

$$f(x^*) = \min f(x) \qquad A = Hf\binom{0}{0} \quad b = \nabla f\binom{0}{0} \quad c = f\binom{0}{0}$$

$$(\Leftrightarrow) \qquad \nabla f(x^*) = 0 \; (\Leftrightarrow) \quad x^* = \;?$$

# Conjugate gradient method : motivation

has been developped for quadratic functions in
high dimension. $f(x) = \frac{1}{2} \langle Ax, x \rangle + \langle b, x \rangle + c$

$A \in S^n_{++}$  $b \in \mathbb{R}^n$

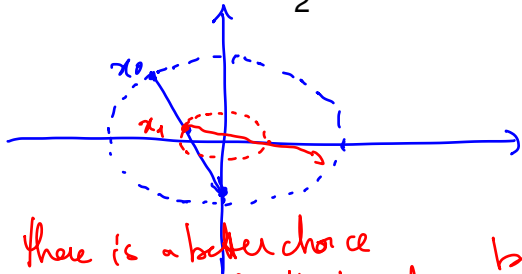$$f(x^*) = \inf_{x \in \mathbb{R}^n} f(x)$$

$\alpha_1 = 1 \quad \alpha_2 = 2$

$$f(x) = \frac{1}{2}(\alpha_1 x_1^2 + \alpha_2 x^2), \quad \text{with } 0 < \alpha_1 < \alpha_2$$

$$= \frac{1}{2}\langle Ax, x \rangle \quad \text{with } A = \begin{pmatrix} \alpha_1 & 0 \\ 0 & \alpha_2 \end{pmatrix}$$

$\nabla f(x) = \begin{pmatrix} \alpha_1 x_1 \\ \alpha_2 x_2 \end{pmatrix}$

$\nabla f\begin{pmatrix} -1 \\ 1 \end{pmatrix} = \begin{pmatrix} -1 \\ 2 \end{pmatrix}$



$x_0$

$x_1$

$x_1 = x_0 - \alpha_0 \nabla f(x_0)$

$x_k \xrightarrow[k \to \infty]{} x^* = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$

but $x_k \neq x^*$ $\forall k$

there is a better choice
than $-\nabla f(x)$ for the descent

# A-conjugate directions $f(x) = \frac{1}{2}\langle Ax, x \rangle + \langle b, x \rangle$

*Definition :* Let $A \in S_{++}^n$.

- ▶ 2 non zero vectors $v$, $w$ are called $A-$conjugate iff $\langle Av, w \rangle = 0$. $= \langle v, Aw \rangle$
- ▶ A family of non zero vectors $(v_i)_{i=1,\ldots m}$, is called $A-$conjugate iff $\langle Av_i, v_j \rangle = 0$ for all $i = 1, \ldots, m$, $j = 1, \ldots, m$, $i \neq j$.

*Property :* $A-$conjugate vectors are linearly independent. If $m = n$ a $A-$conjugate family is a basis of $\mathbb{R}^n$.

*Definition :* a conjugate descent method is a method where the successive descent directions form a $A-$conjugate family

Expression of the minimum of *f* in a *A*−conjugate basis the minimum of $f(x)$ satisfies $\nabla f(x^\star) = 0$
$Ax^\star + b = 0$

$$f(x) = \frac{1}{2}\langle Ax, x \rangle + \langle b, x \rangle$$

Suppose we have a basis $(d_i)_{i=1,\ldots n}$, such that $\langle Ad_i, d_j \rangle = 0$ for $j \neq i$

$$x^\star = \sum_{i=1}^{n} \alpha_i d_i, \text{ and } Ax^\star + b = 0,$$

therefore $Ax^\star = -b = \sum_{i=1}^{n} \alpha_i Ad_i$, then for any $j = 1, \ldots, n$

$$-\langle b, d_j \rangle = \sum_{i=1}^{n} \alpha_i \langle Ad_i, d_j \rangle = \alpha_j \langle Ad_j, d_j \rangle$$

$$\alpha_j = \frac{-\langle b, d_j \rangle}{\langle Ad_j, d_j \rangle}$$

181

# Construction of the $A-$conjugate basis *iteratively*

*choose initial point $x_0$   $g_0 = Ax_0 + b$*

Let $g_k = \nabla f(x_k) = Ax_k + b$ be the gradient at step $k$

Choose $d_0 = -g_0$ (The first step is a standard gradient descent step)

Then $d_k = -g_k + \beta_{k-1}d_{k-1}$ satisfying:    $\beta_{k-1} \in \mathbb{R}$

$\qquad d_1 = -g_1 + \beta_0 d_0$

(CG1)  $\langle Ad_k, d_j \rangle = 0$ for $j = 0, \ldots, k-1$ and    $\langle Ad_1, d_0 \rangle = 0$

(CG2)  $\langle g_k, d_j \rangle = 0$ for $j = 0, \ldots, k-1$    $\langle g_1, d_0 \rangle = 0$

Update at step $k$ : $x_{k+1} = x_k + \alpha_k d_k$

Next gradient $g_{k+1} = Ax_{k+1} + b = g_k + \alpha_k Ad_k$

*Property :* For all initial guess $x_0$ there exists $(\alpha_k)_k$ and $(\beta_k)_k$ such that (CG1) and (CG2) are satisfied.

*Property :* (CG1) and (CG2) $\Rightarrow \langle g_k, g_j \rangle = 0$ for $j \neq k$

$$Ax_1 + b = g_0 + \alpha_0 Ad_0$$
$$\langle g_1, d_0 \rangle = 0 = \langle g_0, d_0 \rangle + \alpha_0 \langle Ad_0, d_0 \rangle = 0 \quad \longrightarrow \quad \alpha_0 = \frac{-\langle g_0, d_0 \rangle}{\langle Ad_0, d_0 \rangle}$$

# Convergence of a conjugate method

*Property :* A conjugate descent method using directions satisfying conditions (CG1) and (CG2) converges in at most *n* steps.

*Property :* $\beta_k = -\dfrac{\langle Ad_{k-1}, g_k \rangle}{\langle Ad_{k-1}, d_{k-1} \rangle} = \dfrac{\|g_k\|^2}{\|g_{k-1}\|^2}$.

*Property :* $\alpha_k = -\dfrac{\langle g^k, d^k \rangle}{\langle Ad^k, d^k \rangle}$

$d_k = -g_k + \beta_{k-1} d_{k-1}$

$x_{k+1} = x_k + \alpha_k d_k$

# Conjugate gradient algorithm

**Data:** Matrix $A$, vector $b$, tolerance $\varepsilon$
**Result:** $x^\star$ such that $f(x^\star) = \min_x f(x)$
**Initialisation** : $k = 0$,
Initial guess for solution $x^0 \in \mathbb{R}^n$
$g^0 = Ax^0 + b$
$d^0 = -g^0$
**while** $\|g^k\| > \varepsilon$ **do**

- Compute directionnal minimum :
  $v^k = Ad^k$
  $\alpha_k = -\dfrac{\langle g^k, d^k \rangle}{\langle v^k, d^k \rangle}$
  $x^{k+1} = x^k + \alpha_k d^k$

- Update gradient :
  $g^{k+1} = g^k + \alpha_k v^k$

*same as the standard gradient descent method with optimal step*

- Compute new direction :
  $\beta_{k+1} = \dfrac{\langle g^{k+1}, g^{k+1} \rangle}{\langle g^k, g^k \rangle}$
  $d^{k+1} = -g^{k+1} + \beta_{k+1} d^k$

  $k \leftarrow k + 1$

**end**
$x^\star \leftarrow x^k$

# Monotonicity of the conjugate gradient algorithm

*Property :* If $d_k \neq 0$ and $\alpha_{k+1} \neq 0$ then $f(x_{k+1}) < f(x_k)$.
If $\alpha_{k+1} = 0$, $x_k$ is the minimizer of $f$ and $Ax_k + b = 0$

# Polak-Ribière method $\quad$ for non quadratic $f : \mathbb{R}^n \to \mathbb{R}$

**Data:** Function $f$, gradient $\nabla f$, tolerance $\varepsilon$
**Result:** $x^\star$ such that $f(x^\star) = \min_x f(x)$
**Initialisation** : $k = 0$,
Initial guess for $x^0 \in \mathbb{R}^n$
$g^0 = \nabla f(x^0)$
$d^0 = -g^0$
**while** $\|g^k\| > \varepsilon$ *and* $k < k_{\max}$ **do**
- Compute the step in direction $d_k$. *an admissible step in the direction*
  $f(x^k + \alpha_k d^k) \phantom{xxxxxx} < f(x^k)$ ~~for all~~
- Compute new position :
  $x^{k+1} = x^k + \alpha_k d^k$
- Compute new direction :
  $g^{k+1} = \nabla f(x^{k+1})$
  $c_{k+1} = \dfrac{\langle g^{k+1} - g^k, g^{k+1} \rangle}{\langle g^k, g^k \rangle}$
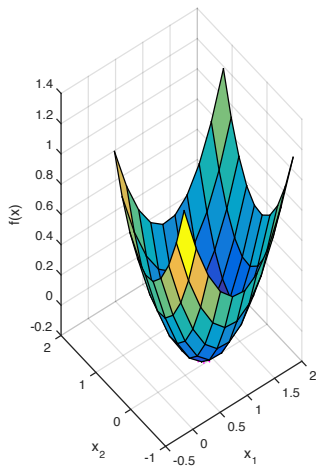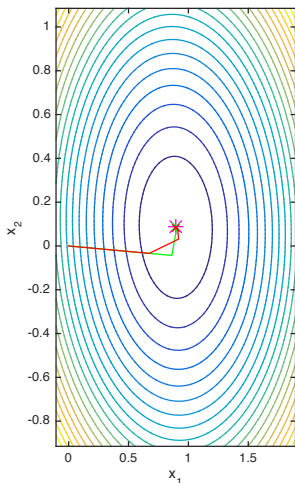  $d^{k+1} = -g^{k+1} + c_{k+1} d^k$

$\quad k \leftarrow k + 1$
**end**
$x^\star \leftarrow x^k$

# Comparaison of conjugate Gradient (green, 4 steps)and Polak-Ribière (red, 8 steps) methods.

$f$ quadratic function in $\mathbb{R}^5$. Projection on $(0, x_1, x_2)$.

# Linear regression

$m >> n$ if at least $n$ points are different $\text{rg}(X) = n$

$\theta \in \mathbb{R}^n$     $h_\theta : \mathbb{R}^n \longrightarrow \mathbb{R}$

Find $\theta$ defining a linear model

$$\hat{y} = h_\theta(x) = \theta^T.x$$

Let $m$ measurements $(x_i, y_i)$, $i = 1, \ldots, m$, where explaining variables are in $\mathbb{R}^n$ $(x_i = (x_i^j)_{j=1,\ldots,n}$. $\theta$ is found by minimizing the least squared error

$$Y = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix}$$

$$E(\theta) = \frac{1}{m} \sum_{i=1}^m \left( \theta^T.x_i - y_i \right)^2 = \frac{1}{m} \| X\theta - Y \|^2$$

The normal equation gives the best solution

$$X = \begin{pmatrix} x_1^1 \cdots x_1^n \\ \vdots \\ x_m^1 \cdots x_m^n \end{pmatrix}$$

$f(\theta) = \| X\theta - Y \|^2$

$\nabla f(\theta) = 0$

$\nabla f(\theta) = 2 X^T (X\theta - Y)$

$$\hat{\theta} = (X^T.X)^{-1}.X^T.y$$

complexity in $O(n^3)$ and $O(m)$.

198

# Outline

# Nonlinear least squares

$$f : \begin{cases} \mathbb{R}^P & \to & \mathbb{R}^Q \\ x = (x_1, \ldots, x_P)^t & \mapsto & (f_1(x), \ldots, f_Q(x))^t \end{cases}$$

for $Q > P$ we seek a solution to the problem $f(x) = 0$.

even if $Q = P$ $f(x) = 0$ is difficult

$\Downarrow$

Newton

# Examples

▶ Find a line that passes through $Q$ points with $Q > 2$

if the points
are not aligned the Pb has no solution
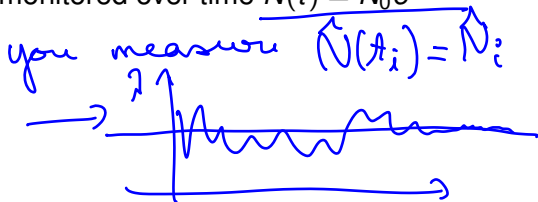
$$g(x) = \|f(x)\|^2 \qquad \text{minimise} \quad g(x)$$

# Examples

▶ Find the parameters $N_0$ and $\lambda$ of a radioactive material whose emissions are monitored over time $N(t) = N_0 e^{-\lambda t}$

each month you measure $\widehat{N}(t_i) = \widehat{N}_i$

$$\min_{N_0, \lambda} \sum_{i=1}^{m} \| N_0 - \bar{e}^{-\lambda t_i} - N_i \|^2$$

Toy example

$Q$ is the number of months were you measured $N_i$

$N_0 = x_1$    $\lambda = x_2$

$Q$ measurements $N_i \sim N(t_i) = N_0 e^{-\lambda t_i}$    $i = 1, \ldots, Q$

$f : \mathbb{R}^2 \to R^Q$ with $Q$ large
$(N_i)_{i=1,\ldots,Q}$ radioactivity measurements at times $(t_i)_{i=1,\ldots,Q}$

$$f(x) = \begin{pmatrix} f_1(x) \\ f_2(x) \\ \vdots \\ f_Q(x) \end{pmatrix} = \begin{pmatrix} x_1 e^{-x_2 t_1} - N_1 \\ x_1 e^{-x_2 t_2} - N_2 \\ \vdots \\ x_1 e^{-x_2 t_Q} - N_Q \end{pmatrix}.$$

$f(x) = 0$

Calculate the Jacobian matrix $Jf(x)$

## Toy example

$f : \mathbb{R}^2 \to R^Q$ with $Q$ large
$(N_i)_{i=1,\ldots,Q}$ radioactivity measurements at times $(t_i)_{i=1,\ldots,Q}$

$$f(x) = \begin{pmatrix} f_1(x) \\ f_2(x) \\ \vdots \\ f_Q(x) \end{pmatrix} = \begin{pmatrix} x_1 e^{-x_2 t_1} - N_1 \\ x_1 e^{-x_2 t_2} - N_2 \\ \vdots \\ x_1 e^{-x_2 t_Q} - N_Q \end{pmatrix}.$$

Calculate the Jacobian matrix $Jf(x)$

$$Jf(x) = \begin{pmatrix} e^{-x_2 t_1} & -x_1 t_1 e^{-x_2 t_1} \\ e^{-x_2 t_2} & -x_1 t_2 e^{-x_2 t_2} \\ \vdots \\ e^{-x_2 t_Q} & -x_1 t_Q e^{-x_2 t_Q} \end{pmatrix}.$$

## Reminders: linear least squares

$Ax = b$ for $b \in \mathbb{R}^Q$ and $A \in \mathcal{M}_{Q,P}(\mathbb{R})$ with $Q > P$ and
$\mathrm{rg}(A) = P$.
The problem: find $x \in \mathbb{R}^P$ such that

$$\|Ax - b\|^2 = \min_{y \in \mathbb{R}^P} \|Ay - b\|^2$$

admits a unique solution given by the normal equation

$$A^t A x = A^t b.$$

Nonlinear case

instead of looking for $f(x^*) = 0$
(which doest not exists)

us minimize $g(x) = \|f(x)\|^2 = \sum_{i=1}^{Q} f_i(x)^2$

$g(x^*)$ $\begin{cases} \text{Find } x^* \in \mathbb{R}^P \text{ such that} \\ \|f(x^*)\|^2 = \min_{x \in \mathbb{R}^P} \|f(x)\|^2 \quad \left( \|f(x)\|^2 = \sum_{k=1}^{Q} (f_k(x))^2 \right), \end{cases}$

We suppose that :

$$\forall x \in \mathbb{R}^P, \qquad J_f(x) \in \mathcal{M}_{Q,P}(\mathbb{R}) \text{ has rank } P.$$

In particular, we will have $(J_f(x))^t J_f(x)$ symmetric defined positive.

Nonlinear case (continued)

$f : \mathbb{R}^2 \to \mathbb{R}^Q \qquad 3 \in \mathbb{R}^2$

with $N(y) = \|y\|^2$

$g(x) = \|f(x)\|^2 = N \circ f$

$Dg(x)h = DN(f(x)) Df(x)h$

$DN(y)h = \langle 2y, h \rangle$

$Df(x)3 = Jf(x)3$

$= \langle 2f(x), Jf(x)h \rangle$

We notice

$$g : \left\{ \begin{array}{ccc} \mathbb{R}^P & \to & \mathbb{R} \\ x & \mapsto & \|f(x)\|^2 \end{array} \right.$$

If $g$ is strictly convex and coercive then the problem $g(x^*) = \min_x g(x)$ admits a unique solution $x^*$

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

$Dg(x)h = \langle 2 Jf(x)^T f(x), h \rangle$

$\nabla g(x) = 2 Jf(x)^T f(x)$

$Jf(x) \in \mathcal{M}_{Q \times 2}(\mathbb{R})$

# Calculating the gradient of *g*

$g = N \text{o} f$ composition of
$N : \mathbb{R}^Q \to \mathbb{R}$, $N(y) = ||y||^2$ and $f : \mathbb{R}^P \to \mathbb{R}^Q$.
The rule for differentiating a composite function gives
$Dg(x) = DN(f(x))Df(x)$
For $y, \delta \in \mathbb{R}^Q$, $DN(y)\delta = \langle 2y, \delta \rangle$
for $x, h \in \mathbb{R}^P$, $Df(x)h = Jf(x)h \in \mathbb{R}^Q$

$$h, x \in \mathbb{R}^P, \quad Dg(x)h = \langle 2f(x), Jf(x)h \rangle = \langle 2Jf(x)^T f(x), h \rangle$$

$$\nabla g(x) = 2Jf(x)^T f(x).$$

# Find the zeros of $\nabla g$ or the zeros of $f(x)$

*P = 2 for our example*

- $\nabla g(x) = 2Jf(x)^T f(x)$ Newton method requires $Hf(x)$
- If $f(x)$ is a function of $\mathbb{R}^P$ in $\mathbb{R}^P$ we find the zeros by Newton's algorithm

$$x_{k+1} = x_k + d_k$$

$$\text{with } Jf(x_k)d_k = -f(x_k).$$

- Here $f(x)$ is a function of $\mathbb{R}^P$ in $\mathbb{R}^Q$ so the system $Jf(x_k)d_k = -f(x_k)$ of size $Q \times P$ is solved in the least squares sense

*given the best solution for $f(x_k) = 0$*

$$Jf(x_k)^T Jf(x_k)d_k = -Jf(x_k)^T f(x_k)$$

$$\Leftrightarrow \quad d_k = -(Jf(x_k)^T Jf(x_k))^{-1} Jf(x_k)^T f(x_k).$$

# Gauss Newton method *for N.L. systems of equations*

- Initialize $x_0 \in \mathbb{R}^P$    $\| d_k \| > \varepsilon$
- While ~~$\| \text{...} \|$~~ and $k < k_{\max}$
  - Solve $(Jf(x_k)^T Jf(x_k))d_k = -Jf(x_k)^T f(x_k)$    $\Leftarrow$ LSQ system solution
  - Update $x_{k+1} = x_k + d_k$
  - Update $k \to k + 1$

$$Jf(x)\, d = -f(x)$$

# Convergence of the Gauss Newton method

We recall that $Jf(x)$ of rank $P$ and $g(x)$ is strictly convex coercive

- Let $x_k \in \mathbb{R}^P$, then the direction
  $d_k = -(Jf(x_k)^T Jf(x_k))^{-1} Jf(x_k)^T f(x_k)$ satisfies

  $$\langle \nabla g(x_k), d_k \rangle \leq 0.$$

  *descent direction*

  If $x_k \neq x^*$ then

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

  So $d_k$ is a descent direction for $g$ at $x_k$.

- If the sequence $(x_k)_k$ converges, then its limit is $x^*$.

## Exercice

$$f(x) = 3x_1^2 + x_2^2 - 2x_1 x_2 + x_1 + x_2 + 1$$

write in the quadratic form    $A, b, c$

$f : \mathbb{R}^2 \to \mathbb{R}$

$f(x_1, x_2) = 4(x_1^2 + x_2^2) - (x_1^2 + x_2^2)^2$     $X = (x_1, x_2)$

1) Compute $\nabla f(X)$ and $Hf(X)$

2) Compute the set of points $\{\nabla f(X) = 0\} = S$

3) Say if $X \in S$ is minimum or maximum or what?

$$\nabla f(X) = \begin{pmatrix} 8x_1 - 4x_1(x_1^2 + x_2^2) \\[2mm] 8x_2 - 4x_2(x_1^2 + x_2^2) \end{pmatrix} \qquad Hf(x) = \begin{pmatrix} 8 - 12x_1^2 - 4x_2^2 & -8x_1x_2 \\[2mm] -8x_1x_2 & 8 - 12x_2^2 - 4x_1^2 \end{pmatrix}$$

$$\begin{pmatrix} 2x_1 = x_1(x_1^2 + x_2^2) \\[2mm] 2x_2 = x_2(x_1^2 + x_2^2) \end{pmatrix} \Longrightarrow S = \{(0,0)\} \cup \left\{ \begin{array}{l} x_1^2 + x_2^2 = 2 \\ \text{circle of center } (0,0) \\ \text{radius } \sqrt{2} \end{array} \right\}$$

question 3 for tomorow