1/ About the EXAM :

 written exam week from 14-18 December at the university. 13:30 $\longrightarrow$ 15:30 2 Hours

 without documents.

For the 3/4 of you who cannot be present, we will organize an oral exam.

- Gradient direction: $\nabla f(x)$

- Newton direction: $-\left[\nabla^2 f(x)\right]^{-1} \nabla f(x)$

- $f(x) = \frac{1}{2} x^T A x \qquad x \in \mathbb{R}^2, \quad A = \begin{pmatrix} 9 & 0 \\ 0 & 1 \end{pmatrix}$

Plot $\nabla f(x)$, $\left[\nabla^2 f(x)\right]^{-1} \nabla f(x)$ and level set of $f$.
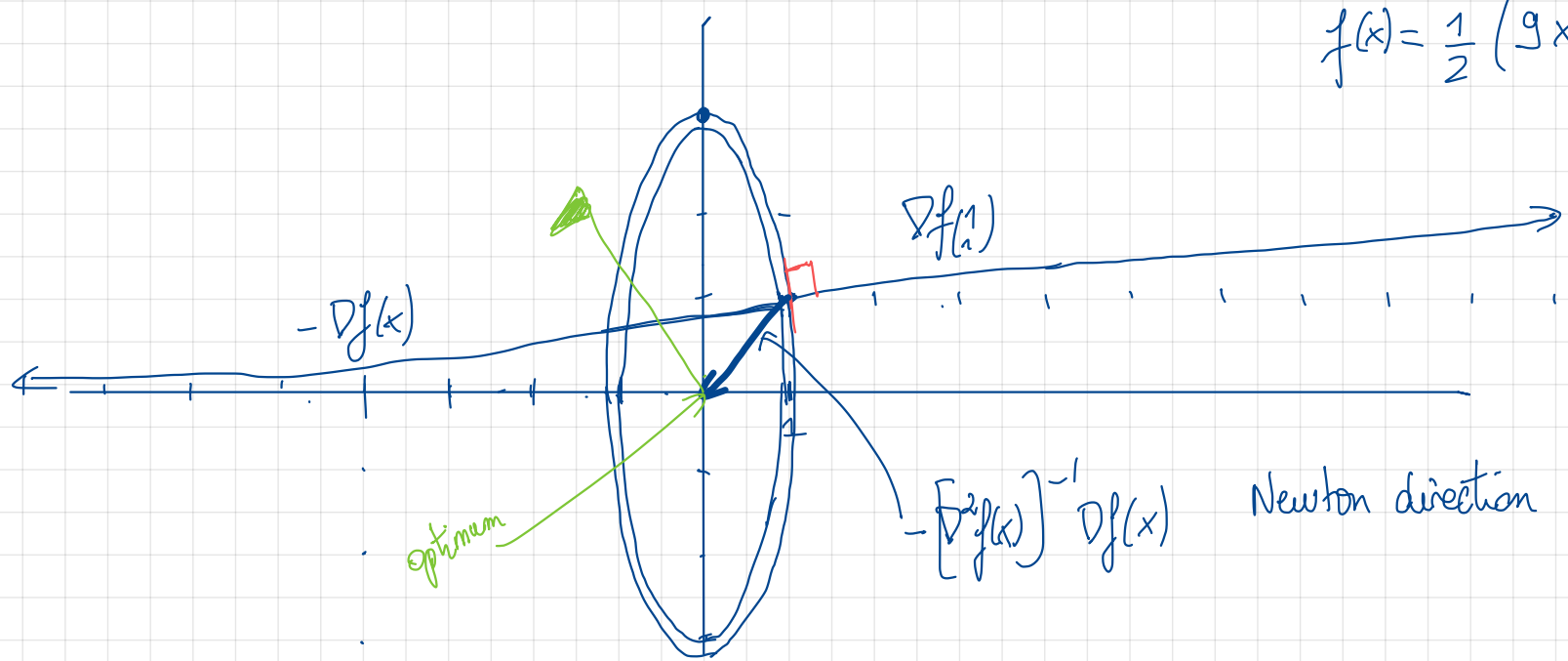
$$\nabla f(x) = \begin{pmatrix} 9x_1 \\ x_2 \end{pmatrix}$$

$$\nabla^2 f(x) = A = \begin{pmatrix} 9 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\left[\nabla^2 f(x)\right]^{-1} = \begin{pmatrix} \frac{1}{9} & 0 \\ 0 & 1 \end{pmatrix} \; ; \; \left[\nabla^2 f(x)\right]^{-1} \nabla f(x) = \begin{pmatrix} \frac{1}{9} & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 9x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$
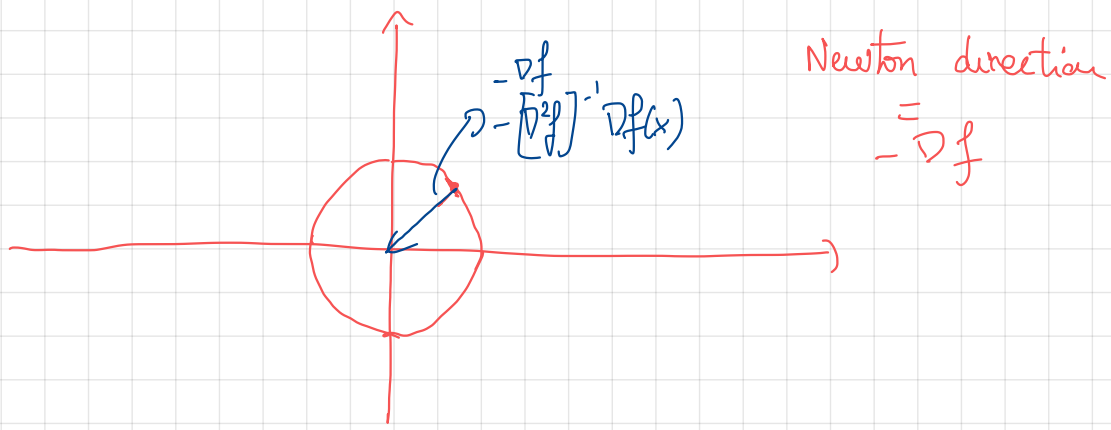
$$= x$$

Newton direction: $-x$

$$f(x) = \frac{1}{2}\left(9x_1^2 + x_2^2\right)$$

$$Df\begin{pmatrix}1\\1\end{pmatrix}$$

$$-Df(x)$$

$$-\left[\nabla^2 f(x)\right]^{-1} Df(x) \quad \text{Newton direction}$$

optimum

At $x = \begin{pmatrix} -4 \\ -5 \end{pmatrix}$  $Df(x) = \begin{pmatrix} -9.4 \\ -5 \end{pmatrix}$

What if $f(x) = \frac{1}{2}\left(x_1^2 + x_2^2\right)$  $Df^2 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

What about the Newton and $-Df$ in this case?

$$-Df$$
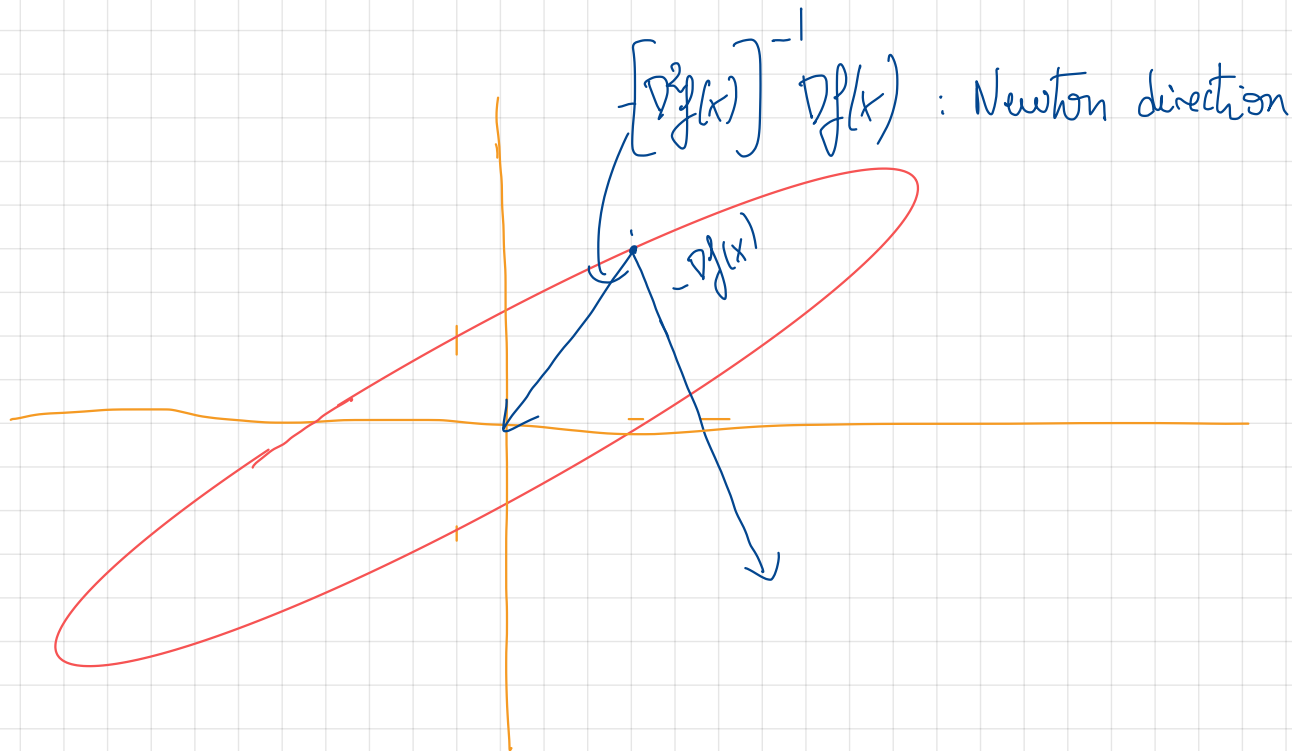$$-\left[\nabla^2 f\right]^{-1} Df(x)$$

Newton direction
$$= -Df$$

We observe that the Newton direction points towards the optimum independently of the condition number of the Hessian matrix.

whereas $-Df(x)$ points towards the optimum, if and only if at $x = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ and only if $\nabla^2 f(x) = Id$ and the condition number equal to $1$.

If the Hessian matrix is not diagonal anymore:

$$f(x) = \frac{1}{2} x^T A x$$

A positive, definite
A not diagonal

$[\nabla^2 f(x)]^{-1} Df(x)$ : Newton direction



$-Df(x)$

$-Df(x)(h) = -Df(x) \cdot h$

# Optimality conditions:

Assume $f: \mathbb{R} \longrightarrow \mathbb{R}$ is differentiable ( $f'(x)$ exists for all $x$)
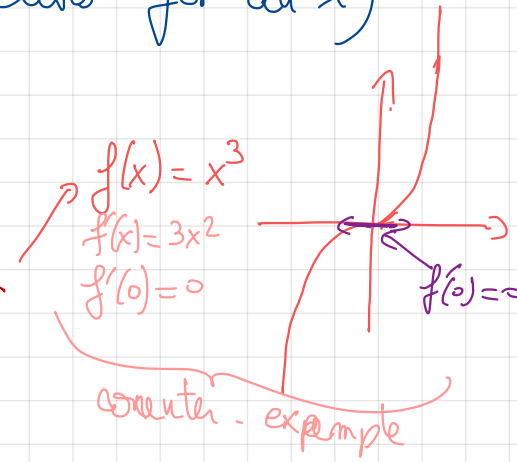
Which one of the following statements are correct:

① $f'(x^*) = 0 \implies x^*$ is a local optimum **WRONG**

② $x^*$ is a local optimum $\implies f'(x^*) = 0$ **CORRECT**

③ $f'(x^*) = 0 \implies x^*$ is a global optimum **WRONG**

④ $x^*$ is a global optimum $\implies f'(x^*) = 0$ **CORRECT**

$f(x) = x^3$
$f(x) = 3x^2$
$f'(0) = 0$

$f'(0) = 0$

counter-example

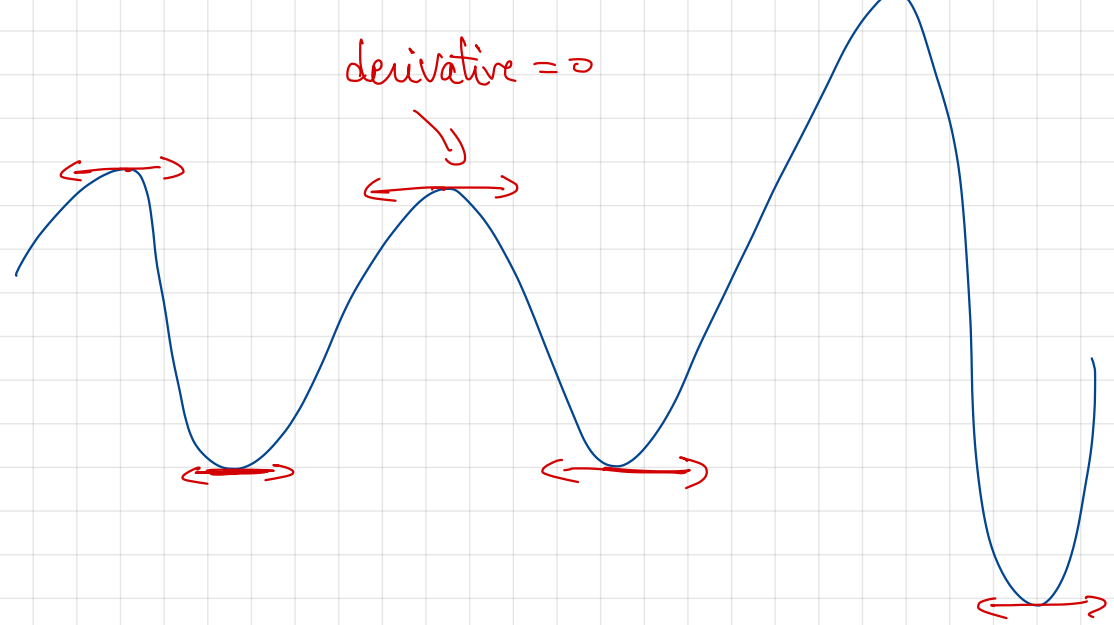② gives a first order necessary condition.

## THEOREM: (first order necessary condition)

Let $f: \mathbb{R}^n \longrightarrow \mathbb{R}$ be a differentiable function. If $x^*$ is a local optimum of $f$ then $Df(x^*) = 0$.

minimum
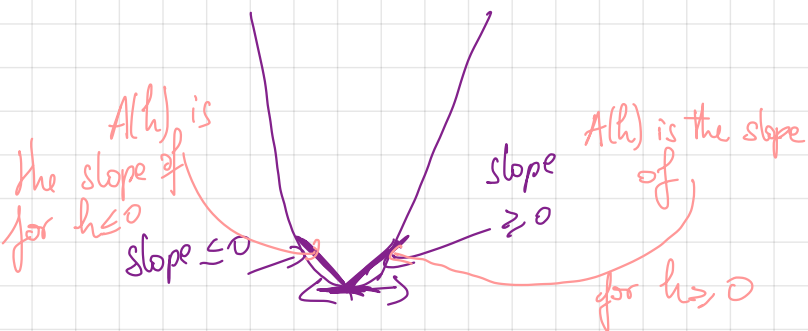or maximum

# Interpretation when $n = 1$:

derivative $= 0$

# Proof for $n = 1$:

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

assume that $x^{\alpha}$ is a local minimum: $f(x^{\alpha}) \leq f(x^{\alpha} + h) \quad \forall h$ small enough

$$A(h) = \frac{\overbrace{f(x^{\alpha} + h) - f(x^{\alpha})}^{\geq 0}}{h}$$

$\rightarrow$ if $h \geq 0 \quad A(h) \geq 0$

if $h \leq 0 \quad A(h) \leq 0$

$A(h)$ is the slope of for $h < 0$

slope $\leq 0$

slope $\geq 0$

$A(h)$ is the slope of for $h \geq 0$

$\left. \begin{array}{l} \lim\limits_{\substack{h \to 0 \\ h \geq 0}} \underbrace{A(h)}_{\geq 0} = f'(x^{\alpha}) \geq 0 \\[2em] \text{if} \quad \lim\limits_{\substack{h \to 0 \\ h \leq 0}} \underbrace{A(h)}_{\leq 0} = f'(x) \leq 0 \end{array} \right\} f'(x) = 0$

# SECOND ORDER NECESSARY AND SUFFICIENT CONDITIONS:
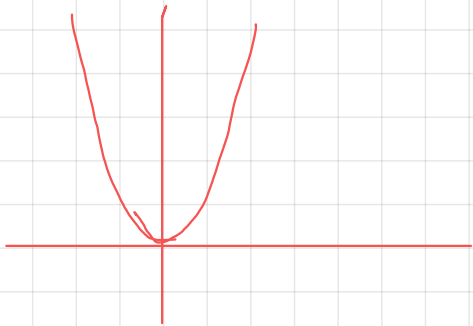
Let assume that $f$ is twice continuously differentiable

NECESSARY CONDITION: If $x^*$ is a local minimum, then $Df(x^*) = 0$
and $D^2f(x)$ is positive semi-definite.

$$\left( \text{if } n = 1, \quad x^* \text{ is a local minimum} \Rightarrow f'(x^*) = 0, \quad f''(x) \geq 0 \right)$$

SUFFICIENT CONDITION: If $x^*$ which satisfies $Df(x^*) = 0$ and $D^2f(x)$ is

positive definite, then $x^*$ is a strict local minimum.

$$\left( \text{if } n = 1, \quad x^* \text{ such that } f'(x^*) = 0 \quad f''(x) > 0 \Rightarrow x^* \text{ is a strict local minimum} \right)$$
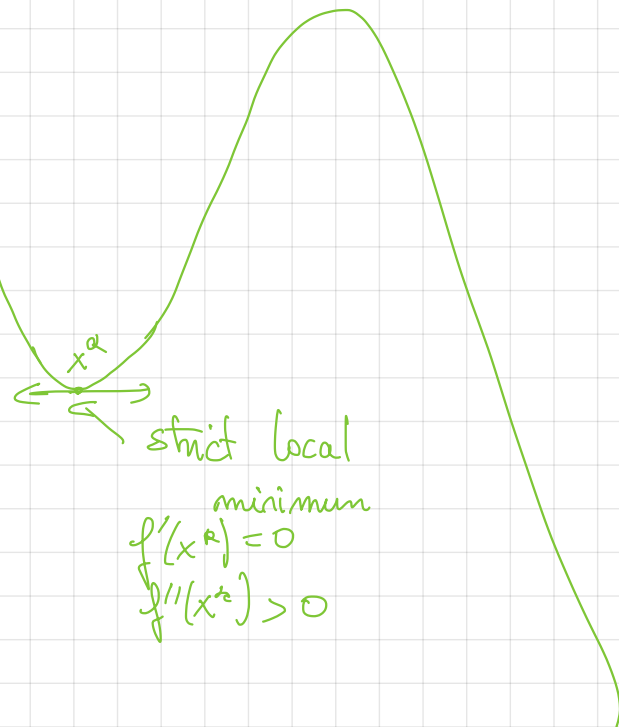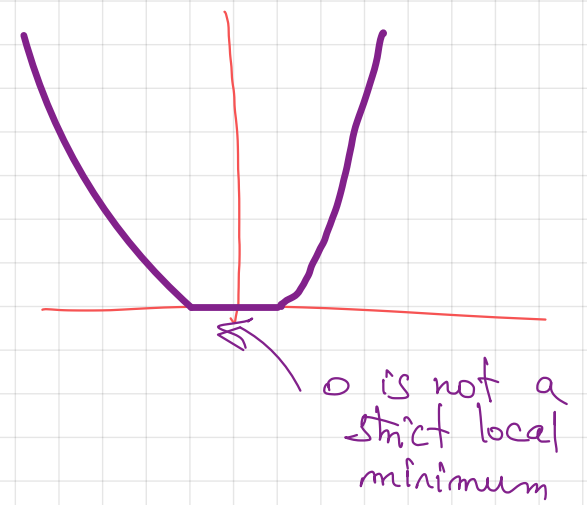
Example: $f(x) = x^2, \quad f'(x) = 2x \quad f''(x) = 2$



0 satisfies that $f'(0) = 2 \times 0 = 0$ and $f''(0) = 2 > 0$

$\Rightarrow$ 0 is a strict local minimum.

strict local minimum:

strict local
minimum

o is not a
strict local
minimum

$x^a$

strict local
minimum
$f'(x^R) = 0$
$f''(x^*) > 0$
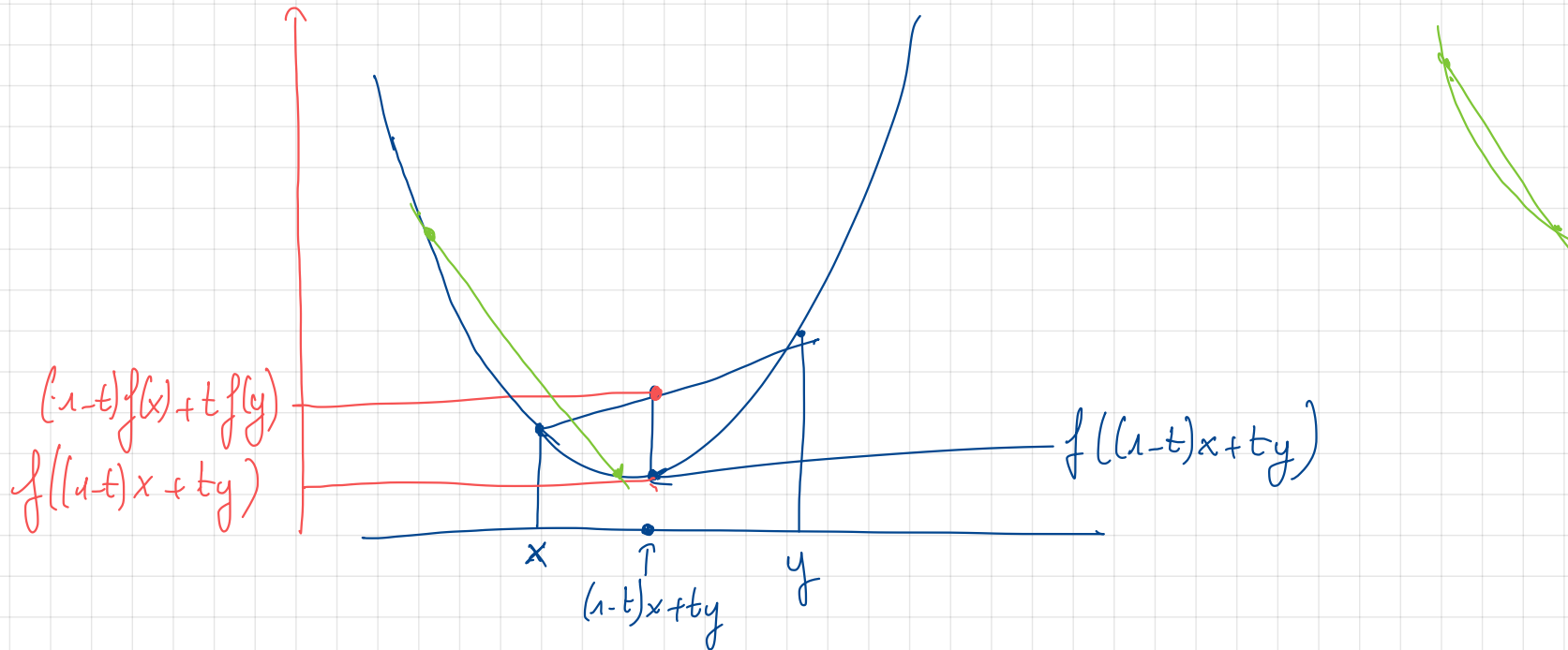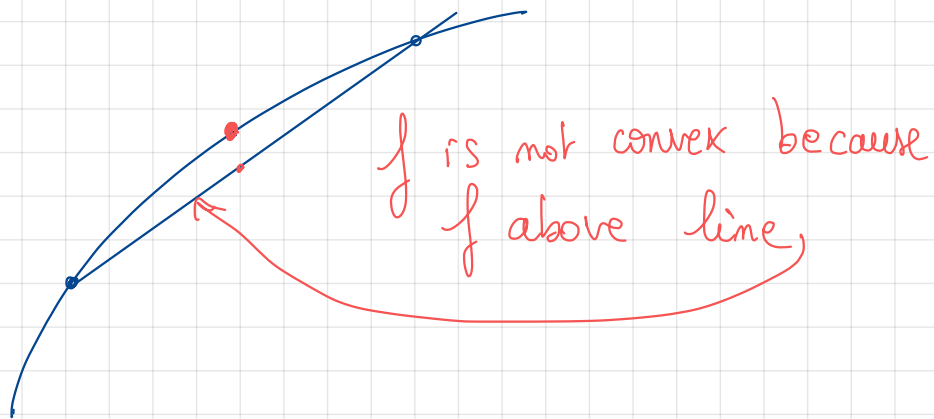
# CONVEX FUNCTIONS

Let $f: U \subset \mathbb{R}^n \longrightarrow \mathbb{R}$. We say that $f$ is convex, if for all $x, y \in U$

$U$ ↑ open convex set

$$\forall \ t \in [0, 1]$$

$$f\big((1-t)x + ty\big) \leq (1-t)\,f(x) + t\,f(y)$$

$(1-t)f(x) + t f(y)$

$f\big((1-t)x + ty\big)$

$f\big((1-t)x + ty\big)$

$x$

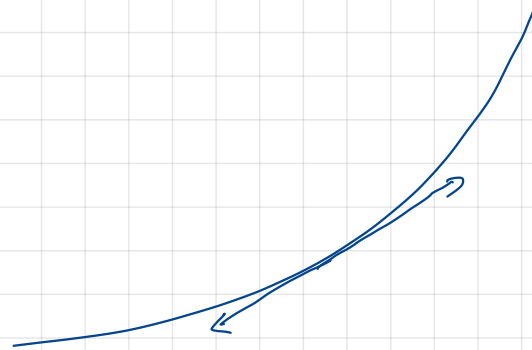$(1-t)x + ty$

$y$

$f$ is not convex because
$f$ above line.

THEOREM: If $f$ is differentiable, then $f$ is convex if and only if

for all $x, y$    $f(y) - f(x) \geq Df(x)^\top (y - x) = Df(x) \cdot (y - x)$

If $n = 1$   $f(y) - f(x) \geq f'(x)(y - x)$

$f$ is convex if and only if the function is above the tangent.

THEOREM: If $f$ is twice continuously differentiable, then $f$ is convex if and only if $D^2 f(x)$ is positive semi-definite for all $x$.

If $n = 1$ $f$ is twice derivable, then $f$ is convex if and only if $f''(x) \geq 0$

Examples: $f(x) = x^2$ is convex (because $f''(x) = 2 \geq 0$)

$f(x) = -x^2$ $(f''(x) = -2 \Rightarrow f$ is not convex $)$

$f(x) = \log(x)$ $(f'(x) = \frac{1}{x}, f''(x) = -\frac{1}{x^2} \leq 0 \Rightarrow f$ is not convex $)$
$\underset{x > 0}{\Gamma}$

$f(x) = x$ $f$ is convex $f''(x) = 0$

## Examples of convex functions:

- $f(x) = \frac{1}{2} x^T A x$ $A$ sym. pos. definite.

- $f(x) = a^T x + b$ $a \in \mathbb{R}^n, b \in \mathbb{R}^n$

- the negative of the entropy: $f(x) = -\sum_{i=1}^{n} x_i \log(x_i)$

**EXERCICE:** Let $f: U \subset \mathbb{R}^n \longrightarrow \mathbb{R}$ be a convex and differentiable function.

Prove that if $Df(x^*) = 0$, then $x^*$ is a global minimum.

If $f$ is convex and differentiable we have: $\forall x, y$
$$f(y) - f(x) \geq Df(x)^\top (y - x)$$

If $x^*$ is such that $Df(x^*) = 0$, then $f(y) - f(x^*) \geq \underbrace{Df(x^*)^\top (y - x^*)}_{= 0}$
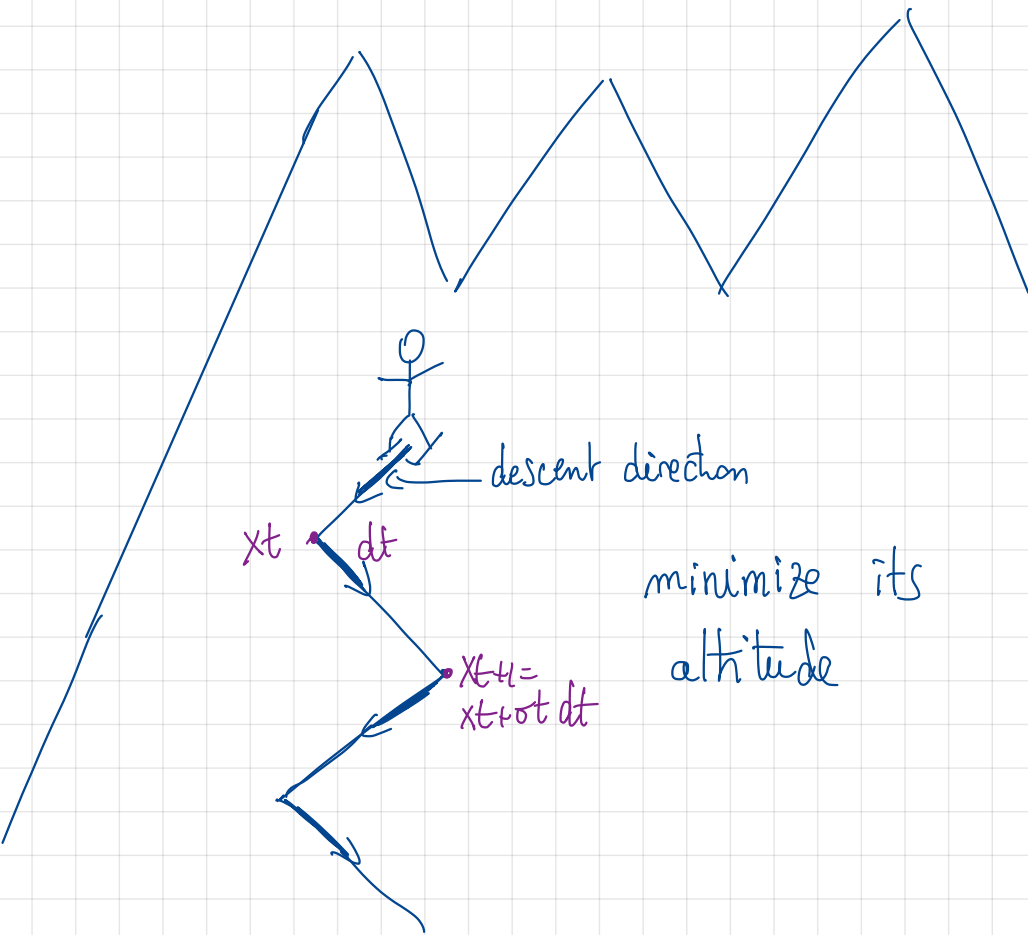
$$f(y) - f(x^*) \geq 0 \quad \forall y$$

then $\forall y \quad f(y) \geq f(x^*)$

which means that $x^*$ is the global minimum of $f$.

The important consequence is that for convex $\overset{\text{differentiable}}{\text{functions}}$ critical points, points where $Df(x) = 0$ are global minima of the functions.

# DESCENT METHODS

OBJECTIVE :
Minimize $f : \mathbb{R}^n \longrightarrow \mathbb{R}$
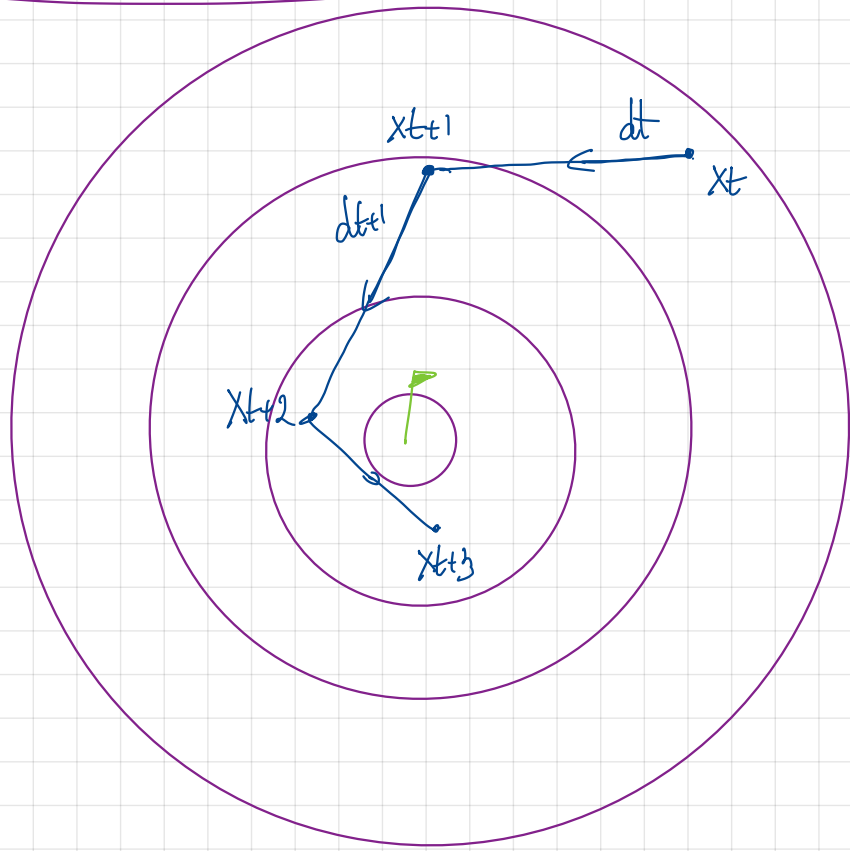
## General principle

1/ choose an initial point $x_0$, $t = 0$

WHILE NOT HAPPY [WHILE $f$ not minimized enough].

- choose a descent direction $dt \neq 0$   $dt \in \mathbb{R}^n$
- line search
  - $\rightarrow$ choose a step-size $\sigma t > 0$
  - $\rightarrow$ set $x_{t+1} = x_t + \sigma t \, dt$
- set $t = t+1$

## Remaining questions :

- how to choose $dt$ ?

- how to choose $\sigma t$ ?

descent direction

$x_t$  $dt$

$x_{t+1} = x_t + \sigma t \, dt$

minimize its altitude

# Picture with level sets

We can choose for $dt = -Df(xt)$

this is a descent direction:

if $f$ is differentiable and if $\sigma$
is small enough then

$$f(xt - \sigma Df(xt)) \approx f(xt) - \sigma Df(xt)^T Df(xt)$$
$$= f(xt) - \sigma \| Df(xt) \|^2$$
$$< f(xt)$$

$\sigma$ small enough

$\hookrightarrow -Df(xt)$ is a descent direction

from Taylor formula:

$$f(x+h) = f(x) + Df(x)^T h + o(\|h\|)$$

$h$ small $\quad f(x+h) \simeq f(x) + Df(x)^T h$

$\hookrightarrow f(xt - \underbrace{\sigma Df(xt)}_{h})) \simeq f(xt) + Df(xt)^T(-\sigma Df(xt)) = f(xt) - \sigma Df(xt)^T Df(xt) = f(xt) - \sigma \| Df(xt) \|^2$

# Choice of the step-size ?

optimal step-size: $\sigma t = \underset{\sigma \geq 0}{\arg\min} \; f(xt - \sigma Df(xt))$

$$\sigma \overset{g}{\mapsto} f(xt - \sigma Df(xt))$$
$$\sigma t = \underset{\sigma}{\arg\min} \; g(\sigma)$$

Typically too expensive to do those 1D optimization perfectly

There exists different techniques. One widely used one is Armijo rule.

# When do we stop the overall algorithm

$\rightarrow$ We can track $f(xt+1) - f(xt)$   (stop when it's small)

$\rightarrow$ We can stop $\wedge \| Df(xt) \|$ is small.
      when