# L1. Minimization and saddle point problem

- Non-convex minimization

$$\min_{z \in \mathbb{R}^d} f(z)$$

$O(\varepsilon^{-2})$ steps to find $\varepsilon$-stationary pt

assump. on $f$ : "nothing but $C^1$"

Reference: Section 1.2.3
[N2018]Nesterov. *Introductory lectures on convex optimization*. Vol. 137. Berlin: Springer International Publishing, 2018.

(A1) $f$ is $C^1$ on $\mathbb{R}^d$, $l$-smooth

   ie. $\nabla f$ is $l$-Lipschitz

(A2) $\min_z f(z) = f^*$ finite, ie. $f$ is lower bounded
   below on $\mathbb{R}^d$.

### GD (Gradient Descent)

$$z_{k+1} = x_k - s \nabla f(x_k)$$

No other information, $f$ may have local minima, and generally GD only finds a stationary pt.

Rk. Without (A2), GD may diverge

   eg. $f(x,y) = xy$, along $\begin{bmatrix} 1 \\ -1 \end{bmatrix}$, $f \to -\infty$.

define   $\Delta_f := f(x_0) - f^* \geq 0$.

Thm  If $0 < s < \frac{2}{l}$, then for $N = O(\frac{1}{\varepsilon^2})$, $\exists k \leq N$

   s.t.   $\| \nabla f(x_k) \| \leq \varepsilon$.

pf.  idea: make $f(x_k) \downarrow$, then bcz $f$ has a lower bound, it can not always make proper.

$$f(x_{k+1}) \leq f(x_k) + \langle \nabla f(x_k), x_{k+1} - x_k \rangle + \frac{\ell}{2} \| x_{k+1} - x_k \|$$

by Taylor expansion and $\ell$-smoothness of $f$

$$x_{k+1} - x_k = -s \nabla f(x_k)$$

$$\Rightarrow f(x_{k+1}) - f(x_k) \leq -s \| \nabla f(x_k) \|^2 + \frac{\ell}{2} s^2 \| \nabla f(x_k) \|_2^2$$

$$= -s(1 - \frac{\ell s}{2}) \| \nabla f(x_k) \|^2$$

$$\underbrace{\qquad\qquad} > 0 \text{ if } \frac{\ell s}{2} < 1$$

suppose $s = \frac{2}{\ell} \alpha$, $0 < \alpha < 1$, then

$$f(x_{k+1}) - f(x_k) \leq -\frac{2\alpha(1-\alpha)}{\ell} \| \nabla f(x_k) \|^2$$

telescopic sum $k = 0, \cdots, N$

$$\frac{2\alpha(1-\alpha)}{\ell} \sum_{k=0}^{N} \| \nabla f(x_k) \|^2 \leq f(x_0) - f(x_{N+1}) \leq \Delta_f$$
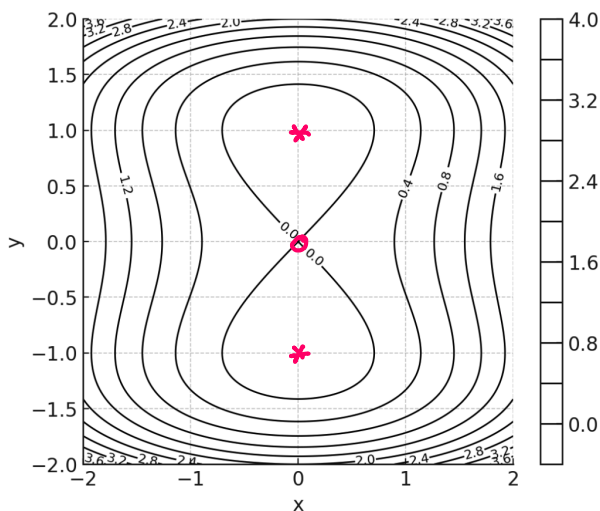
$\Rightarrow \exists \; k \leq N$ s.t.

$$\| \nabla f(x_k) \|^2 \leq \frac{\ell \Delta_f}{2\alpha(1-\alpha)} \cdot \frac{1}{N+1} \sim \frac{1}{N} \qquad \#$$

Rk. The $O(\varepsilon^2)$ iteration complexity is tight, to find $\varepsilon$-stationary pt is $\Omega(\varepsilon^{-2})$ grad evaluations.

[CDHS2017] Carmon, Duchi, Hinder, Sidford. Lower bounds for finding stationary points. arXiv:1710.11606, 2017.

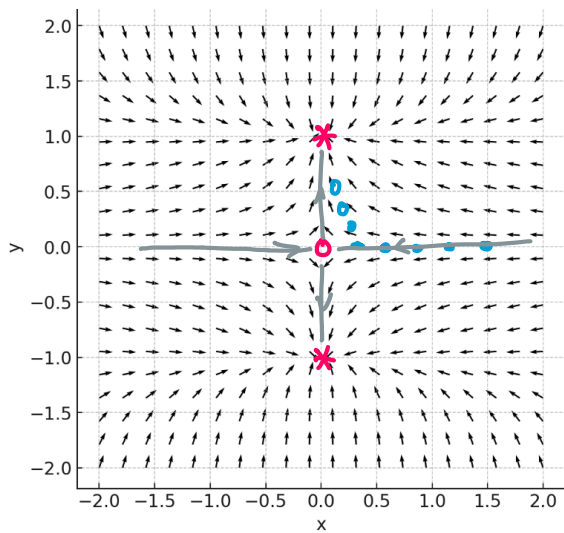Eg. $f(x,y) = \frac{1}{2} x^2 + \frac{1}{4} y^4 - \frac{1}{2} y^2$



Gradient Descent Directions for f(x,y)

3 stationary pts

$(0, \pm 1)$ local min

$(0, 0)$ saddle pt

Starting GD from

$z_0 = (1, 0)$

in theory $x_k \to (0, 0)$

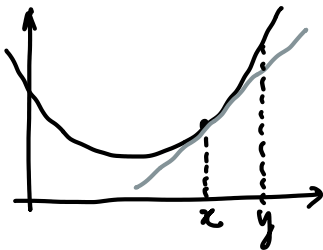RK. finding saddle pt is not easy in practical computation

- Convex minimization

$$\min_x f(x) \qquad f \text{ is } C^1 \text{ on } \mathbb{R}^d$$

$f(x)$ is convex on $\mathbb{R}^d$

fact If $f$ is $C^1$ and convex, then

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle, \quad \forall x, y \in \mathbb{R}^d.$$



"above the tangent line"

lemma If $f$ is $C^1$ and convex, $\nabla f(x^*) = 0$, then

$$f(x^*) = \min_x f(x),$$

i.e. $x^*$ is a global minimum. (may not unique)

pf. $\forall x$, $f(x) \geq f(x^*) + \langle \nabla f(x^*), x - x^* \rangle$ ♯.

Prop $f$ is $C^1$, then $f$ is convex iff.

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq 0, \quad \forall x, y.$$

Rk. The field $F(x) = \nabla f(x)$ is called **monotone**.

pf. "$\Rightarrow$" $\left. \begin{array}{l} f(x) - f(y) \geq \langle \nabla f(y), x-y \rangle \\ f(y) - f(x) \geq \langle \nabla f(x), y-x \rangle \end{array} \right\} \Rightarrow$ adding the two.

"$\Leftarrow$" ☺

$$F(t) := f(x + t(y-x)), \quad F(1) = f(y), \quad F(0) = f(x)$$

$$f(y) - f(x) = \int_0^1 F'(t)\, dt$$

$$= \int_0^1 \langle \nabla f(\underbrace{x + t(y-x)}_{\substack{\parallel \\ \nabla f(x) + \nabla f(x(t)) - \nabla f(x(0))}}^{\overbrace{\quad}^{x(t)}}), y-x \rangle\, dt$$

$$= \langle \nabla f(x), y-x \rangle$$

$$+ \int_0^1 \langle \underbrace{\nabla f(x(t)) - \nabla f(x(0)), \frac{x(t) - x(0)}{t}}_{\geq 0} \rangle\, dt$$

\#

Prop. $f$ is $C^2$ on $\mathbb{R}^d$, $f$ is convex iff $\nabla^2 f \succeq 0$.

  pf (Ex)                                    i.e. $\nabla^2 f(x) \succeq 0, \forall x$.

• strongly convexity

Def $f$ is $C^1$ on $\mathbb{R}^d$, $f$ is $\mu$-strongly convex if ( $\mu > 0$, $\mu$-S.C. )

$$f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle + \frac{\mu}{2} \|y-x\|^2, \quad \forall x, y.$$

Prop. $f$ is $C^1$ and $\mu$-S.C., then

$$\langle \nabla f(x) - \nabla f(y), x-y \rangle \geq \mu \|x-y\|^2, \quad \forall x, y.$$

Rk. $F(x) = \nabla f(x)$ is $\mu$-strongly monotone.

Prop $f$ is $C^2$ on $\mathbb{R}^d$, $f$ is $\mu$-S.C. iff $\nabla^2 f \succeq \mu I_d$.

Lemma $f$ is $C^1$ and $\mu$-S.C. on $\mathbb{R}^d$, then $\min_x f(x)$ is

attained at a single (global) minimizer.

pf. First use "quadratic growth" to show that
Q.G.
level set is compact, then $\min_x f(x)$ is attained within a compact set.

QG: $f(y) \geq f(x_0) + \langle \nabla f(x_0), y-x_0 \rangle + \frac{\mu}{2} \|y-x_0\|^2$

$$\sim \|y-x_0\|^2 \quad \text{as} \quad \|y\| \to \infty.$$

Then use $\mu$-S.C. to prove uniqueness of $x^*$ #

Thm If $f$ is $C^1$, $\ell$-smooth, $\mu$-S.C., then GD converges
(Rk. $\ell \geq \mu$)
exponentially fast (or "linear convergence"), ie

1) $f(x_k) \leq f^* + \varepsilon$ (function value convergence)

2) $\|\nabla f(x_k)\| \leq \varepsilon$ (first-order optimality)

3) $\|x_k - x^*\| \leq \varepsilon$ (variable convergence)

"identification of $x^*$ parameters"

within $k = O(\log \frac{1}{\varepsilon})$ steps, assuming $s < \frac{2\mu}{\ell^2}$

$\frac{1}{k} \frac{2}{\ell}$

Rk $\kappa := \frac{\ell}{\mu}$ is called <u>condition number</u>, $\kappa \geq 1$.

pf. $x^* \exists!$ by S.C.

3) $\Rightarrow$ 2): $\|\nabla f(x_k) - \nabla f(x^*)\| \leq \ell \|x_k - x^*\|$
$\overset{O}{\nearrow}$
by $\ell$-smoothness.

3) $\Rightarrow$ 1): $f(x^*) \geq f(x_k) + \underline{\langle \nabla f(x_k), x^* - x_k \rangle}$
$\underset{f^*}{\uparrow}$
$\leq \underbrace{\|\nabla f(x_k)\|}_{\text{by 2)} \to \leq \varepsilon} \underbrace{\|x_k - x^*\|}_{\leq \varepsilon} \leq \varepsilon^2$

function value actually achieves $O(\varepsilon^2)$ approximation.

We thus only prove 3).

Recall that $\quad x_{k+1} = x_k - s\nabla f(x_k)$.

$x_{k+1} - x^* = x_k - x^* - s\nabla f(x_k)$

$\Rightarrow \|x_{k+1} - x^*\|^2 = \|x_k - x^*\|^2 - 2s\langle \nabla f(x_k), x_k - x^*\rangle$

$$+ s^2 \underbrace{\|\nabla f(x_k)\|^2}_{\text{bdd.}} \qquad \text{by monotonicity}$$

$\langle \nabla f(x_k) - \underset{0}{\nabla f(x^*)}, x_k - x^*\rangle \geq \mu\|x_k - x^*\|^2$

$\|\nabla f(x_k) - \underset{0}{\nabla f(x^*)}\| \leq \ell\|x_k - x^*\|$

$\leq \|x_k - x^*\|^2 - 2s\mu\|x_k - x^*\|^2 + s^2\ell^2\|x_k - x^*\|^2$

$= \underbrace{(1 - 2s\mu + s^2\ell^2)}_{< 1 \,?} \|x_k - x^*\|^2$

$<1\,?\quad$ when $\quad 0 < 2s\mu - s^2\ell^2$

$\qquad\qquad\qquad \Leftrightarrow \; s < \dfrac{2\mu}{\ell^2} \qquad$ #

<u>Rk.</u> The proof only uses that $\nabla f(x)$ is strongly-monotone

(in addition to $\ell$-smoothness of $f$). See L2.