$$\min_{x} \max_{y} L(x,y)$$

NC-$\underline{SC}$ :  $\overset{\text{weakly convex}}{\text{nonconvex in } x}$, strongly concave in $y$

↑
more generally can be PL

- GDA with properly set (two scales of) time steps

  $\underline{\text{converges}}$ in $O(\frac{1}{\varepsilon^2})$ steps
  ↓
  find a "stionary pt" with $O(\varepsilon)$ gradient

Q. why $O(\frac{1}{\varepsilon^2})$ ?

Recall that this is the iteration complexity of vanilla GD

on nonconvex minimization, and it is tight (see L1)

Thus, suppose we use that $\max_{y}$ is SC and can be

solved fast (in $O(\log\frac{1}{\varepsilon})$ inner-loop steps) to "eliminate"

the variable, then we have the minimization problem

$$\min_{x} \bar{L}(x), \quad \bar{L}(x) := \max_{y} L(x,y),$$

and this still needs $O(\frac{1}{\varepsilon^2})$ outer-loop steps.

In this sense, the $O(\frac{1}{\varepsilon^2})$ iteration of GDA for the

minimax problem is also "tight".

- Below, we first prove the $O(\frac{1}{\varepsilon^2})$ convergence of GDA for

  the NC-PL type. The variants :

  { A-GDA      Alternative-GDA

1 S-GDA       Smoothed-GDA (with proximal step)

[LJJ 2020] Lin, Jin, Jordan. On gradient descent ascent for nonconvex-concave minimax problems. ICML 2020.
[YOLH 2022] Yang, Orvieto, Lucchi, He. Faster single-loop algorithms for minimax optimization without strong concavity. AISTATS 2022.

- PL condition (Polyak-Łojasiewicz)

$$\min_x h(x) = h^* \quad . \quad h \ C^1 \ \text{on} \ \mathbb{R}^d.$$

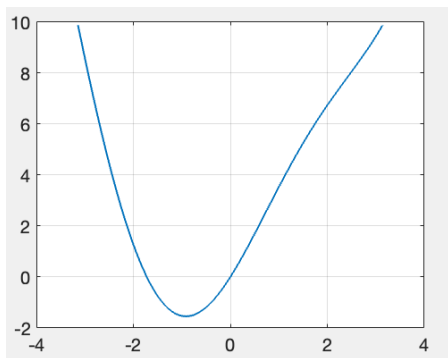def $h$ is $\mu$-PL if

$$\frac{1}{2} \|\nabla h(x)\|^2 \geq \mu (h(x) - h^*), \quad \forall x \in \mathbb{R}^d.$$

lemma $h$ is $\mu$-strongly convex $\Rightarrow \mu$-PL
                    S.C.

fact. $h$ can be $\mu$-PL without $\mu$-S.C.
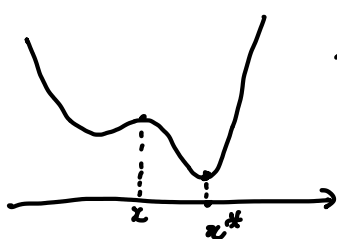
Eg. non-convex but PL.

$$h(x) = x^2 + 3 \sin x$$



$h' = 2x + 3\cos x$
$\quad x^* \approx -0.9148$

$h'' = x - 3 \sin x$, can $< 0$.

$\dfrac{\frac{1}{2} |h'(x)|^2}{h(x) - h^*} \geq \mu \approx 0.416$

Q. what does NOT satisfy PL?

$$h'(x) = 0, \quad h(x) - h^* > 0$$



Eg. Moreau envelope

$$\ell(x) = -\frac{1}{2} x_2^-, \qquad x = [x_i]$$

$$\nabla^2 \ell = \begin{bmatrix} 0 & 0 \\ 0 & -\rho \end{bmatrix} \succeq -\rho I. \qquad \rho\text{-weakly convex}$$

$$\ell_\eta(x) = \max_z \underbrace{\ell(z) + \frac{\eta}{2}\|z - x\|^2}_{L_{\eta,x}(z) \;=\; L(z)}$$

**Q** Is $L_{\eta,x}(z)$ PL?

- if $\eta > \rho$, $L(z)$ is $(\eta - \rho)$- S.C.

$$L(z) = -\frac{\rho}{2} z_2^2 + \frac{\eta}{2}(z_1 - x_1)^2 + \frac{\eta}{2}(z_2 - x_2)^2$$

- if $\eta < \rho$, no minimiser on $\mathbb{R}^2$

$\eta = \rho$, let $x_2 = 0$,

$$L(z) = \frac{\eta}{2}(z_1 - x_1)^2, \qquad \nabla^2 L = \begin{bmatrix} \eta & 0 \\ 0 & 0 \end{bmatrix}$$
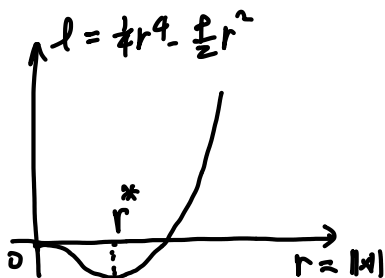
convex but not S.C.

$$\nabla L(z) = \begin{bmatrix} \eta(z_1 - x_1) \\ 0 \end{bmatrix}, \quad L^* = 0$$

$$\frac{\frac{1}{2}\|\nabla L(z)\|^2}{L(z) - L^*} = \frac{\frac{1}{2}\eta^2(z_1 - x_1)^2}{\frac{\eta}{2}(z_1 - x_1)^2} = \eta > 0. \quad \forall z$$

Thus, $L$ is convex, non-S.C., but $\mu$-PL.

**Another example:**

$$\ell(x) = \frac{1}{4}\|x\|^4 - \frac{\rho}{2}\|x\|^2 \qquad \rho\text{-weakly convex}$$

$$\ell = \frac{1}{4}r^4 - \frac{\rho}{2}r^2$$



$$L(z) = \frac{1}{4}\|z\|^4 - \frac{\rho}{2}\|z\|^2 + \frac{\eta}{2}\|z - x\|^2, \qquad x \text{ fixed}$$

$$\nabla L(z) = (\|z\|^2 - \rho)\, z + \eta\,(z - x)$$

suppose $x = 0$, then $L(z) = \frac{1}{4}\|z\|^4 - \frac{1}{2}(\rho - \eta)\|z\|^2$.

$$L(z) - L^* = \frac{1}{4}\left(\|z\|^2 - (\rho - \eta)\right)^2.$$

$$\nabla L(z) = \left(\|z\|^2 - (\rho - \eta)\right)\|z\|$$

$$\Rightarrow \frac{\frac{1}{2}\|\nabla L(z)\|^2}{L(z) - L^*} = 2\|z\|^2 \to 0\dagger \text{ when } z \to 0$$

locally PL (on a nbh) but not globally

- PL $\Rightarrow$ linear convergence in minimization

[KNS 2016] Karimi, Nutini, Schmidt. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition. (2016)

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

(PL) $\quad \frac{1}{2}\|\nabla h(x)\|^2 \geq \mu\,(h(x) - h^*),\ \forall x$

(EB) $\quad \|\nabla h(x)\| \geq \mu\|x - x^*\|,\ \forall x$
error bound $\qquad\qquad \uparrow$
$\qquad$ generally $\|x - x_p\|$, $x_p$ is the projection
$\qquad$ of $x$ to the minimizer set $\mathcal{X}^*$

(QG) $\quad h(x) - h^* \geq \frac{\mu}{2}\|x - x^*\|^2,\ \forall x$
quadratic growth

Thm. If $h$ is $C^1$ and $L$-smooth on $\mathbb{R}^d$, then

$$(SC) \Rightarrow (EB) \equiv (PL) \Rightarrow (QG)$$

If furtherly $h$ is convex, then $(EB) \equiv (PL) \equiv (QG)$

ref. Thm 2 of [KNS 2016].

Thm Suppose $h$ is $C^1$ and $L$-smooth, $h^* = h(x^*)$, $x^* \exists !$,

$h$ satisfies (PL) with $\mu > 0$ then GD

$$z_{k+1} \leftarrow z_k - s\nabla h(z_k)$$

with $s < \frac{2}{\ell}$ converges linearly. i.e.

$$h(z_k) - h^* \leq \varepsilon^2 \quad \text{in} \quad k = O(\log\tfrac{1}{\varepsilon}) \text{ steps.}$$

Rk. By (QG), $h(z_k) - h^* \geq \frac{\mu}{2}\|z_k - z^*\|^2$, so objective

value convergence to $O(\varepsilon^2)$ implies $\|z_k - z^*\| = O(\varepsilon)$,

also $\|\nabla h(z_k) - \nabla h(z^*)\| \leq \ell\|z_k - z^*\| = O(\varepsilon)$.

pf. descent of $h(z_k)$: by $\ell$-smoothness of $h$,

$$h(z_{k+1}) \leq h(z_k) + \langle \nabla h(z_k), \underbrace{z_{k+1} - z_k}_{\substack{\| \\ -s\nabla h(z_k)}}\rangle + \frac{\ell}{2}\|z_{k+1} - z_k\|^2$$

$$= h(z_k) - s\|\nabla h(z_k)\|^2 + \frac{\ell}{2}s^2\|\nabla h(z_k)\|^2$$

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

$\underline{\text{set } s = \frac{1}{\ell}} \qquad \underbrace{0 < 1 - \frac{\ell}{2}s}_{} < 1 \text{ if } \frac{\ell}{2}s < 1$

$$\downarrow\ = h(z_k) - \frac{s}{2}\|\nabla h(z_k)\|^2$$

This is $\quad h(z_k) - h(z_{k+1}) \geq \frac{s}{2}\|\nabla h(z_k)\|^2$

so far no difference from the "NC case" proof in L1.

By PL, $\quad \|\nabla h(z_k)\|^2 \geq 2\mu(h(z_k) - h^*)$,

$\Rightarrow \quad h(z_{k+1}) \leq h(z_k) - s\mu(h(z_k) - h^*)$

$\Rightarrow \quad (h(z_{k+1}) - h^*) \leq (1 - s\mu)(h(z_k) - h^*)$

$\underbrace{\qquad}_{\overset{\|}{1 - \frac{\mu}{\ell} = 1 - \frac{1}{\kappa}}}$

$\kappa := \frac{\ell}{\mu}$ condition number, $\kappa \geq 1$

This shows $\quad h(z_k) - h^* \leq (1 - \frac{\mu}{\ell})^k (h(z_0) - h^*)$ $\dots$

Rk. (PL) is a relaxed assump. from (SC) but still ensures

exponential convergence of GD, or GD-like schemes.

• NC-PL minimax problem

$$\min_x \max_y L(x,y)$$

(A1)   $L$ is $C^1$ on $\mathbb{R}^n \times \mathbb{R}^m$, $l$-smooth

(A2)   $L(x, \cdot)$ is $\mu$-PL in $y$ and $y^*(x) \exists !$ , $\forall x$

Intuitively, $y$ is the "fast variable" to be solved.

A-GDA ($\tau, \eta$)        two-timescale

$$\begin{cases} x_{k+1} \leftarrow x_k - \tau \nabla_x L(x_k, y_k) \\ y_{k+1} \leftarrow y_k + \eta \nabla_y L(x_{k+1}, y_k) \end{cases}$$

$\uparrow x_k$ if without alternating
(just GDA)

We will set   $\tau \sim \frac{1}{l} \frac{1}{\kappa^2}$, $\eta \sim \frac{1}{l}$, $\kappa := \frac{l}{\mu} \geq 1$

$\tau$ is smaller, so the update of $x_k$ is slower.