

# Optimization

$$x^* \in \operatorname{argmin}_{x \in X} f(x)$$

$$\iff \delta_{x^*} \in \operatorname{argmin}_{\mu \text{ prob. distrib.}} \int_X f(x) d\mu(x)$$

$$f^*(\omega) \stackrel{\text{def.}}{=} \sup_x \langle x, \omega \rangle - f(x)$$

$$(f \diamond g)(x) \stackrel{\text{def.}}{=} \inf_y f(y) + g(x - y)$$

*Theorem:*  $(f \diamond g)^* = f^* + g^*$

$$\hat{f}(\omega) \stackrel{\text{def.}}{=} \int f(x) e^{-i\langle \omega, x \rangle} dx$$

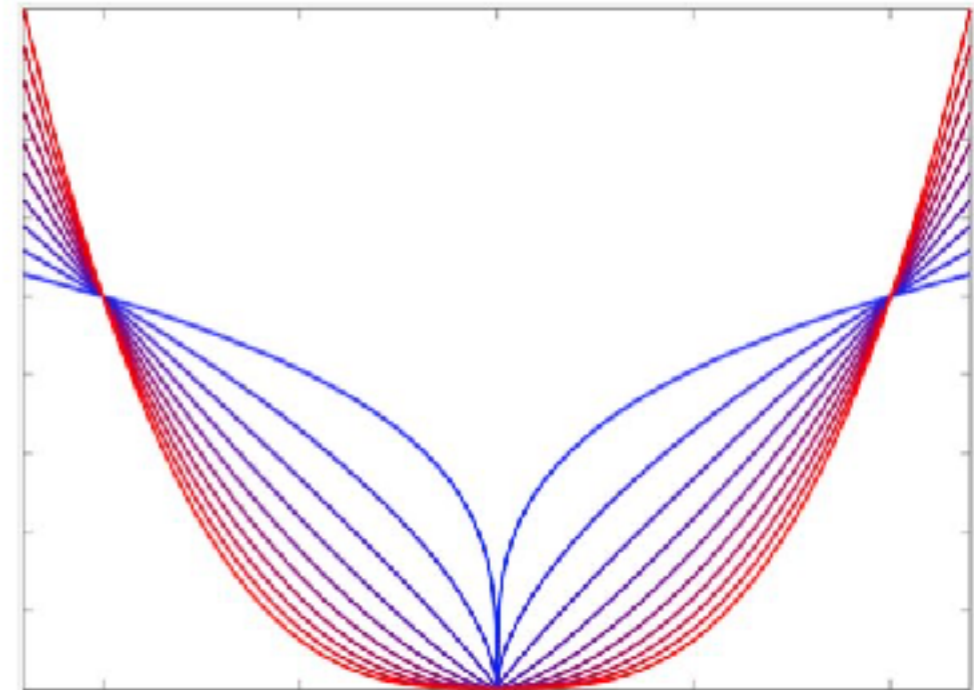
$$(f \star g)(x) \stackrel{\text{def.}}{=} \int f(y) g(x - y) dy$$

*Theorem:*  $\widehat{f \star g} = \hat{f} \cdot \hat{g}$

Kurdyka-Łojasiewicz:  
(at minimum  $f(0) = 0$ )

$$\exists \varphi, \|\nabla(\varphi \circ f)\| \geq 1$$

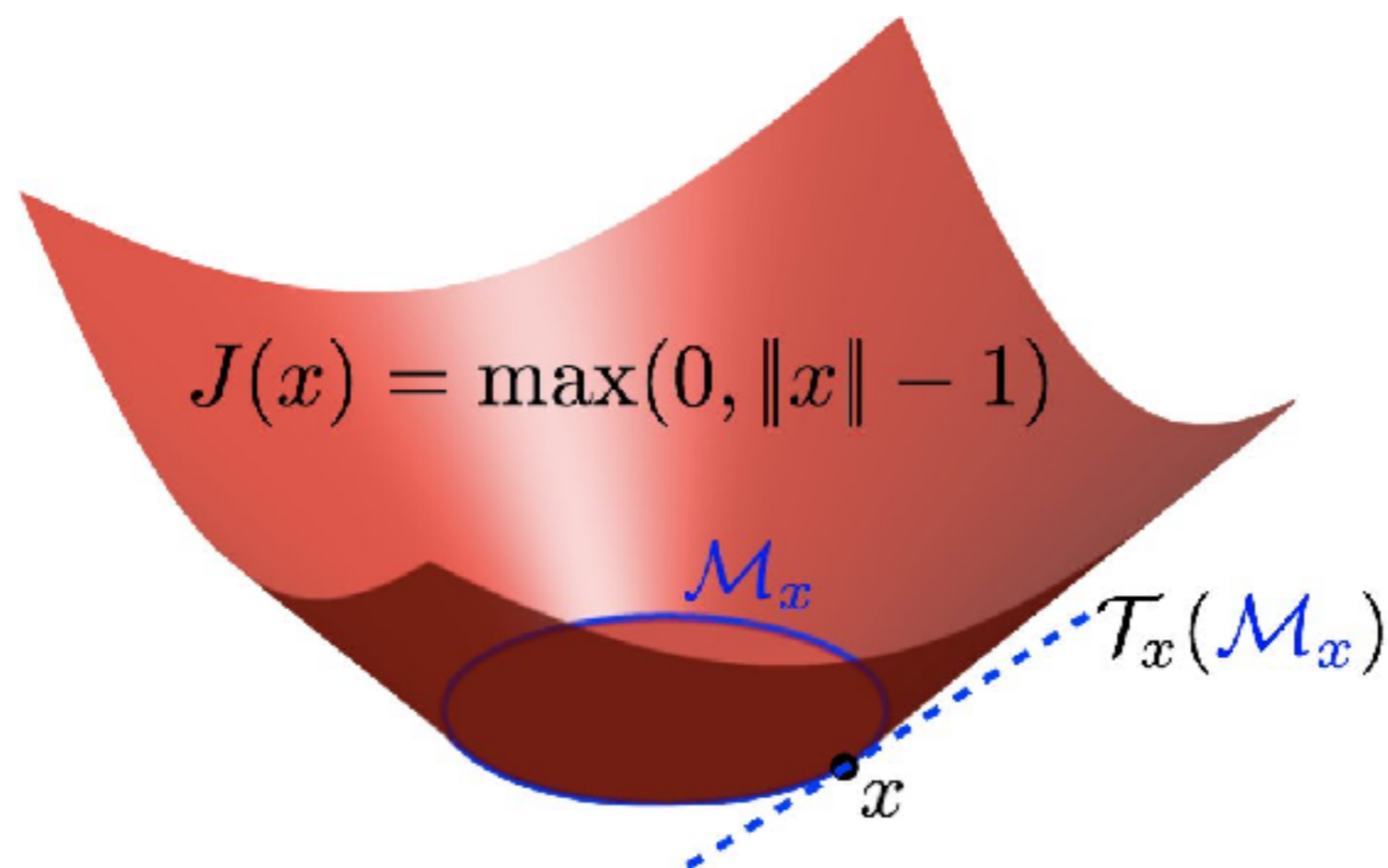
$$\iff \exists (K, \tau), \|\nabla f(x)\| \geq K|f(x)|^\tau$$

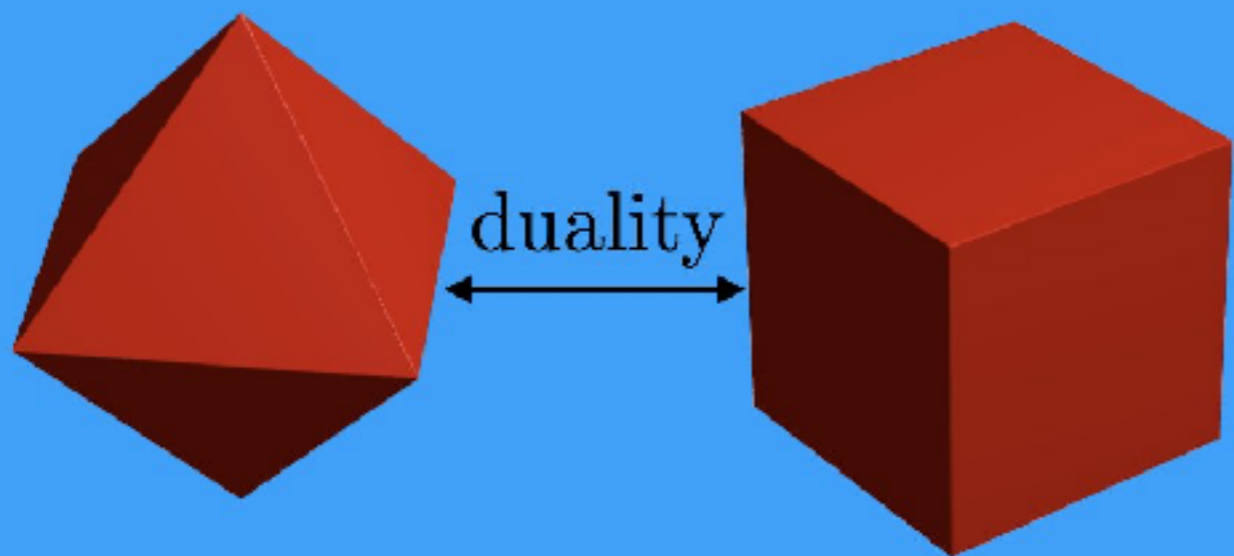


$J : \mathbb{R}^N \rightarrow \mathbb{R}$  is partly smooth at  $x$  for a manifold  $\mathcal{M}_x$

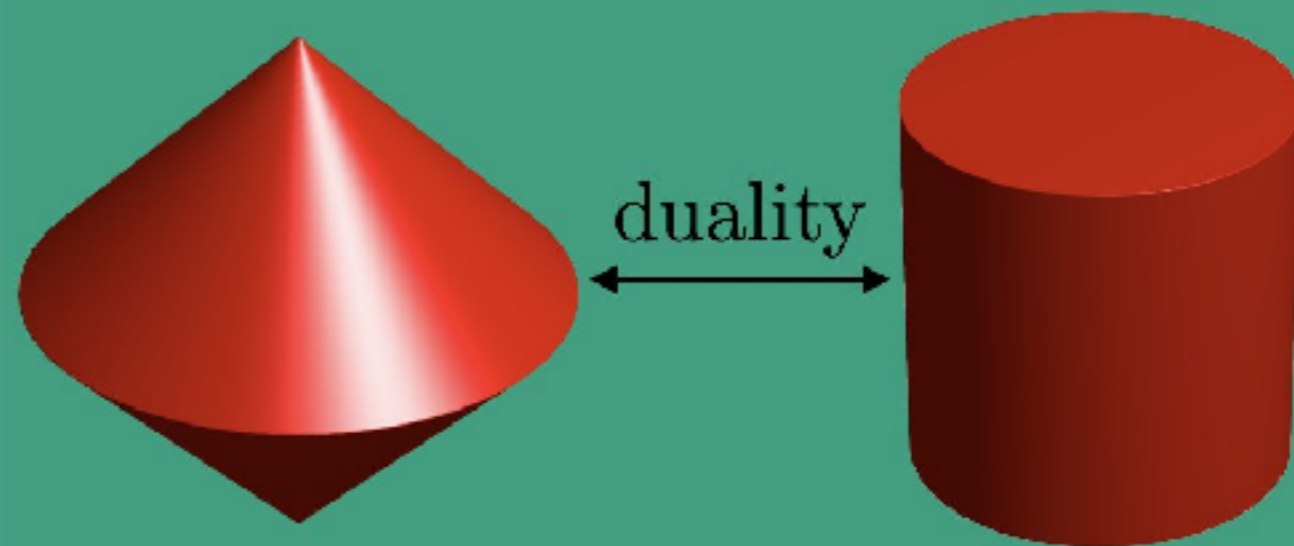
- (i)  $J$  is  $C^2$  along  $\mathcal{M}_x$  around  $x$  ;
- (ii)  $\forall h \in \mathcal{T}_x(\mathcal{M}_x)^\perp, t \mapsto J(x + th)$  non-smooth at  $t = 0$ .
- (iii)  $\partial J$  is continuous on  $\mathcal{M}_x$  around  $x$ .

[Lewis 2003]

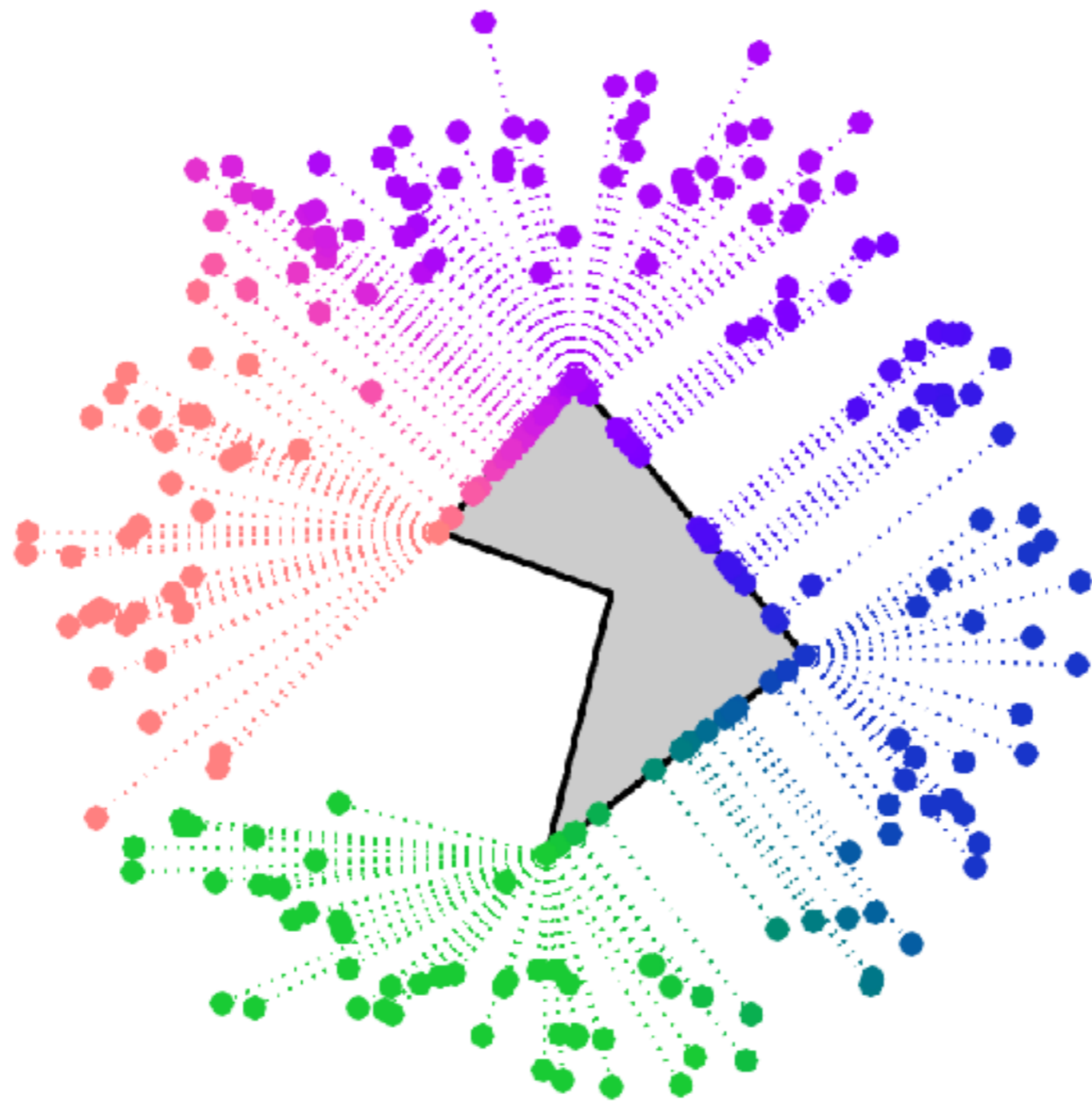




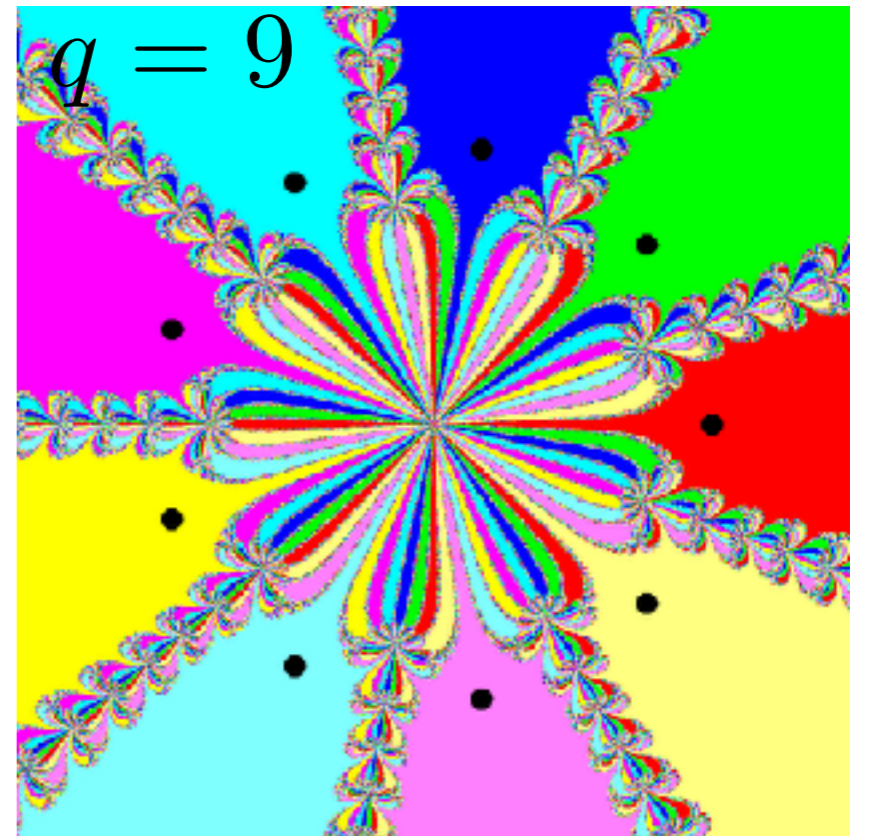
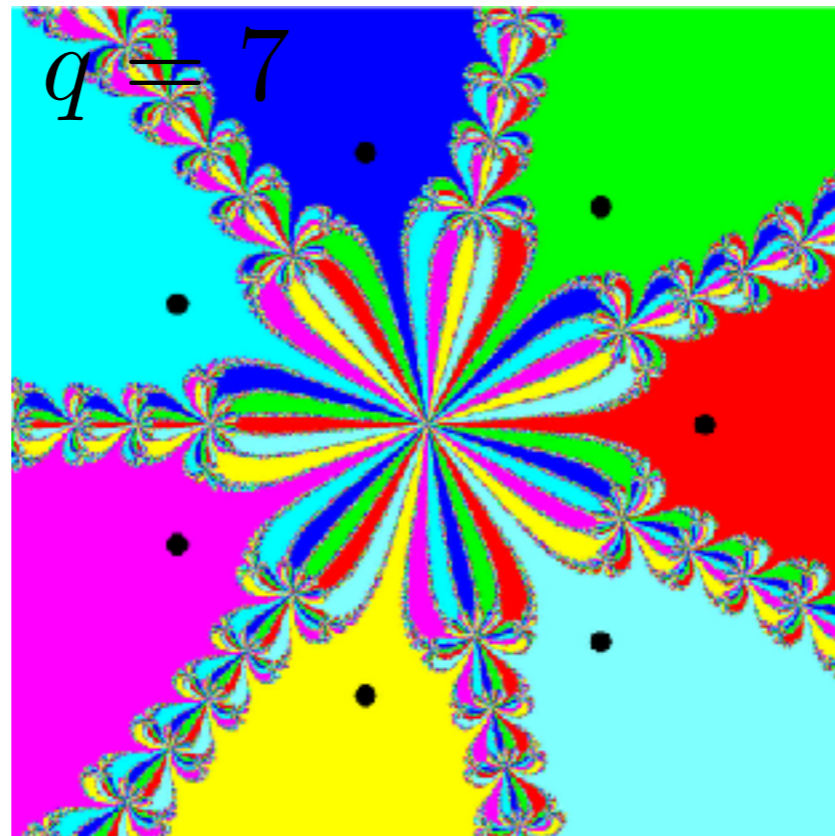
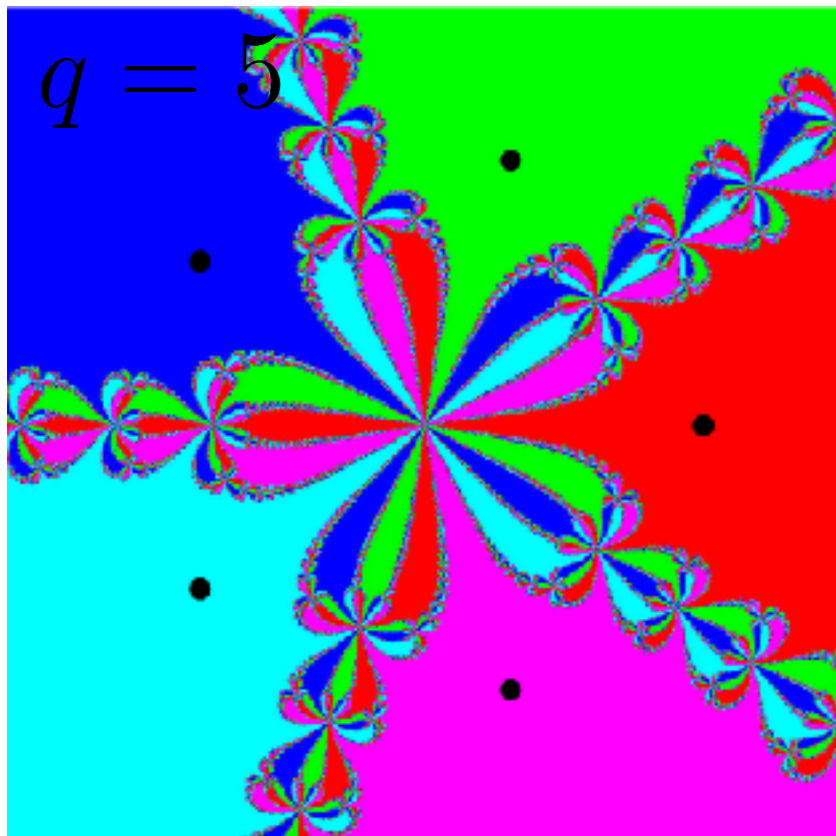
Polyhedral



Non-polyhedral



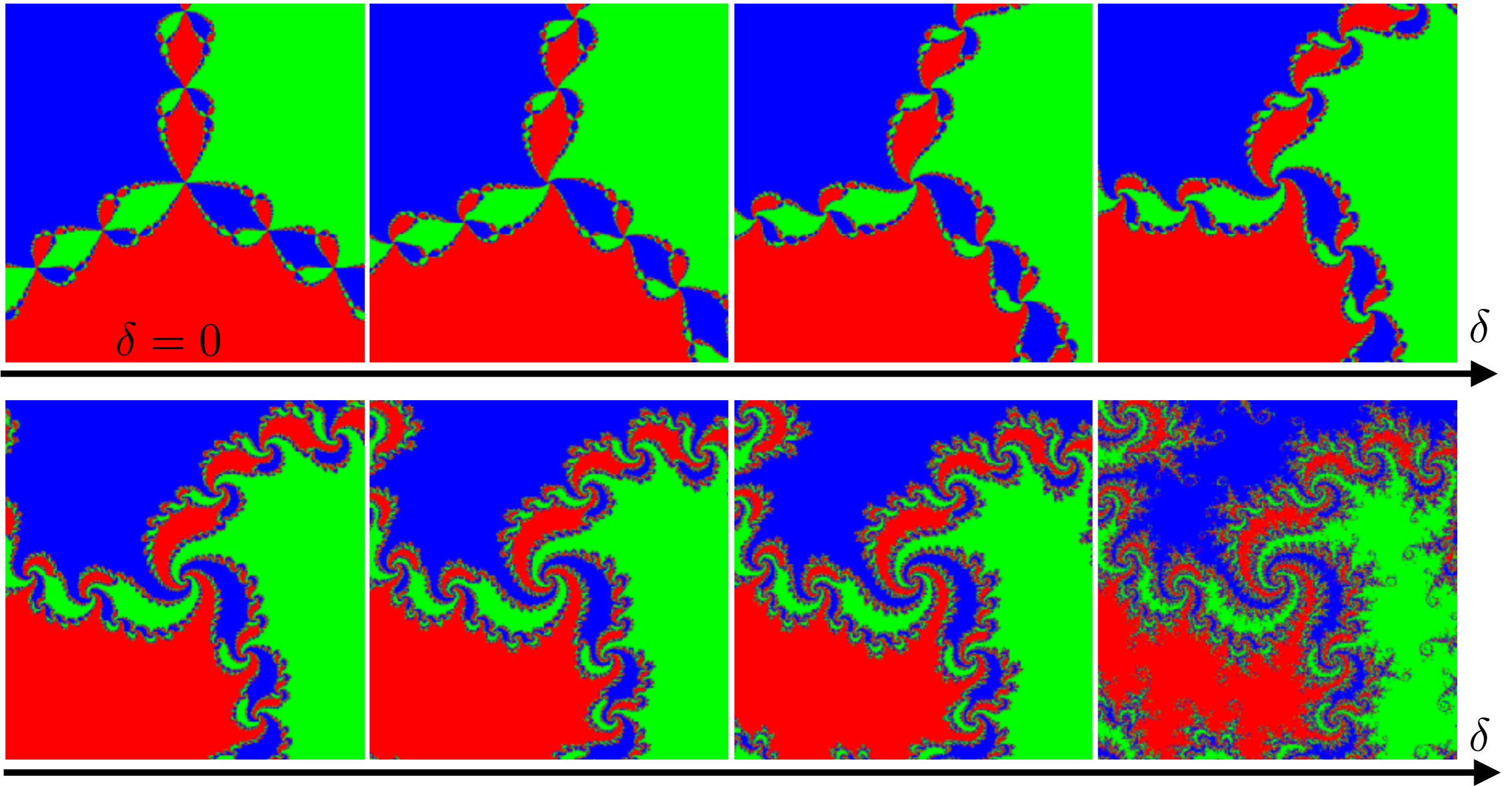
Newton method:  $z_{k+1} = z_k - \frac{f(z_k)}{f'(z_k)}$



Attraction bassins for  $f(z) = z^q - 1$



“Twisted” Newton:  $z_{k+1} = z_k - (1 + \delta e^{i\theta}) \frac{f(z_k)}{f'(z_k)}$        $f(z) = z^3 - 1$

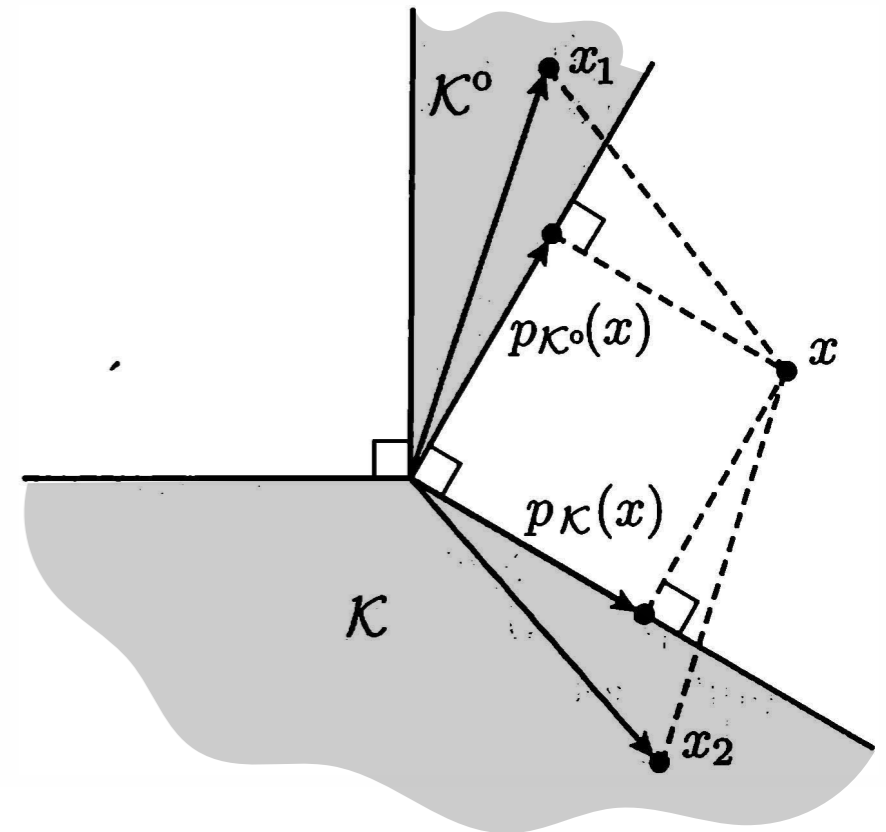


*Polar cone:*

$$\mathcal{K}^\circ \stackrel{\text{def.}}{=} \{x ; \forall y \in \mathcal{K}, \langle x, y \rangle \leq 0\}$$

*Theorem:* (Moreau's decomposition)

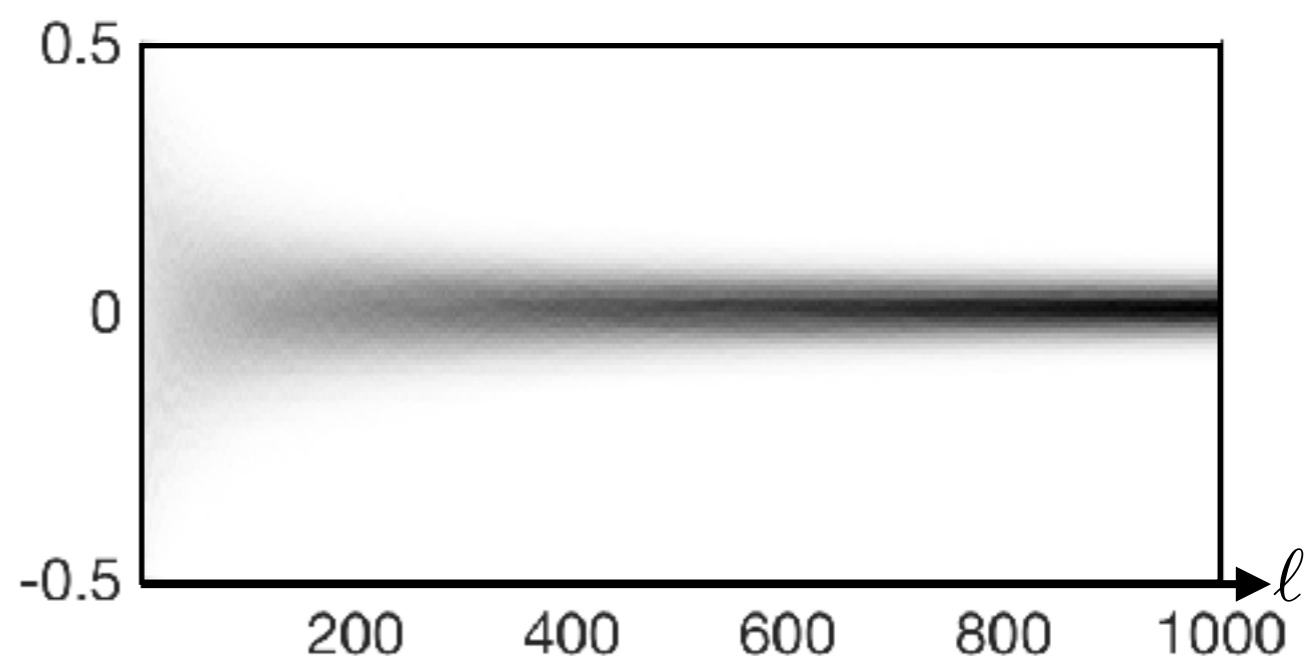
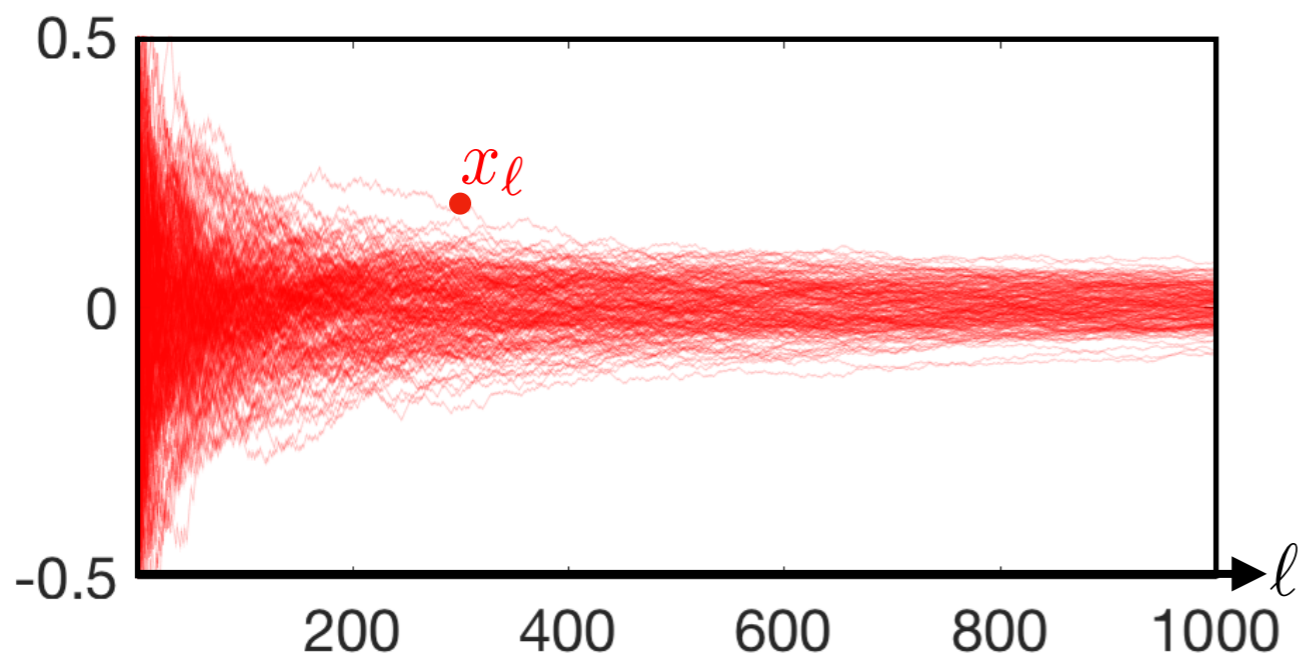
$$x = \text{Proj}_{\mathcal{K}}(x) +^\perp \text{Proj}_{\mathcal{K}^\circ}(x)$$



$$\min_{x \in \mathbb{R}} (x+1)^2 + (x-1)^2$$

$$= f_1(x) \quad = f_2(x)$$

$$x_{\ell+1} \stackrel{\text{def.}}{=} \begin{cases} x_{\ell} - \frac{1}{\ell} \nabla f_1(x_{\ell}) & \text{with proba } \frac{1}{2} \\ x_{\ell} - \frac{1}{\ell} \nabla f_2(x_{\ell}) & \text{with proba } \frac{1}{2} \end{cases}$$



Gradient descent dynamic:

$$x(t) = -\nabla f(x(t))$$

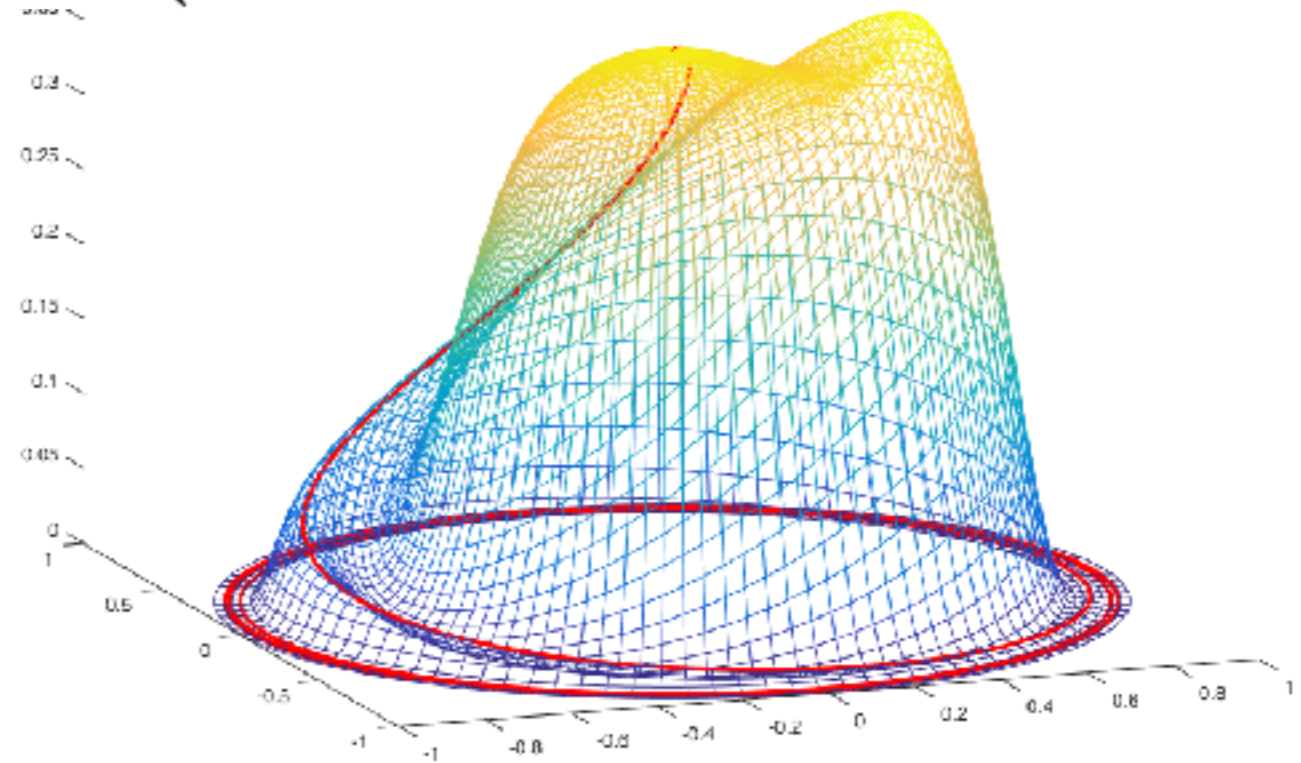
On the “mexican hat”:

$$x(t) = r(t)(\cos(\theta(t)), \sin(\theta(t)))$$

$$\theta(t) = (1 - r(t))^{-2}$$

→ Length( $x$ ) =  $+\infty$

$$f(r, \theta) := \begin{cases} e^{-\frac{1}{1-r^2}} \left[ 1 - \frac{4r^4}{4r^4 + (1-r^2)^4} \sin\left(\theta - \frac{1}{1-r^2}\right) \right] & \text{if } r < 1, \\ 0 & \text{if } r \geq 1, \end{cases}$$



$$f : z \stackrel{\text{def.}}{=} x + iy \in \mathbb{C} \longmapsto f(z) \in \mathbb{C}$$

$$\frac{\partial}{\partial z} \stackrel{\text{def.}}{=} \frac{1}{2} \left( \frac{\partial}{\partial x} - i \frac{\partial}{\partial y} \right)$$

$$\frac{\partial}{\partial \bar{z}} \stackrel{\text{def.}}{=} \frac{1}{2} \left( \frac{\partial}{\partial x} + i \frac{\partial}{\partial y} \right)$$

*Proposition:*  $f$  holomorphic  $\Leftrightarrow \frac{\partial f}{\partial \bar{z}} = 0$ .

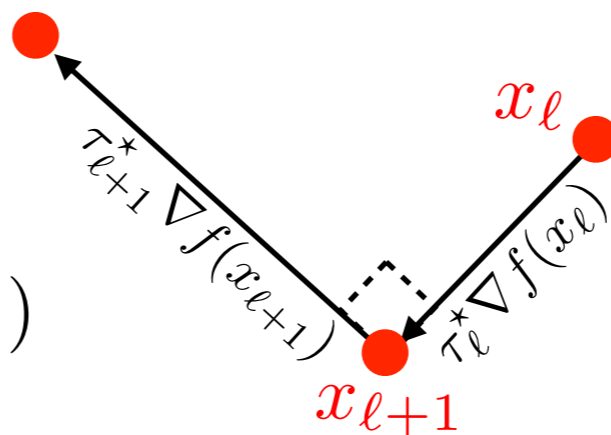
*Examples:*

$$\frac{\partial z}{\partial z} = 1 \quad \frac{\partial \bar{z}}{\partial z} = 0 \quad \frac{\partial z}{\partial \bar{z}} = 0 \quad \frac{\partial \bar{z}}{\partial \bar{z}} = 1$$

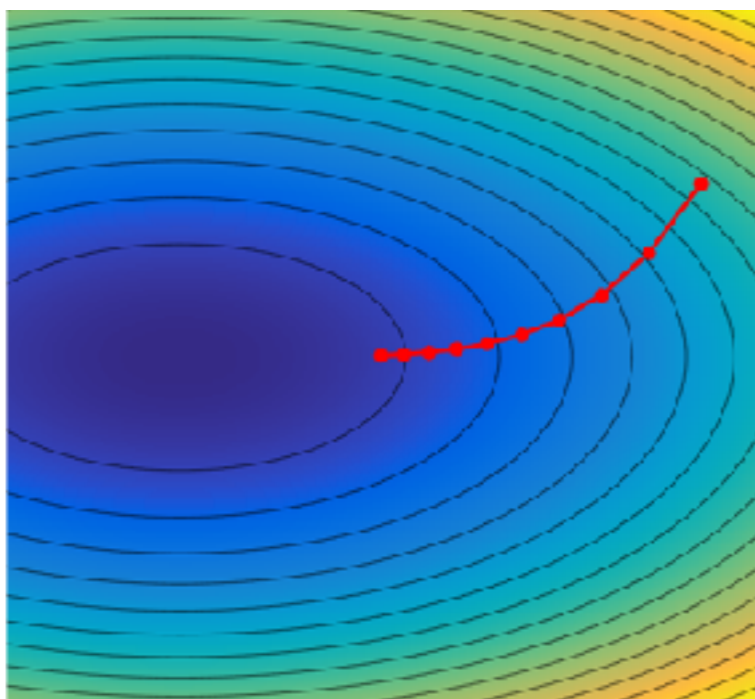
$$\frac{\partial cz}{\partial z} = c \quad \frac{\partial |z|^2}{\partial z} = \frac{\partial (z\bar{z})}{\partial z} = \bar{z}$$

$$x_{l+1} = x_l - \tau_l \nabla f(x_l)$$

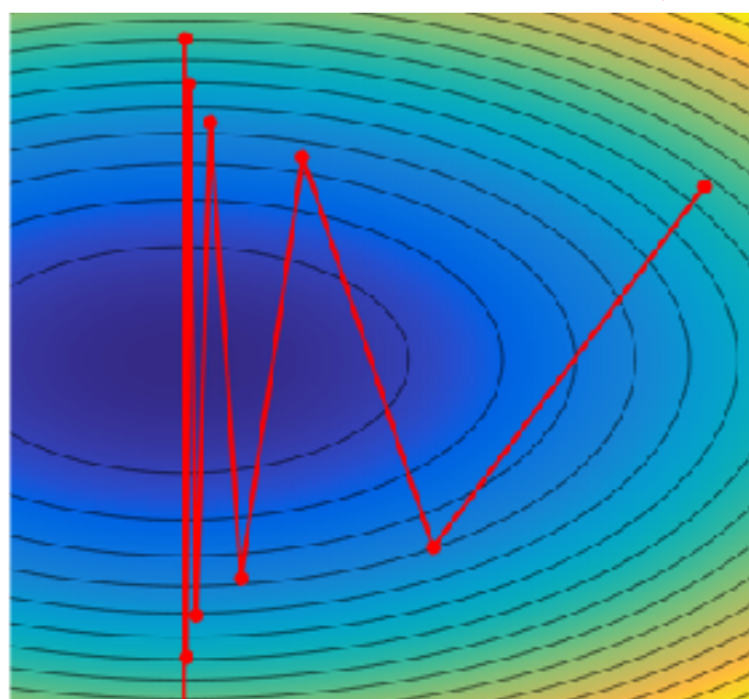
$$\tau_l^* = \operatorname{argmin}_{\tau} f(x_l - \tau \nabla f(x_l))$$



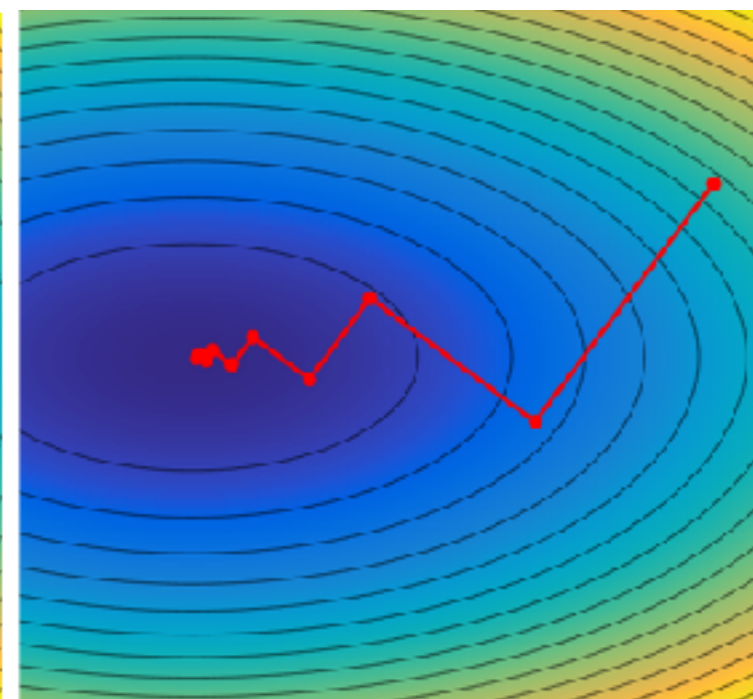
$$\nabla f(x_l) \perp \nabla f(x_{l+1})$$



Small  $\tau_l$

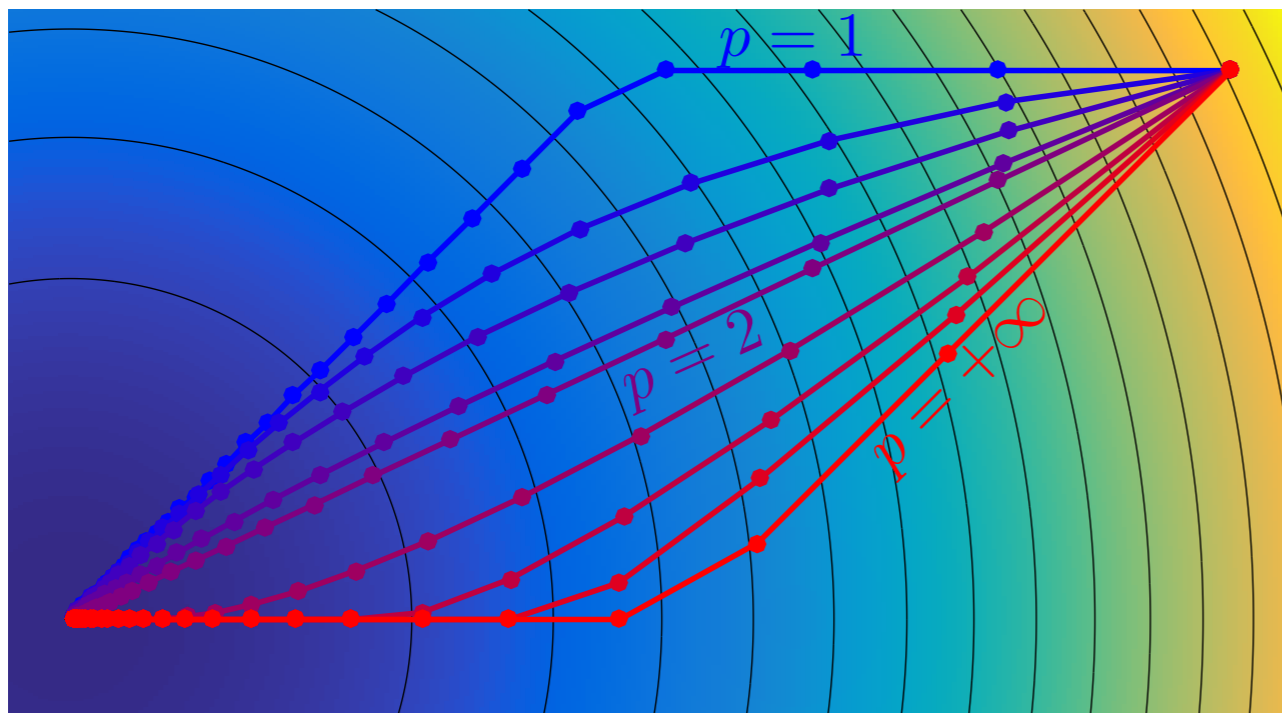


Large  $\tau_l$



Optimal  $\tau_l = \tau_l^*$

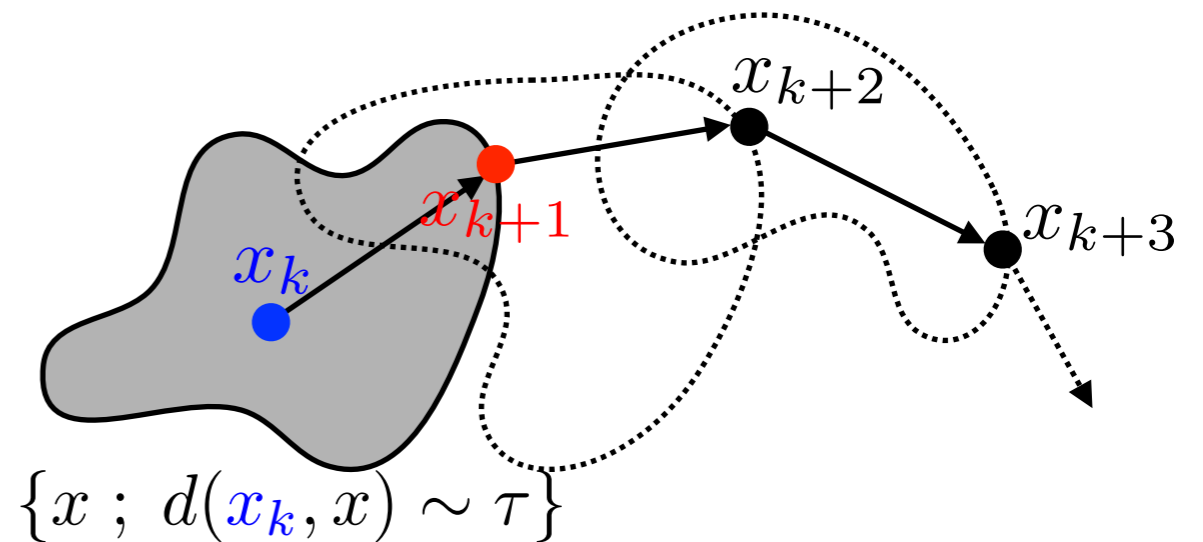
Metric space  $(\mathcal{X}, d)$ , minimize  $F(x)$  on  $\mathcal{X}$ .



$$F(x) = \|x\|^2 \text{ on } (\mathcal{X} = \mathbb{R}^2, \|\cdot\|_p)$$

Implicit Euler step:

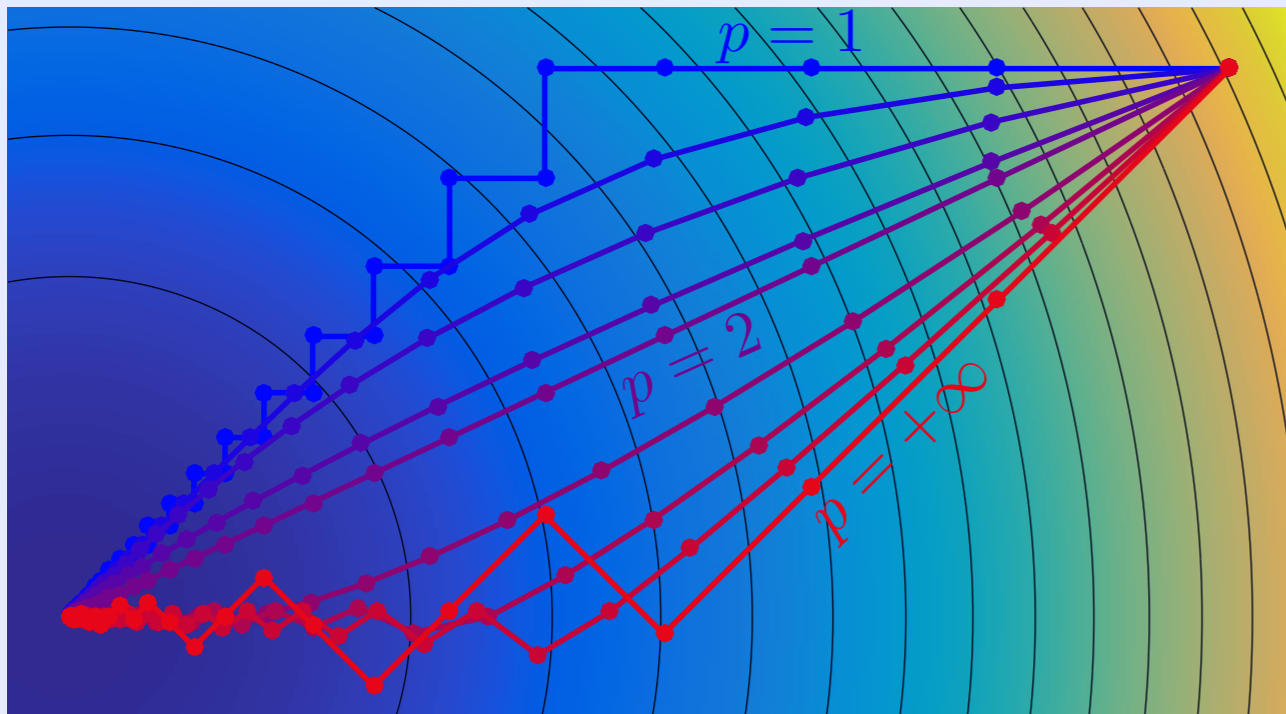
$$x_{k+1} \stackrel{\text{def.}}{=} \operatorname{argmin}_{x \in \mathcal{X}} d(x_k, x)^2 + \tau F(x)$$



Metric space  $(\mathcal{X}, d)$ , minimize  $F(x)$  on  $\mathcal{X}$ .

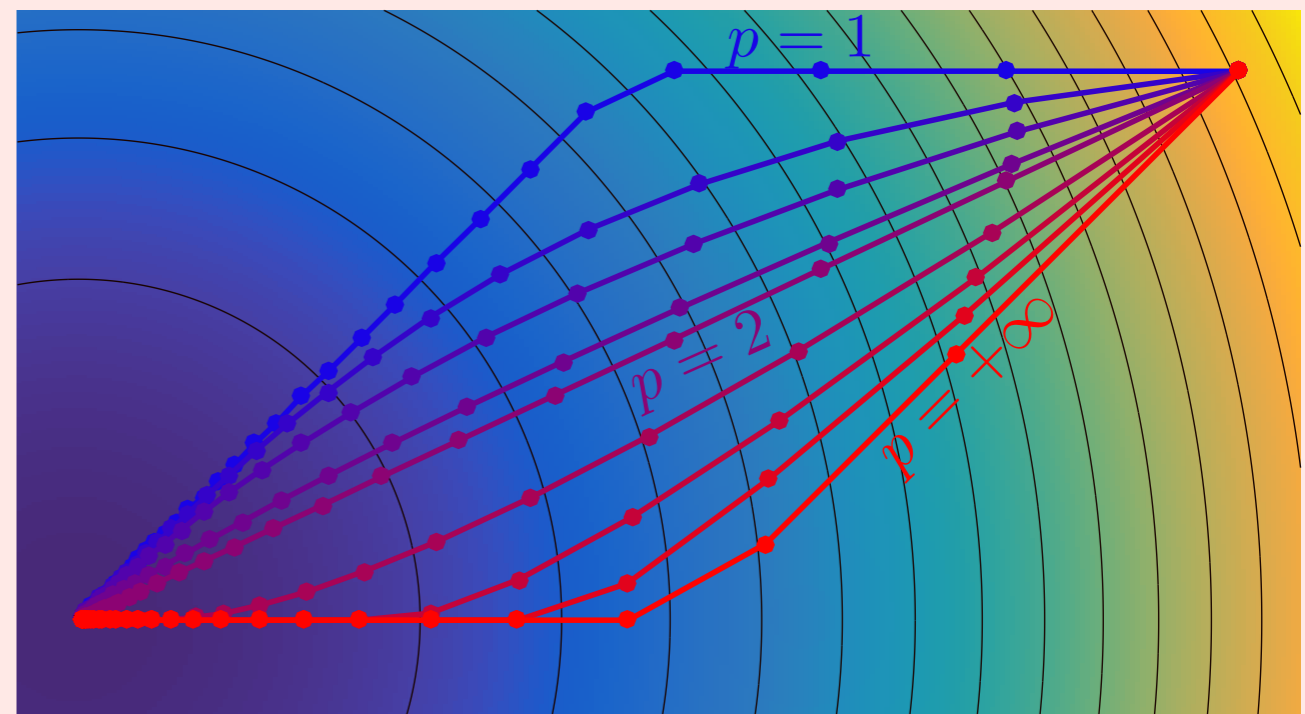
Explicit

$$x_{k+1} = \operatorname{argmin}_{x \in \mathcal{X}} d(x_k, x)^2 + \tau \langle \nabla F(x_k), x \rangle$$



Implicit

$$x_{k+1} = \operatorname{argmin}_{x \in \mathcal{X}} d(x_k, x)^2 + \tau F(x)$$



$$F(x) = \|x\|^2 \text{ on } (\mathcal{X} = \mathbb{R}^2, \|\cdot\|_p)$$

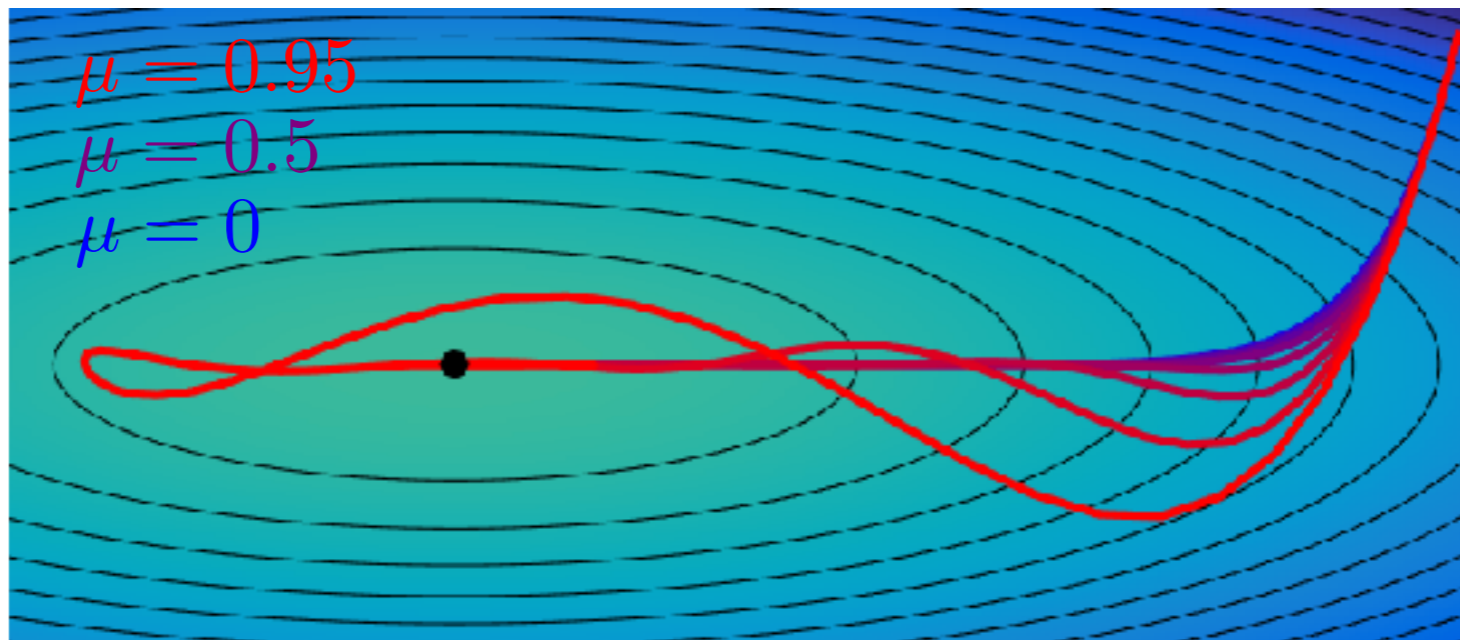


$$x_{k+1} = x_k + p_k$$

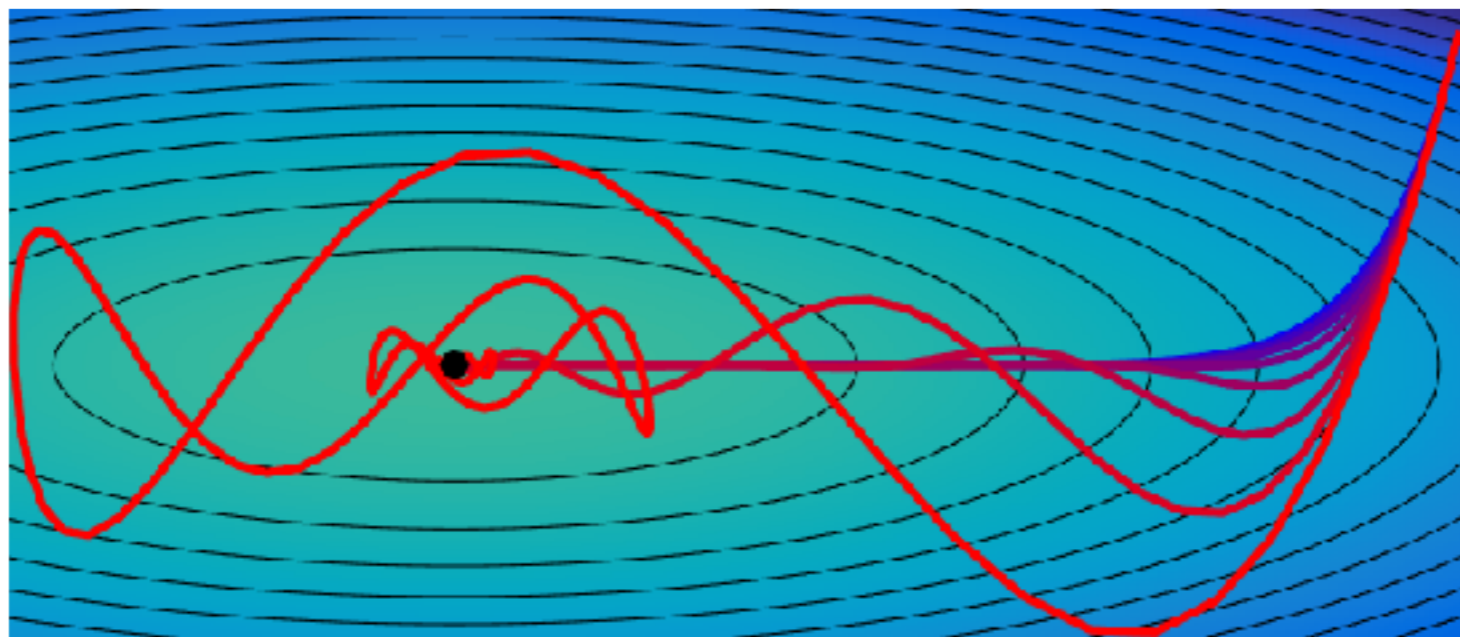
$$p_{k+1} = \mu p_k - \tau \begin{cases} \nabla f(x_k) & \text{Polyak} \\ \nabla f(x_k + \mu p_k) & \text{Nesterov} \end{cases}$$



Yurii  
Nesterov



Boris  
Polyak



## Gradient descent

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

$$\tau \rightarrow 0 \quad \downarrow \quad k\tau \rightarrow t$$

$$\frac{dx(t)}{dt} = -\nabla f(x(t))$$

## Nesterov's acceleration

$$x_{k+1} = y_k - \tau \nabla f(y_k)$$
$$y_{k+1} = x_{k+1} + \frac{k}{k+3} (x_{k+1} - x_k)$$

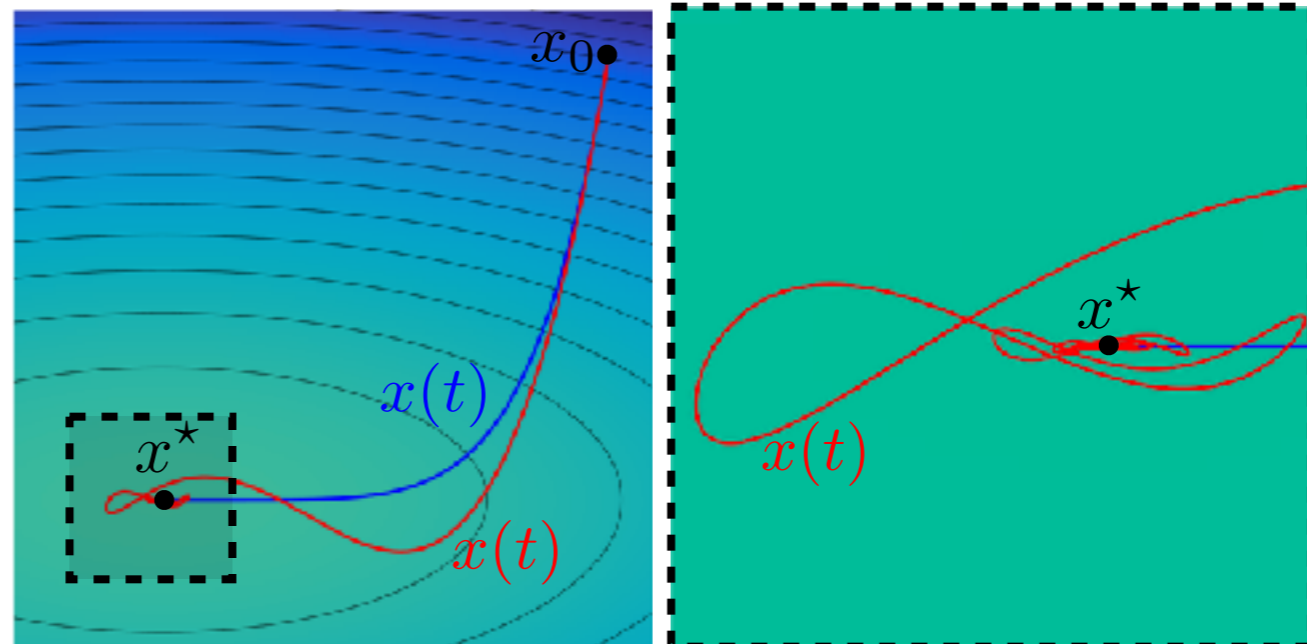
$$\tau \rightarrow 0 \quad \downarrow \quad k\sqrt{\tau} \rightarrow t$$

$$\frac{d^2 x(t)}{dt^2} + \frac{3}{t} \frac{dx(t)}{dt} = -\nabla f(x(t))$$

*Theorem:*

$$f(x_k) - f(x^*) = O(1/k)$$

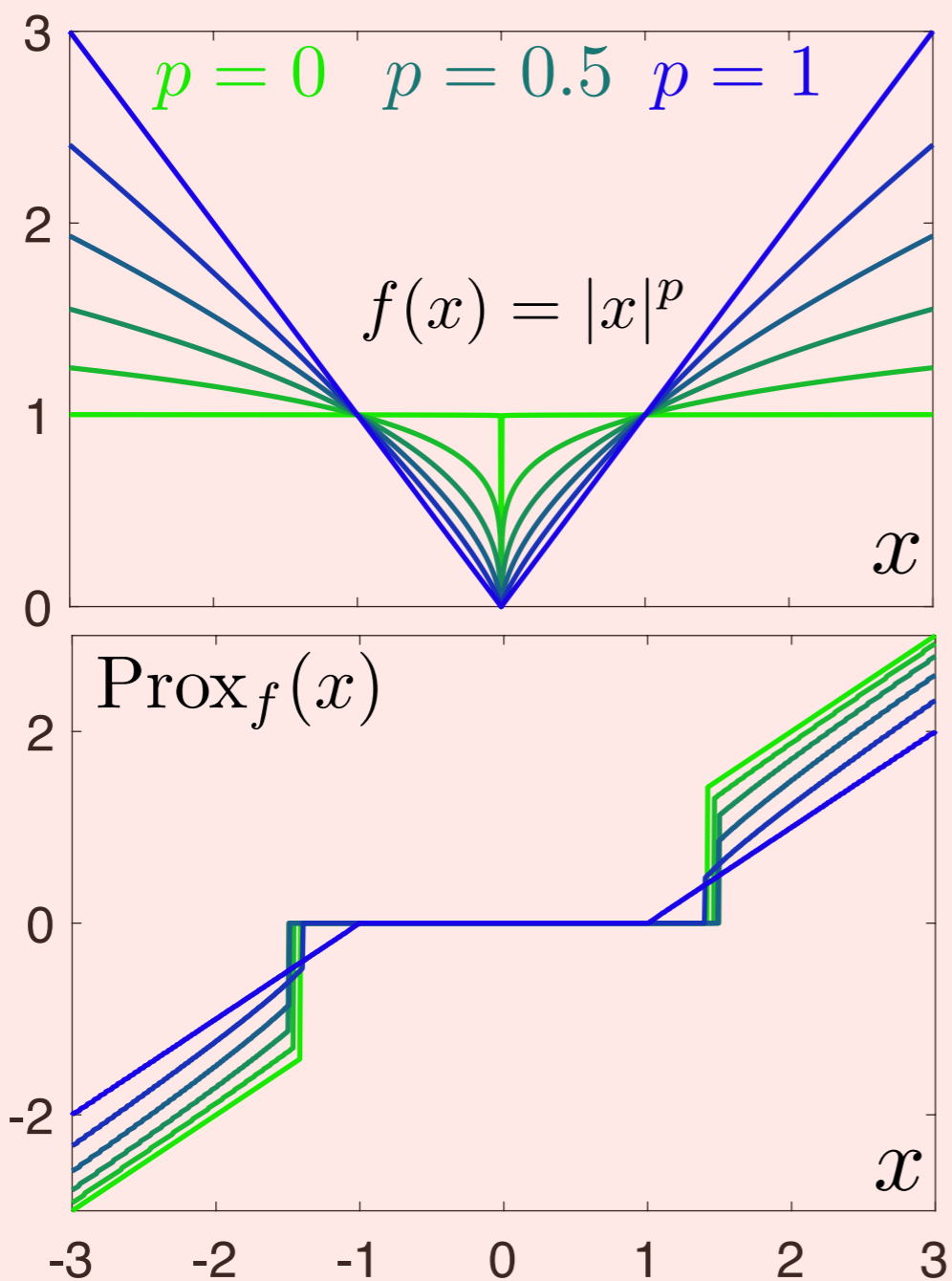
$$f(x_k) - f(x^*) = O(1/k^2)$$



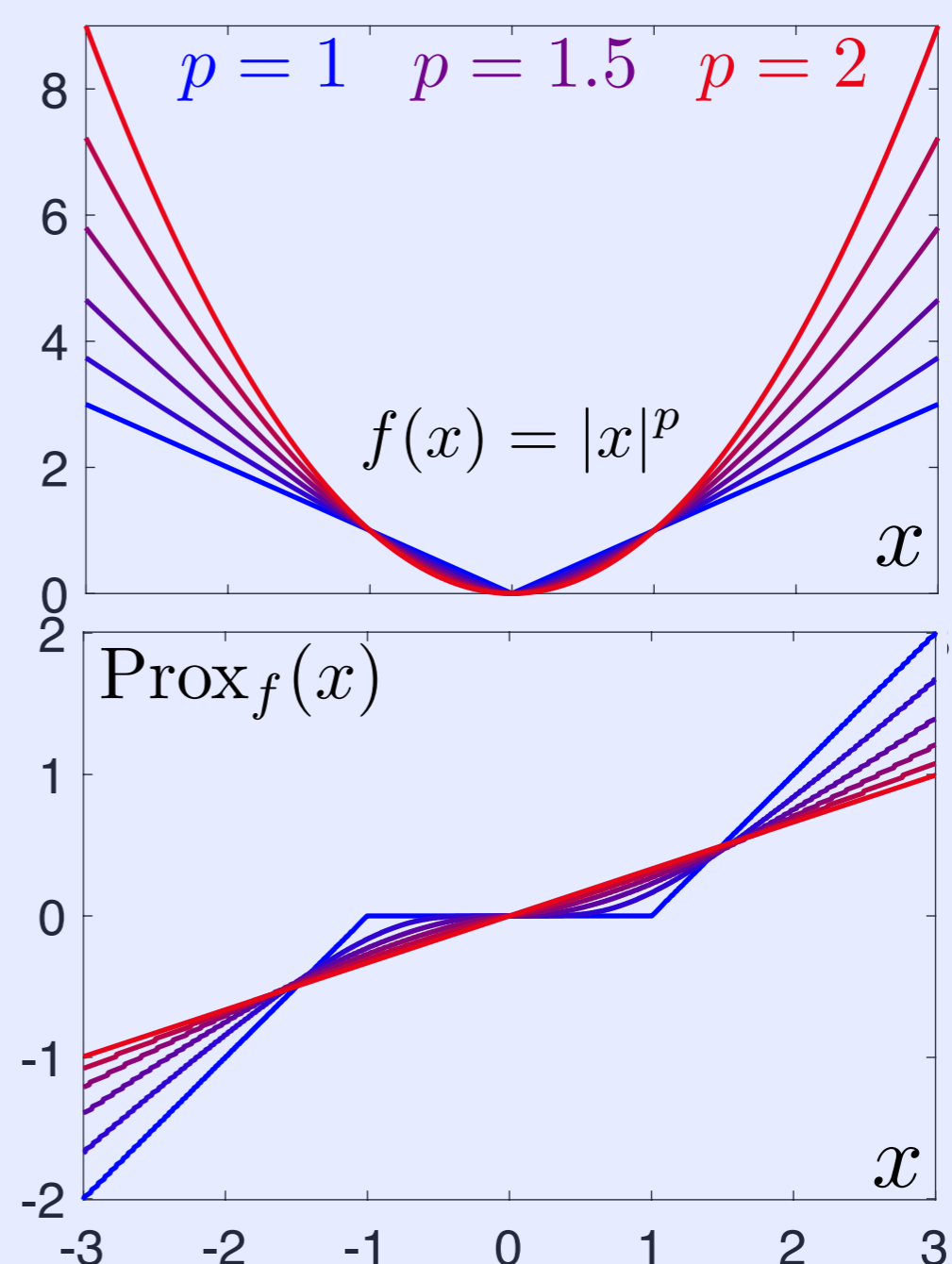
Yurii  
Nesterov

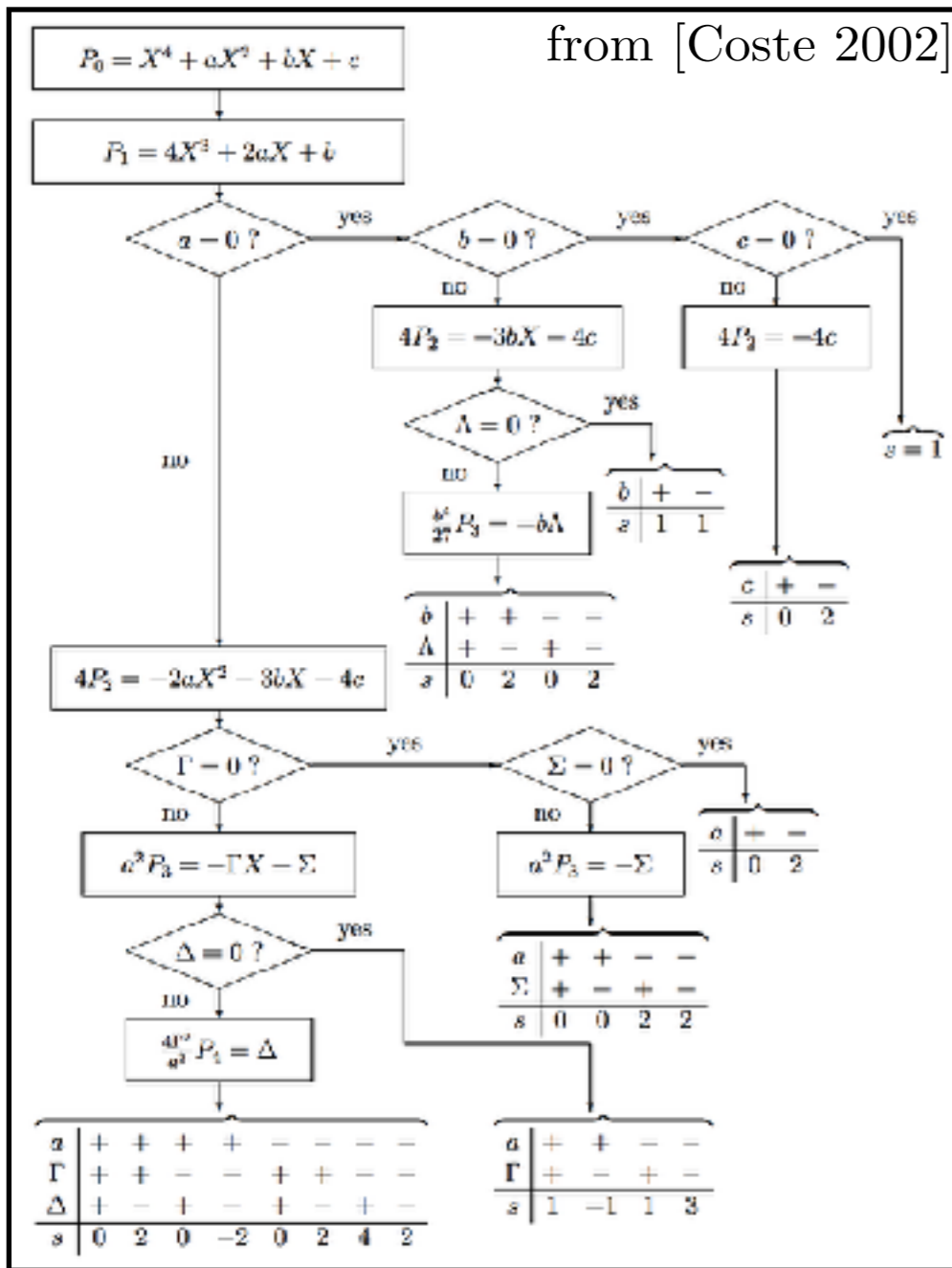
$$\text{Prox}_f(x) = \underset{x'}{\text{argmin}} \frac{1}{2} \|x - x'\|^2 + f(x')$$

Non-convex



Convex





algebraic set

$$\mathcal{X} \stackrel{\text{def.}}{=} \{(a, b, c, X) ; X^4 + aX^2 + bX + c = 0\}$$

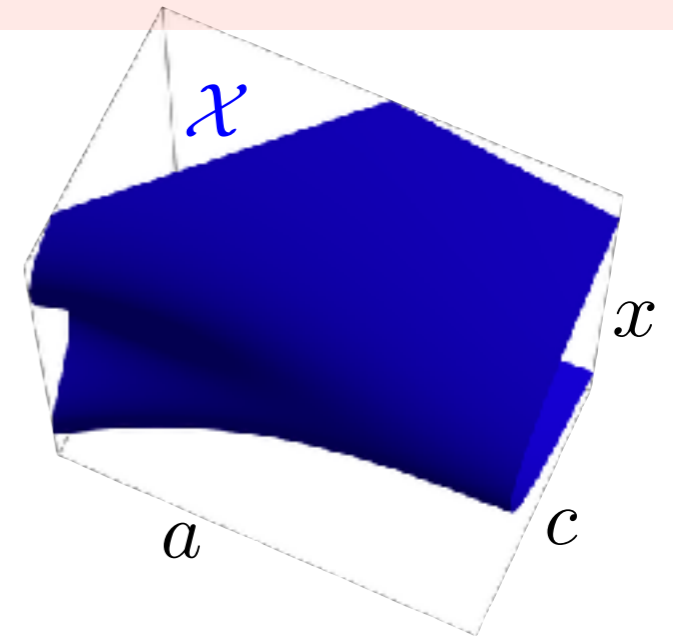
projection

$$(a, b, c, X) \mapsto (a, b, c)$$

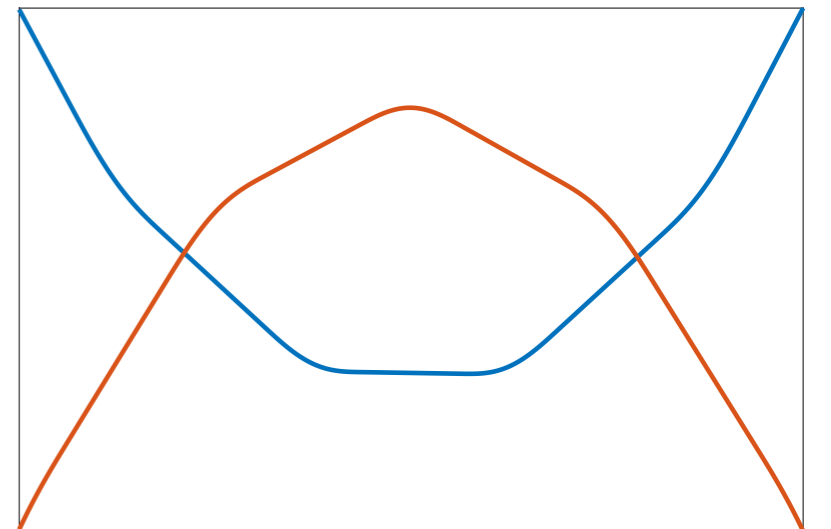
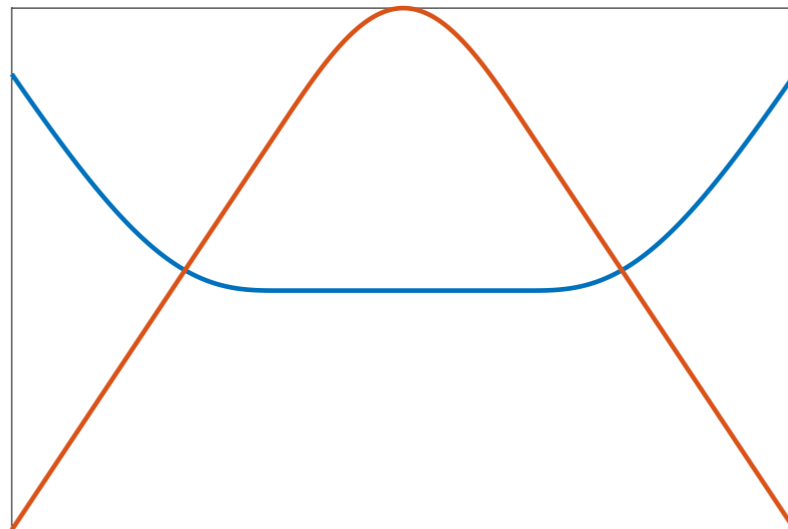
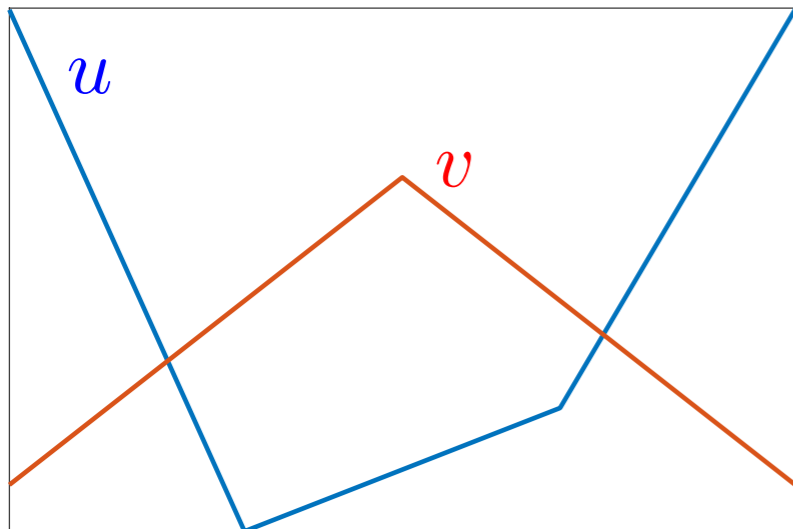
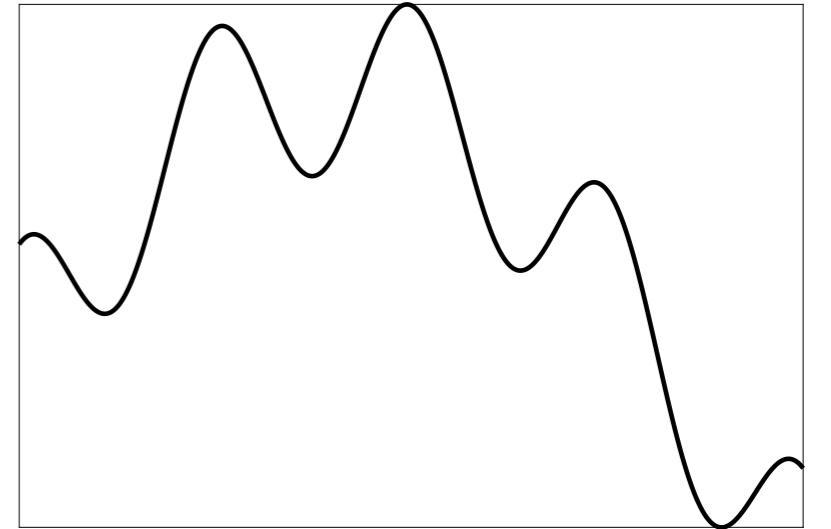
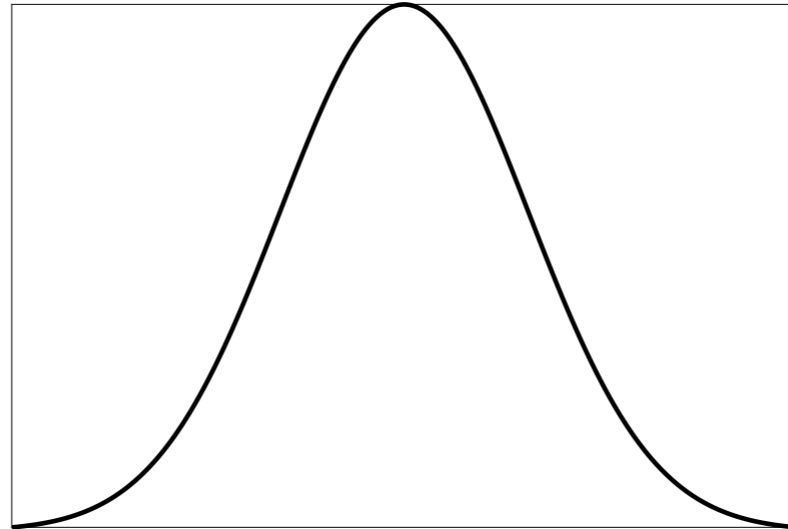
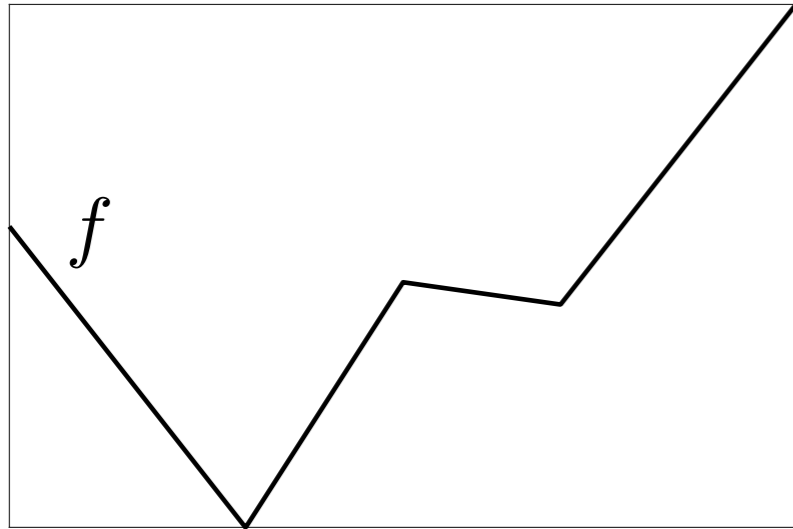
$$\{(a, b, c) ; \exists X \in \mathbb{R}, X^4 + aX^2 + bX + c = 0\}$$

semi-algebraic set

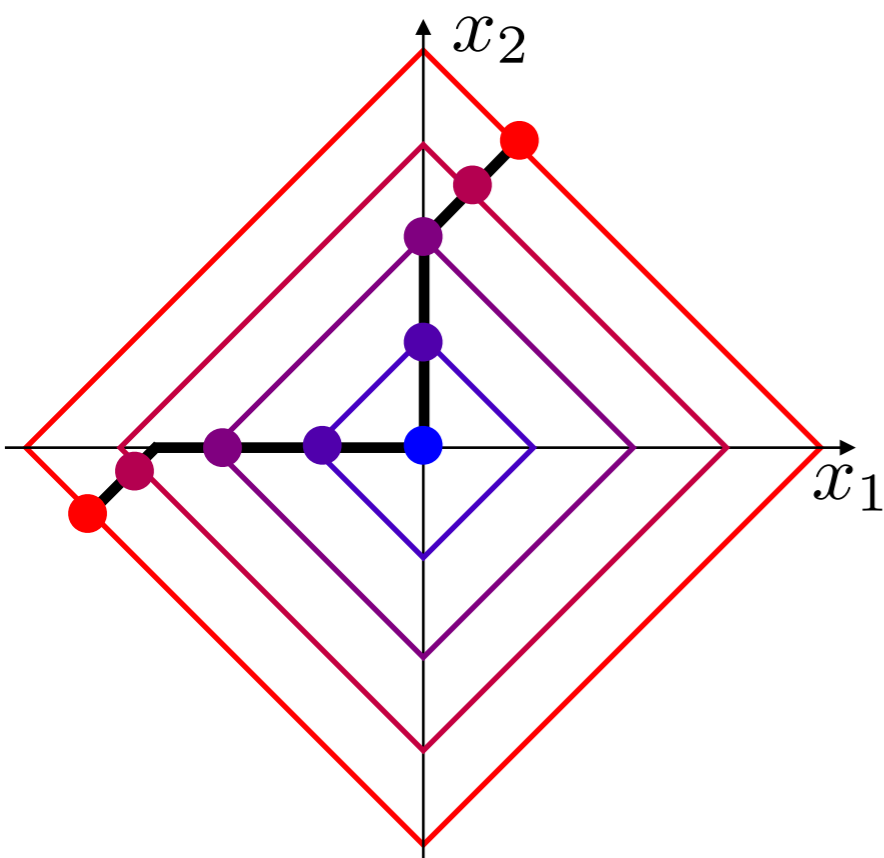
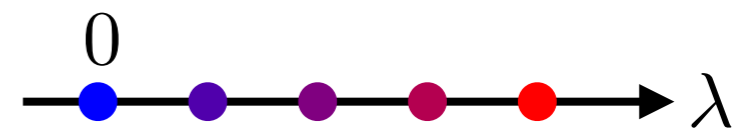
$$b = 0$$



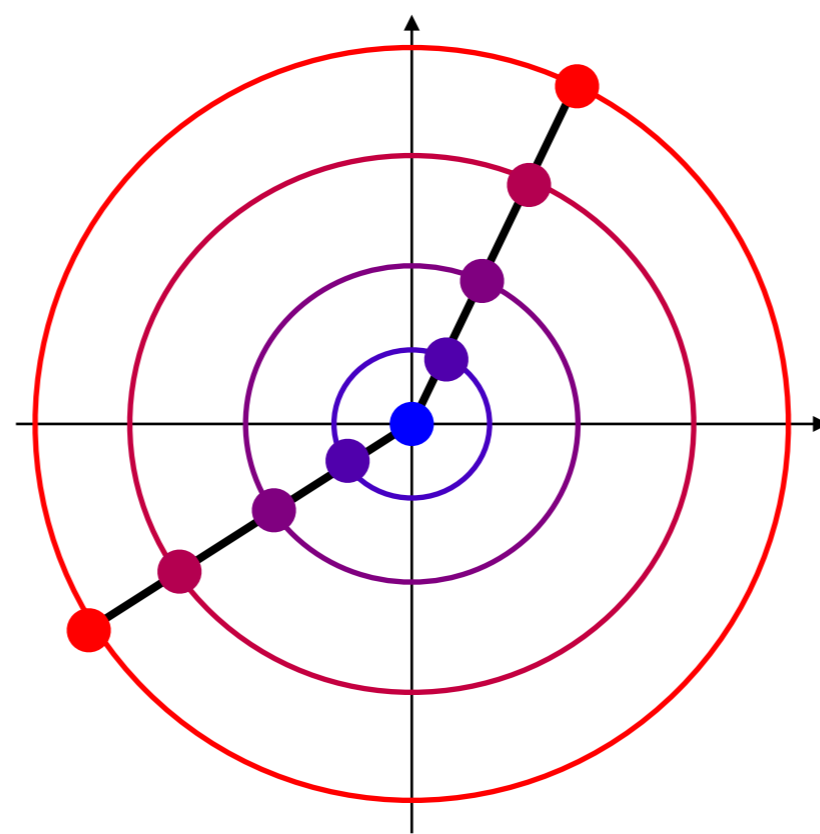
$$f = u + v = \text{convex} + \text{concave}$$



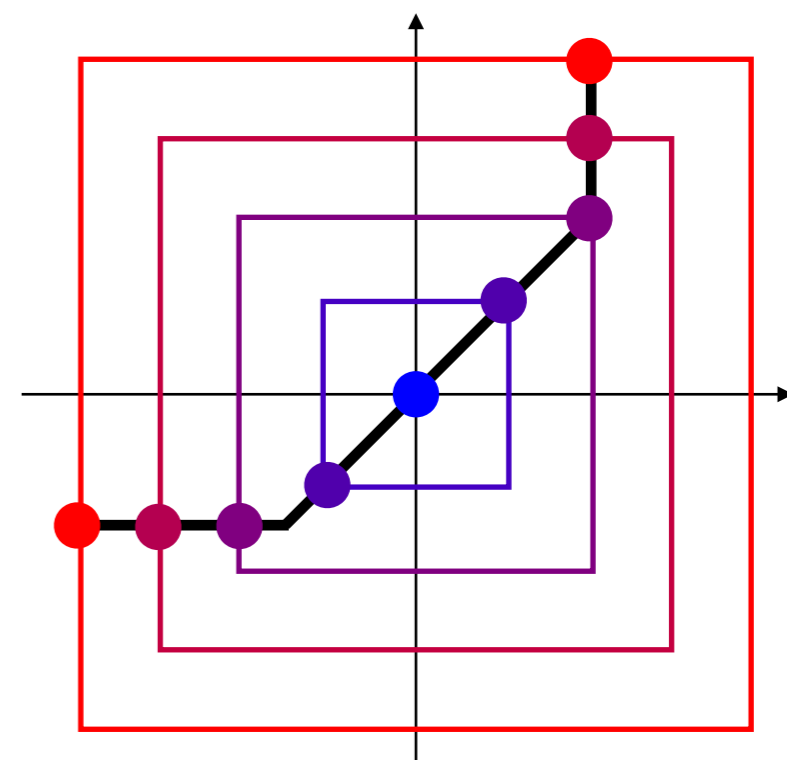
$$\text{Prox}_{\lambda f}(x) = \underset{x'}{\operatorname{argmin}} \frac{1}{2} \|x - x'\|^2 + \lambda f(x')$$



$$f(x) = |x_1| + |x_2|$$



$$f(x) = \sqrt{|x_1|^2 + |x_2|^2}$$



$$f(x) = \max(|x_1|, |x_2|)$$



John Toland

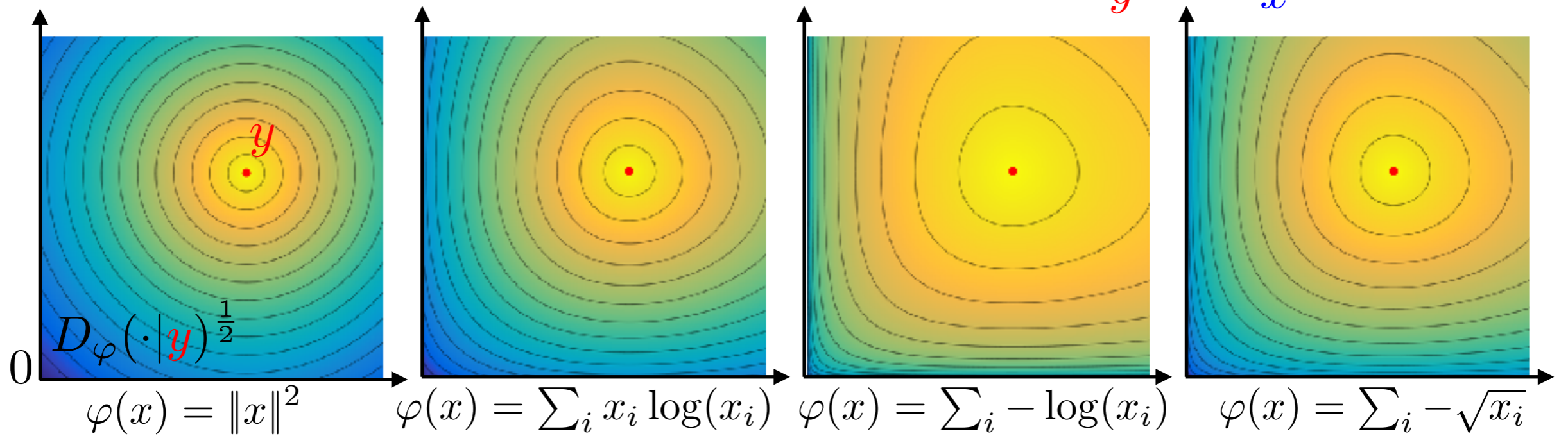
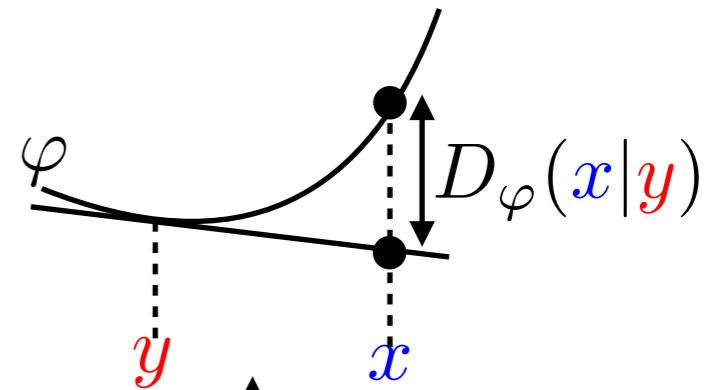
$(f, g)$  convex functions.

$$f^*(p) \stackrel{\text{def.}}{=} \sup_x \langle p, x \rangle - f(x)$$

*Toland's duality:*  $\inf f - g = \inf g^* - f^*$

Bregman divergence:

$$D_\varphi(x|y) \stackrel{\text{def.}}{=} \varphi(x) - \varphi(y) - \langle x - y, \nabla\varphi(y) \rangle$$

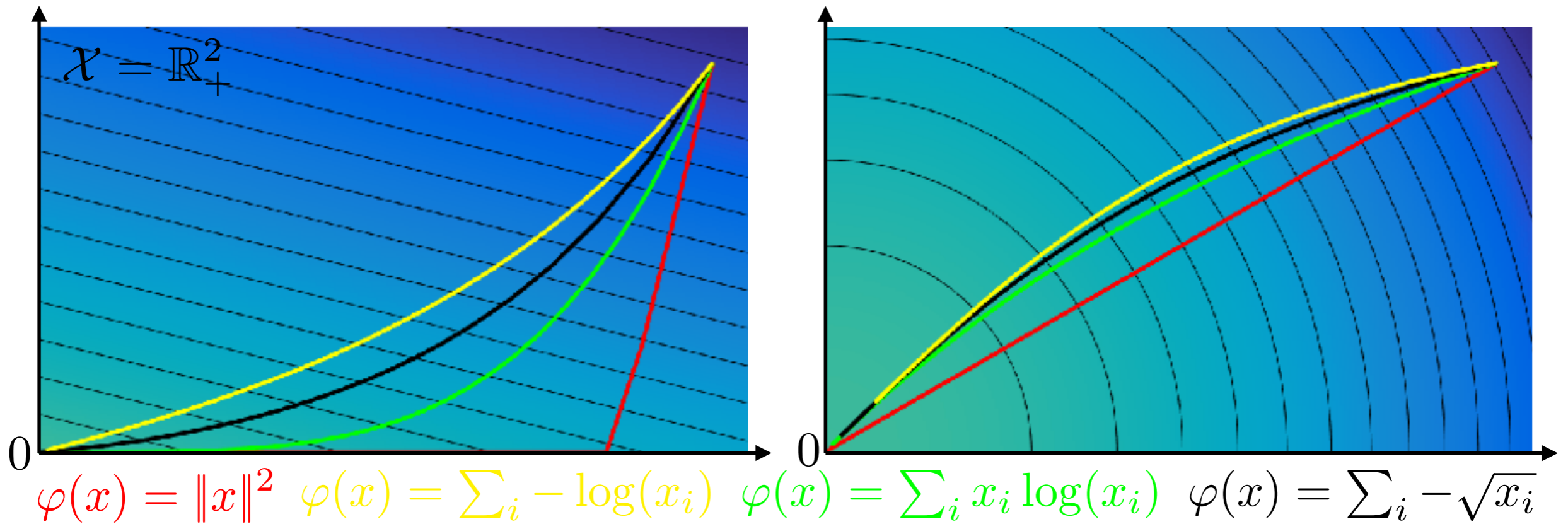


$$\left. \begin{array}{l} D_\varphi(x|x + \varepsilon) \\ D_\varphi(x + \varepsilon|x) \end{array} \right\} = \frac{1}{2} \langle \partial^2 \varphi(x) \varepsilon, \varepsilon \rangle + o(\|\varepsilon\|^2)$$



Bregman divergence:  $D_\varphi(x|y) \stackrel{\text{def.}}{=} \varphi(x) - \varphi(y) - \langle x - y, \nabla\varphi(y) \rangle$

Mirror descent:  $x_{k+1} = \operatorname{argmin}_{x \in \mathcal{X}} D_\varphi(x|x_k) + \tau \langle \nabla f(x_k), x \rangle$   
 $= (\nabla\varphi)^{-1} (\nabla\varphi(x_k) - \tau \nabla f(x_k))$



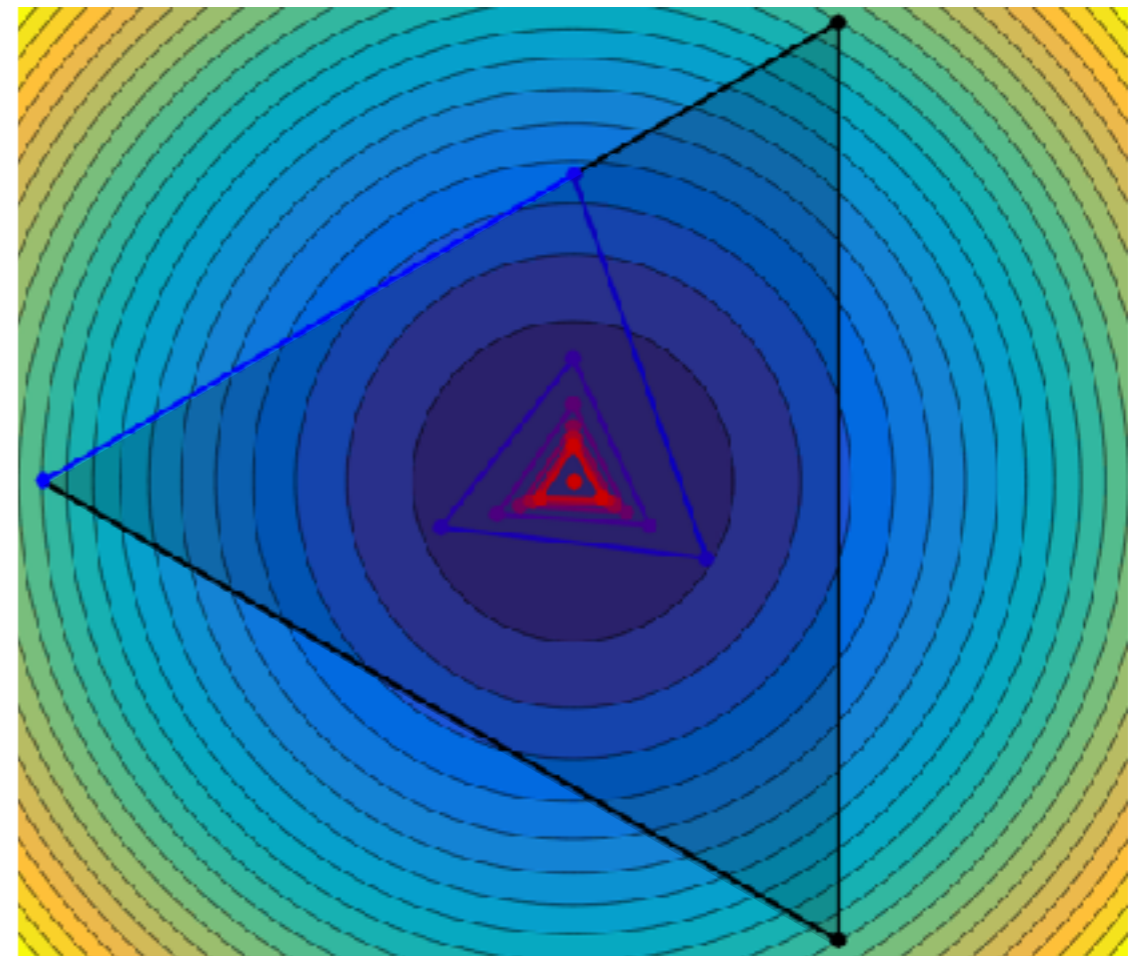
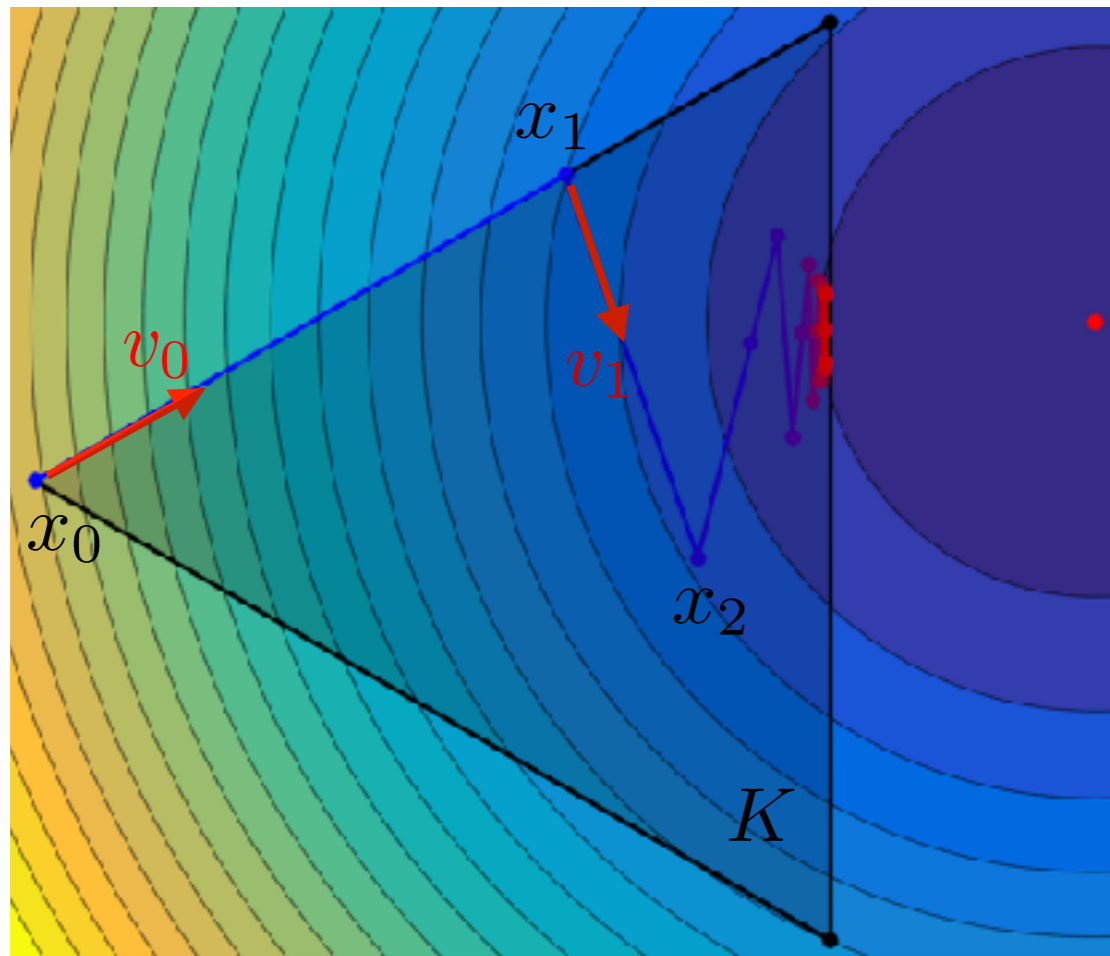
Directional derivative:  $D_v f(x) \stackrel{\text{def.}}{=} \lim_{t \rightarrow 0} \frac{f(x + th) - f(x)}{t}$

$$\min_{x \in K} f(x)$$

Frank-Wolfe

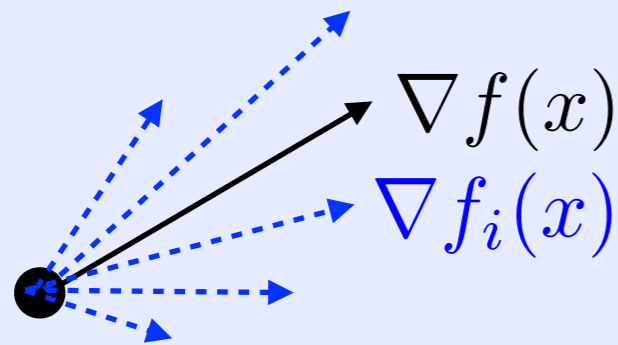
$$v_\ell \stackrel{\text{def.}}{=} \operatorname{argmin}_{v \in K} D_v f(x_\ell)$$

$$x_{\ell+1} \stackrel{\text{def.}}{=} x_\ell + \frac{2}{2 + \ell} (v_\ell - x_\ell)$$



## Finite sums

$$f(x) \stackrel{\text{def.}}{=} \frac{1}{n} \sum_{i=1}^n f_i(x)$$
$$\nabla f(x) = \frac{1}{n} \sum_i \nabla f_i(x)$$

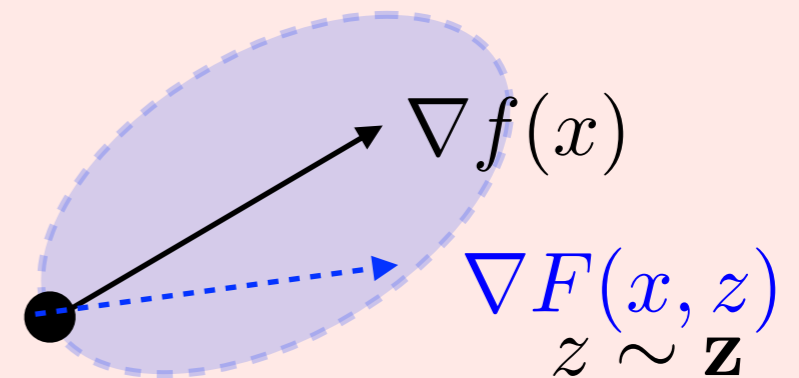


Draw  $i \in \{1, \dots, n\}$  uniformly.

$$x_{k+1} = x_k - \tau_k \nabla f_i(x_k)$$

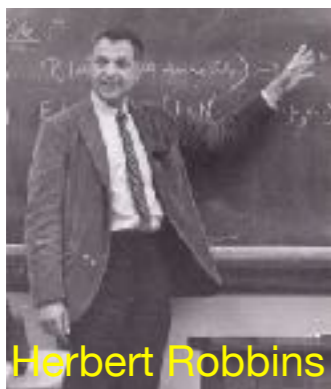
## Expectation

$$f(x) \stackrel{\text{def.}}{=} \mathbb{E}_{\mathbf{z}}(f(x, \mathbf{z}))$$
$$\nabla f(x) = \mathbb{E}_{\mathbf{z}}(\nabla F(x, \mathbf{z}))$$



Draw  $z \sim \mathbf{z}$

$$x_{k+1} = x_k - \tau_k \nabla F(x, z)$$

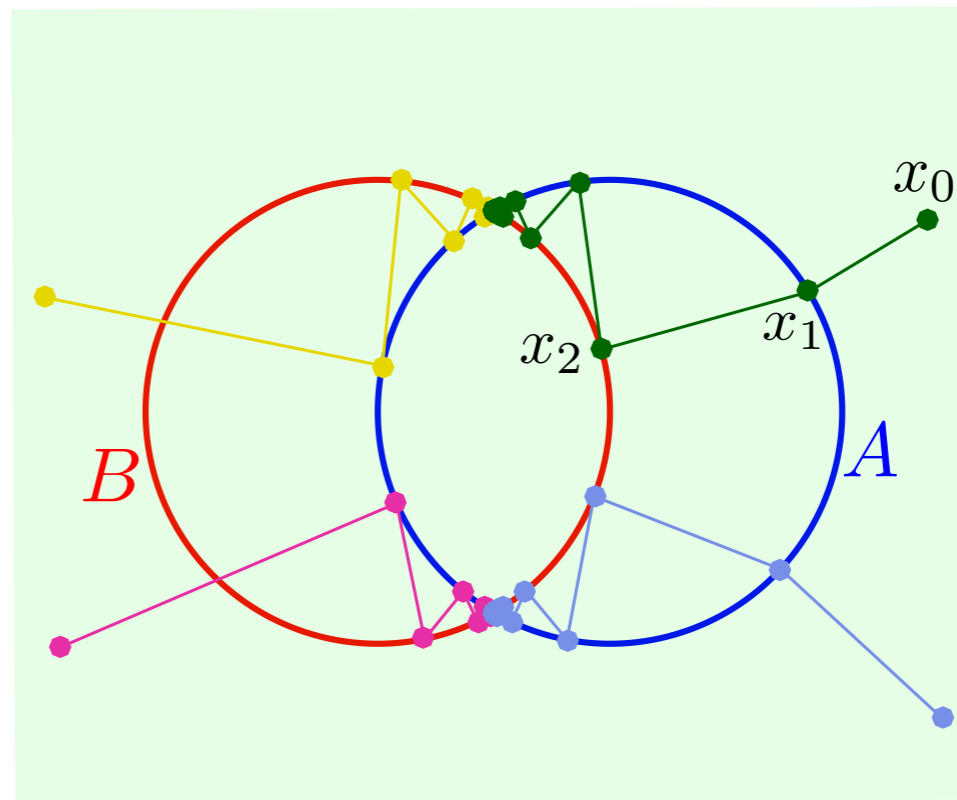
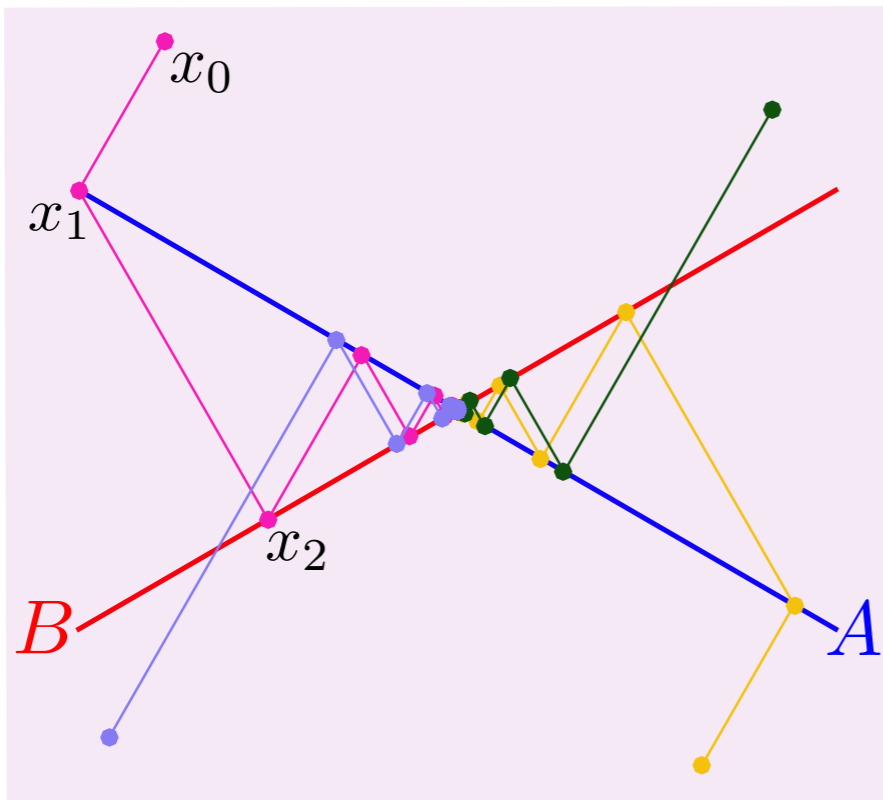


Herbert Robbins

*Theorem:* If  $f$  is strongly convex and  $\tau_k \sim 1/k$ ,  
 $\mathbb{E}(\|x_k - x^*\|^2) = O(1/k)$

Iterative projections: 
$$\begin{cases} x_{2k+1} = \text{Proj}_A(x_{2k}) \\ x_{2k+2} = \text{Proj}_B(x_{2k+1}) \end{cases}$$

*Theorem:* if  $(A, B)$  convex,  $x_k \xrightarrow{k \rightarrow +\infty} A \cap B$

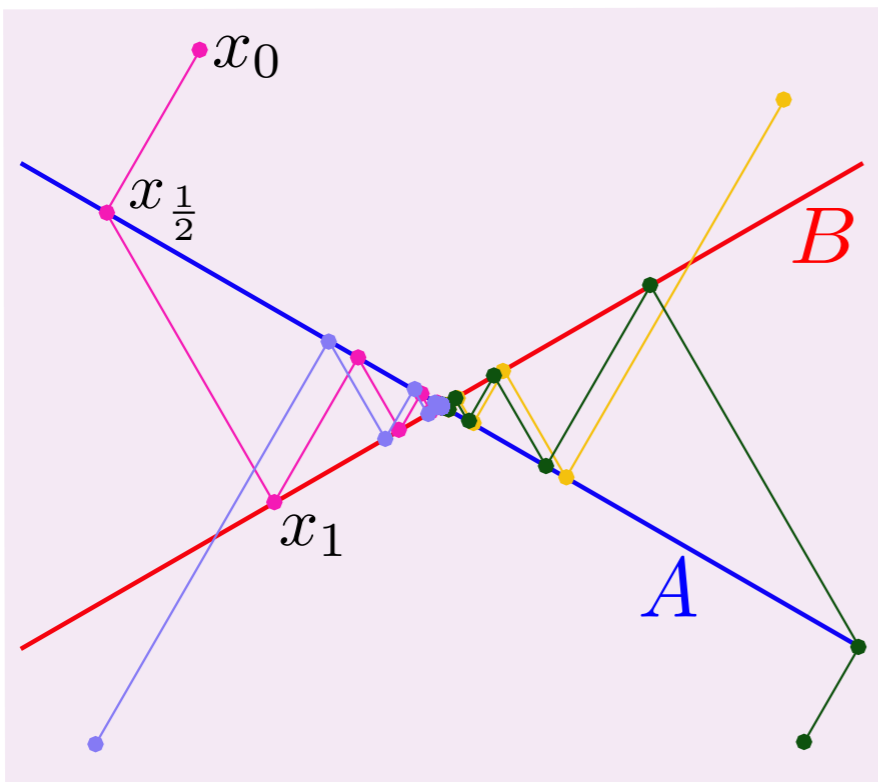


# Iterative Projections

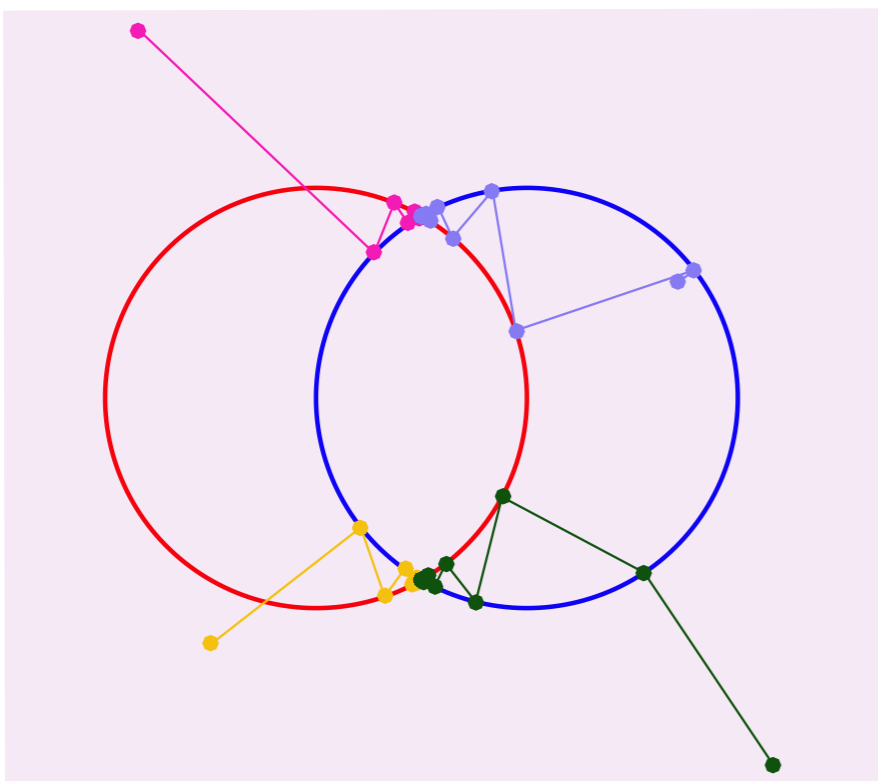
$$x_{k+1} = P_B(P_A x_k)$$

$$P_A \stackrel{\text{def.}}{=} \text{Proj}_A$$

Convex



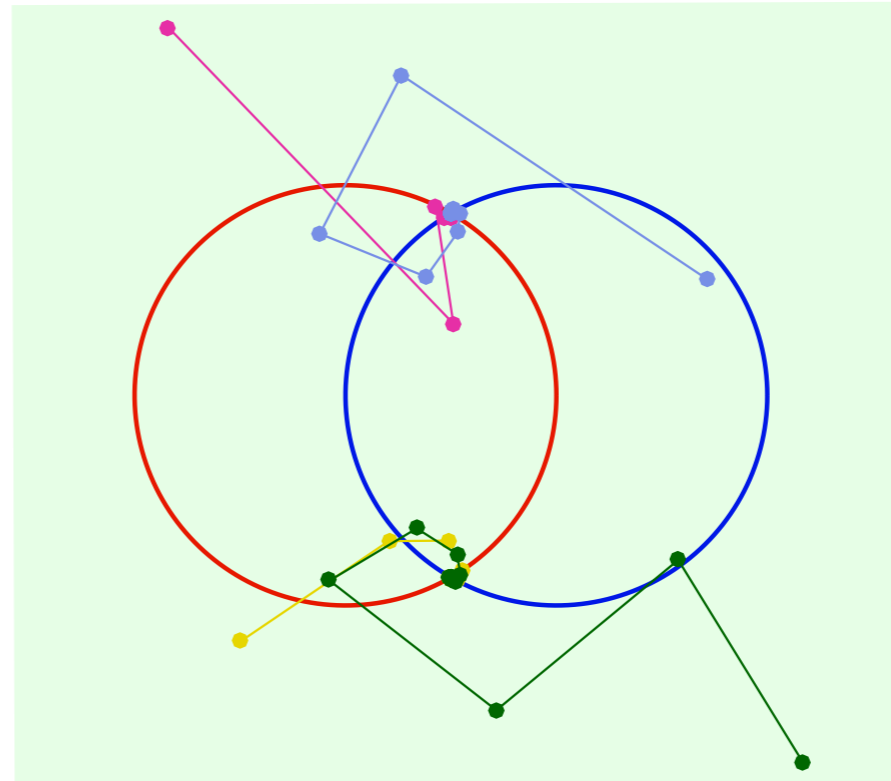
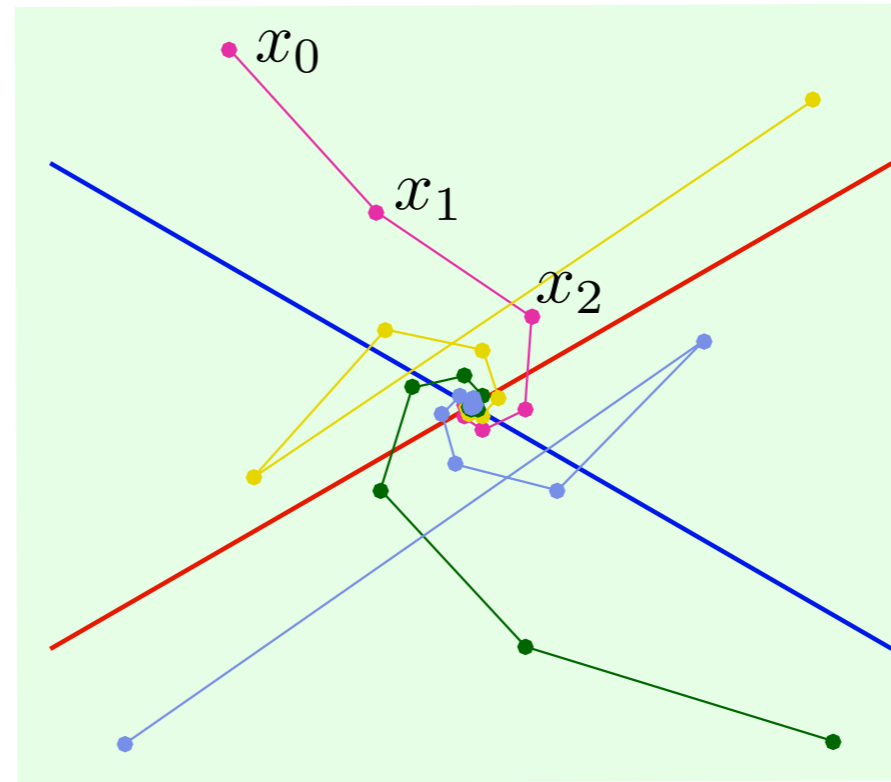
Non-convex



# Douglas-Rachford

$$x_k = \bar{P}_A(y_k) \stackrel{\text{def.}}{=} 2P_A(y_k) - y_k$$

$$y_{k+1} = \frac{1}{2}y_k + \frac{1}{2}\bar{P}_B(x_k)$$



Jim Douglas



Henry Rachford



Pierre-Louis Lions

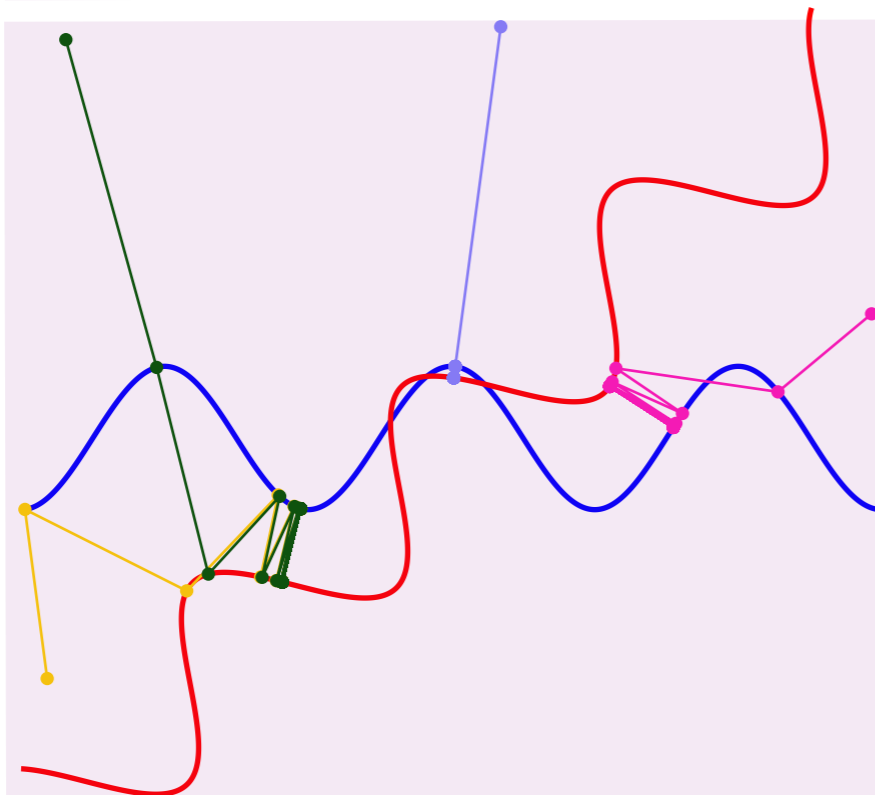


Bertrand Mercier

## Iterative Projections

$$x_{k+1} = P_B(P_A x_k)$$

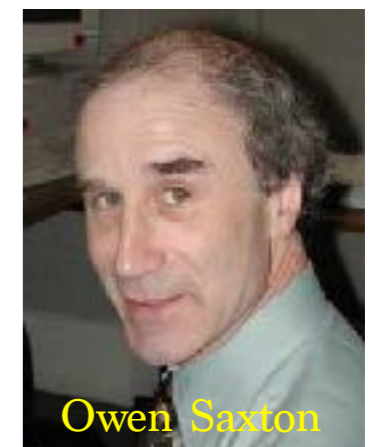
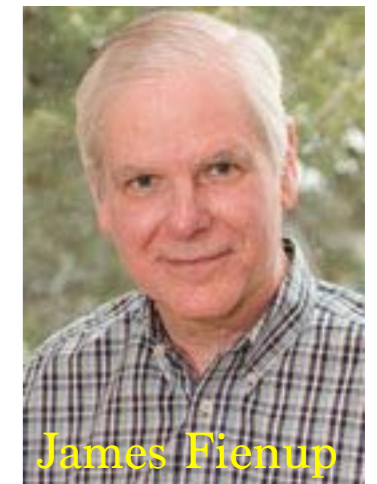
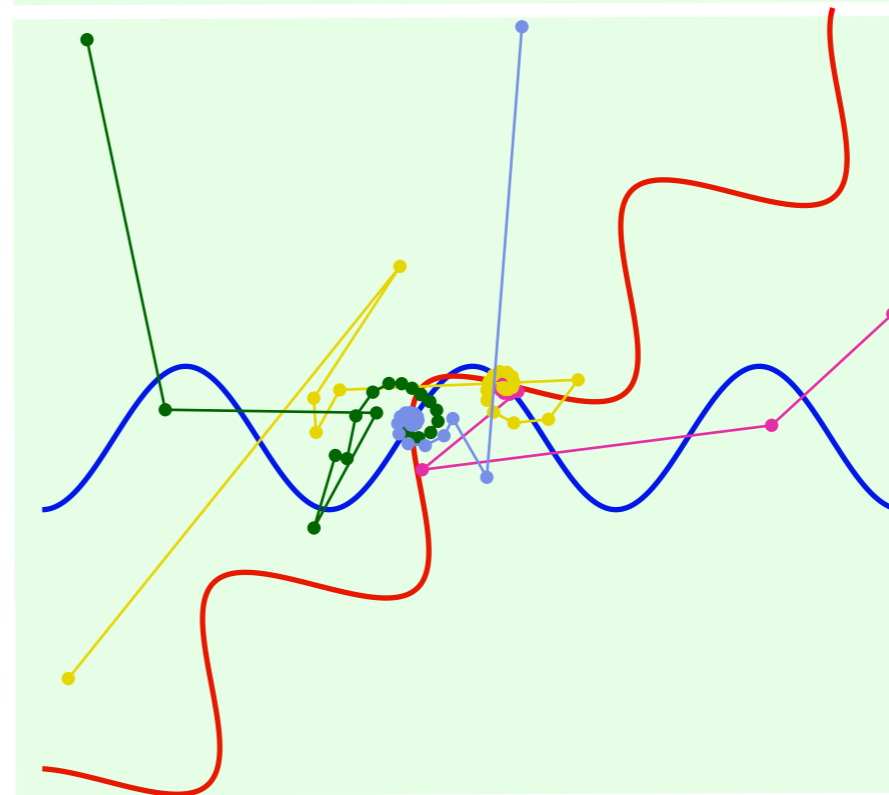
$$P_A \stackrel{\text{def.}}{=} \text{Proj}_A$$



## Douglas-Rachford

$$x_k = \bar{P}_A(y_k) \stackrel{\text{def.}}{=} 2P_A(y_k) - y_k$$

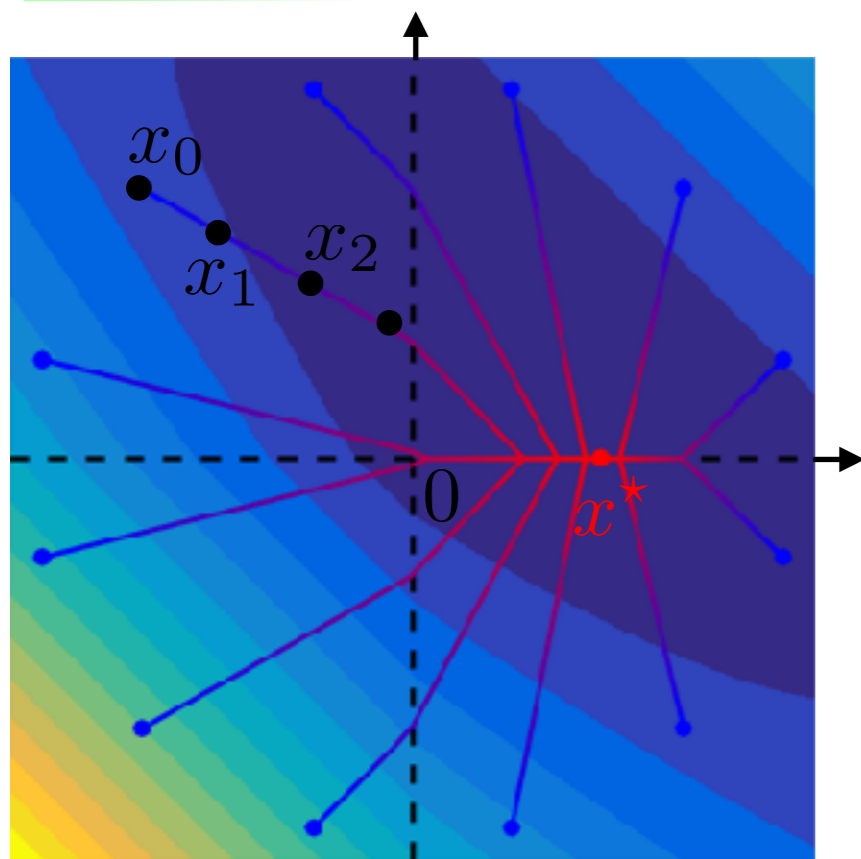
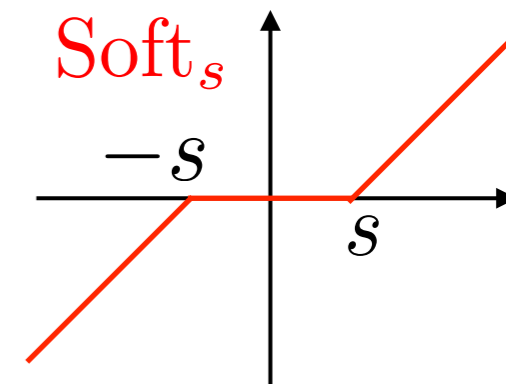
$$y_{k+1} = \frac{1}{2}y_k + \frac{1}{2}\bar{P}_B(x_k)$$



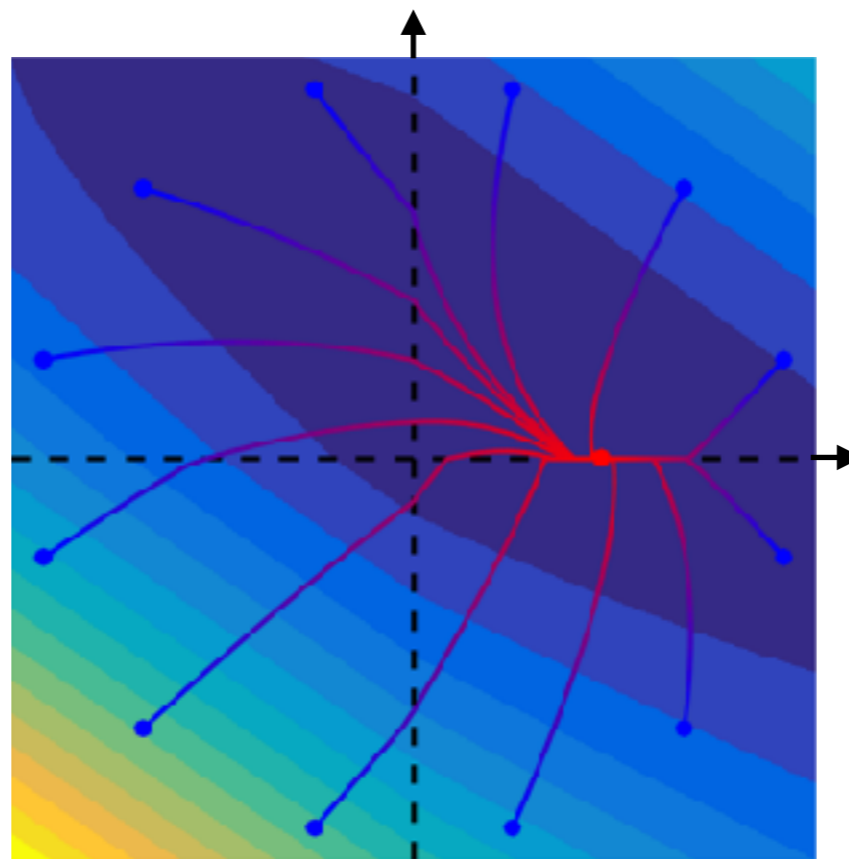
Lasso: 
$$\min_x \frac{1}{2} \|y - Ax\|^2 + \lambda \|x\|_1$$

Fwd-Bwd (ISTA): 
$$x_{k+1} \stackrel{\text{def.}}{=} \text{Soft}_{\tau\lambda}(x_k - \tau A^\top (Ax_k - y))$$

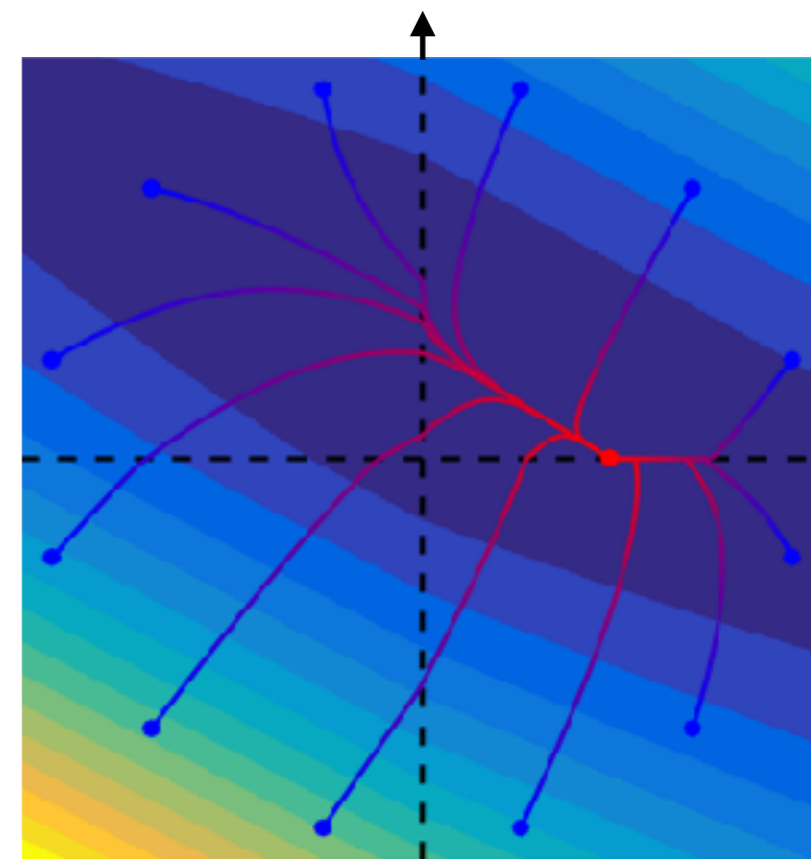
*Theorem:* if  $0 < \tau < 2/\|A\|^2$ ,  $x_k \rightarrow x^*$  solution of Lasso.



$$A = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad y = A \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad \lambda = 0.3$$



$$A = \begin{pmatrix} 1 & 0 \\ 0.5 & 1 \end{pmatrix}$$



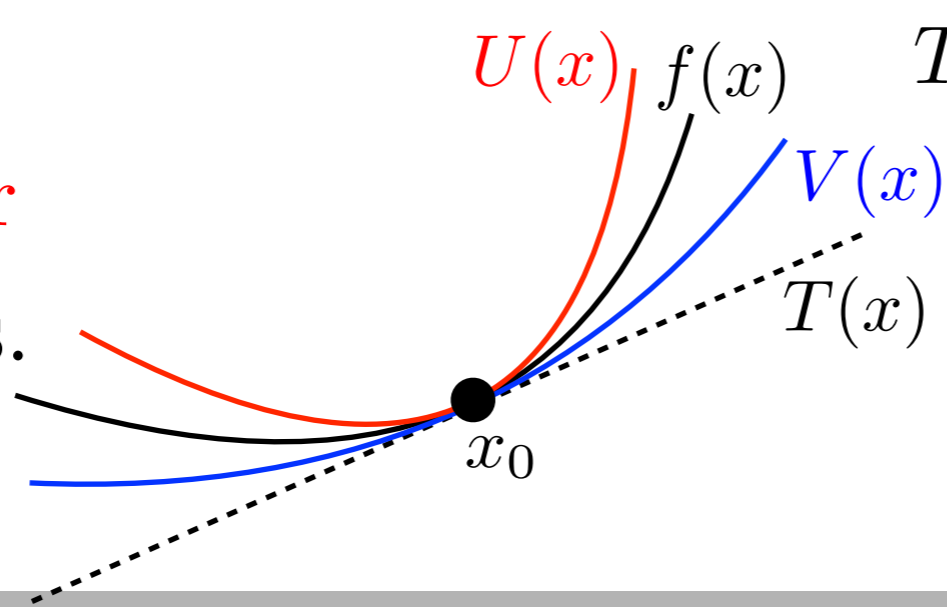
$$A = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}$$

Hypotheses:  $\mu \text{Id}_n \preceq \partial^2 f(x) \preceq L \text{Id}_n$   
 strong convexity                      smoothness

Conditionning:

$$\varepsilon \stackrel{\text{def.}}{=} \frac{\mu}{L} \leq 1$$

Quadratic  
 lower / upper  
 approximants.



$$T(x) \stackrel{\text{def.}}{=} f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle$$

$$U(x) \stackrel{\text{def.}}{=} T(x) + \frac{L}{2} \|x - x_0\|^2$$

$$V(x) \stackrel{\text{def.}}{=} T(x) + \frac{\mu}{2} \|x - x_0\|^2$$

Gradient descent:  $x_{k+1} = x_k - \tau_k \nabla f(x_k)$

*Theorem:*

If  $L < +\infty$ ,  $0 < \tau < \frac{2}{L}$

$$f(x_k) - f(x^*) \leq \frac{C}{k+1}$$

If  $\mu > 0$ ,  $L < +\infty$ ,  $0 < \tau < \frac{2}{L}$

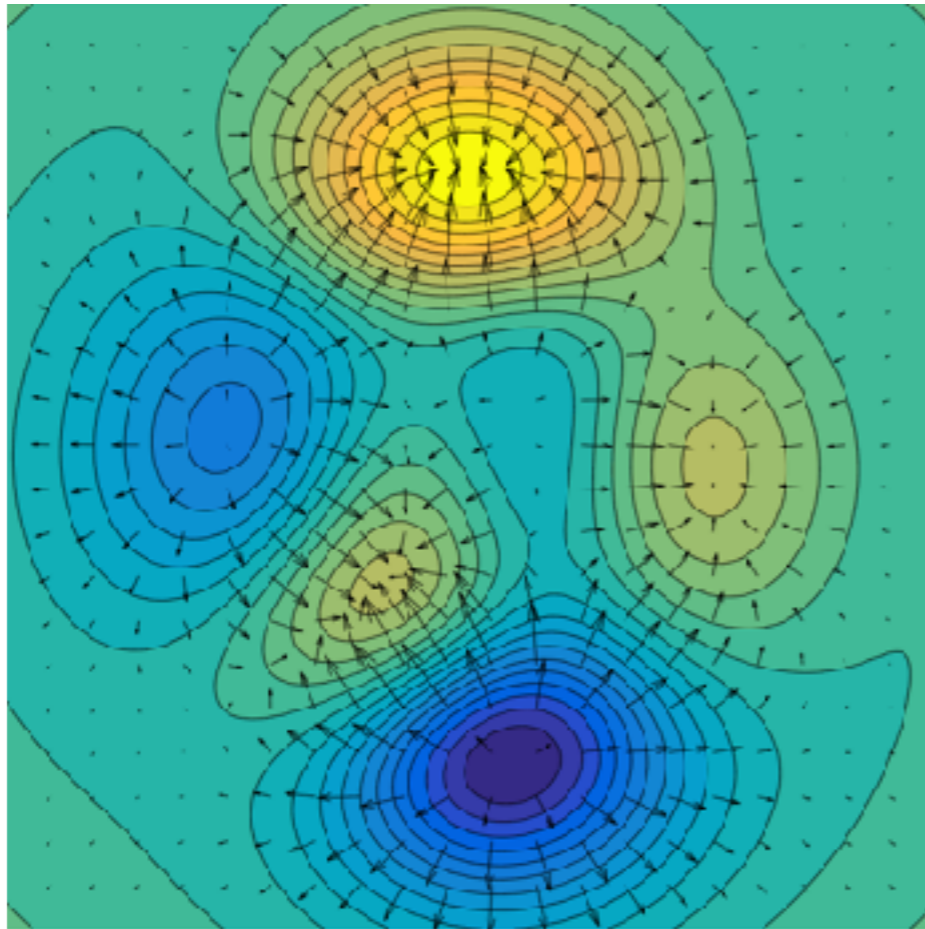
$$\|x_k - x^*\| \leq \rho^k \|x_0 - x^*\|$$

$$\rho = (1 + \varepsilon)^{-\frac{1}{2}} < 1$$

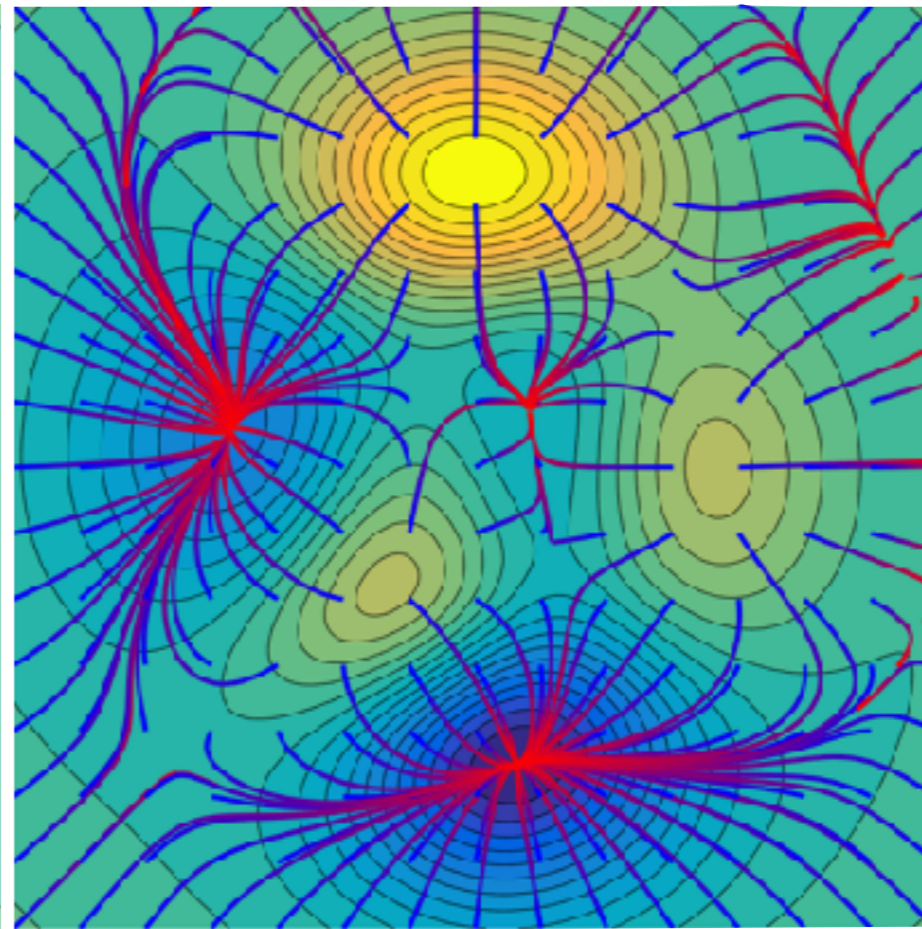


Gradient field:  $f(x + \varepsilon) = f(x) + \varepsilon \langle \varepsilon, \nabla f(x) \rangle + o(\varepsilon)$

Gradient flow:  $x'(t) = -\nabla f(x(t))$



Gradient field  $\nabla f$



Gradient flows  $x(t)$   
 $t = 0$   $t$  medium  $t$  large

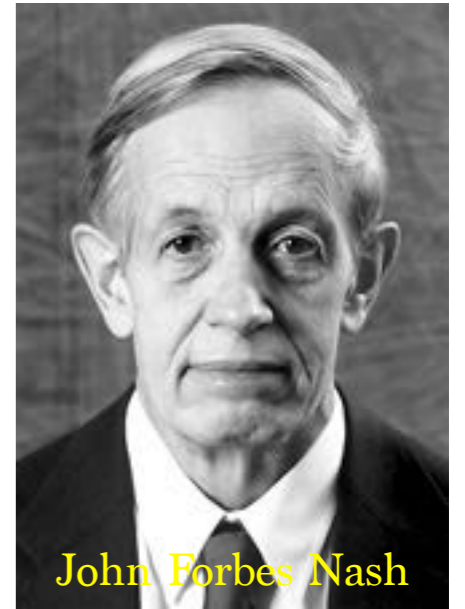
Min-max game:  $\min_x \max_y f(x, y) \geq \max_y \min_x f(x, y)$

convex concave  
↑ ↑

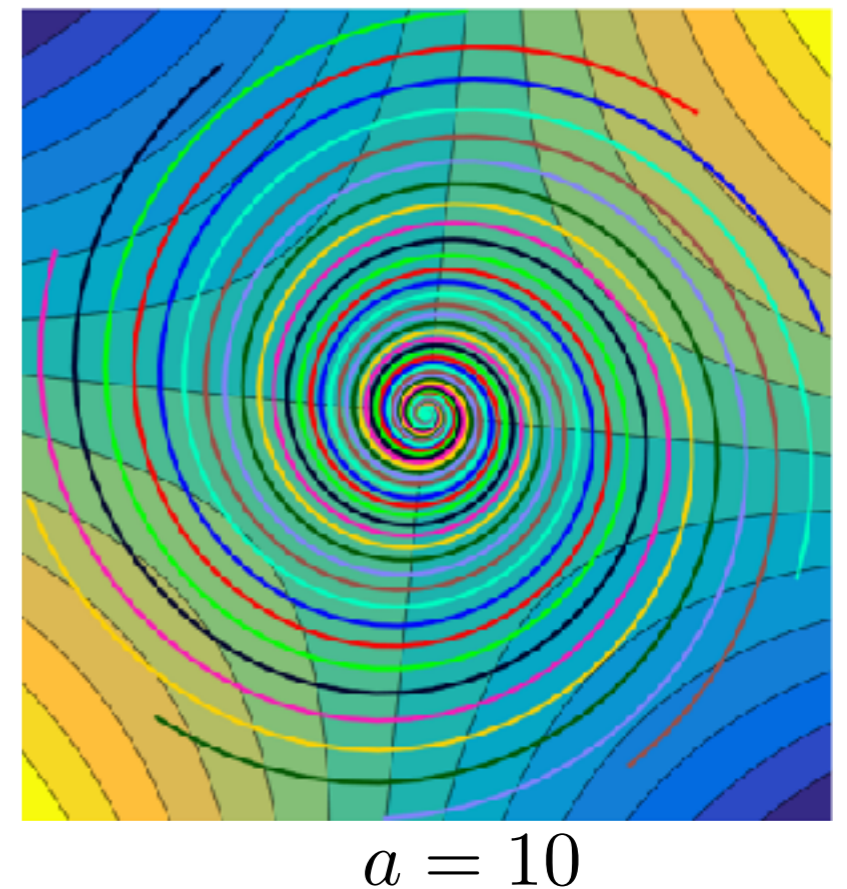
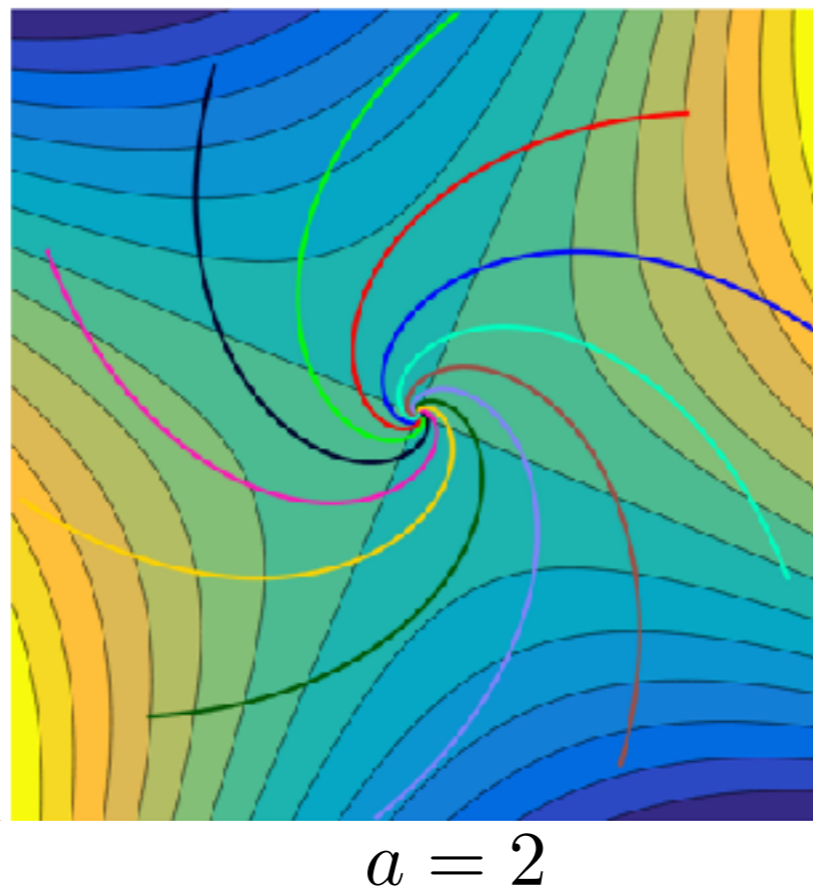
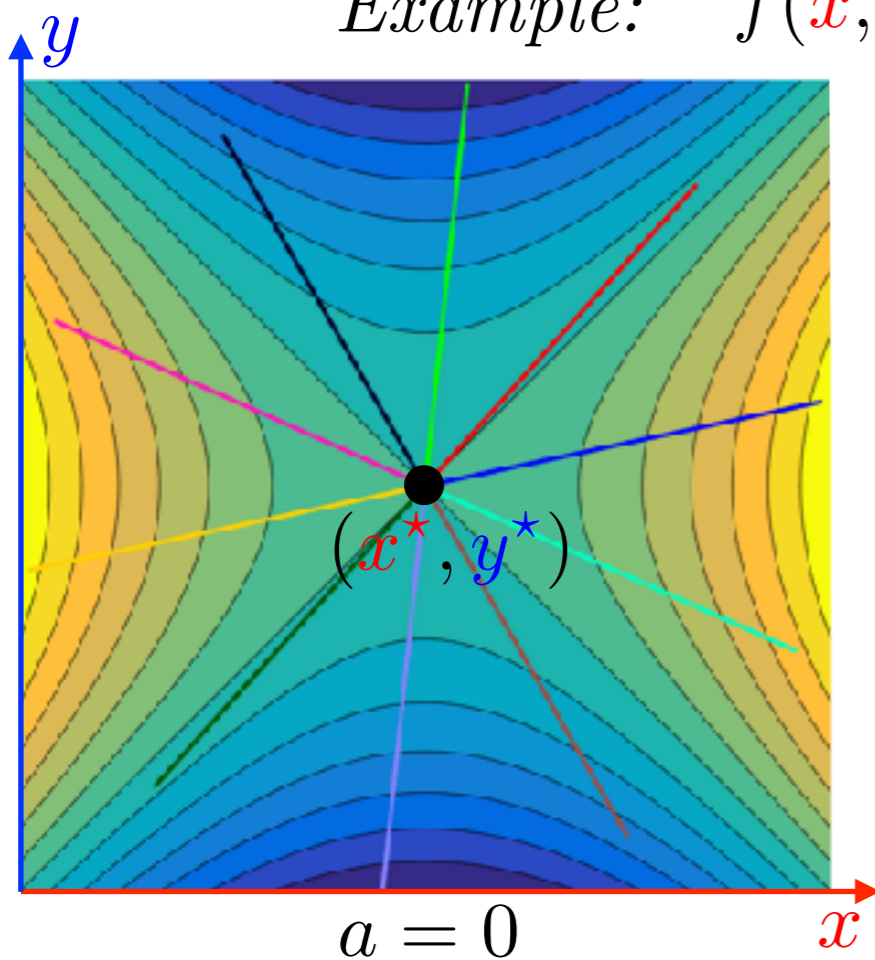
Saddle point  $(x^*, y^*)$ :  $f(x^*, y) \leq f(x^*, y^*) \leq f(x, y^*)$

→ Strong duality:  $\min_x \max_y f(x, y) = \max_y \min_x f(x, y)$

Gradient descent: 
$$\begin{cases} x_{k+1} = x_k - \tau \nabla_x f(x_k, y_k) \\ y_{k+1} = y_k + \tau \nabla_y f(x_k, y_k) \end{cases}$$



Example:  $f(x, y) = a xy + x^2 - y^2$        $a = \text{interaction}$



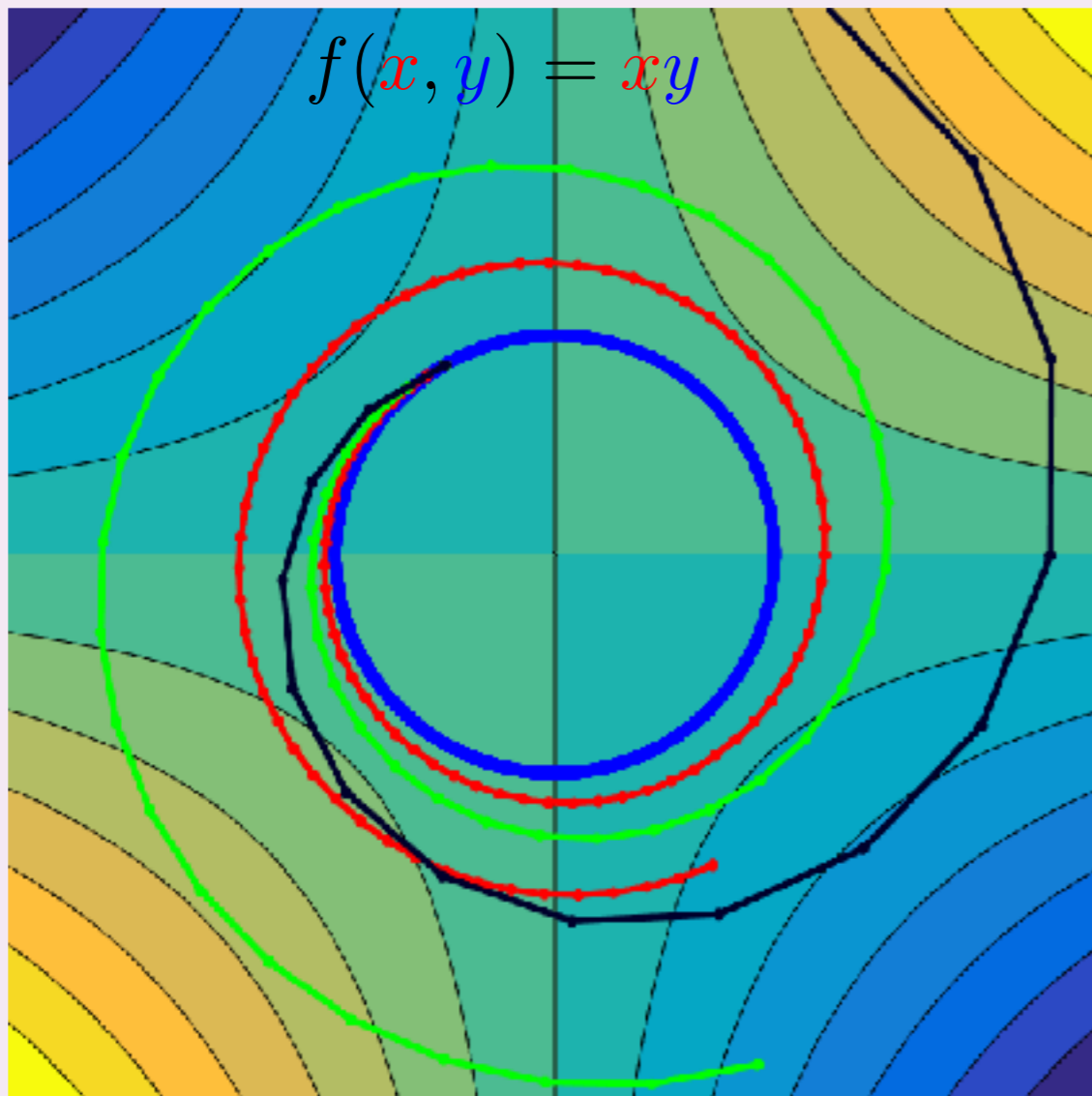


Min-max game:  $\min_x \max_y f(x, y) \geq \max_y \min_x f(x, y)$



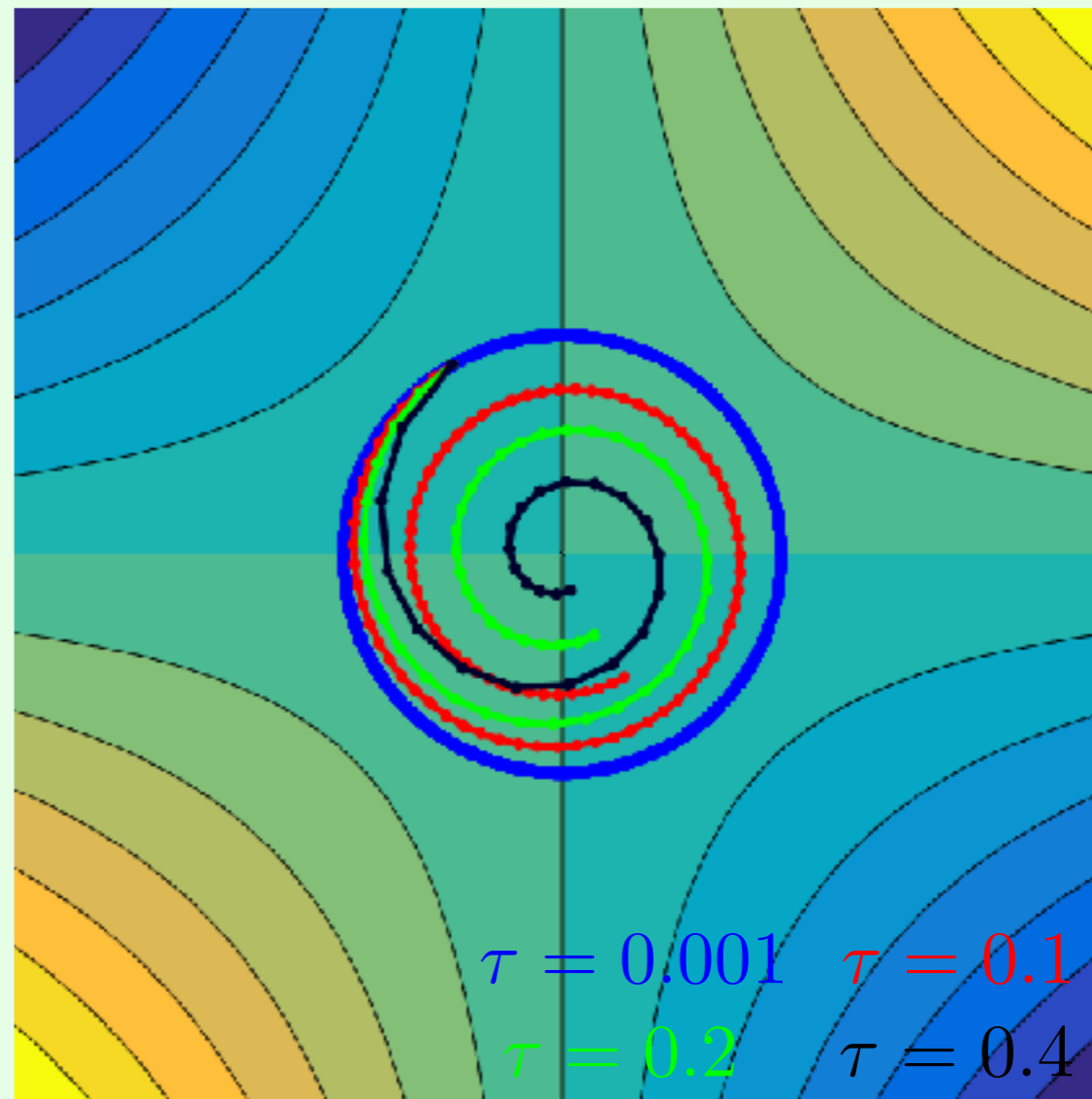
## Explicit

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



## Implicit

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



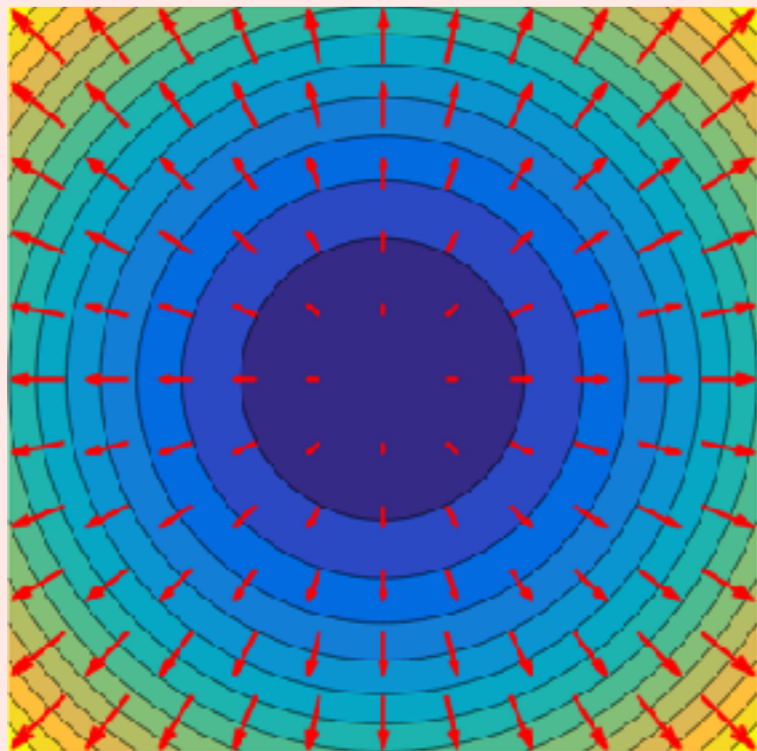
Monotone operator  $V : \mathbb{R}^d \rightarrow \mathbb{R}^d$ :  $\forall (x, y) \in (\mathbb{R}^d)^2, \langle V(x) - V(y), x - y \rangle \geq 0$

### Convex optimization

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

*Theorem:*  $\nabla f$  is monotone.

$f$  quadratic:  $\nabla f$  symmetric.



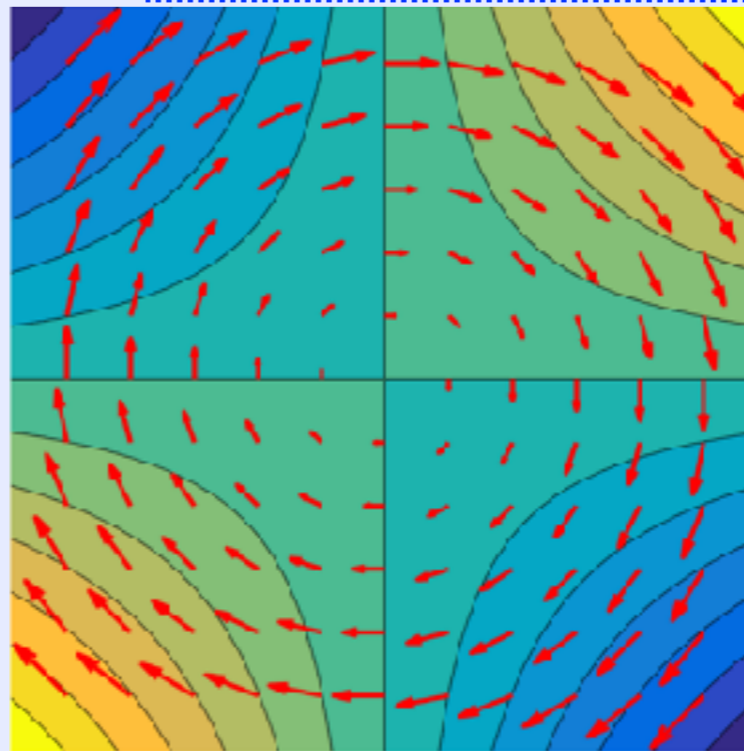
$$f(x) = \|x\|^2/2 \quad V(x) = x$$

### Saddle point

$$\min_{x_1 \in \mathbb{R}^{d_1}} \max_{x_2 \in \mathbb{R}^{d_2}} f(x_1) + \langle Ax_1, x_2 \rangle - g(x_2)$$

*Theorem:*  $\begin{pmatrix} \nabla f & A^* \\ -A & \nabla g \end{pmatrix}$  is monotone.

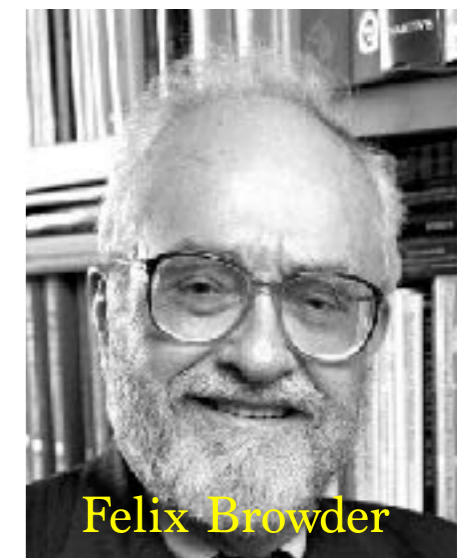
$f = g = 0$ :  $\begin{pmatrix} 0 & A^* \\ -A & 0 \end{pmatrix}$  skew-symmetric.



$$A = 1 \quad V(x) = (x_2, -x_1)^T$$



George Minty



Felix Browder

The two “extremal” cases [Edgar Asplund, 1970]

Legendre-Fenchel transform:

$$f^*(u) \stackrel{\text{def.}}{=} \min_x \langle x, u \rangle - f(x)$$

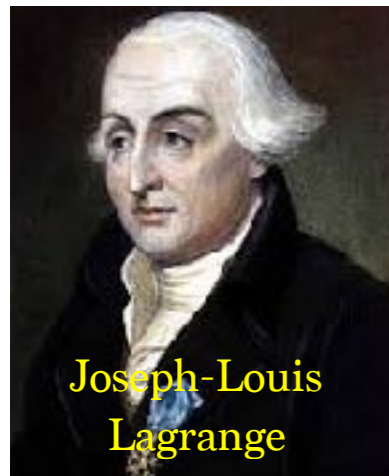


Adrien-Marie Legendre



Werner Fenchel

$$\min_x f(Ax) + g(x) = \min_{Ax=y} f(y) + g(x)$$



Joseph-Louis Lagrange

$$\rightarrow = \min_{x,y} \max_u f(y) + g(x) + \langle Ax - y, u \rangle$$



Ralph Tyrrell Rockafellar

$$= \max_u \left[ \min_y \langle -y, u \rangle + f(y) \right] + \left[ \min_x \langle x, A^*u \rangle + g(x) \right]$$

$$= \max_u -f^*(u) - g^*(-A^*u)$$

Primal-dual relations:

$$g \text{ smooth} \\ \nabla g(x) = -A^*u$$

$\iff$

$$g \text{ strongly convex} \\ x = \nabla g^*(-A^*u)$$

Moreau-Yosida regularization:

$$f_\mu(x) \stackrel{\text{def.}}{=} \min_y f(y) + \frac{1}{2\mu} \|x - y\|^2$$

Proximal operator:

$$\text{Prox}_{\mu f} \stackrel{\text{def.}}{=} \underset{y}{\text{argmin}} f(y) + \frac{1}{2\mu} \|x - y\|^2$$

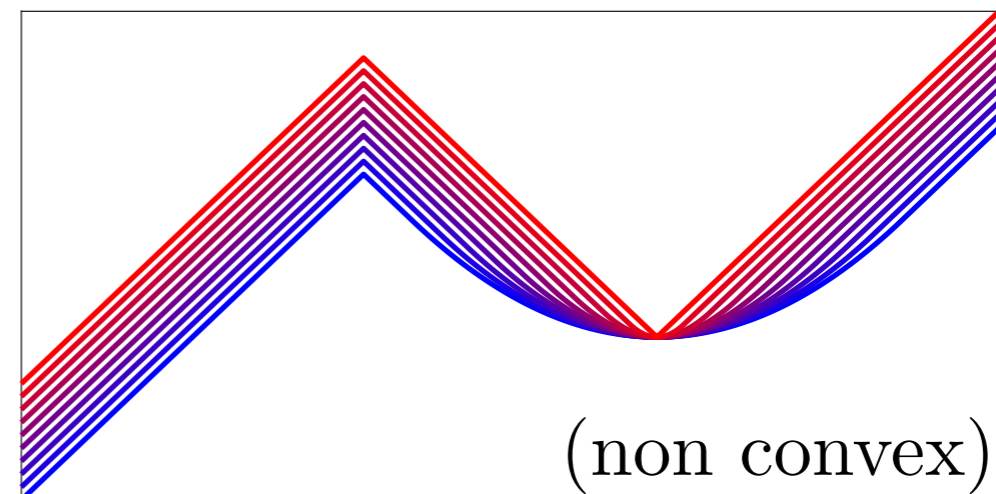
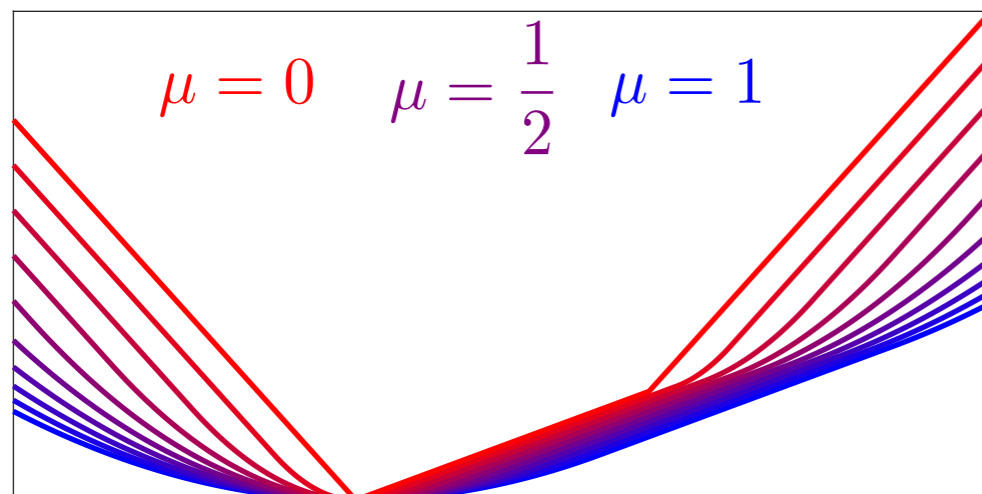
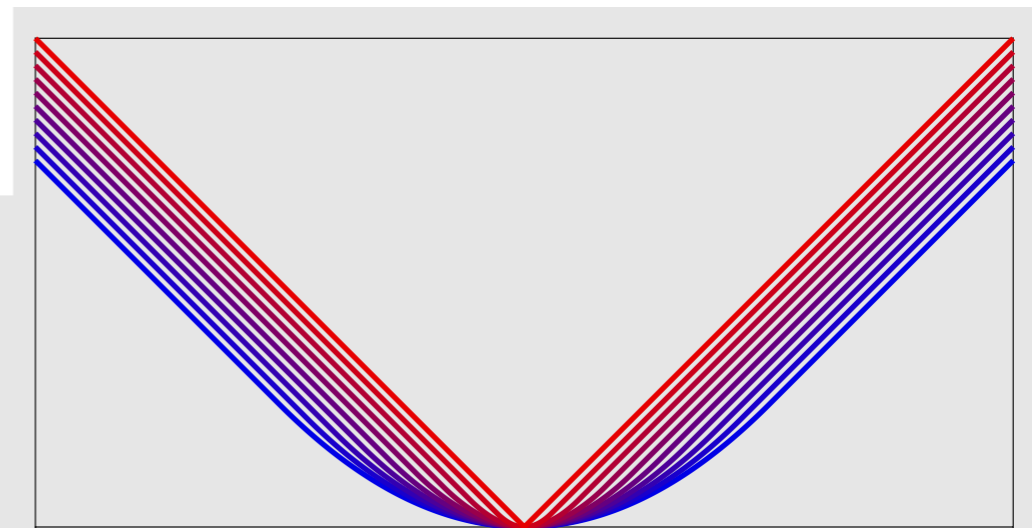


*Prop:*  $\nabla f_\mu$  is  $1/\mu$ -Lipschitz and

$$\mu \nabla f_\mu(x) = x - \text{Prox}_{\mu f}(x)$$

*Huber function:*  $f(x) = |x|$

$$f_\mu(x) = \begin{cases} x^2/(2\mu) & \text{if } |x| \leq \mu \\ |x| - \mu/2 & \text{otherwise.} \end{cases}$$



Fenchel-Legendre transform:

$$f^*(y) = \sup_x \langle x, y \rangle - f(x)$$

Polar of a set:

$$C^\circ = \{y ; \forall x \in C, \langle x, y \rangle \leq 1\}$$

Indicator:

$$\iota_C(x) = \begin{cases} 0 & \text{if } x \in C, \\ +\infty & \text{otherwise.} \end{cases}$$

Gauge:

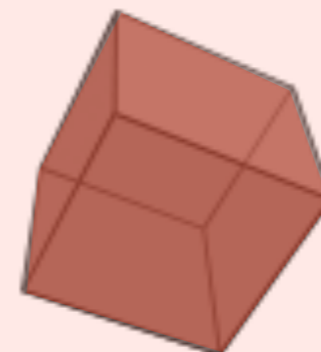
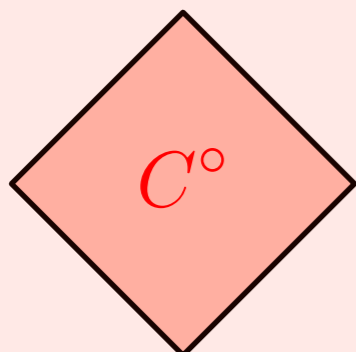
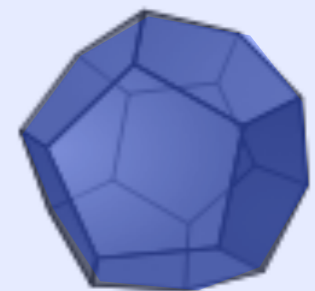
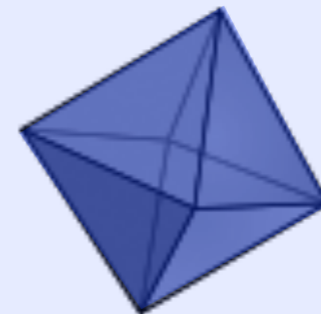
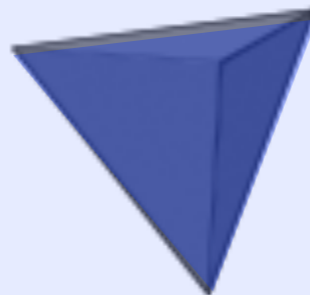
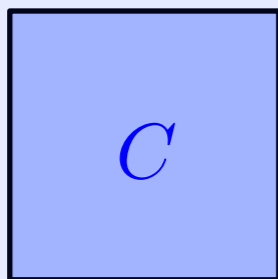
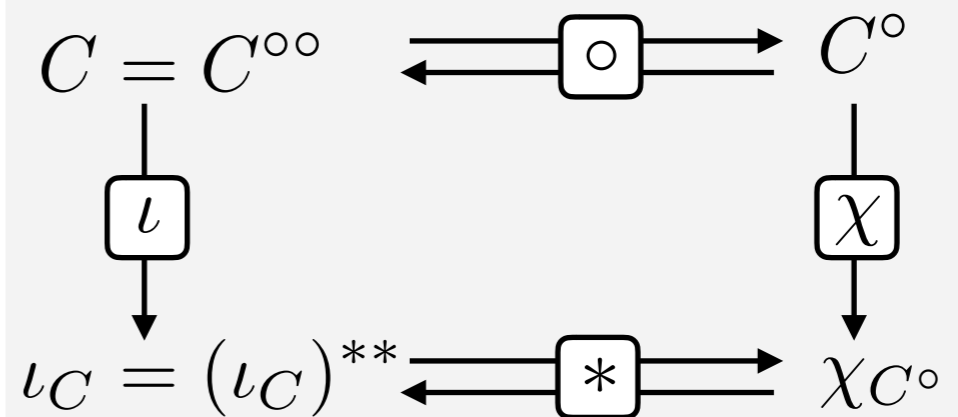
$$\chi_C(x) = \inf \{ \lambda > 0 ; x \in \lambda C \}$$

*Theorem:*

If  $f$  and  $C$  convex, then

$$(C^\circ)^\circ = C \text{ and } (f^*)^* = f.$$

$$(\iota_C)^* = \chi_{C^\circ} \text{ and } (\chi_C)^* = \iota_{C^\circ}.$$

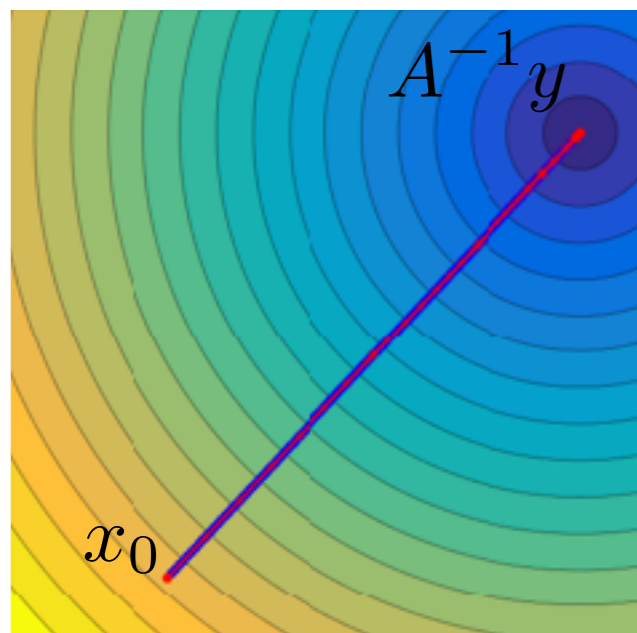


Proximal point:  $x_{k+1} \stackrel{\text{def.}}{=} \text{Prox}_{\tau f}(x_k) = \underset{x}{\text{argmin}} \frac{1}{2} \|x - x_k\|^2 + \tau f(x)$

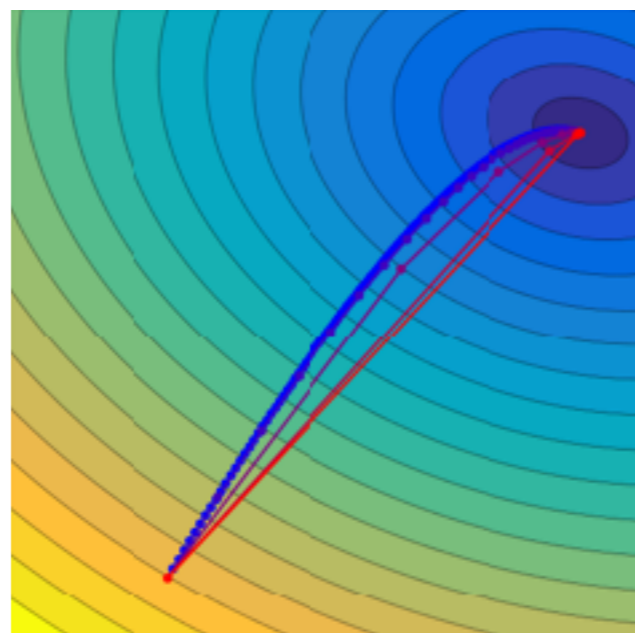
*Theorem:*  $\forall \tau > 0, x_k \xrightarrow{k \rightarrow +\infty} x^* \in \underset{x}{\text{argmin}} f(x)$

*Example:*  $f(x) = \frac{1}{2} \langle Ax, x \rangle - \langle x, y \rangle$

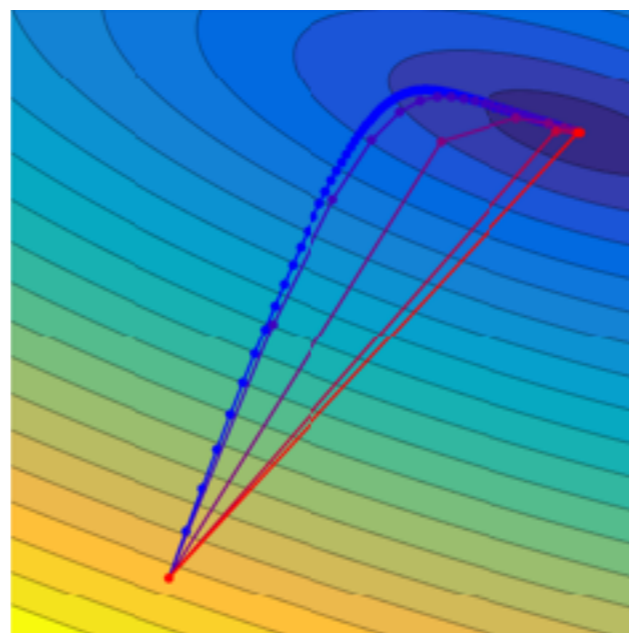
$x_{k+1} = (\text{Id} + \tau A)^{-1}(x_k + \tau y) \xrightarrow{k \rightarrow +\infty} A^{-1}y$



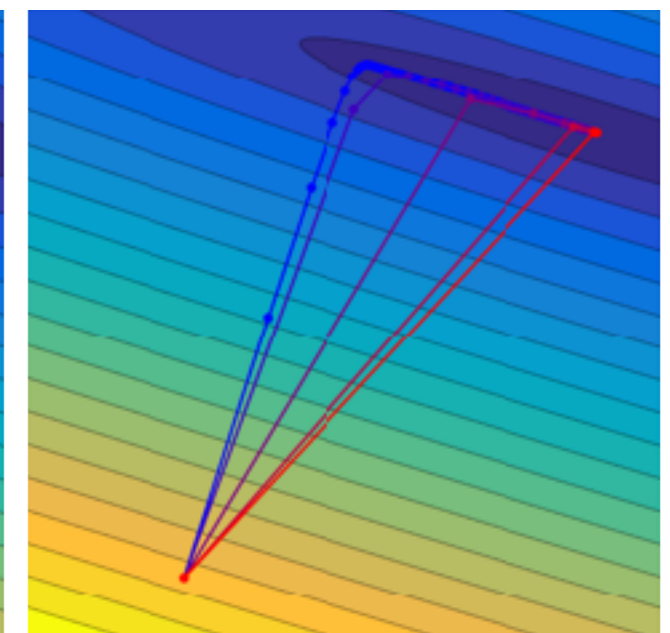
$\tau = 100$



$\tau = 10$



$\tau = 1$



$\tau = 0.1$

$\tau = 0.01$



$f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  convex.

*Subdifferential:*  $\partial f(x) \stackrel{\text{def.}}{=} \{u \in \mathbb{R}^d ; \forall y, f(y) \geq f(x) + \langle u, y - x \rangle\}$

*Theorem:*  $\operatorname{argmin} f = \{x ; 0 \in \partial f(x)\}$

$f$  is differentiable at  $x \Leftrightarrow \partial f(x) = \{\nabla f(x)\}$ .

$\partial f$  is monotone:  $\forall (u, v) \in \partial f(x) \times \partial f(y), \langle u - v, x - y \rangle \geq 0$ .

