

# Optimization

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

Any optimization program is equivalent to a convex (linear) one. #RadonMeasuresRule

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

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

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

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

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

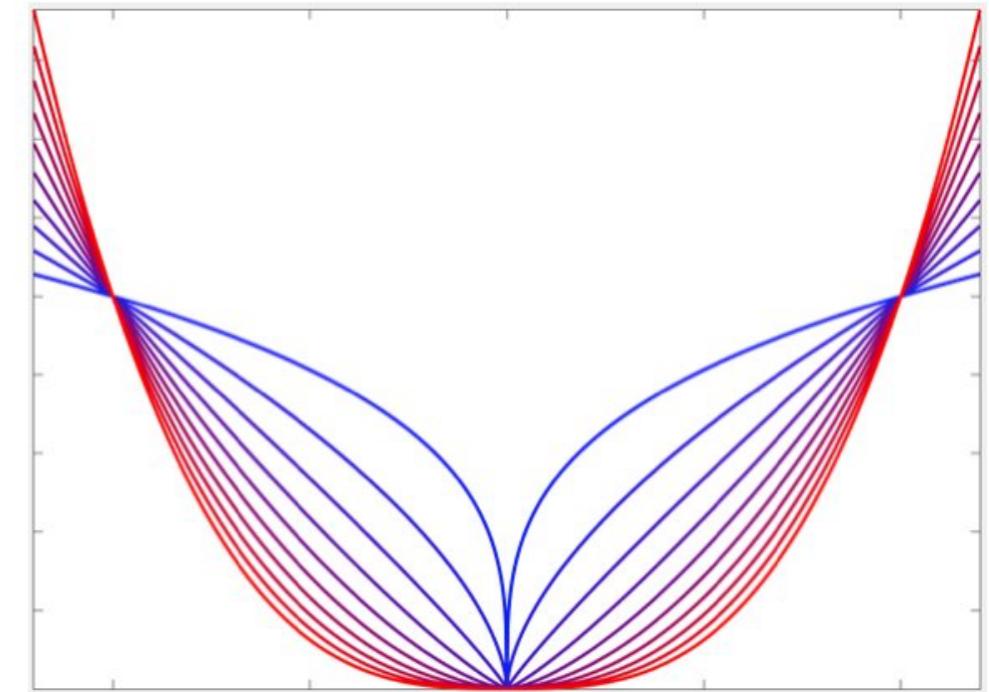
*Theorem:*  $\widehat{f \star g} = \hat{f} \cdot \hat{g}$

Fourier is to convolution what Legendre is to inf-convolution.

[https://en.wikipedia.org/wiki/Fourier\\_transform#Convolution\\_theorem ...](https://en.wikipedia.org/wiki/Fourier_transform#Convolution_theorem ...)

[https://en.wikipedia.org/wiki/Legendre\\_transformation#Infimal\\_convolution ...](https://en.wikipedia.org/wiki/Legendre_transformation#Infimal_convolution ...)

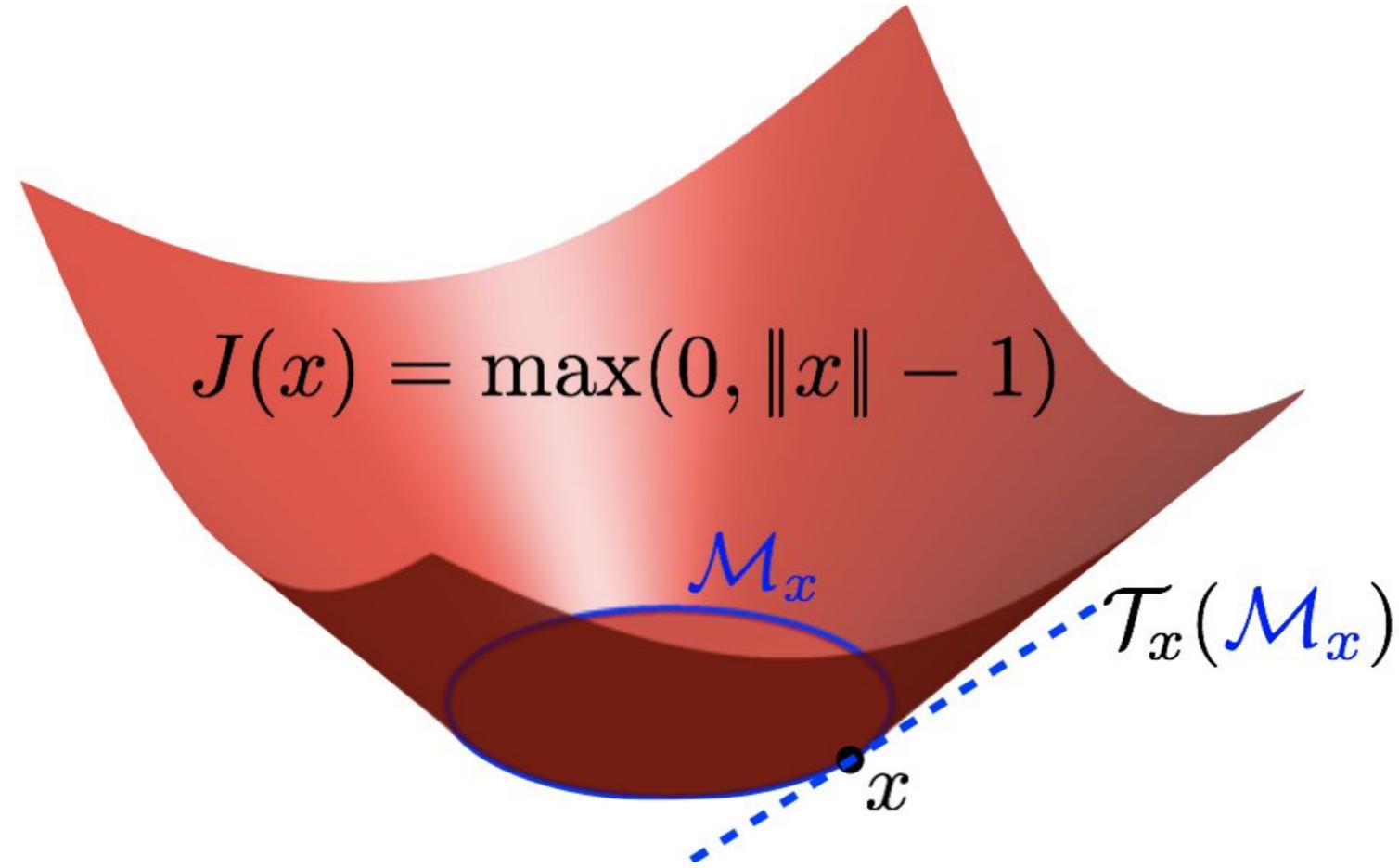
Kurdyka-Łojasiewicz:  
(at minimum  $f(0) = 0$ )  
 $\exists \varphi, \|\nabla(\varphi \circ f)\| \geq 1$   
 $\iff \exists (K, \tau), \|\nabla f(x)\| \geq K|f(x)|^\tau$

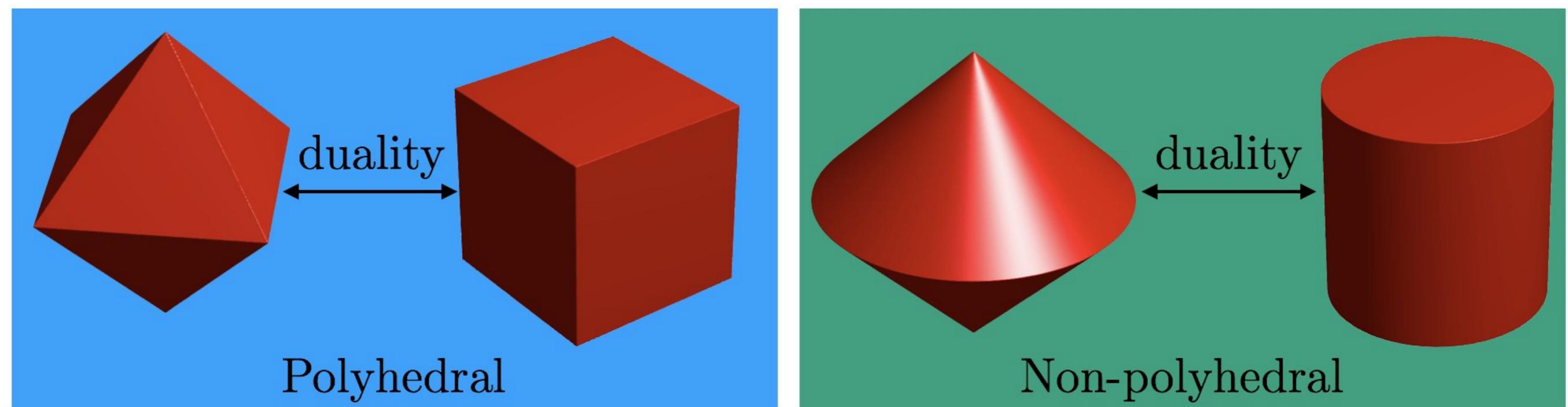


Kurdyka-Łojasiewicz property: fundamental to show convergence of descent optimization methods. tau controls rate. <https://arxiv.org/abs/0802.0826v1>

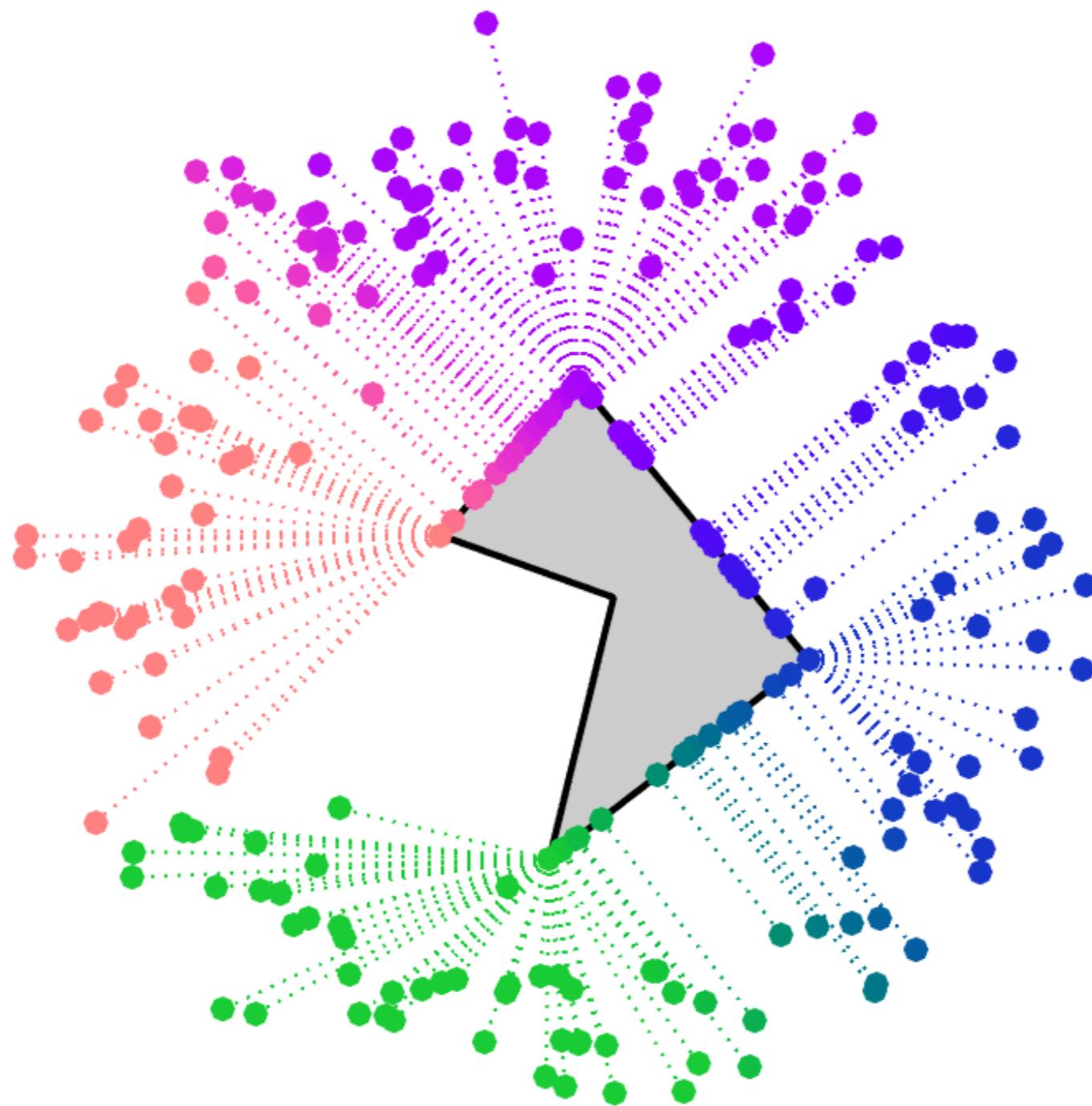
$J : \mathbb{R}^N \rightarrow \mathbb{R}$  is partly smooth at  $x$  for a manifold  $\mathcal{M}_x$

- (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  $x$ .



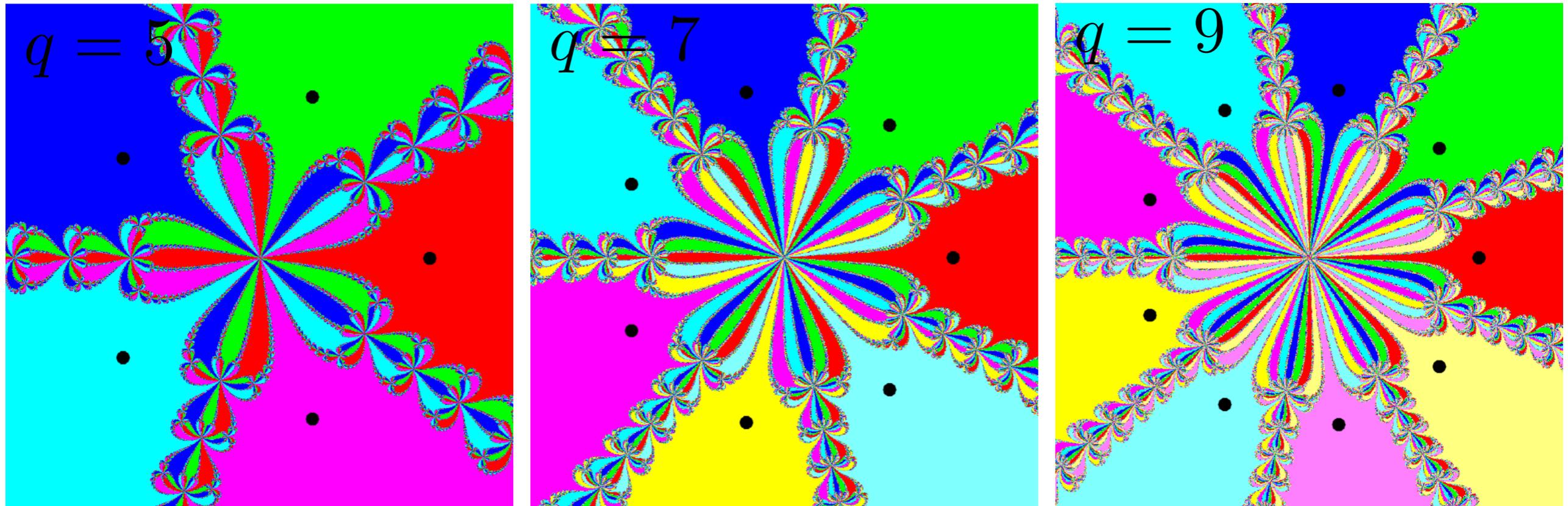


Mirror-stratifiability (Drusvyatskiy/Lewis): generalizes duality of polyhedra. Should become mainstream! <https://arxiv.org/abs/1707.03194>



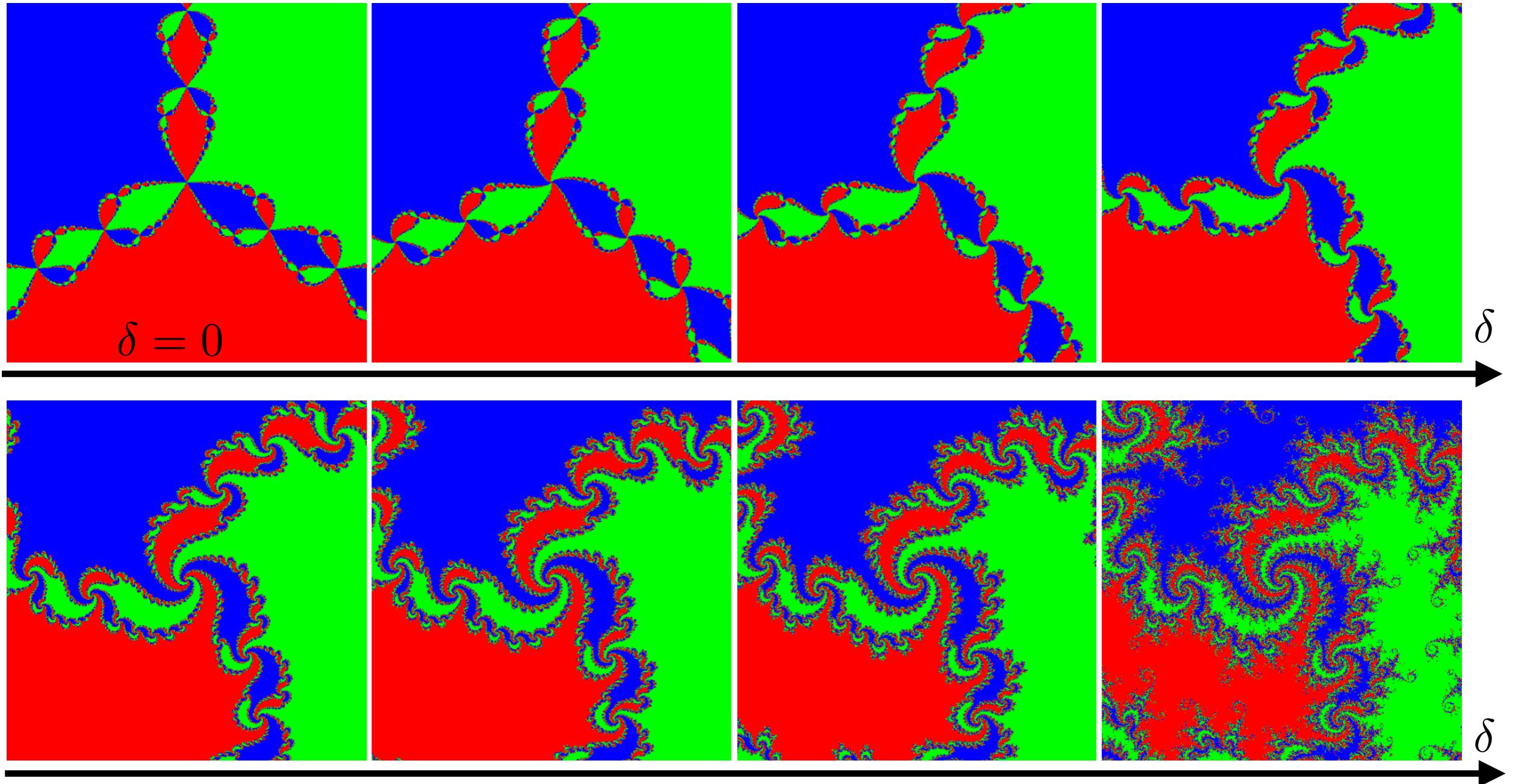
The unreasonable effectiveness of non-smooth optimization: sharp singularities attract solutions of optimization problems.

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



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

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



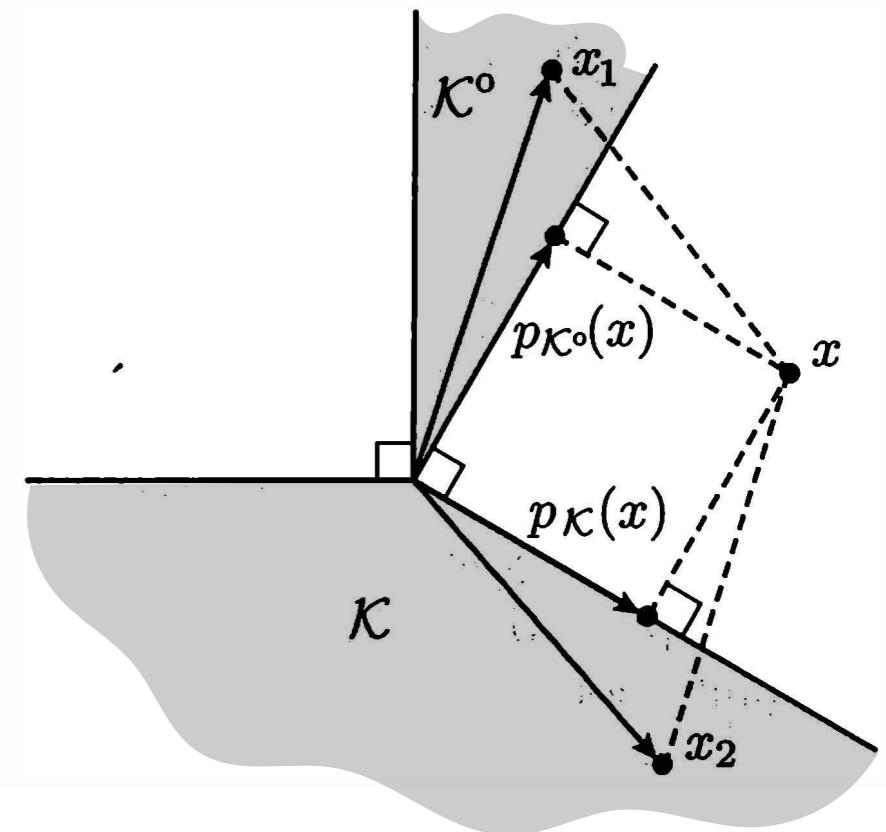
"Twisted" Newton iterations generate nice fractal patterns! Colors display bassin of attraction of cubic roots of unity.

*Polar cone:*

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

*Theorem:* (Moreau's decomposition)

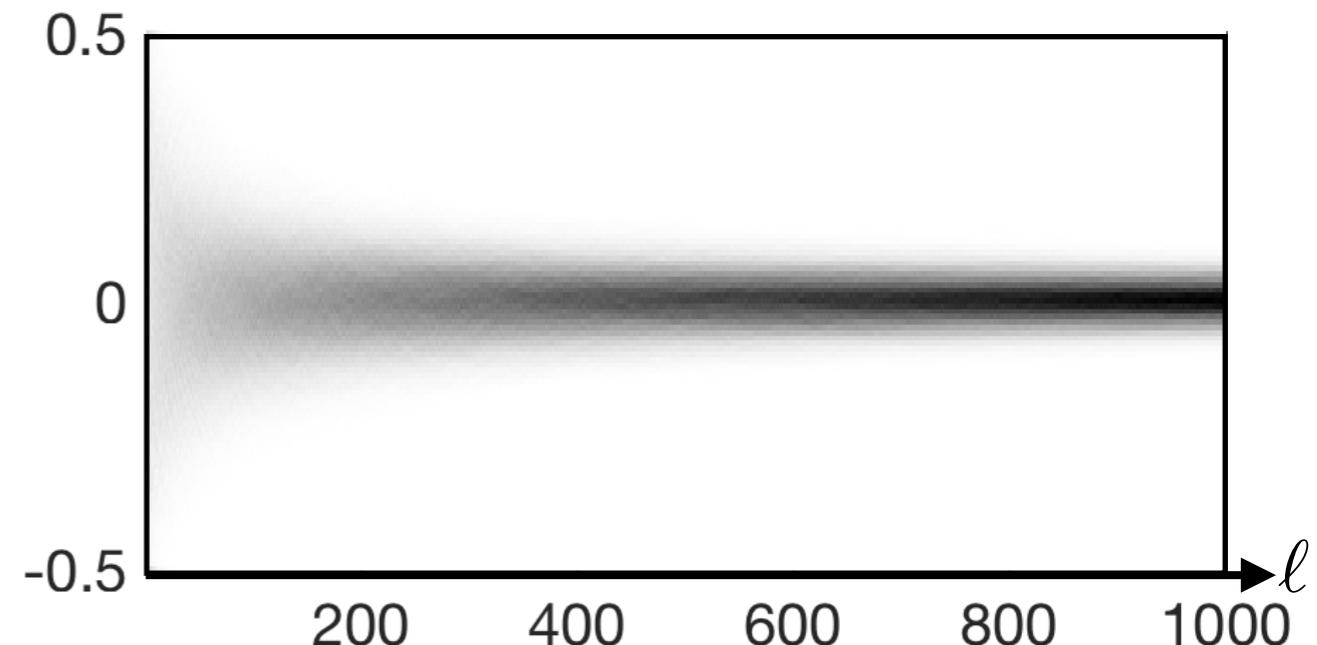
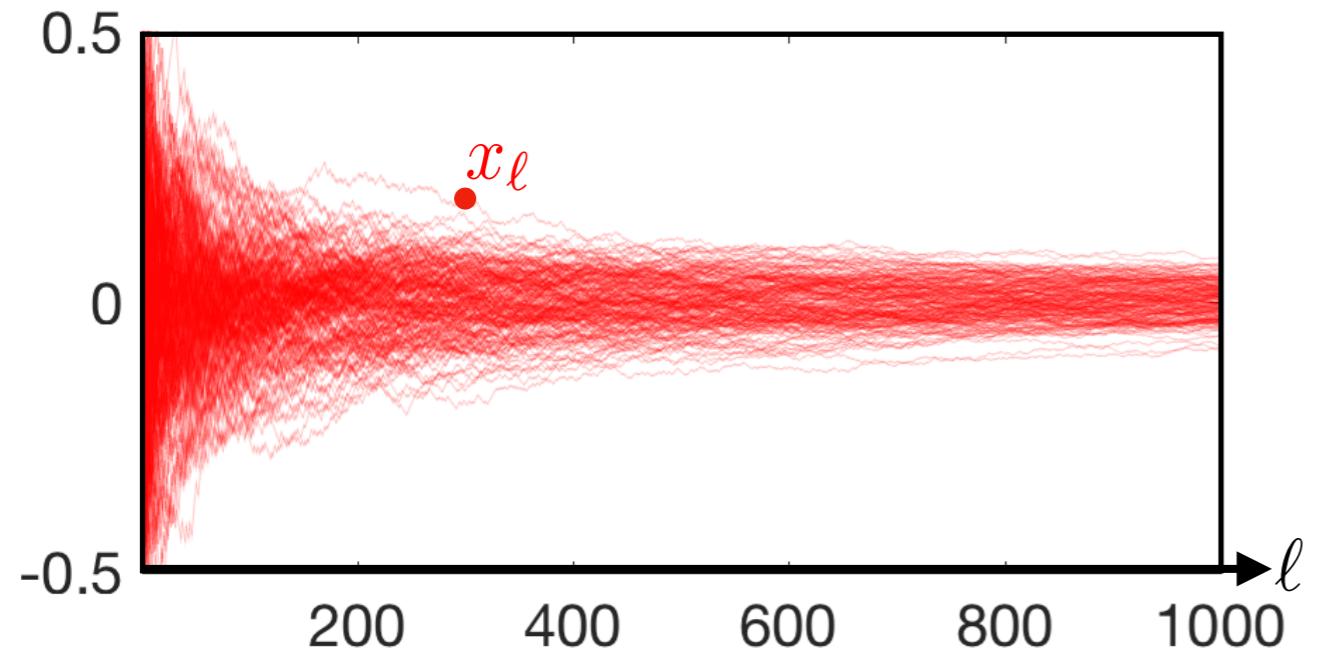
$$x = \text{Proj}_{\mathcal{K}}(x) + {}^\perp \text{Proj}_{\mathcal{K}^\circ}(x)$$



Moreau's decomposition generalizes orthogonal decomposition from linear spaces to convex cones.

$$\min_{x \in \mathbb{R}} (x+1)^2 + (x-1)^2 \\ = f_1(x) \quad \quad = f_2(x)$$

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



Stochastic gradient descent dynamic on a simple 1-D example.

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

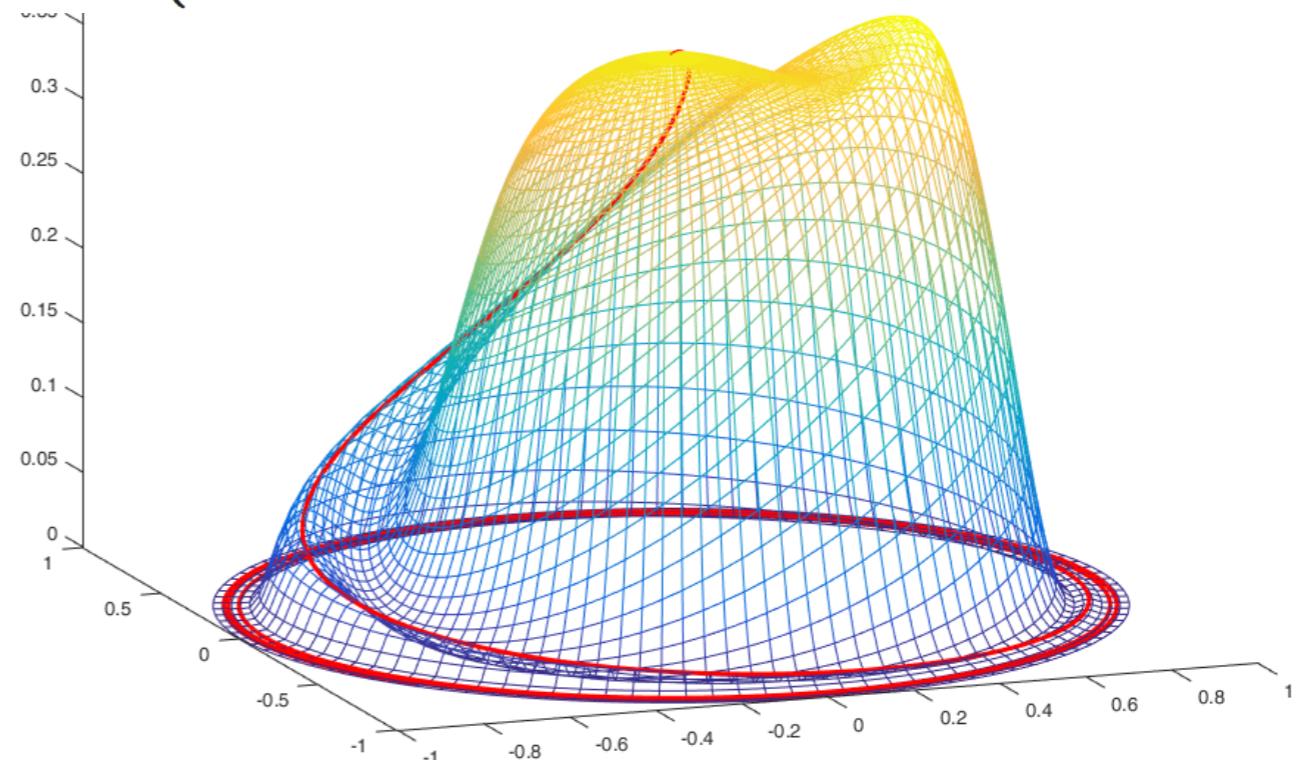
$$x(t) = -\nabla f(x(t)))$$

On the “mexican hat”:

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

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

$$\rightarrow \text{Length}(x) = +\infty$$



Gradient descent does not always converge. Trajectories can have infinite length. Picture & insights from [@GuillaumeG\\_](#) <http://epubs.siam.org/doi/abs/10.1137/040605266>

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

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

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

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

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

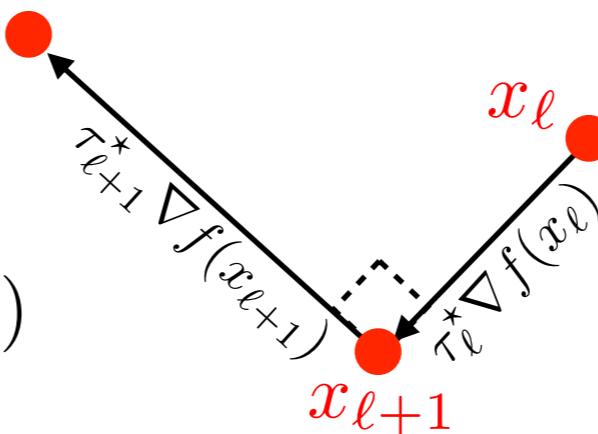
$$\frac{\partial c z}{\partial z} = c \quad \frac{\partial |z|^2}{\partial z} = \frac{\partial (z \bar{z})}{\partial z} = \bar{z}$$

Wirtinger derivatives: treat complex valued functions of complex variables as if they were real functions. Very convenient.

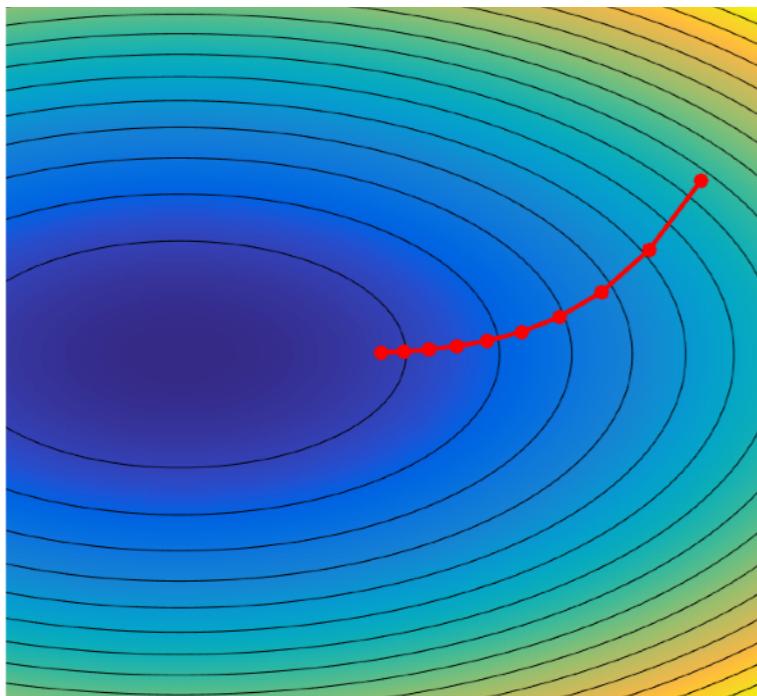
[https://en.wikipedia.org/wiki/Wirtinger\\_derivatives](https://en.wikipedia.org/wiki/Wirtinger_derivatives)

$$x_{\ell+1} = x_\ell - \tau_\ell \nabla f(x_\ell)$$

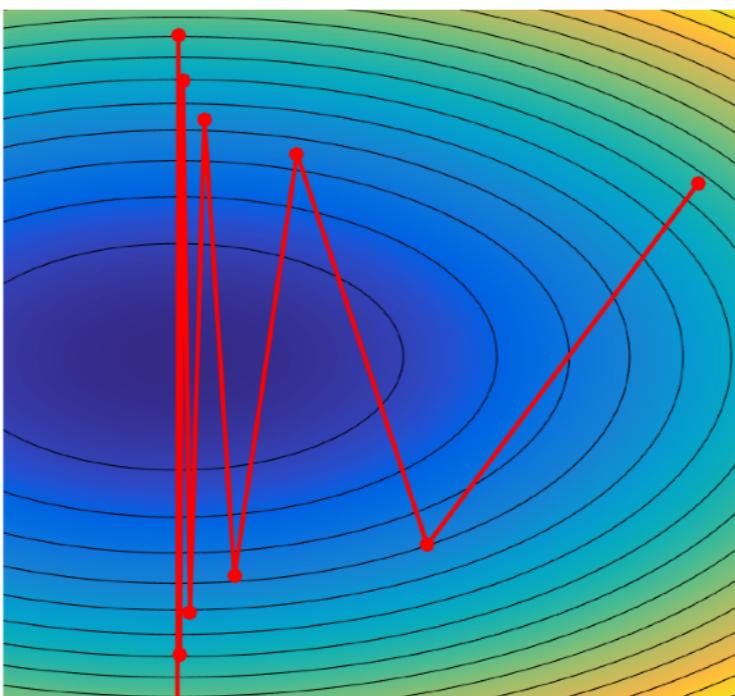
$$\tau_\ell^* = \operatorname{argmin}_\tau f(x_\ell - \tau \nabla f(x_\ell))$$



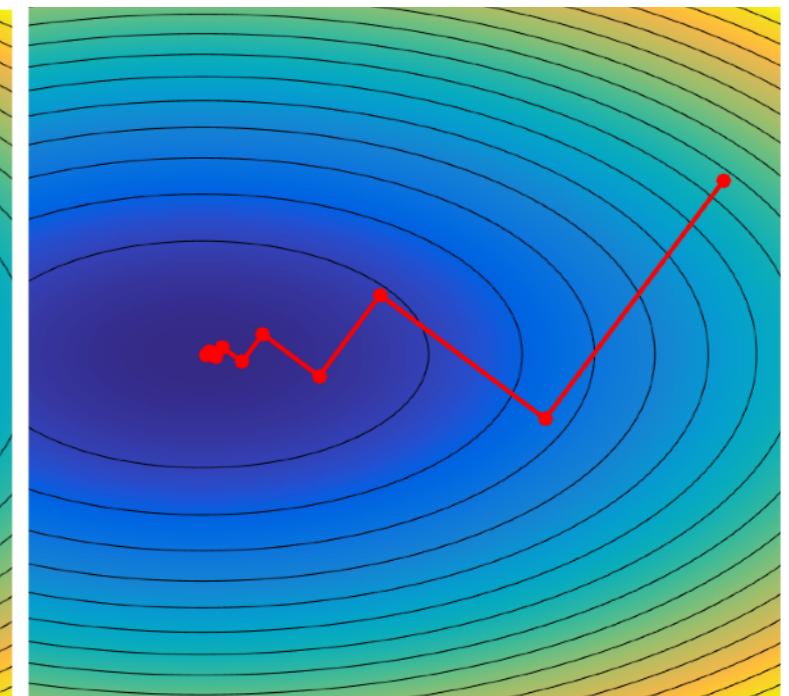
$$\nabla f(x_\ell) \perp \nabla f(x_{\ell+1})$$



Small  $\tau_\ell$



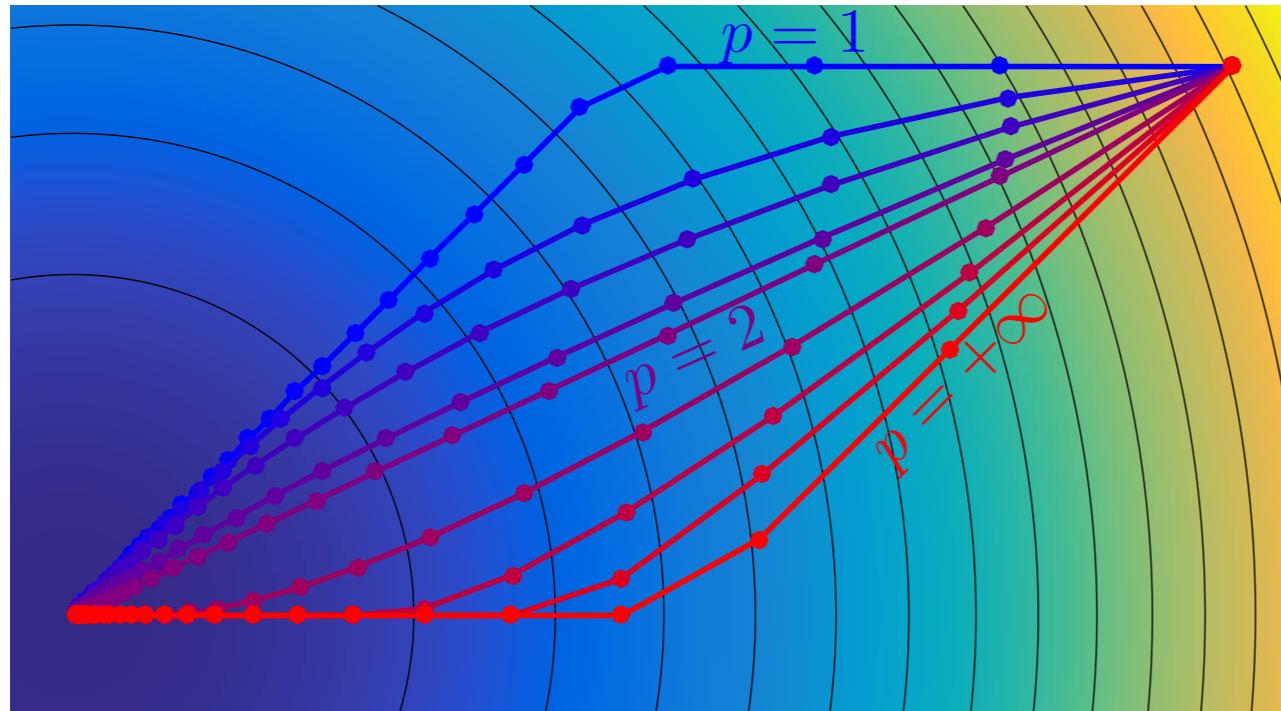
Large  $\tau_\ell$



Optimal  $\tau_\ell = \tau_\ell^*$

Step size selection is important for gradient descent on ill-conditioned functions.

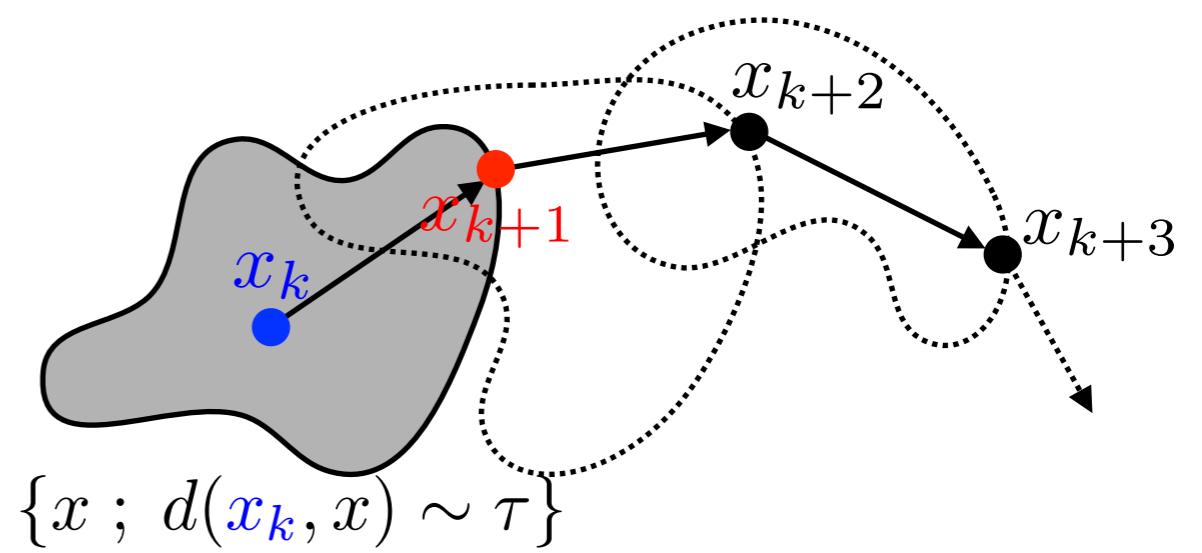
Metric space  $(\mathcal{X}, d)$ , minimize  $F(x)$  on  $\mathcal{X}$ .



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

Implicit Euler step:

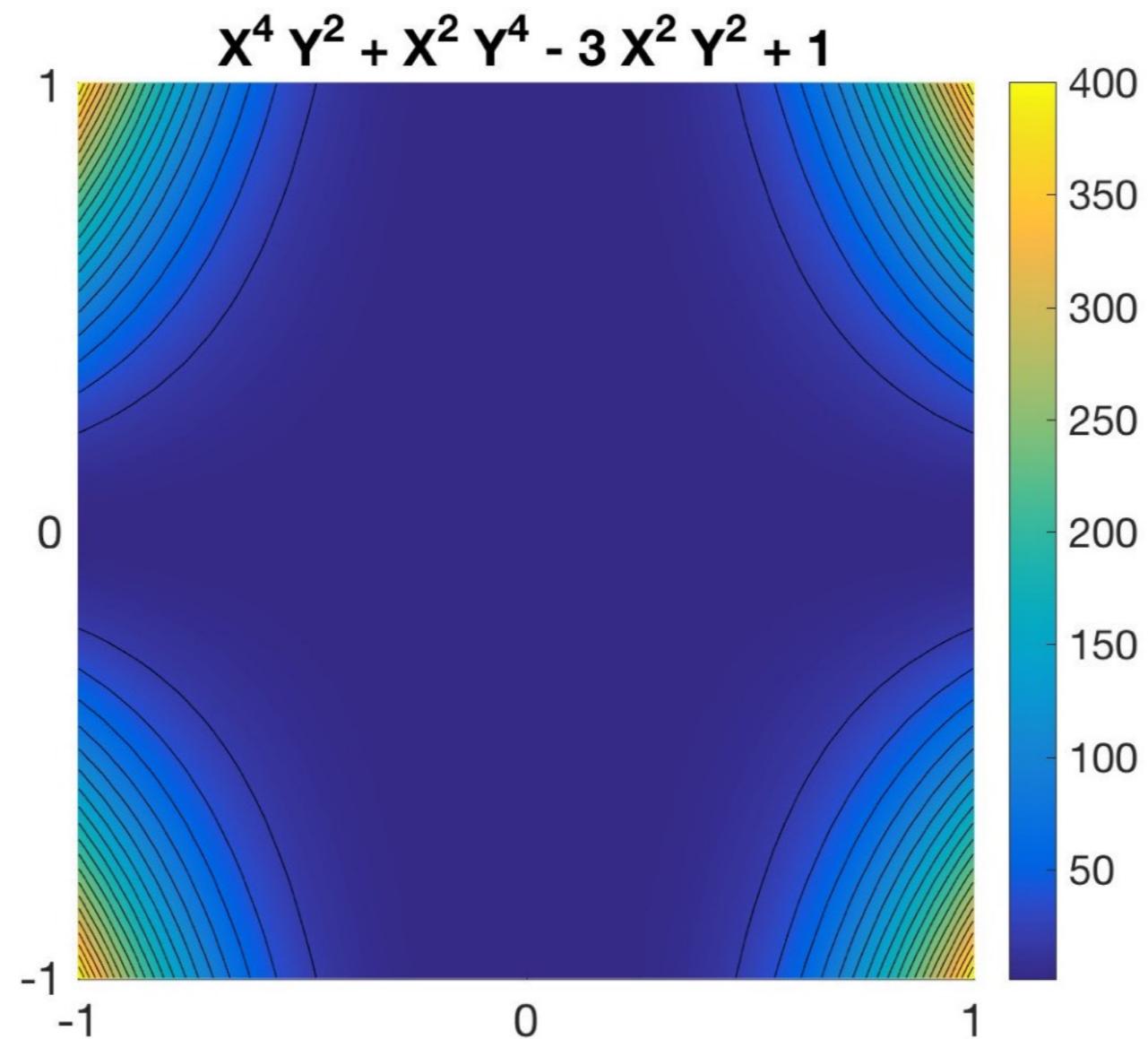
$$x_{k+1} \stackrel{\text{def.}}{=} \operatorname{argmin}_{x \in \mathcal{X}} d(\mathbf{x}_k, x)^2 + \tau F(x)$$



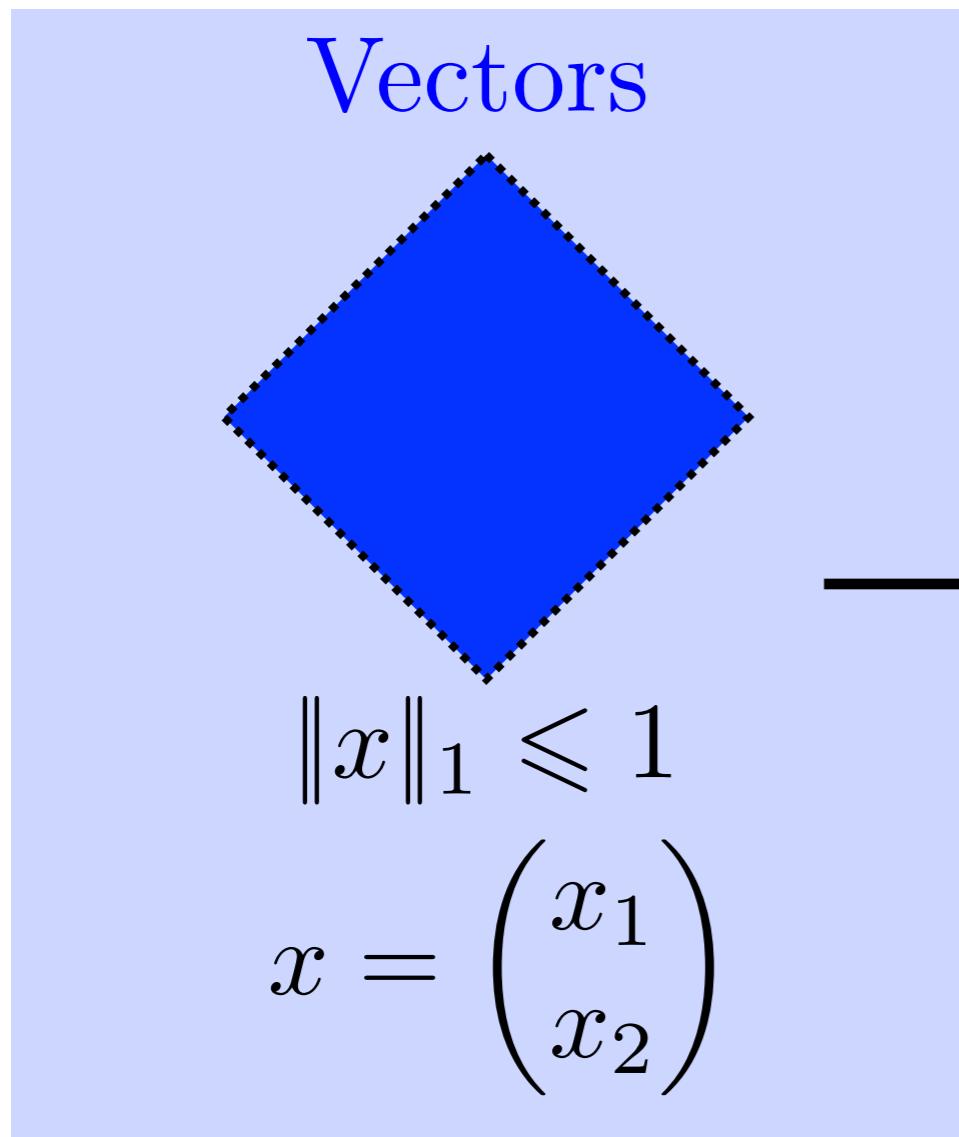
Gradient flows on metric spaces: define descent methods on non-Euclidean domains without using a gradient. A key tool studied in details by Ambrosio, Gigli and Savare.

<http://www.springer.com/la/book/9783764387211>

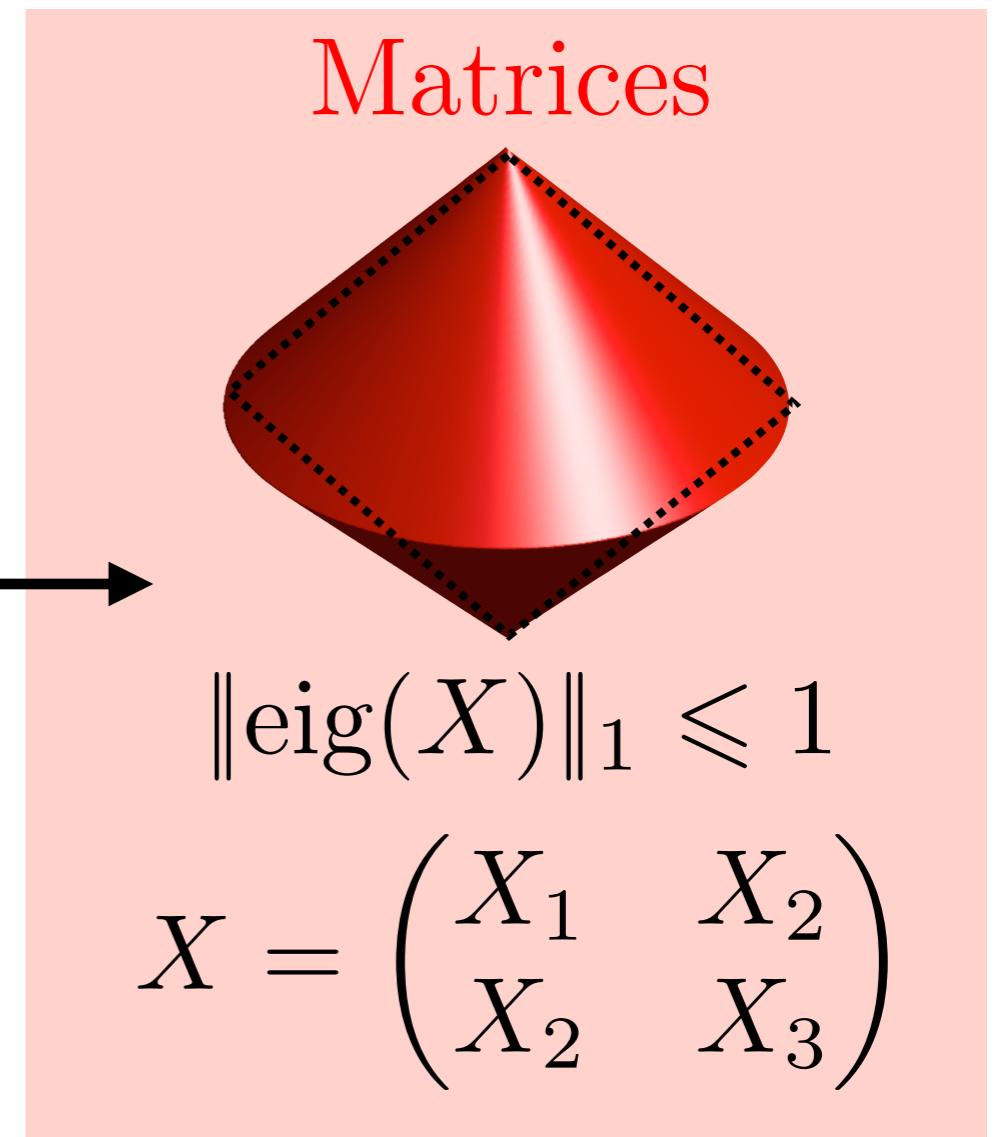
# **Matrix and Polynomials**



Motzkin's polynomial is positive but is not a sum of squares.  
#LifeIsHarderIn2D [https://en.wikipedia.org/wiki/Positive\\_polynomial...](https://en.wikipedia.org/wiki/Positive_polynomial...)



spectral  
lift



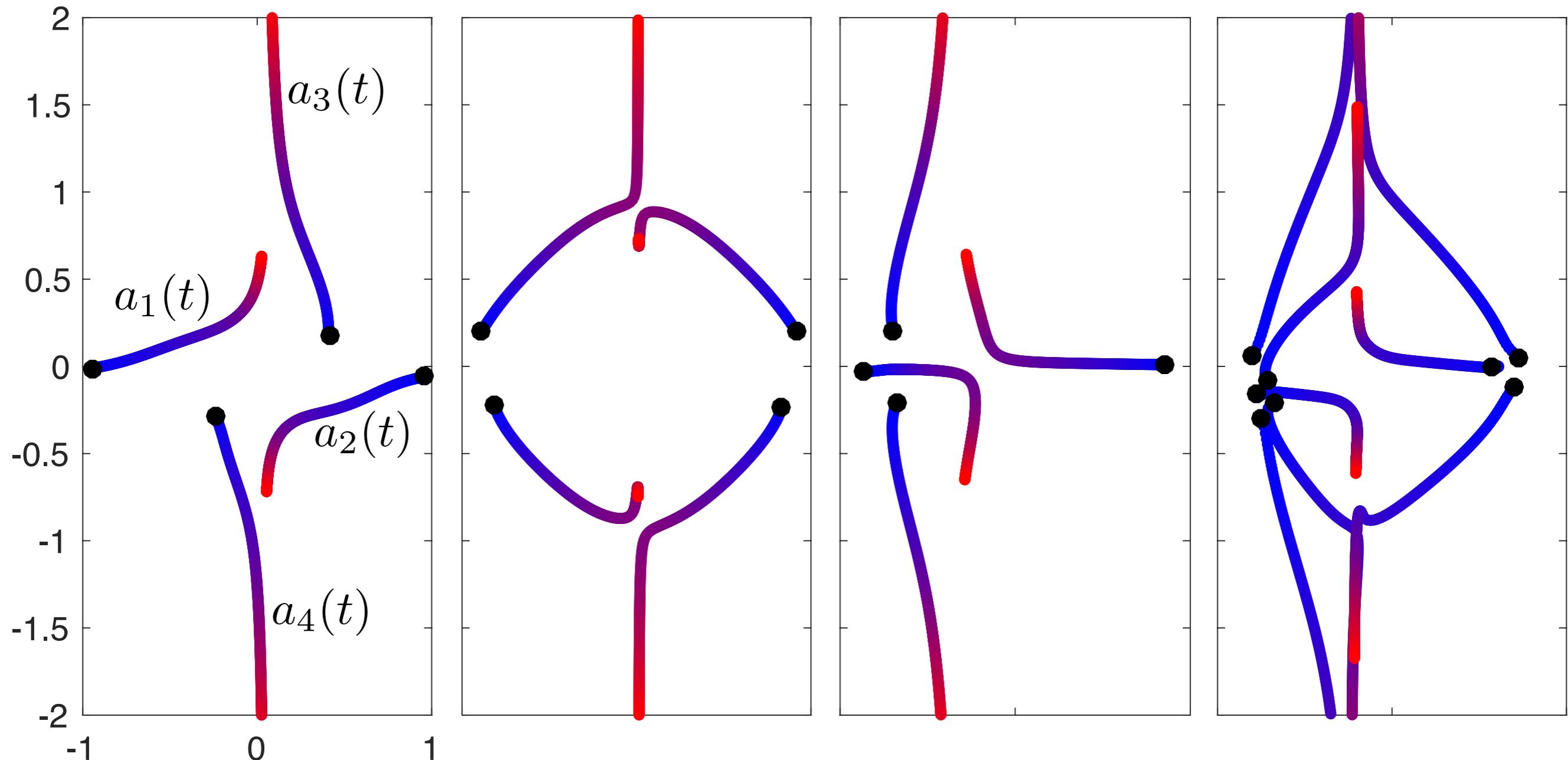
Reasonable properties (convexity, differential, stratifiability, partial-smooth...) lift from vectors to matrices! <https://people.orie.cornell.edu/aslewis/publications/03-mathematics.pdf> ...



The QR algorithm is a gem from numerical matrix analysis.  
Showing convergence toward triangular matrix. [https://en.wikipedia.org/wiki/QR\\_algorithm](https://en.wikipedia.org/wiki/QR_algorithm) ...

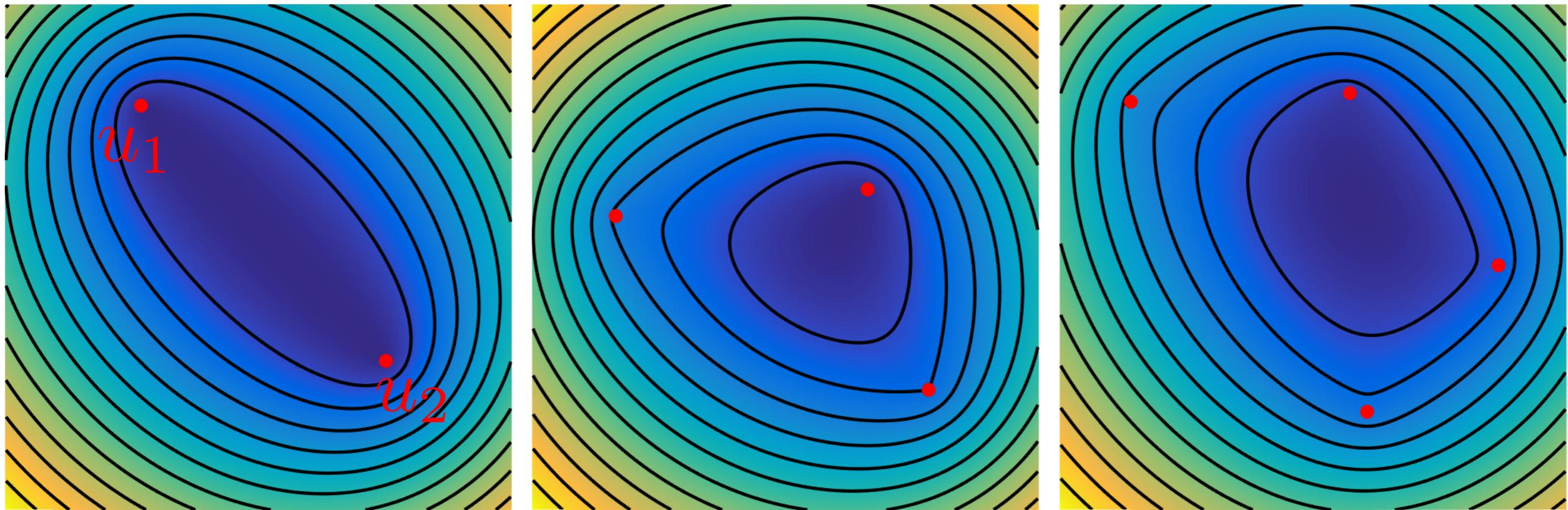
$$\partial_t P_t(x) = \partial_x^2 P_t(x)$$

$$P_t(x) = \prod_i (x - a_i(t))$$



Heat eq on polynomials generates nice roots dynamics. Post by Terry Tao. Code provided!  
[https://terrytao.wordpress.com/2017/10/17/heat-flow-and-zeroes-of-polynomials/ ...](https://terrytao.wordpress.com/2017/10/17/heat-flow-and-zeroes-of-polynomials/)  
<https://goo.gl/hPE5uo>

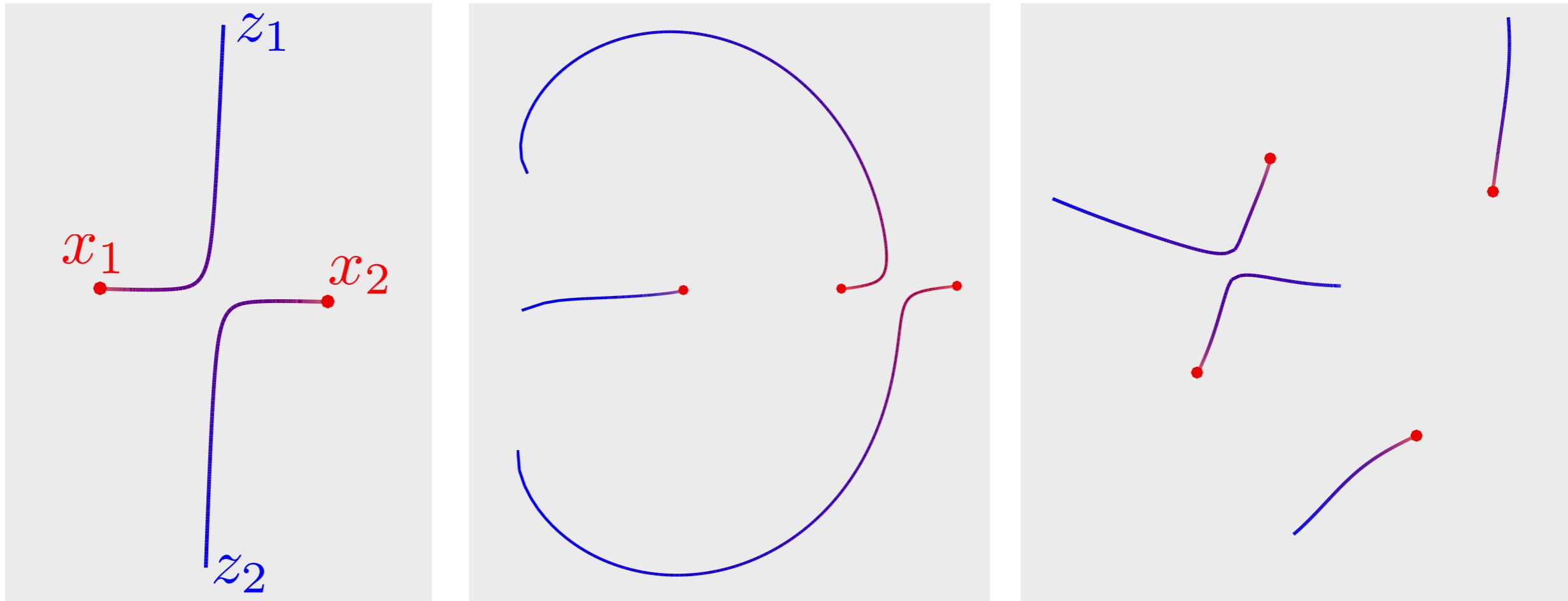
$$\left\{ (x, y) \in \mathbb{R}^2 ; \sum_{i=1}^n \sqrt{(x - u_k)^2 + (y - v_k)^2} \leq t \right\}$$



n-ellipses generalize ellipses with n foci. Spectrahedra (hence convex) algebraic curves of degree less than  $2^n$ . <https://en.wikipedia.org/wiki/N-ellipse>

$$z_i \longleftrightarrow z_i - \frac{P(z_i)}{\prod_{j \neq i} (z_i - z_j)}$$

Theorem:  
 $(z_i)_i \rightarrow \text{Roots}(P)$

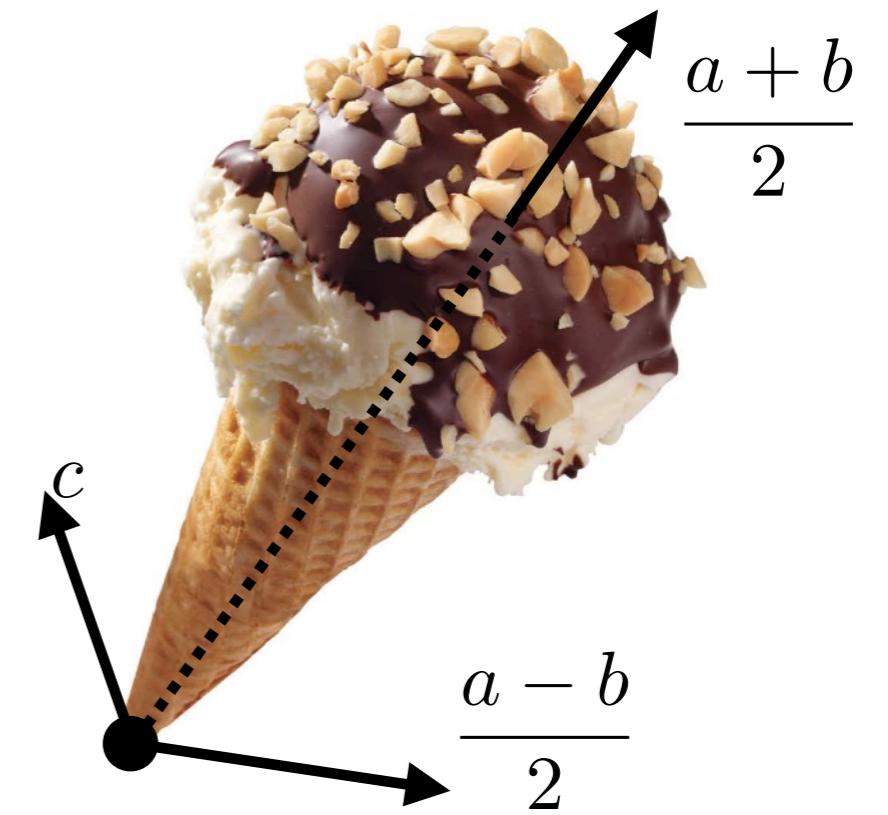


$$P(x) = \prod_i (x - x_i)$$

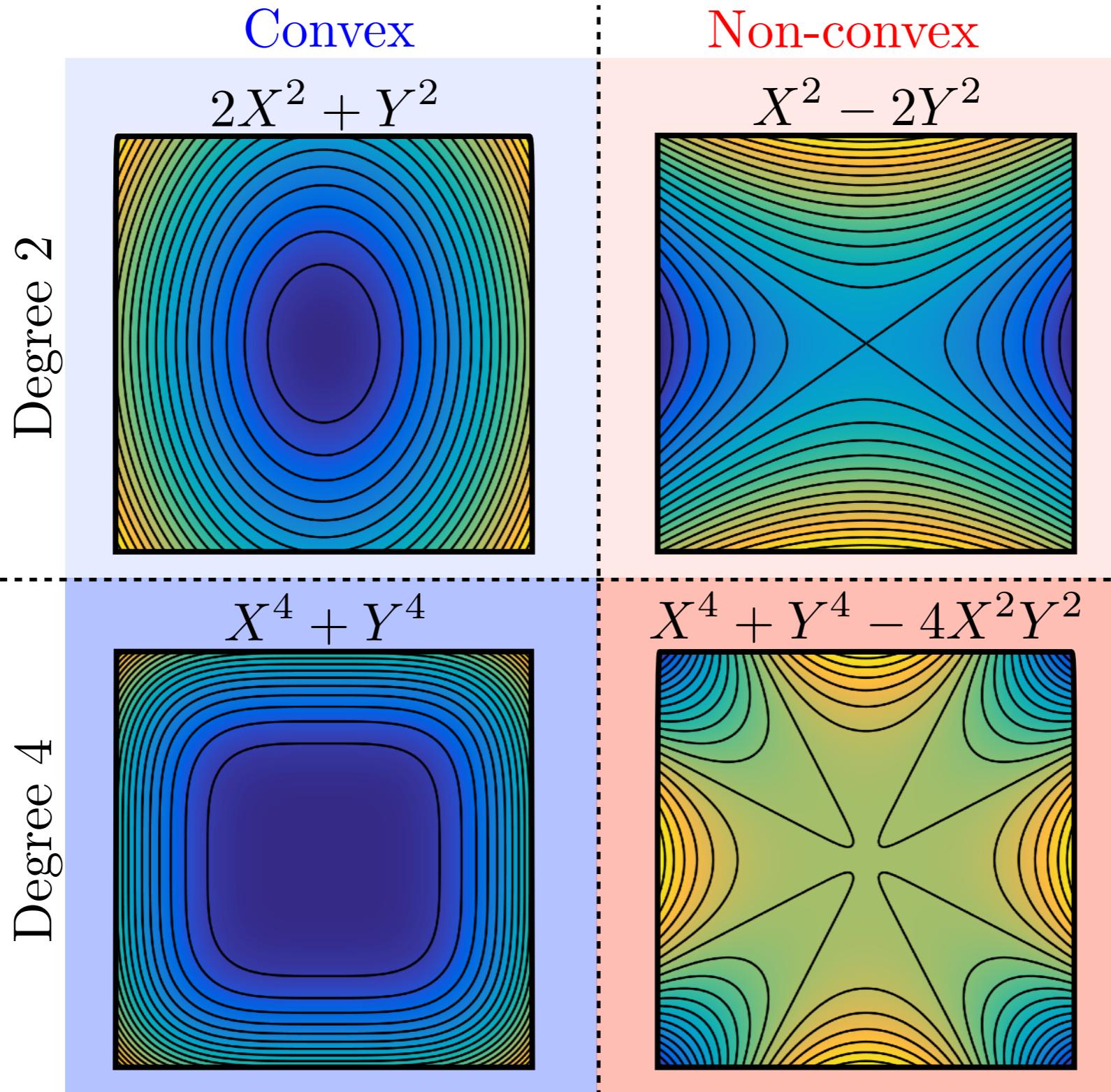
$$\{(a, b, c) ; \forall x, a + 2cx + bx^2 \geq 0\}$$

$$= \left\{ (a, b, c) ; \begin{pmatrix} a & c \\ c & b \end{pmatrix} \succeq 0 \right\}$$

$$= \left\{ (a, b, c) ; \left( \frac{a-b}{2} \right)^2 + c^2 \leq \left( \frac{a+b}{2} \right)^2 \right\}$$



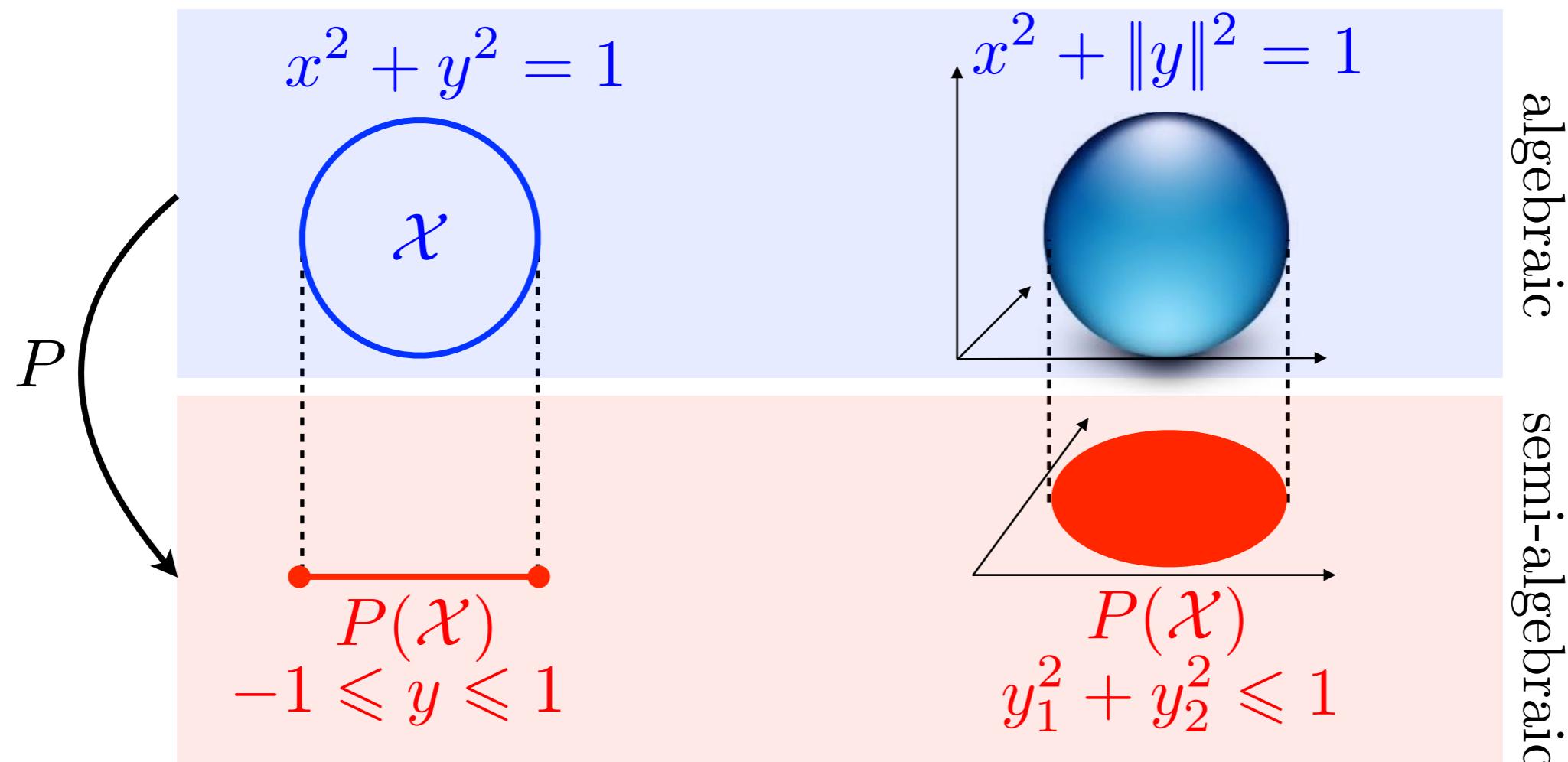
I knew it! Ice-creams, positive polynomials and PSD matrices all have the same flavor!



Checking convexity of degree 2 polynomials is trivial.  
 Checking convexity of degree 4 polynomials is NP-hard.  
[http://web.mit.edu/~a\\_a\\_a/Public/Publications/convexity\\_nphard.pdf](http://web.mit.edu/~a_a_a/Public/Publications/convexity_nphard.pdf)

$\mathcal{X} \subset \mathbb{R}^n$  semi-algebraic  $\implies P(\mathcal{X}) \in \mathbb{R}^{n-1}$  semi-algebraic.

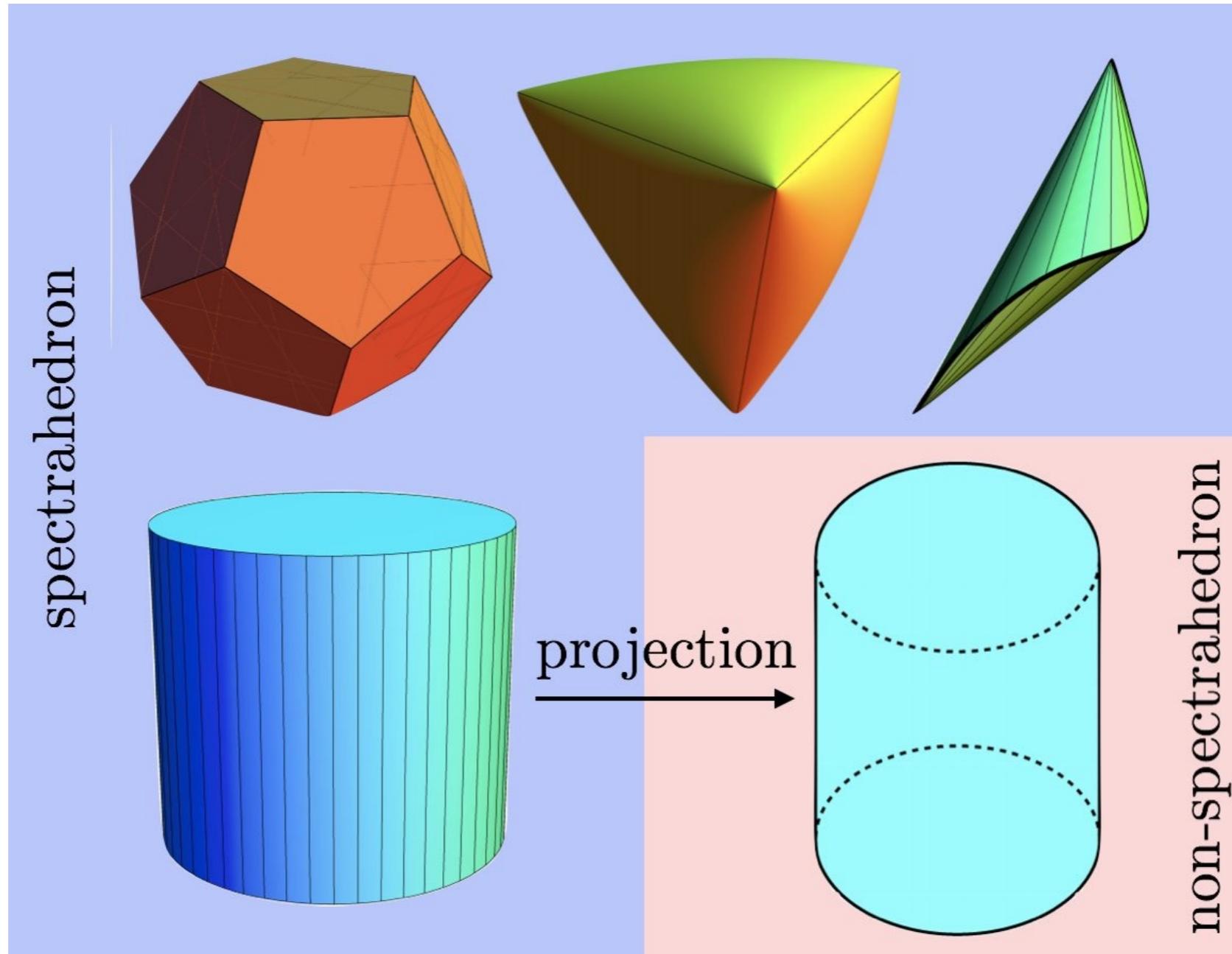
$P : (x, y) \in \mathbb{R} \times \mathbb{R}^{n-1} \mapsto y \in \mathbb{R}^{n-1}$



Tarski–Seidenberg theorem: semi-algebraicity is stable by projection. The most fundamental result in semi-algebraic geometry. Equivalent to the elimination of the existential quantifier.

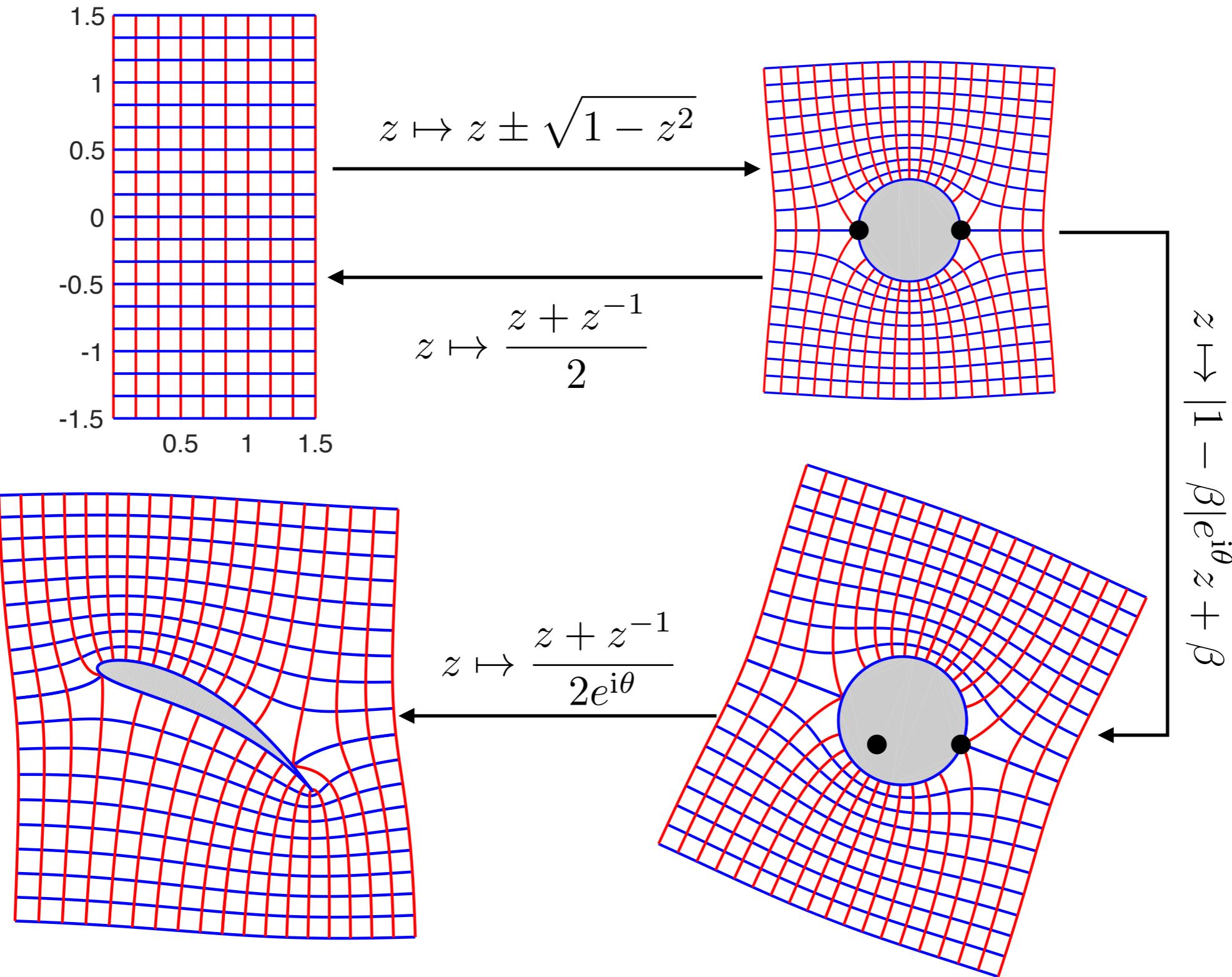
[https://en.wikipedia.org/wiki/Tarski%26Seidenberg\\_theorem](https://en.wikipedia.org/wiki/Tarski%26Seidenberg_theorem)

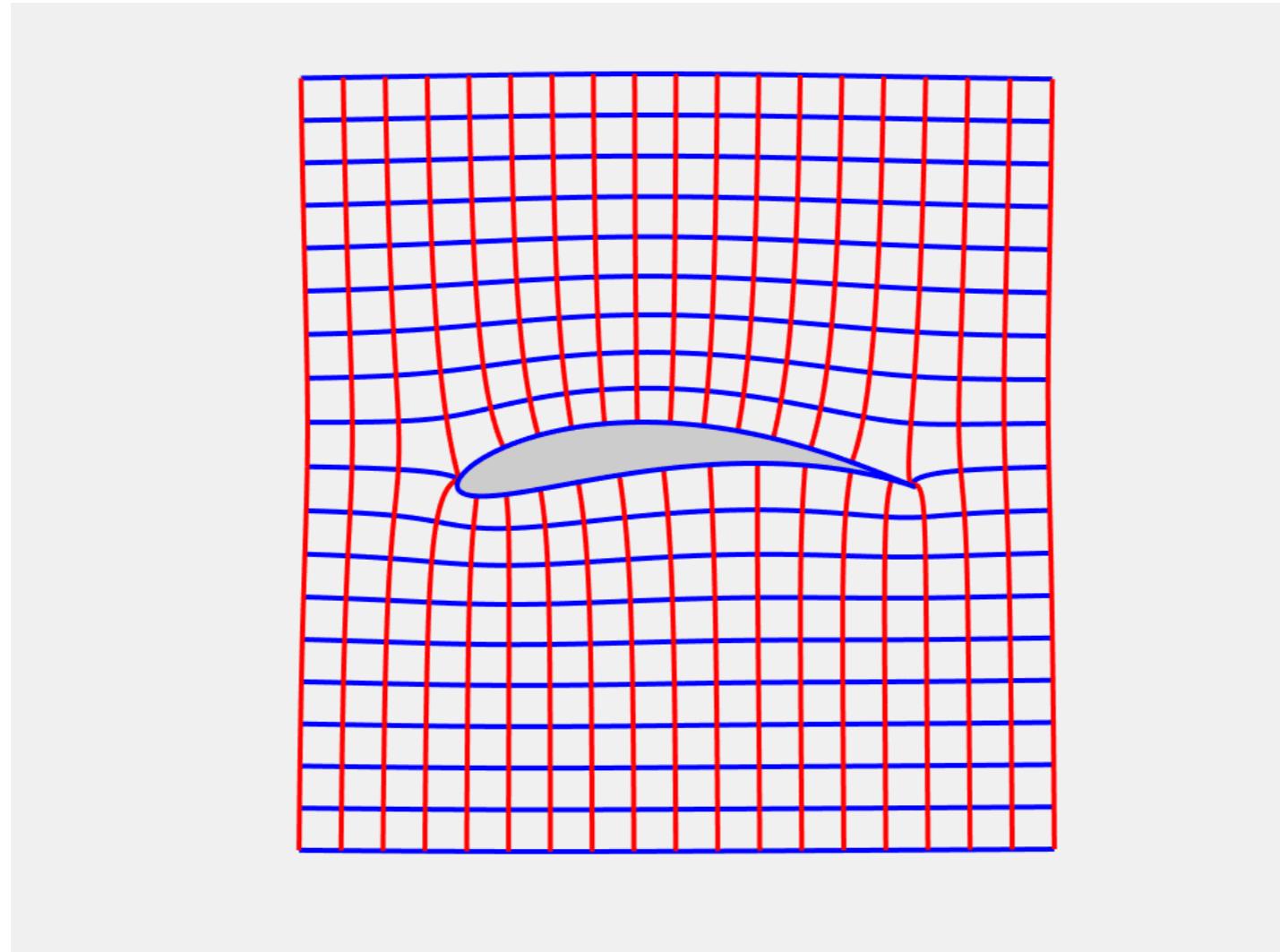
# Geometry



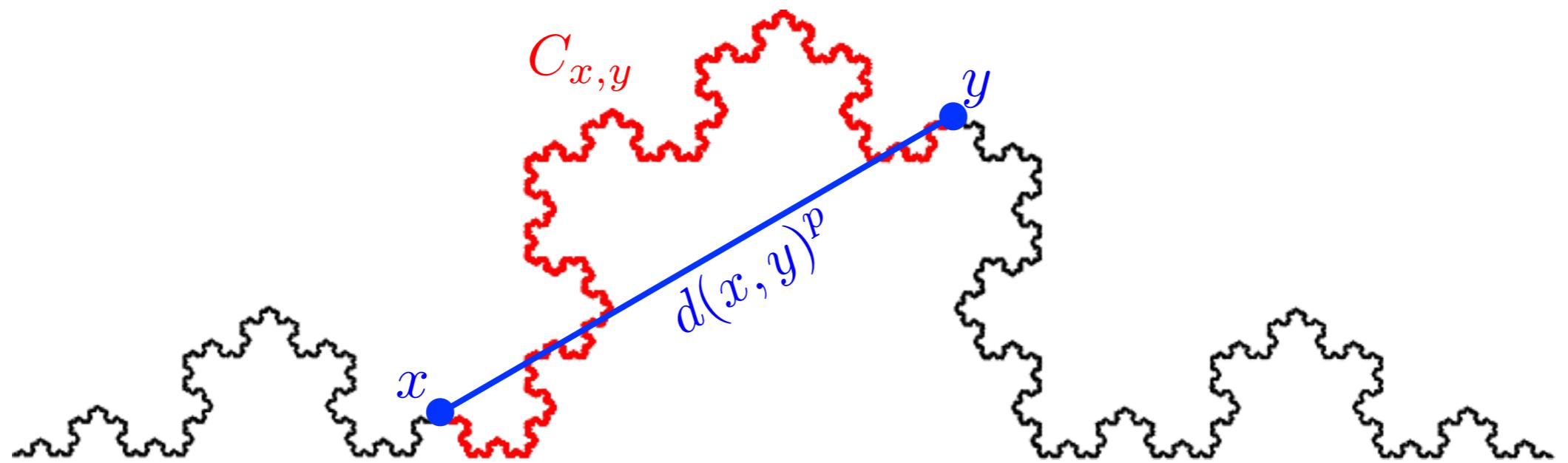
Spectrahedra are matrix generalizations of polyhedra. But the class of spectrahedra is not stable by projection. <https://math.berkeley.edu/~bernd/WhatIsS>

# Joukowski Conformal Mapping





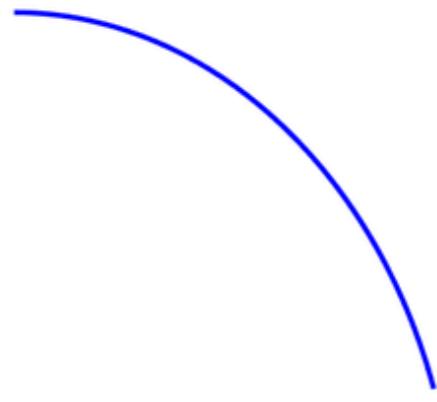
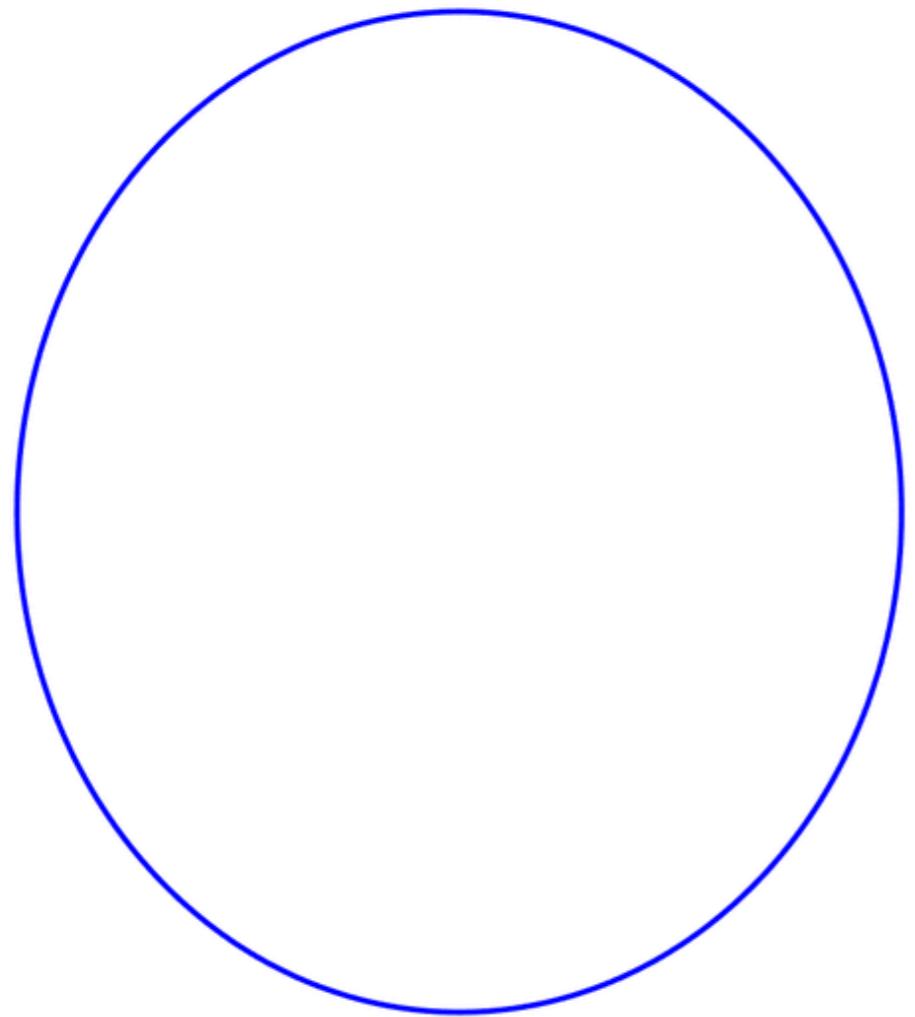
Rotating airfoil using Joukowsky conformal map. Code provided! <https://goo.gl/W6igS5>



$$d(x, y) = \text{Hausdorff}_p(C_{x,y})^{\frac{1}{p}}$$

$$p = \frac{\log(4)}{\log(3)}$$

If  $d(x,y)$  is a distance, then  $d(x,y)^p$  for  $0 < p < 1$  is its snowflake distance. Naming bc of the von Koch fractal. <https://www.emis.de/journals/AASF/Vol30/tynson.pdf> ...



Spirogram (hypotrochoid) vs lissajou curves: messing around with a bunch of cos and sin.

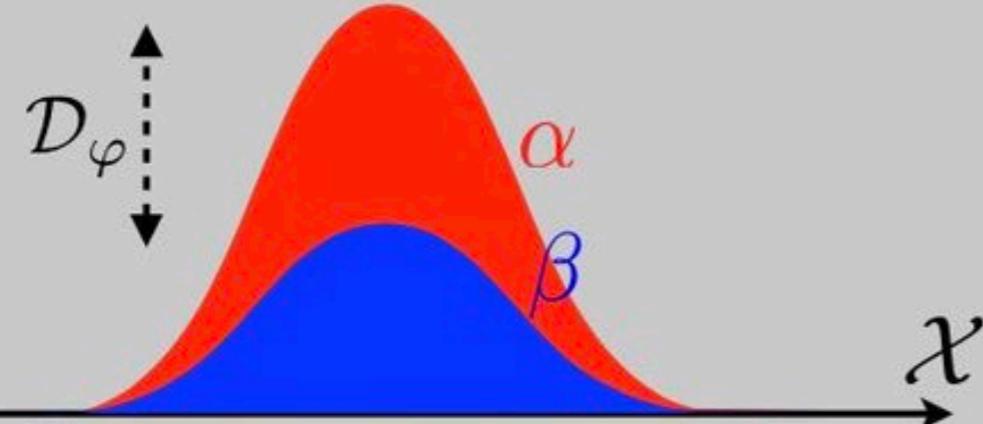
[https://en.wikipedia.org/wiki/Lissajous\\_curve](https://en.wikipedia.org/wiki/Lissajous_curve)

<https://en.wikipedia.org/wiki/Hypotrochoid>

# **Measures and Probability**

Csiszár divergences:

$$\mathcal{D}_\varphi(\alpha|\beta) \stackrel{\text{def.}}{=} \int_{\mathcal{X}} \varphi\left(\frac{d\alpha}{d\beta}\right) d\beta$$

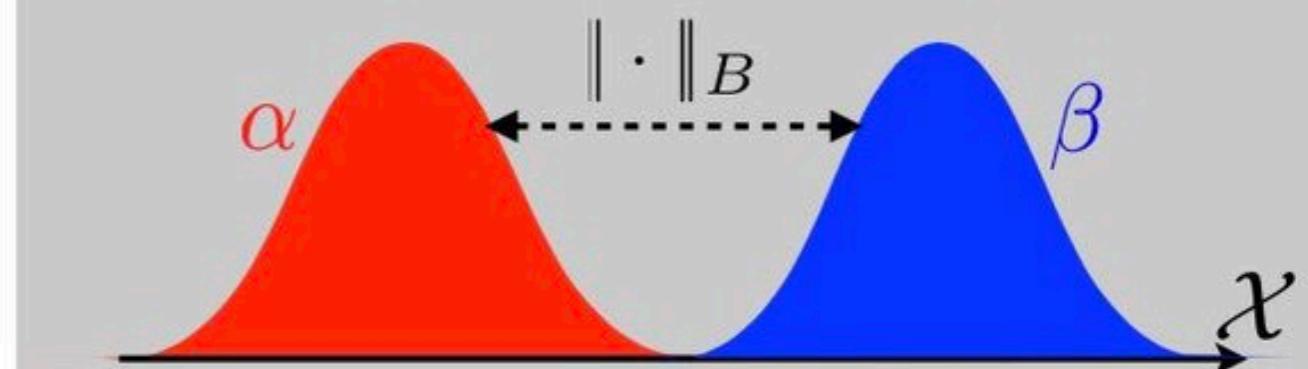


*Strong topology*

→ KL, TV,  $\chi^2$ , Hellinger ...

Dual norms:

$$\|\alpha - \beta\|_B \stackrel{\text{def.}}{=} \max_{f \in B} \int_{\mathcal{X}} f(x)(d\alpha(x) - d\beta(x))$$

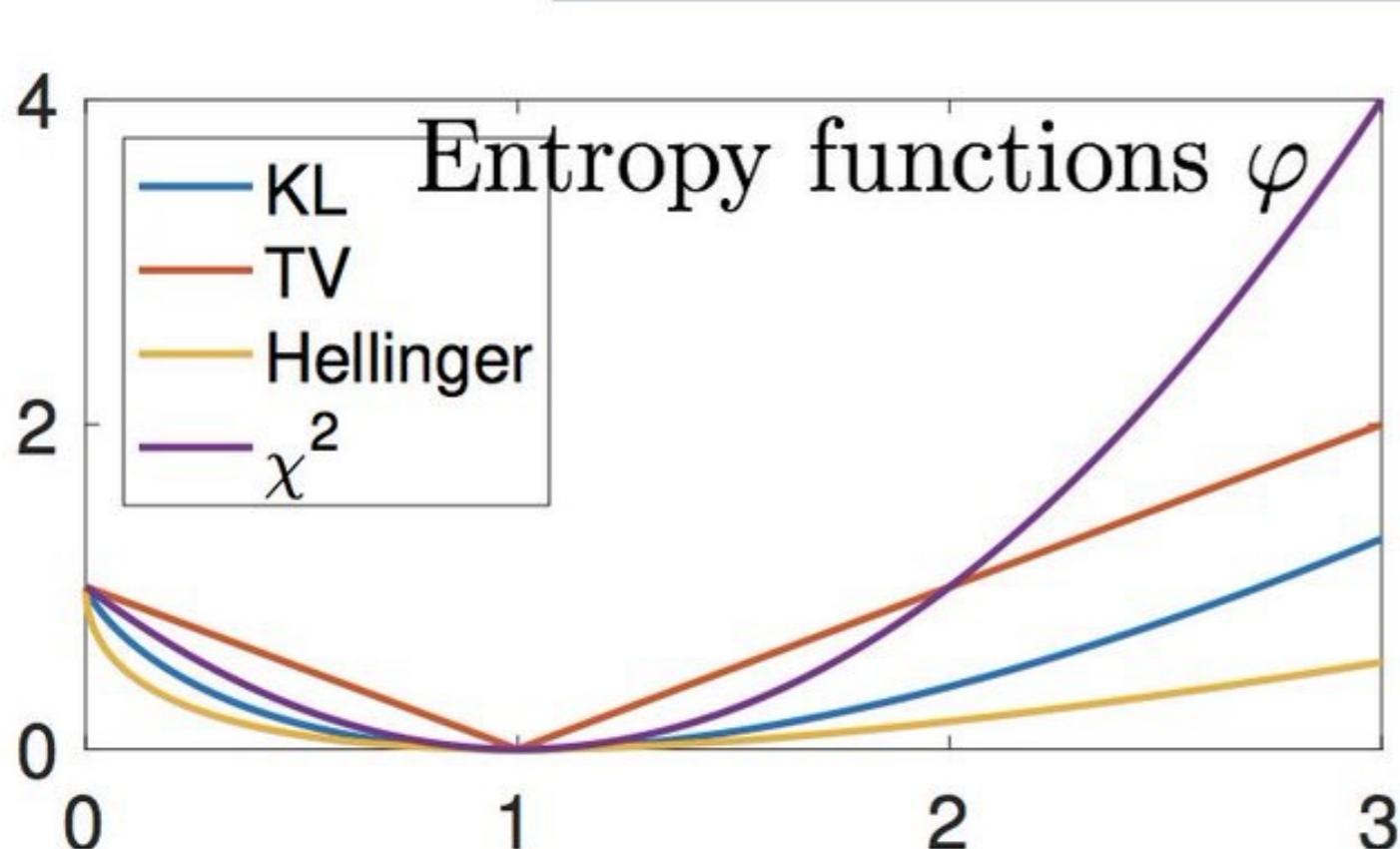


*Weak topology*

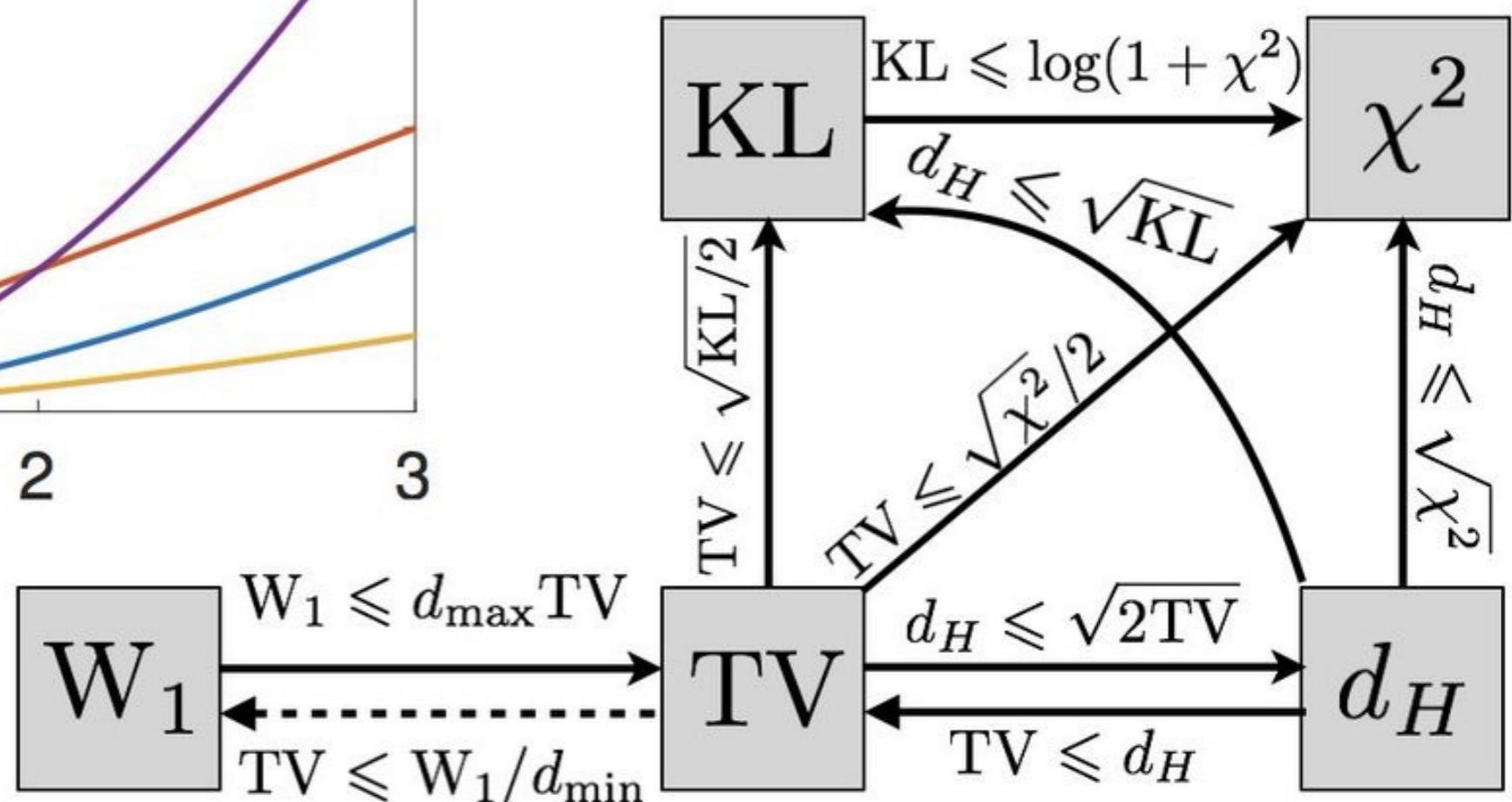
→  $W_1$ , flat, RKHS\*, energy dist, ...

Comparing probability measures: vertical vs horizontal displaceent / Csiszár divergence vs dual norm / strong vs weak topology

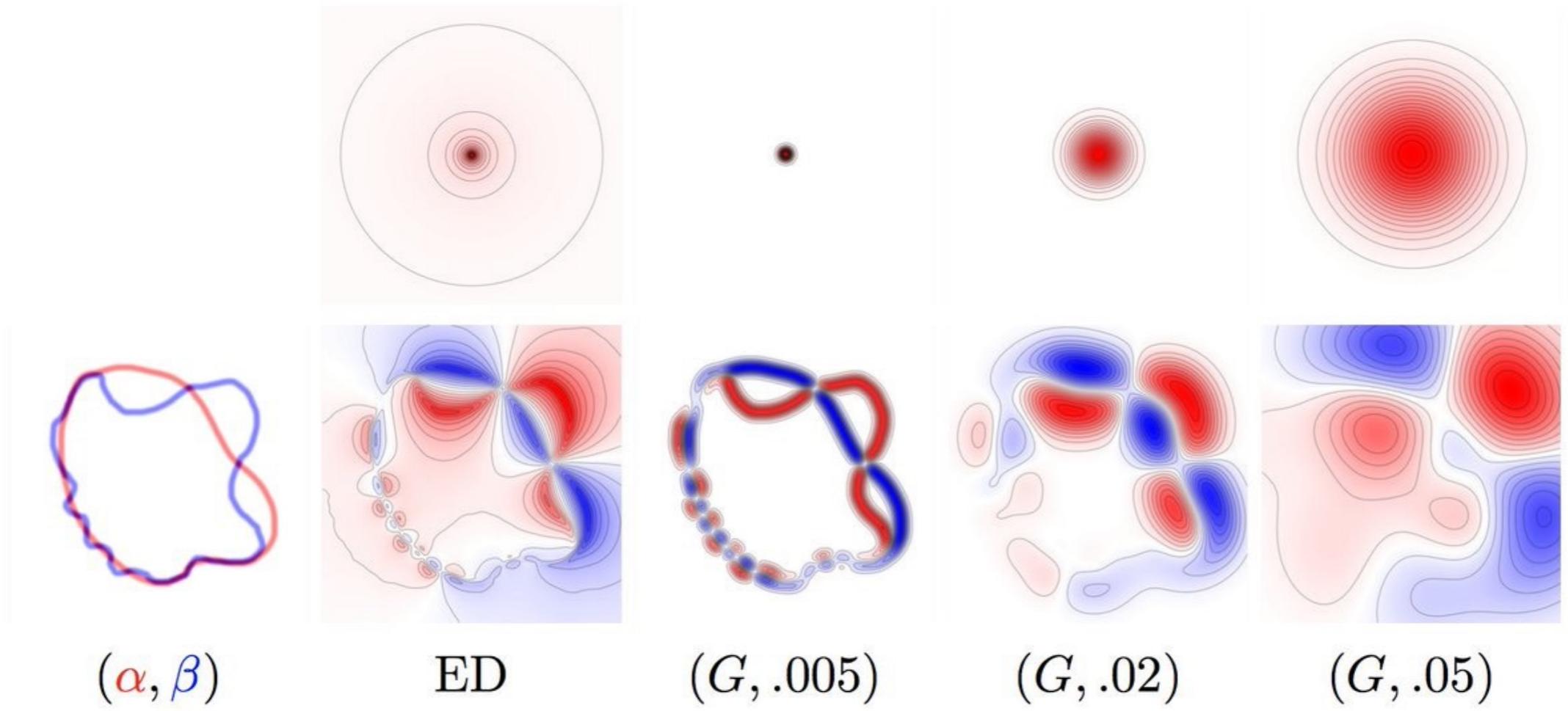
$$\mathcal{D}_\varphi(\alpha|\beta) \stackrel{\text{def.}}{=} \int_{\mathcal{X}} \varphi\left(\frac{d\alpha}{d\beta}\right) d\beta + \varphi'_\infty \alpha^\perp(\mathcal{X})$$



$$\varphi'_\infty = \lim_{x \uparrow +\infty} \varphi(x)/x \in \mathbb{R} \cup \{\infty\}$$

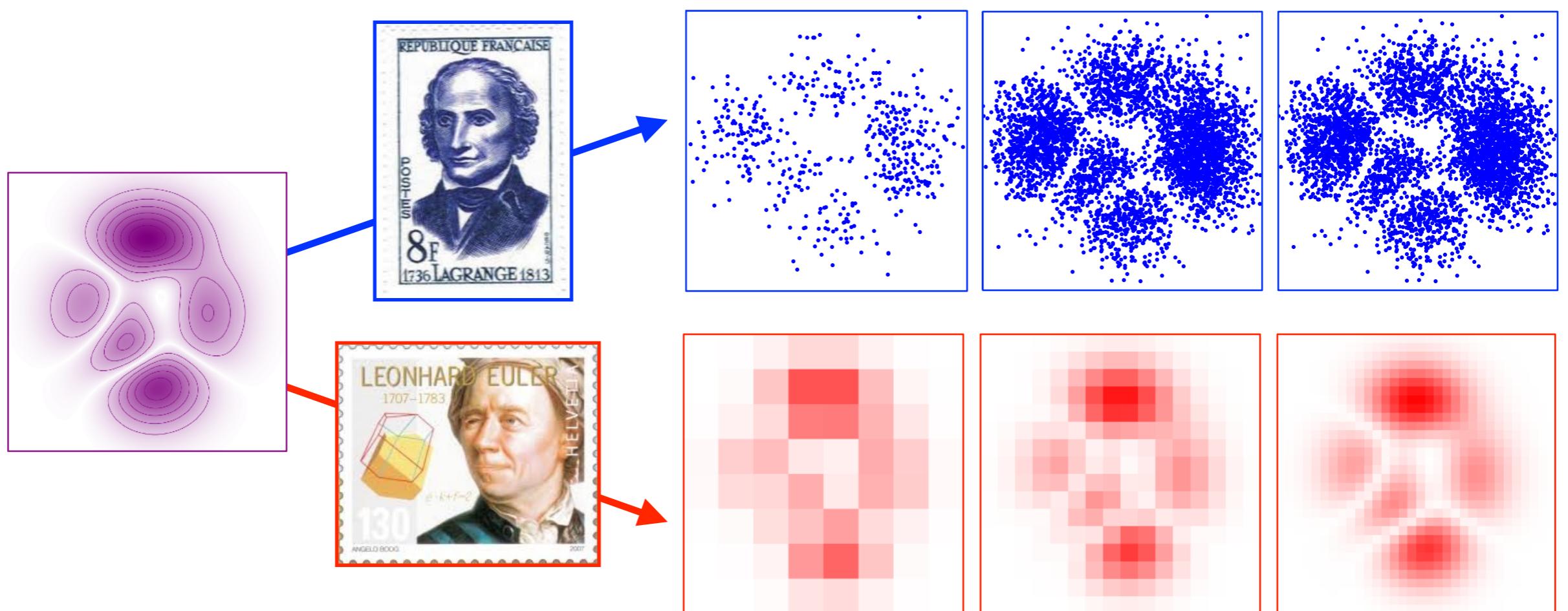


Csiszár divergences, a unifying way to define losses between arbitrary positive measures (discrete & densities). [https://en.wikipedia.org/wiki/F-divergence ...](https://en.wikipedia.org/wiki/F-divergence)

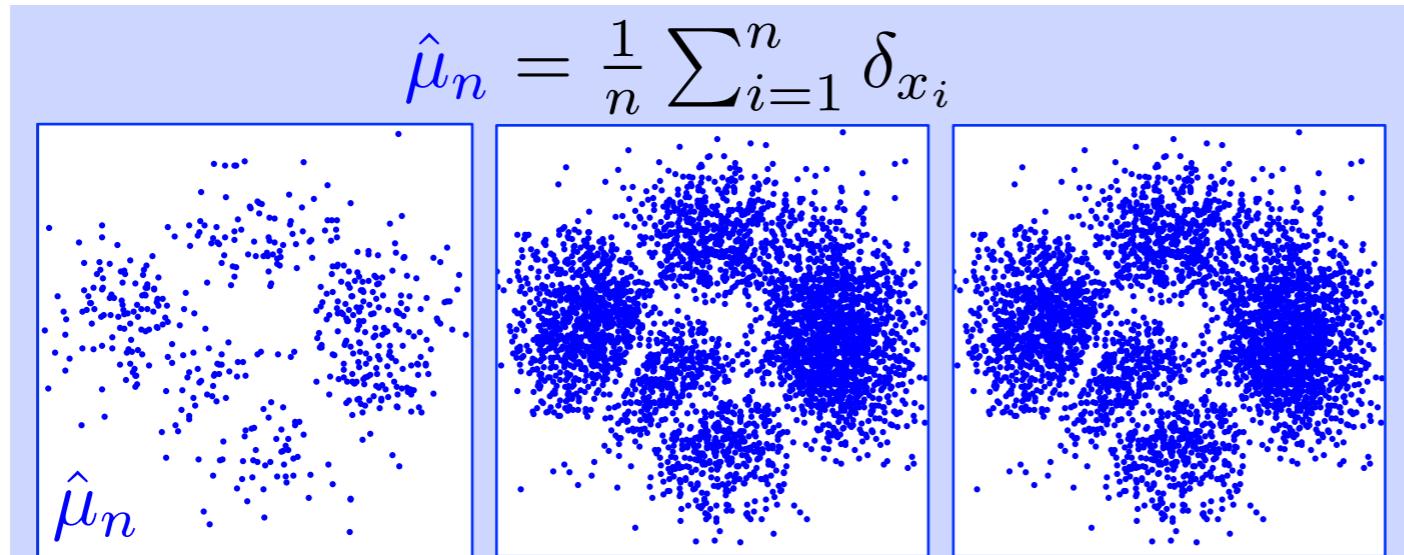


**Figure 8.4:** Top row: display of  $\psi$  such that  $\|\alpha - \beta\|_k = \|\psi \star (\alpha - \beta)\|_{L^2(\mathbb{R}^2)}$ , formally defined over Fourier as  $\hat{\psi}(\omega) = \sqrt{\hat{\varphi}(\omega)}$  where  $k^*(x, x') = \varphi(x - x')$ . Bottom row: display of  $\psi \star (\alpha - \beta)$ .  $(G, \sigma)$  stands for Gaussian kernel of variance  $\sigma^2$  and ED for Energy Distance kernel (in which case  $\psi(x) = 1/\sqrt{\|x\|}$ ).

Comparing distributions using  
 kernels: heavy tail vs. local (Gaussian)  
 kernels.

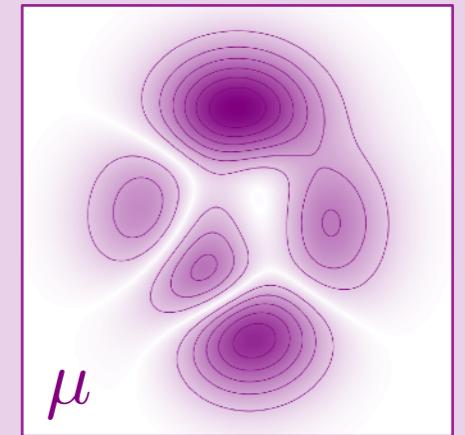


Eulerian vs Lagrangian discretization of a probability distribution. Two philosophies of life.

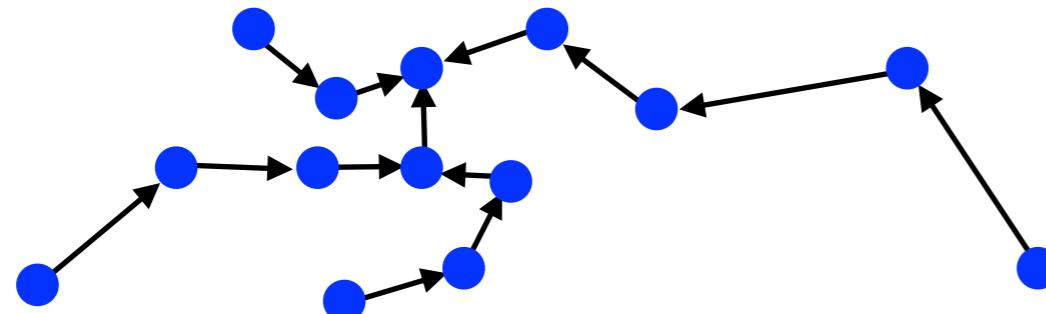


$$\hat{H}(\hat{\mu}_n) \stackrel{\text{def.}}{=} \sum_i \log(\min_{j \neq i} \|x_i - x_j\|)$$

$n \rightarrow +\infty$



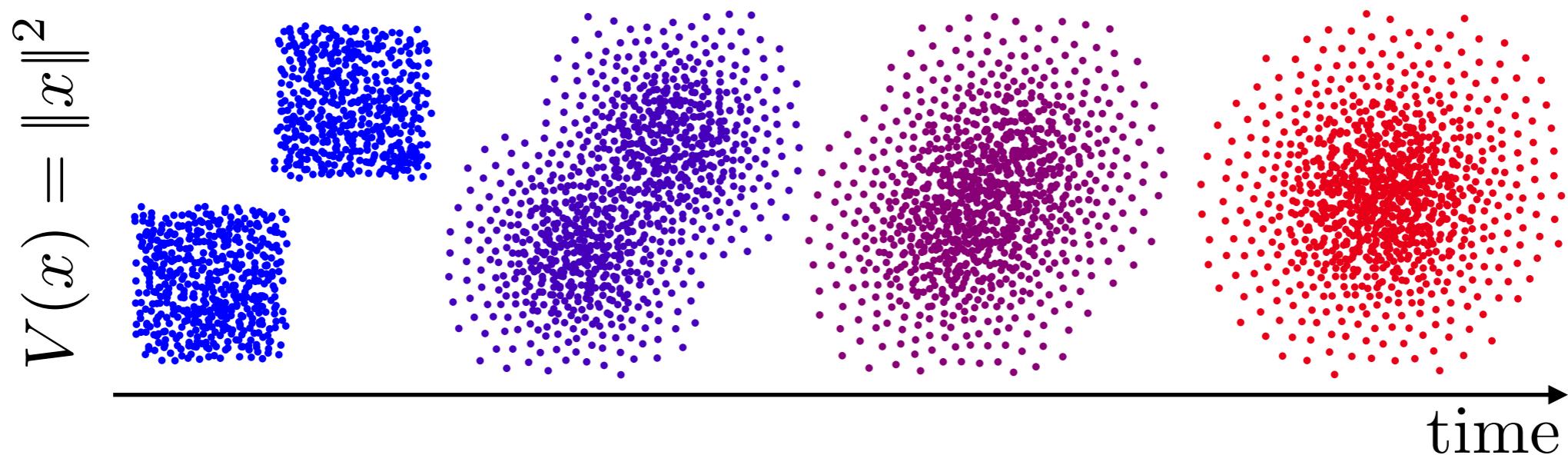
$$H(\mu) \stackrel{\text{def.}}{=} - \int \log\left(\frac{d\mu}{dx}(x)\right) d\mu(x)$$



Summing the log of the nearest neighbors distances: a simple Lagrangian estimator for the entropy.

$$\min_{\rho} E(\rho) \stackrel{\text{def.}}{=} \int V(x)\rho(x)dx + \int \rho(x) \log(\rho(x))dx$$

Wasserstein flow of  $E$ :  $\frac{d\rho_t}{dt} = \Delta\rho_t + \nabla(V\rho_t)$



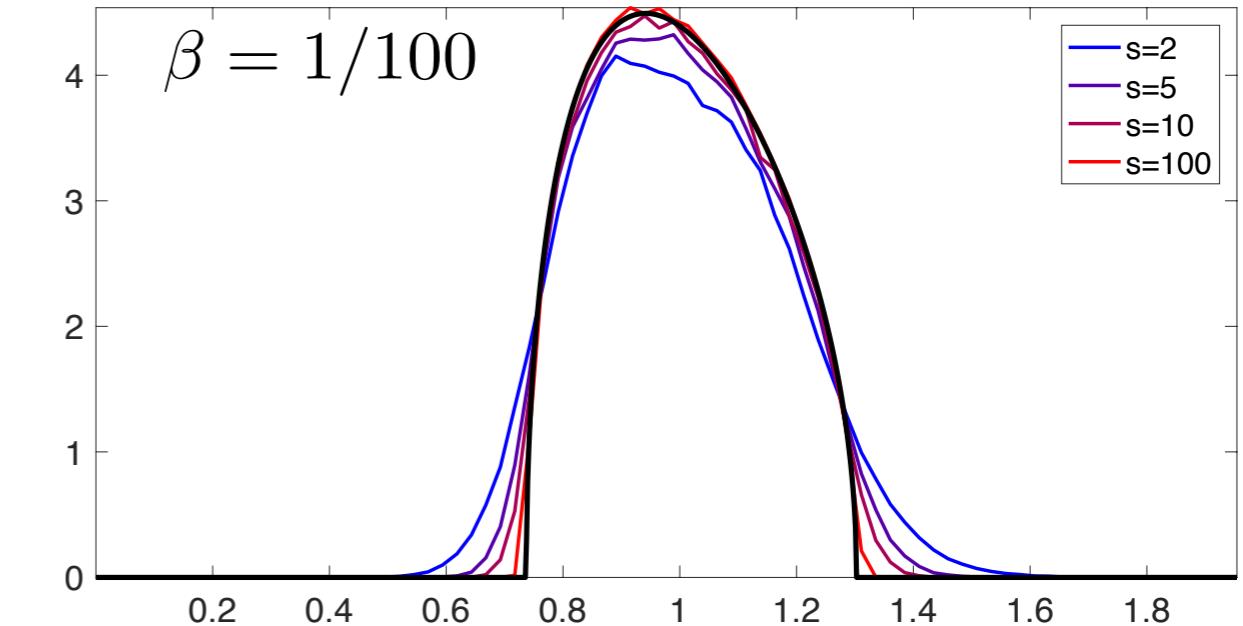
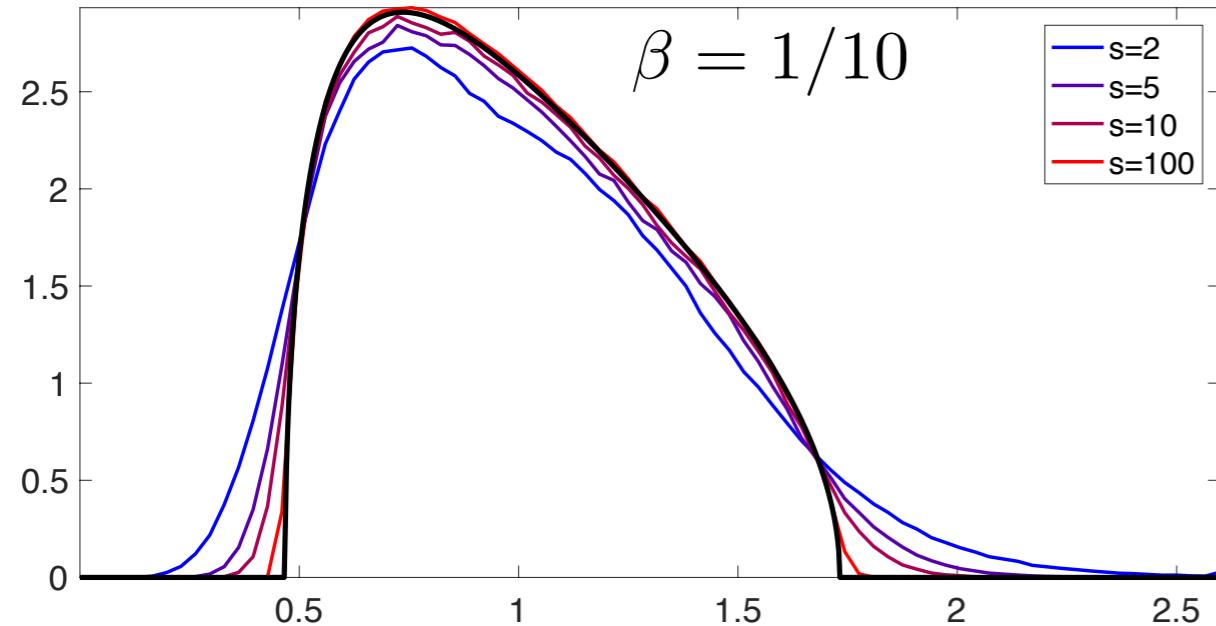
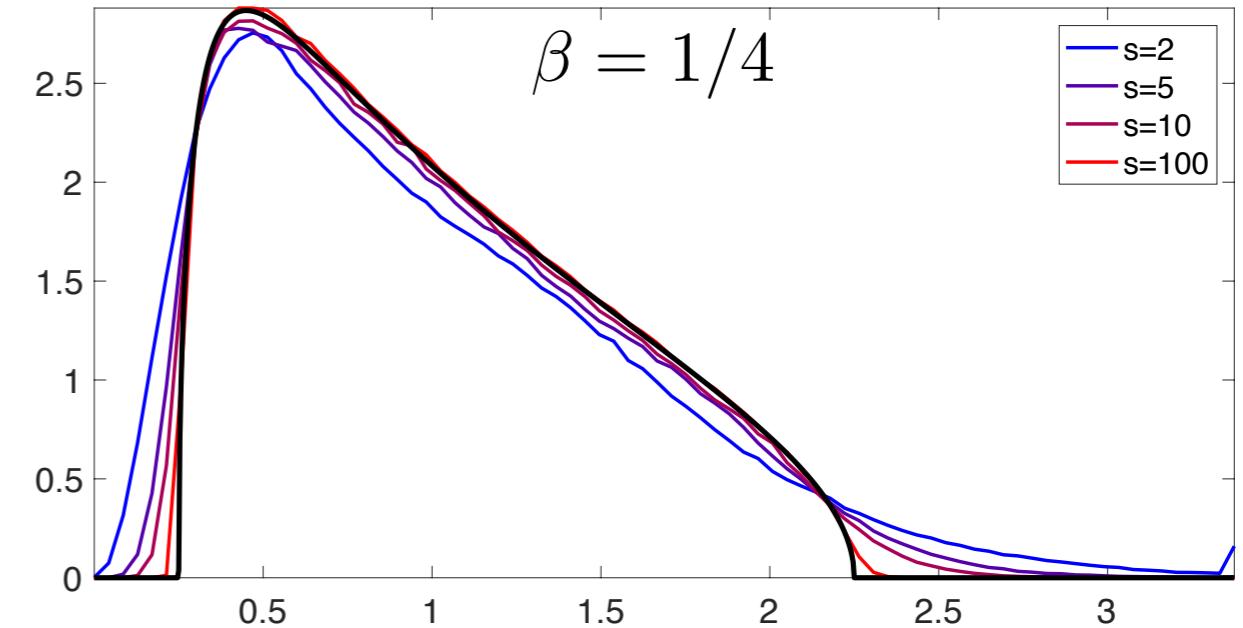
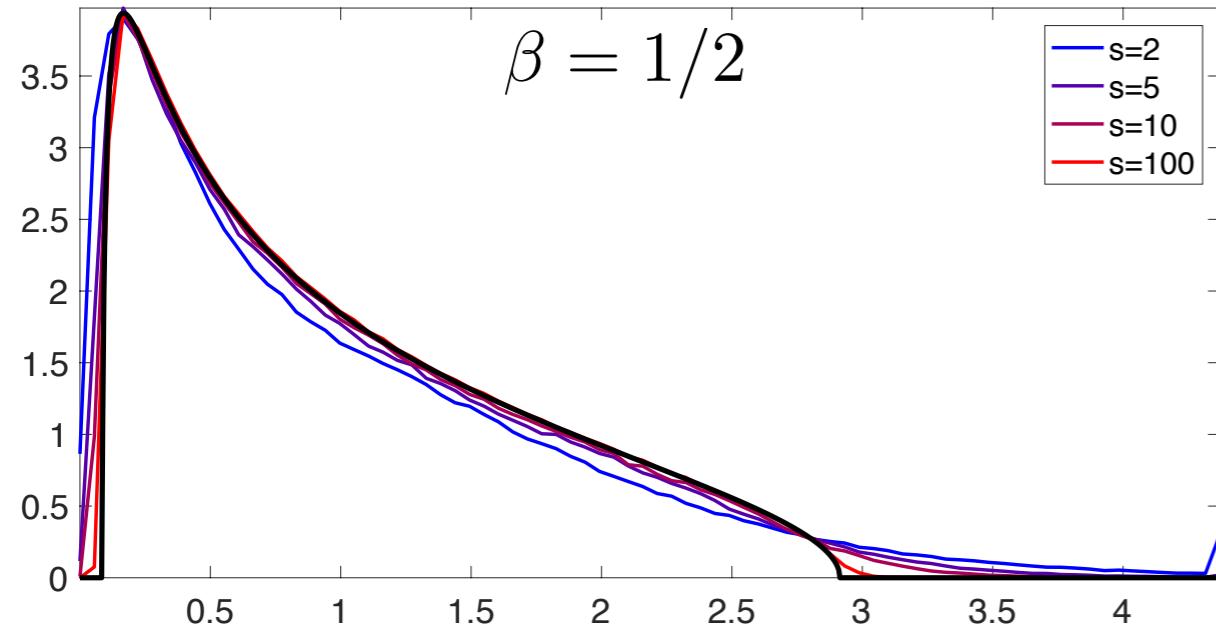
Fokker-Planck equations are Optimal Transport flows of entropies. Nice review by Filippo Santambrogio. <https://arxiv.org/abs/1609.03890>

$$B \sim \frac{1}{\sqrt{P}} \text{randn}(P, s)$$

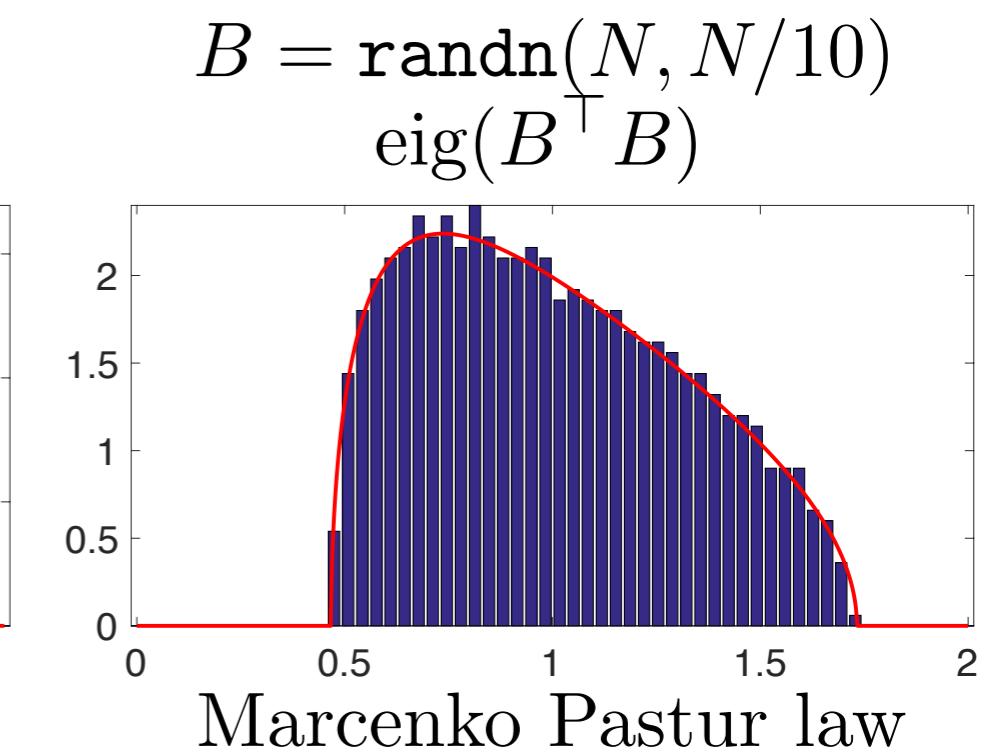
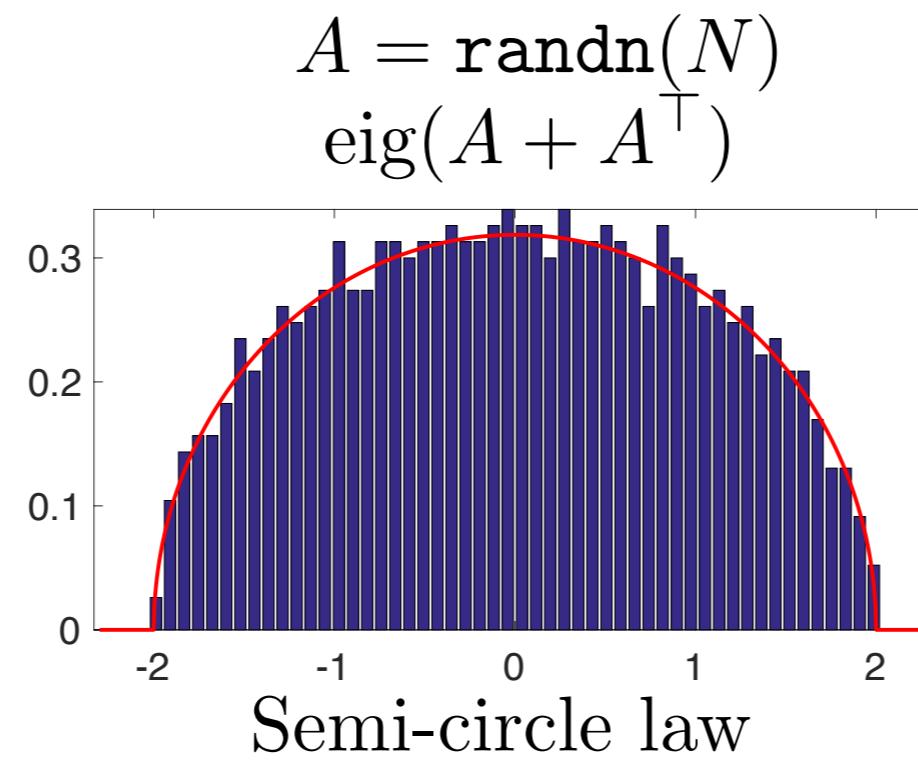
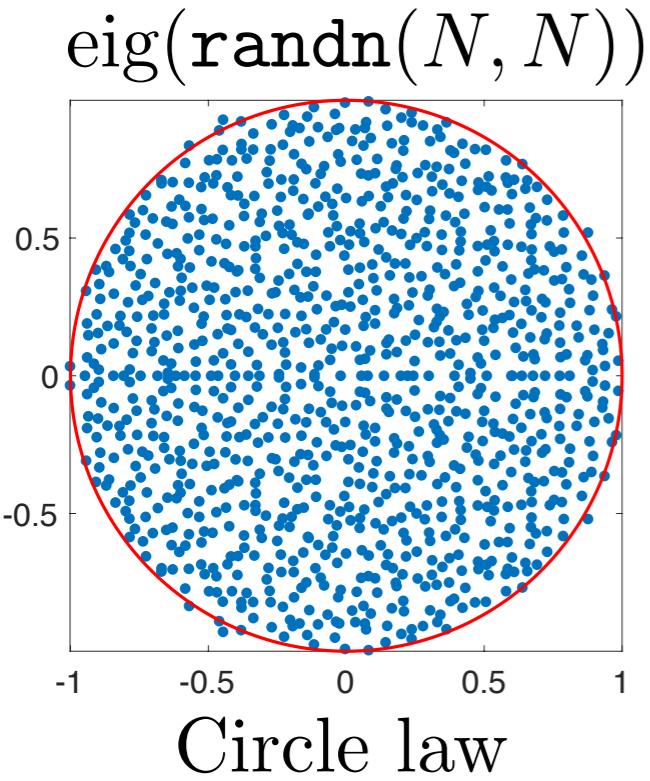
$$\mathbb{P}(\text{eig}(B^\top B) \in [u, v]) \xrightarrow[s \rightarrow +\infty]{\beta \stackrel{\text{def.}}{=} \frac{s}{P}} \int_u^v f_\beta(\lambda) d\lambda$$

$$f_\beta(\lambda) \stackrel{\text{def.}}{=} \frac{1}{2\pi\beta\lambda} \sqrt{(\lambda - \lambda_-)(\lambda_+ - \lambda)} 1_{[\lambda_-, \lambda_+] }(\lambda)$$

[Marcenko-Pastur]  
 $\lambda_\pm \stackrel{\text{def.}}{=} (1 \pm \sqrt{\beta})^2$



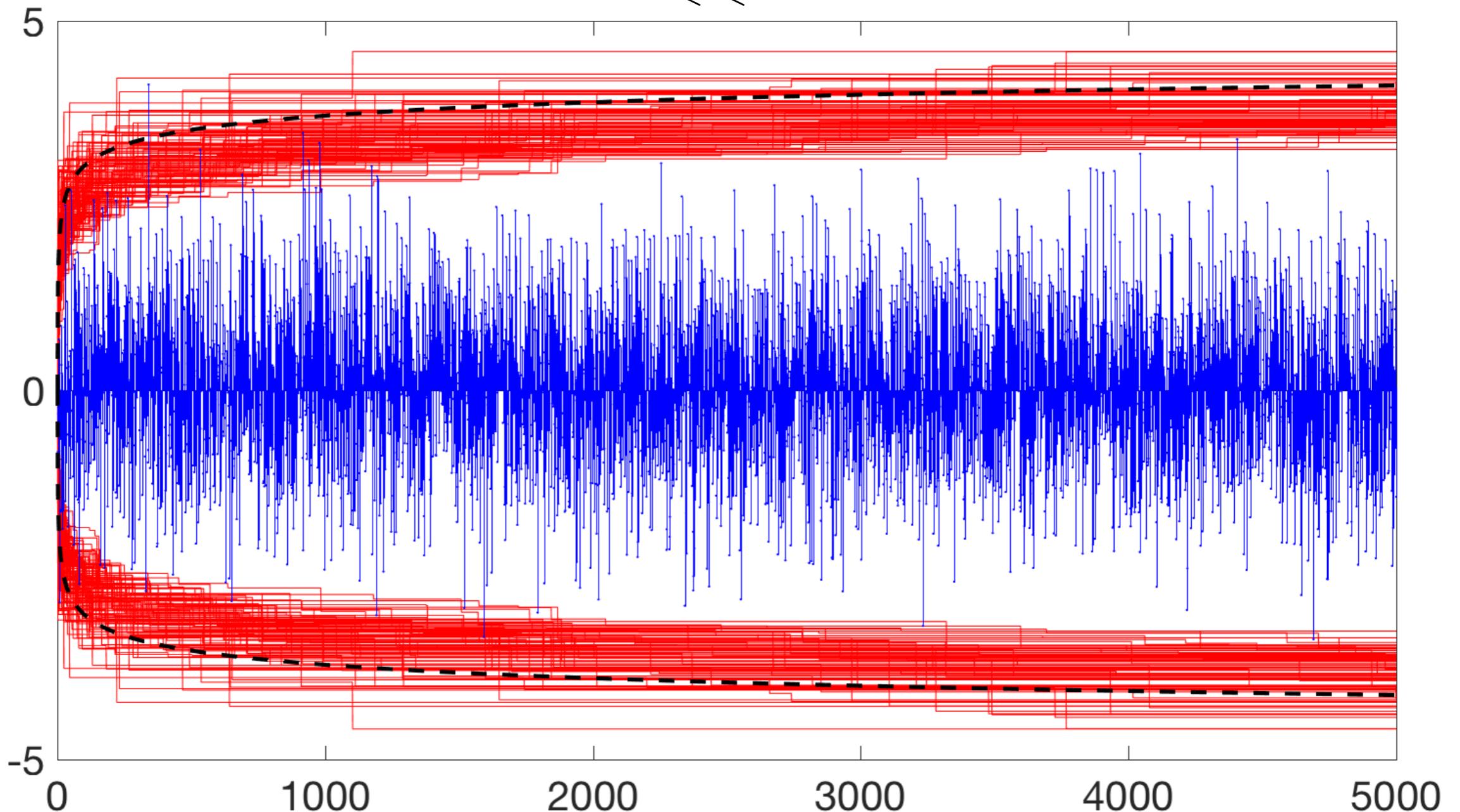
The Marchenko-Pastur law is the distribution limit of eigenvalues of covariance matrices when samples and size grow together. <http://djalil.chafai.net/blog/2011/01/29/the-marchenko-pastur->



The circle, semi-circle and Marčenko–Pastur laws are three simple examples of determinantal processes. [https://terrytao.wordpress.com/2009/08/23/determinantal-processes/ ...](https://terrytao.wordpress.com/2009/08/23/determinantal-processes/)

$$X_i \sim \mathcal{N}(0, 1)$$

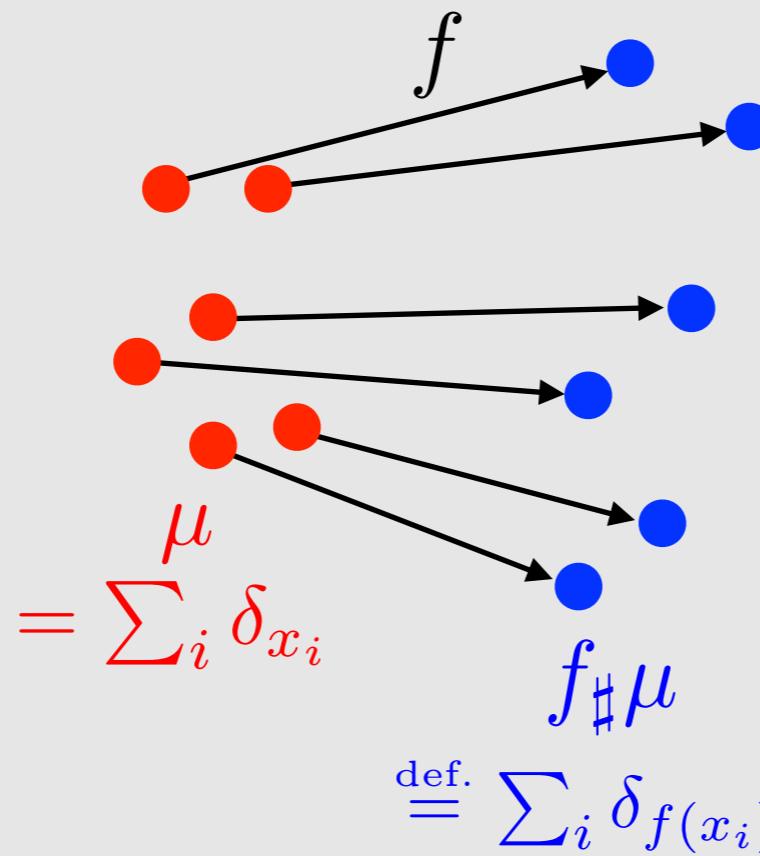
$$\max_{1 \leq i \leq n} |X_i| \sim \sqrt{2 \log(n)}$$



The maximum of  $n$  i.i.d. Gaussians is roughly  $\sqrt{2\log(n)}$ . This is where log factors come from in Compressed Sensing.

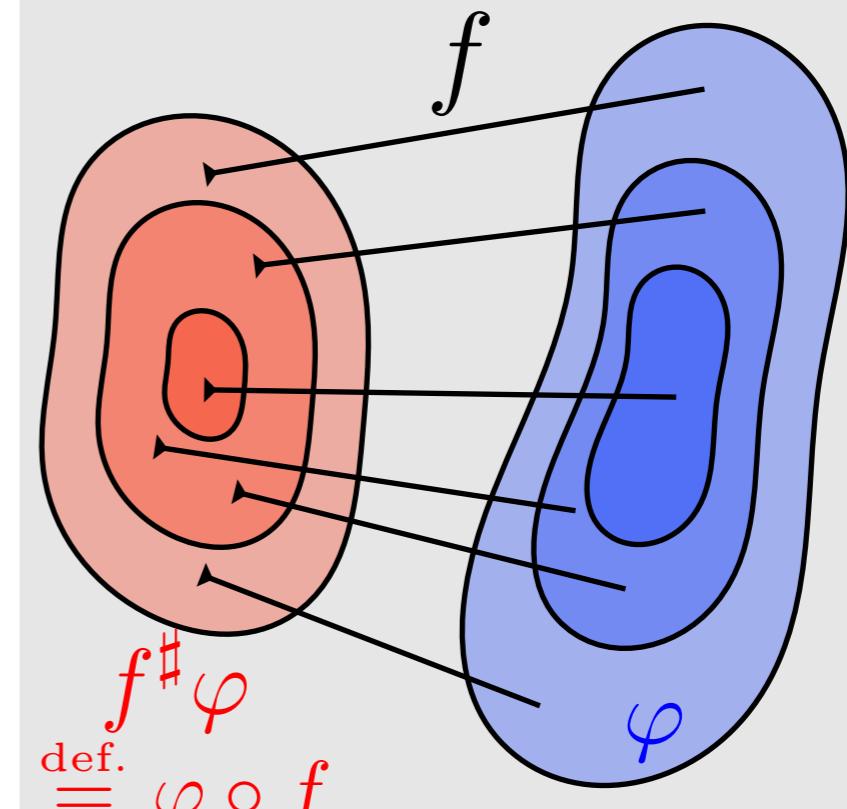
Measures:  
push-forward

$\mathcal{Y}$   
 $\mathcal{X}$   
 $f$



$$f_\# : \mathcal{M}(\mathcal{X}) \rightarrow \mathcal{M}(\mathcal{Y})$$

Functions:  
pull-back



$$f^\# : \mathcal{C}(\mathcal{Y}) \rightarrow \mathcal{C}(\mathcal{X})$$

Remark:  $f^\#$  and  $f_\#$  are adjoints

$$\int_{\mathcal{Y}} \varphi d(f_\# \mu) = \int_{\mathcal{X}} (f^\# \varphi) d\mu$$

Functions should be manipulated by pullback, measures by pushforward. These are adjoint operations.

Random vectors

$$\mathbb{P}(\textcolor{red}{X} \in A)$$

Weak\* convergence:

$\forall$  set  $A$

$$\mathbb{P}(\textcolor{red}{X}_n \in A) \xrightarrow{n \rightarrow +\infty} \mathbb{P}(\textcolor{red}{X} \in A)$$

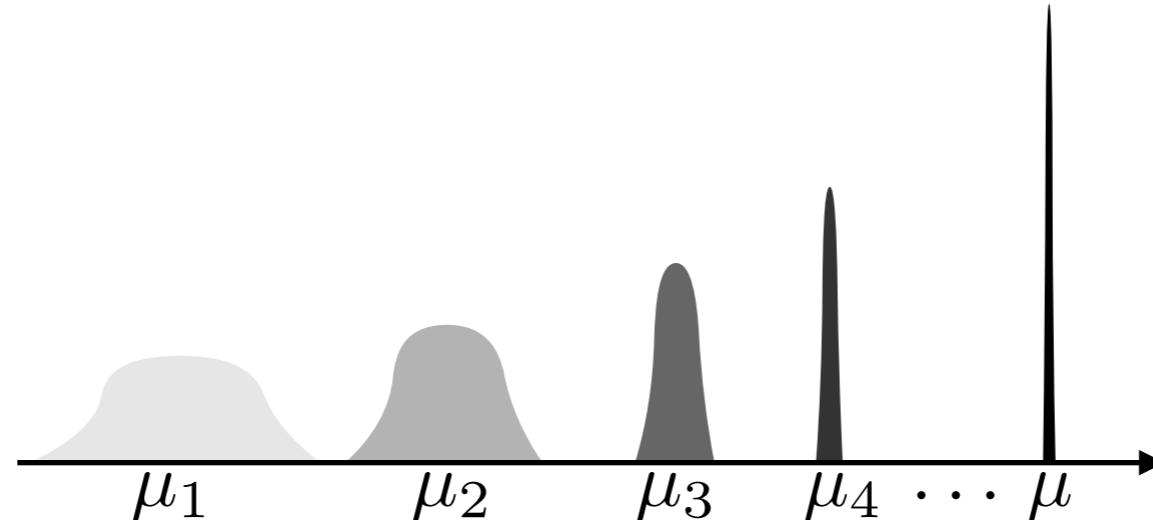
Radon measures

$$\int_A d\mu(x)$$

Convergence in law:

$\forall$  continuous function  $f$

$$\int f d\mu_n \xrightarrow{n \rightarrow +\infty} \int f d\mu$$



Convergence in law of random vectors and  
weak\* convergence of measures are the same. #LetsBeFriends  
[https://en.wikipedia.org/wiki/Convergence\\_of\\_measures](https://en.wikipedia.org/wiki/Convergence_of_measures)  
[https://en.wikipedia.org/wiki/Convergence\\_of\\_random\\_variables](https://en.wikipedia.org/wiki/Convergence_of_random_variables)

In mean

$$\lim_{n \rightarrow +\infty} \mathbb{E}(|X_n - X|^p) = 0$$

Almost sure

$$\mathbb{P}\left(\lim_{n \rightarrow +\infty} X_n = X\right) = 1$$



In probability

$$\forall \varepsilon > 0, \mathbb{P}(|X_n - X| > \varepsilon) \xrightarrow{n \rightarrow +\infty} 0$$



In law

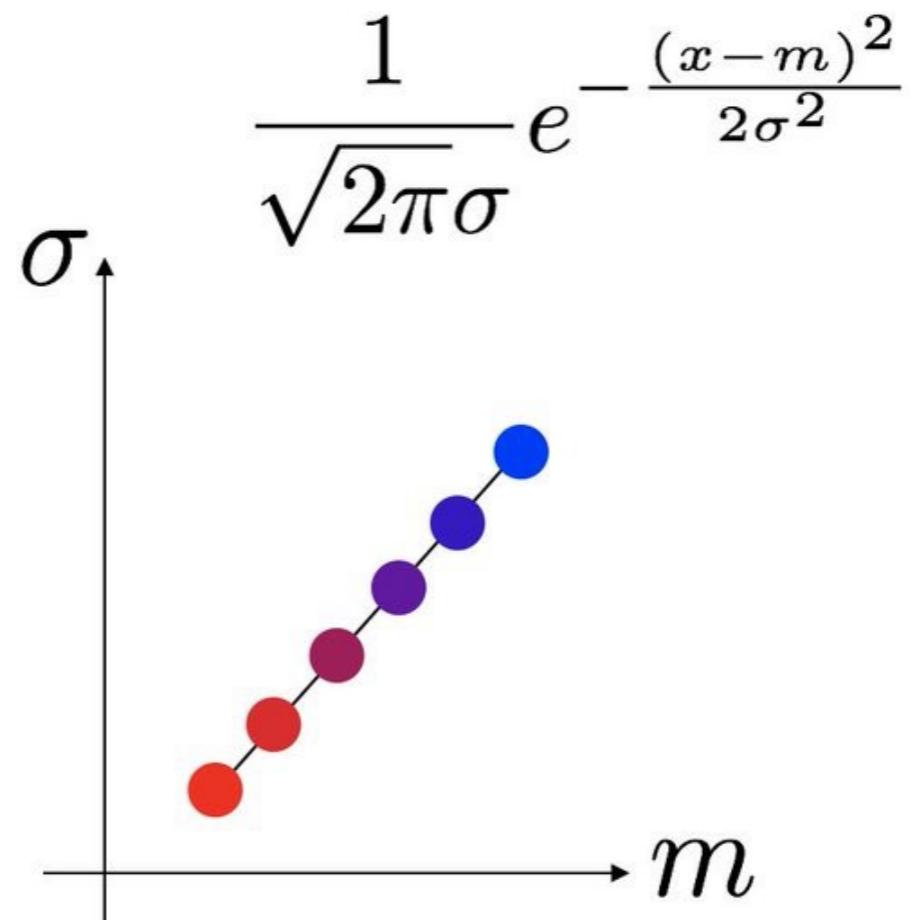
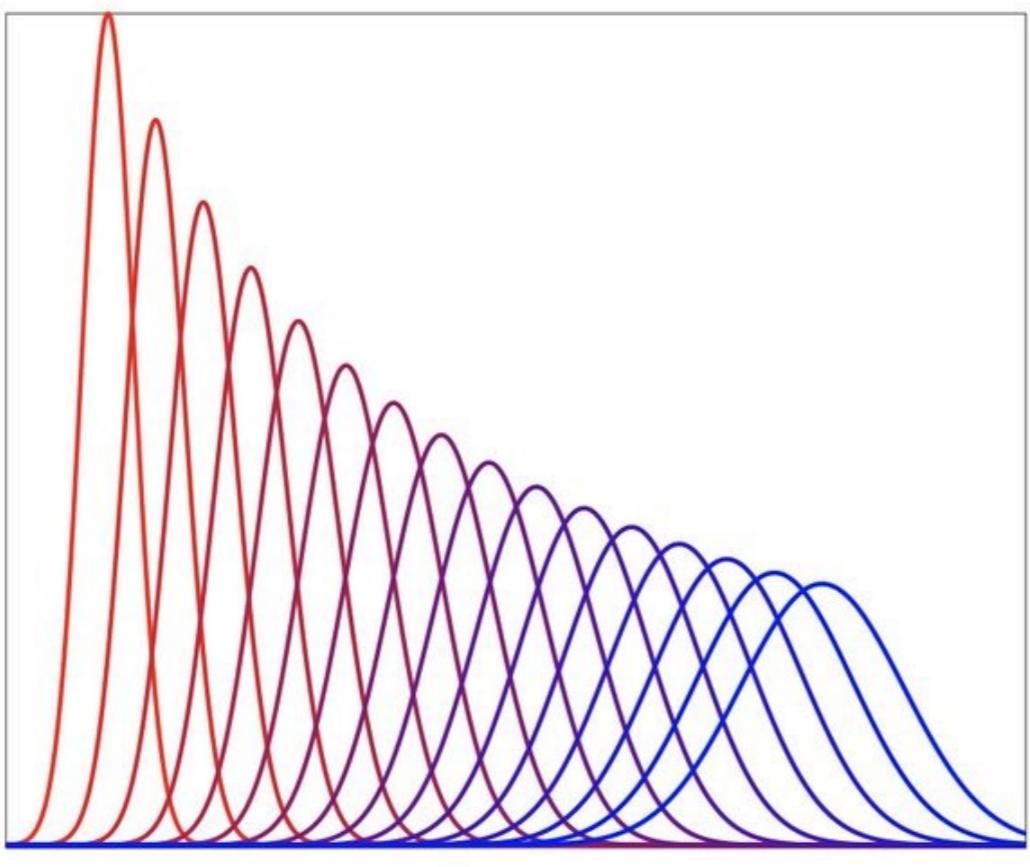
$$\mathbb{P}(X_n \in A) \xrightarrow{n \rightarrow +\infty} \mathbb{P}(X \in A)$$

(the  $X_n$  can be defined on different spaces)

Convergence of random vectors comes with lots of flavors!

[https://en.wikipedia.org/wiki/Convergence\\_of\\_random\\_variables](https://en.wikipedia.org/wiki/Convergence_of_random_variables)

# Optimal Transport



The Optimal Transport geometry of 1-D Gaussians is flat in  
the (mean,std) plane. [#DoesNotWorkIn2D](#)

Dual norms: (aka Integral Probability Metrics)

$$\|\alpha - \beta\|_B \stackrel{\text{def.}}{=} \max \left\{ \int_{\mathcal{X}} f(x)(d\alpha(x) - d\beta(x)) ; f \in B \right\}$$

Wasserstein 1:  $B = \{f ; \|\nabla f\|_\infty \leq 1\}$ .

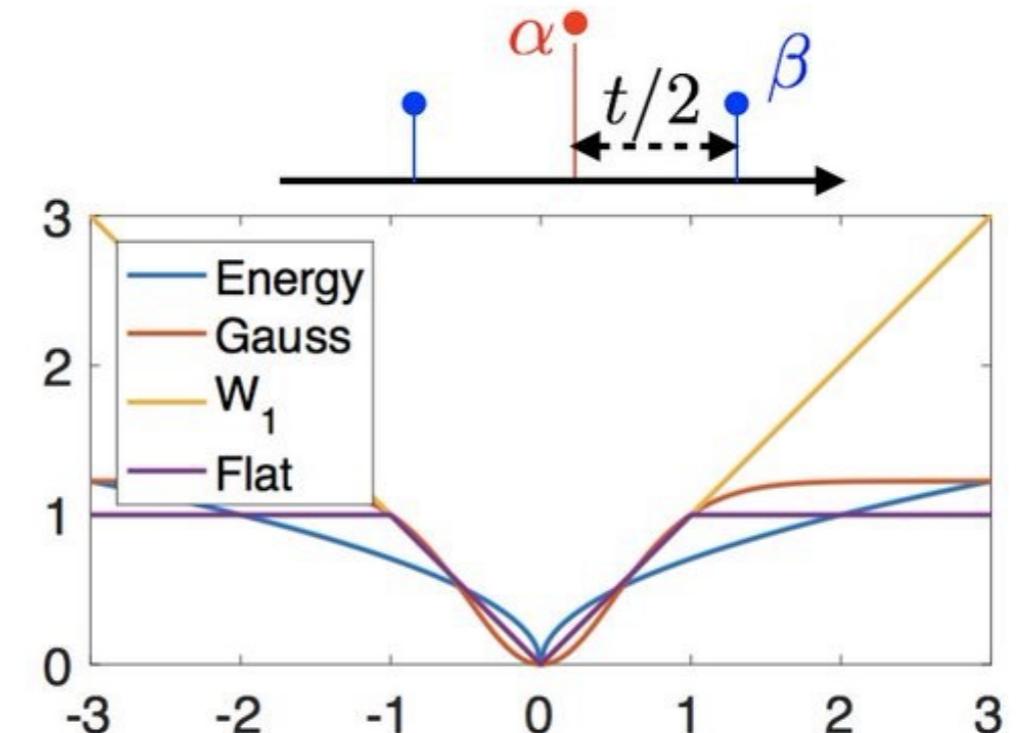
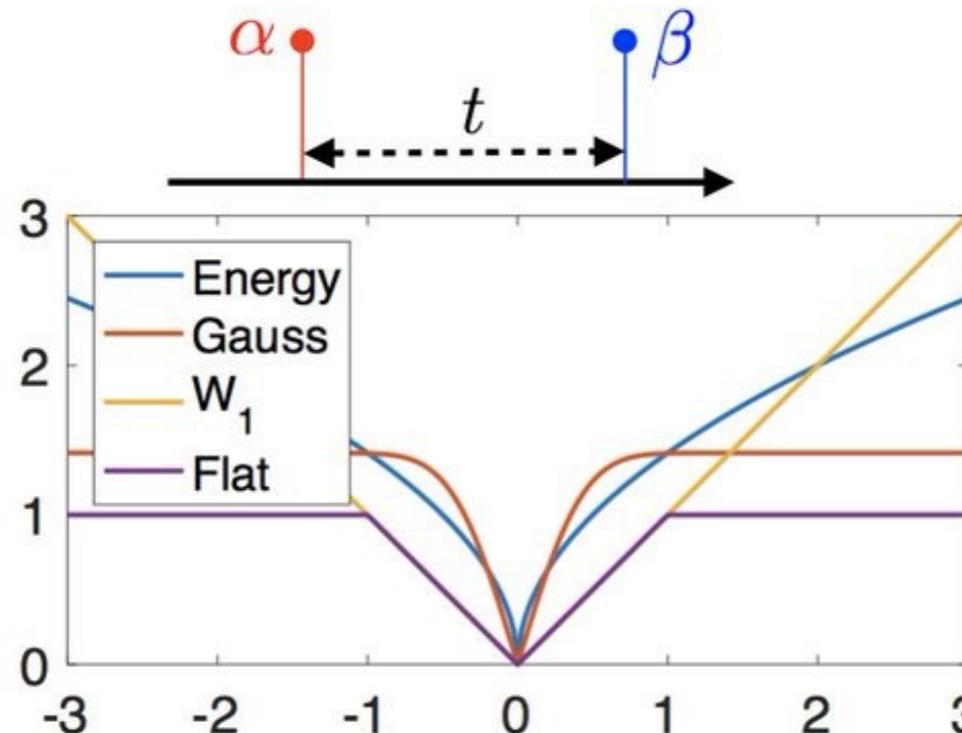
Flat norm:  $B = \{f ; \|f\|_\infty \leq 1, \|\nabla f\|_\infty \leq 1\}$ .

RKHS:  $B = \{f ; \|f\|_k^2 \leq 1\}$ .

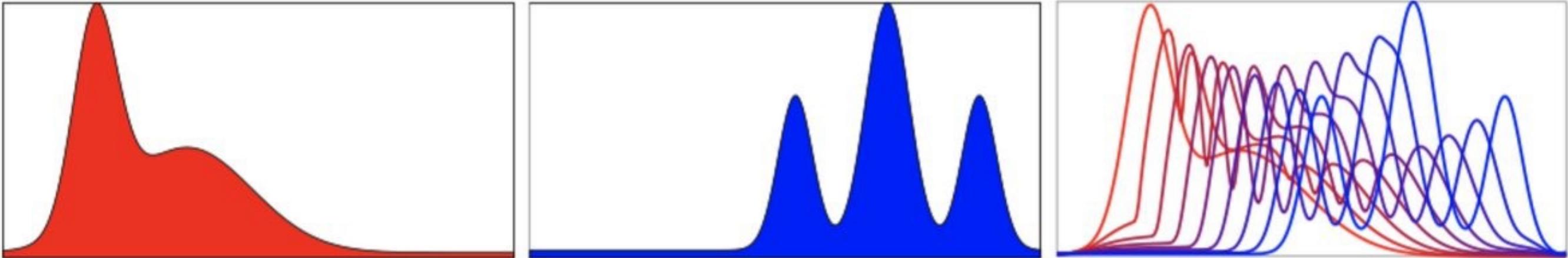
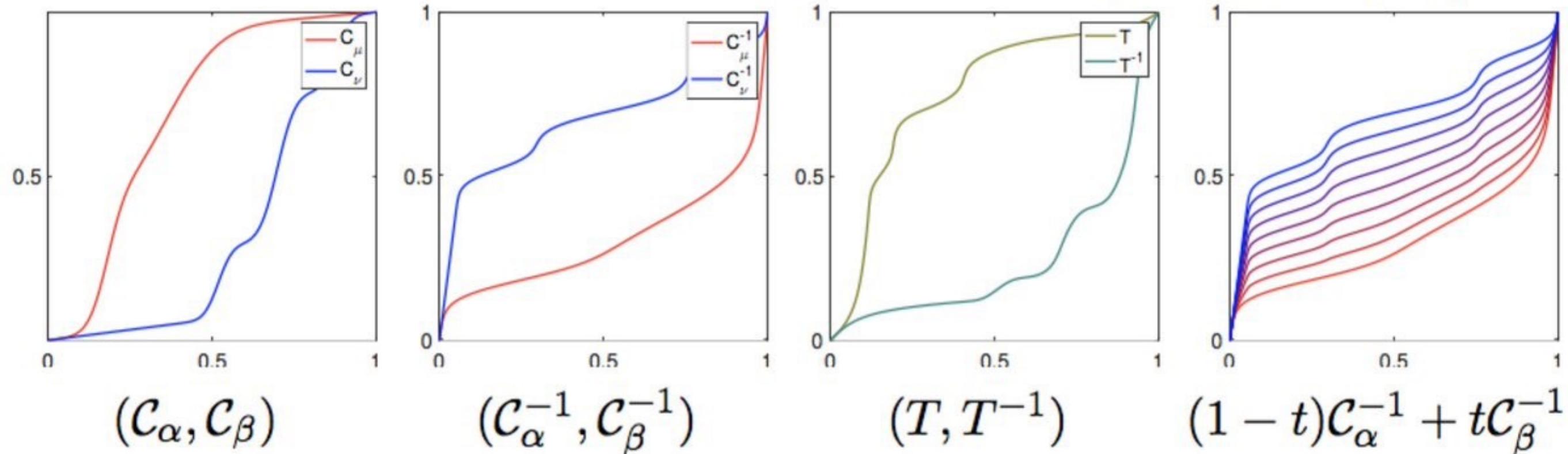
$$\|\alpha - \beta\|_B^2 = \int k(x, x') d\alpha(x) d\alpha(x') + \int k(x, x') d\beta(y) d\beta(y') - 2 \int k(x, y) d\alpha(x) d\beta(y)$$

Energy distance:  $k(x, y) = -\|x - y\|_2^2$

Gaussian:  $k(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}}$



Dual weak norms (aka Integral Probability Metrics): a unifying way to define norms over measures compatible with weak convergence.


 $\alpha$ 
 $\beta$ 
 $((t\text{Id} + (1-t)\text{Id}) \sharp \alpha)$ 


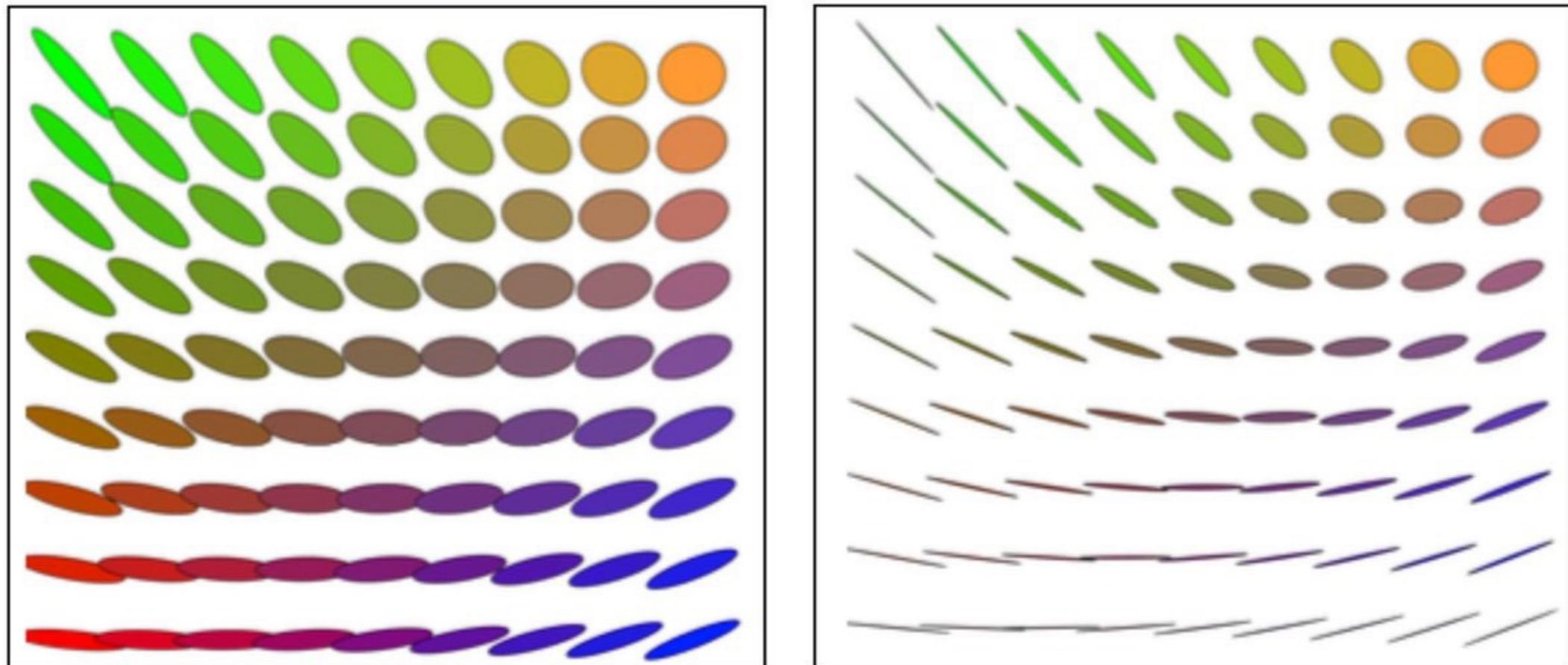
Displacement interpolation aka Optimal Transport in 1D is equivalent to interpolating the \*inverse\* cumulative functions.

**Remark 2.11** (Distance between Gaussians). If  $\alpha = \mathcal{N}(m_\alpha, C_\alpha)$  and  $\beta = \mathcal{N}(m_\beta, C_\beta)$ , then one can show that

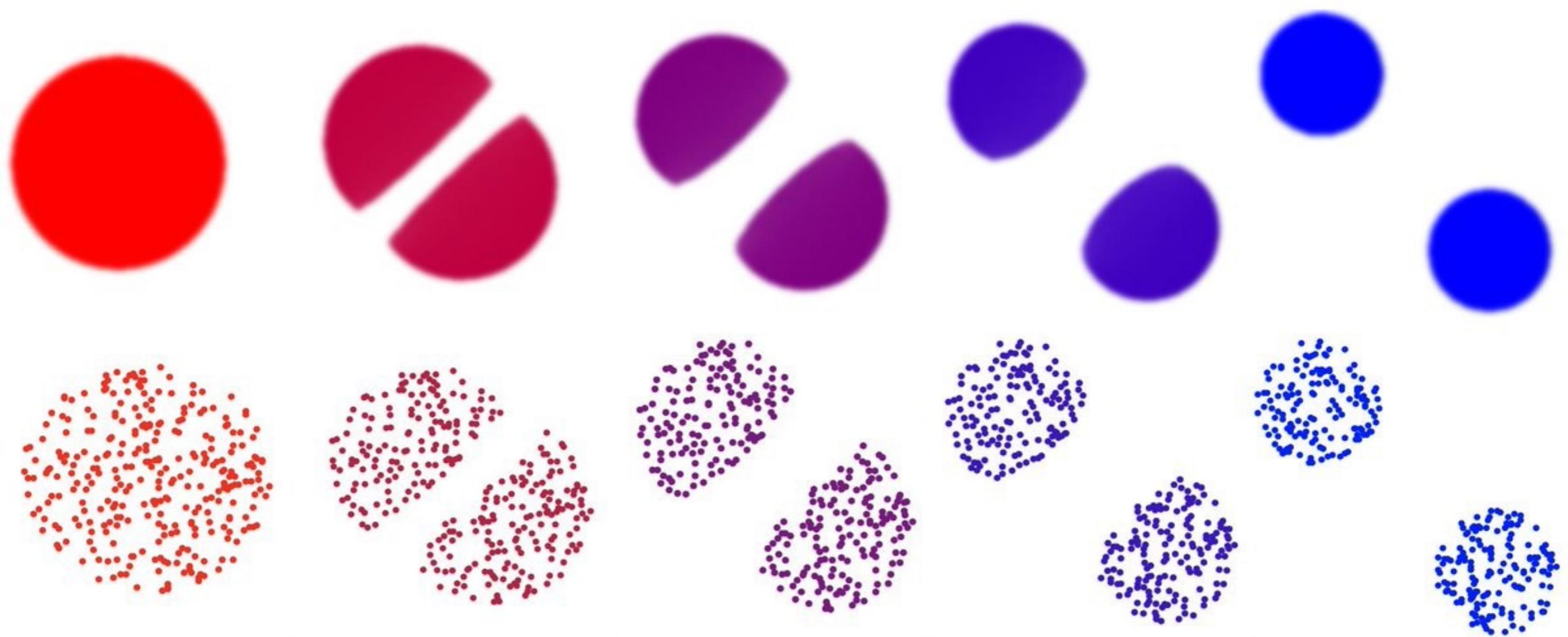
$$\mathcal{W}_2^2(\alpha, \beta) = \|m_\alpha - m_\beta\|^2 + \mathcal{B}(C_\alpha, C_\beta)^2 \quad (2.19)$$

where  $\mathcal{B}$  is the so-called Bures metric

$$\mathcal{B}(C_\alpha, C_\beta)^2 \stackrel{\text{def.}}{=} \text{tr} \left( C_\alpha + C_\beta - 2(C_\alpha^{1/2} C_\beta C_\alpha^{1/2})^{1/2} \right) \quad (2.20)$$

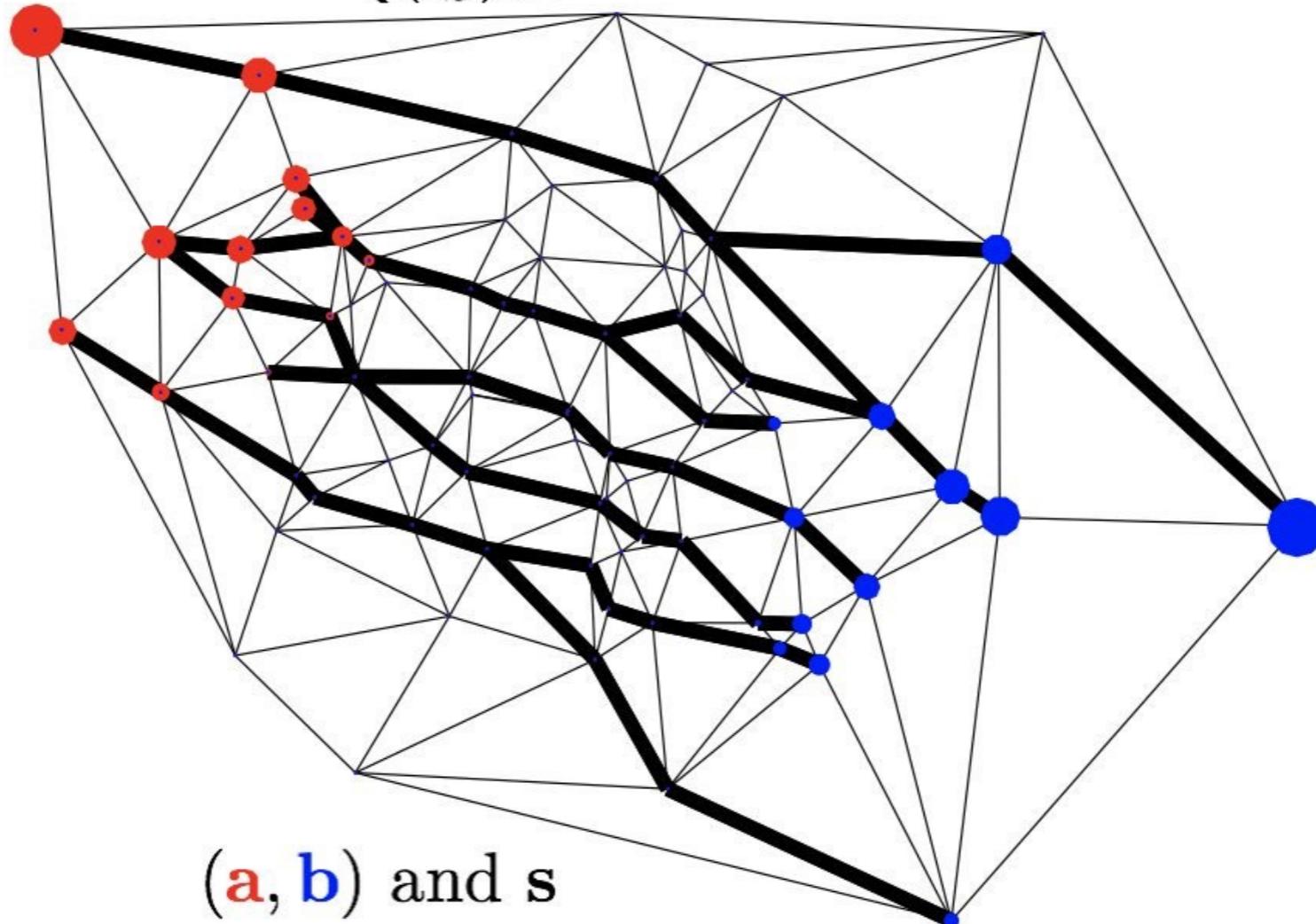


Bures distance on PSD matrices, aka Optimal Transport between Gaussians,  
handles nicely singular matrices [https://en.wikipedia.org/wiki/Bures\\_metric ...](https://en.wikipedia.org/wiki/Bures_metric)

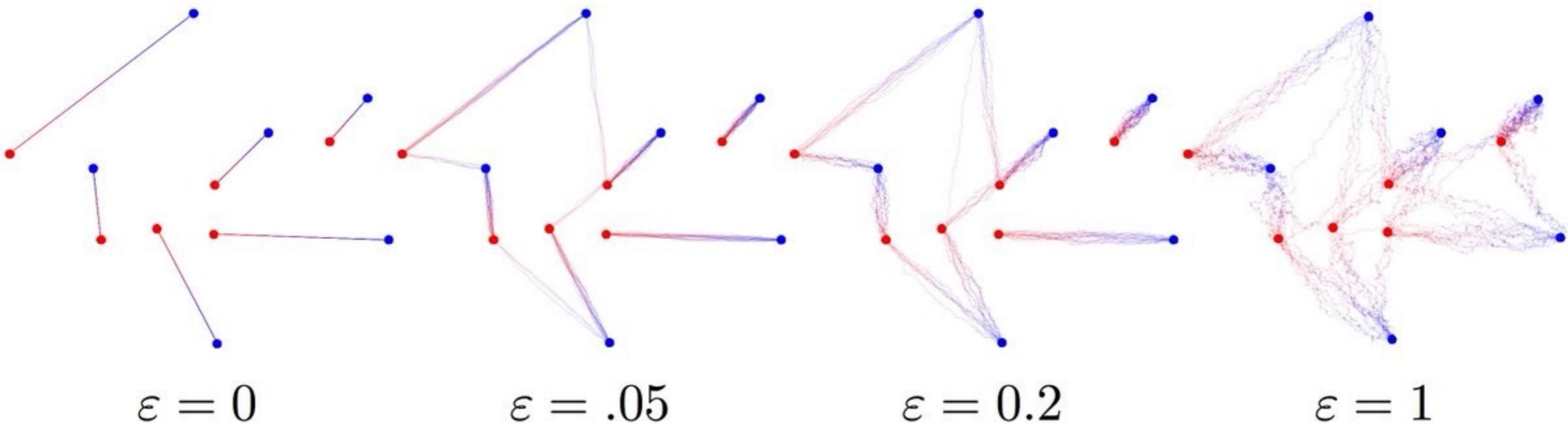


Displacement interpolation aka Optimal Transport: discrete vs continuous.

$$W_1(\mathbf{a}, \mathbf{b}) = \min_{\mathbf{s} \in \mathbb{R}_{+}^{\mathcal{E}}} \left\{ \sum_{(i,j) \in \mathcal{E}} \mathbf{w}_{i,j} s_{i,j} : \text{div}(\mathbf{s}) = \mathbf{a} - \mathbf{b} \right\}$$



Wasserstein-1 distance (norm!) between measures on graphs: equivalent to a min-cost flow. Bold black indicates edges where mass is flowing.



From deterministic to stochastic matching: Schrödinger problem.  
See Leonard's survey for the maths. <https://arxiv.org/abs/1308.0215>

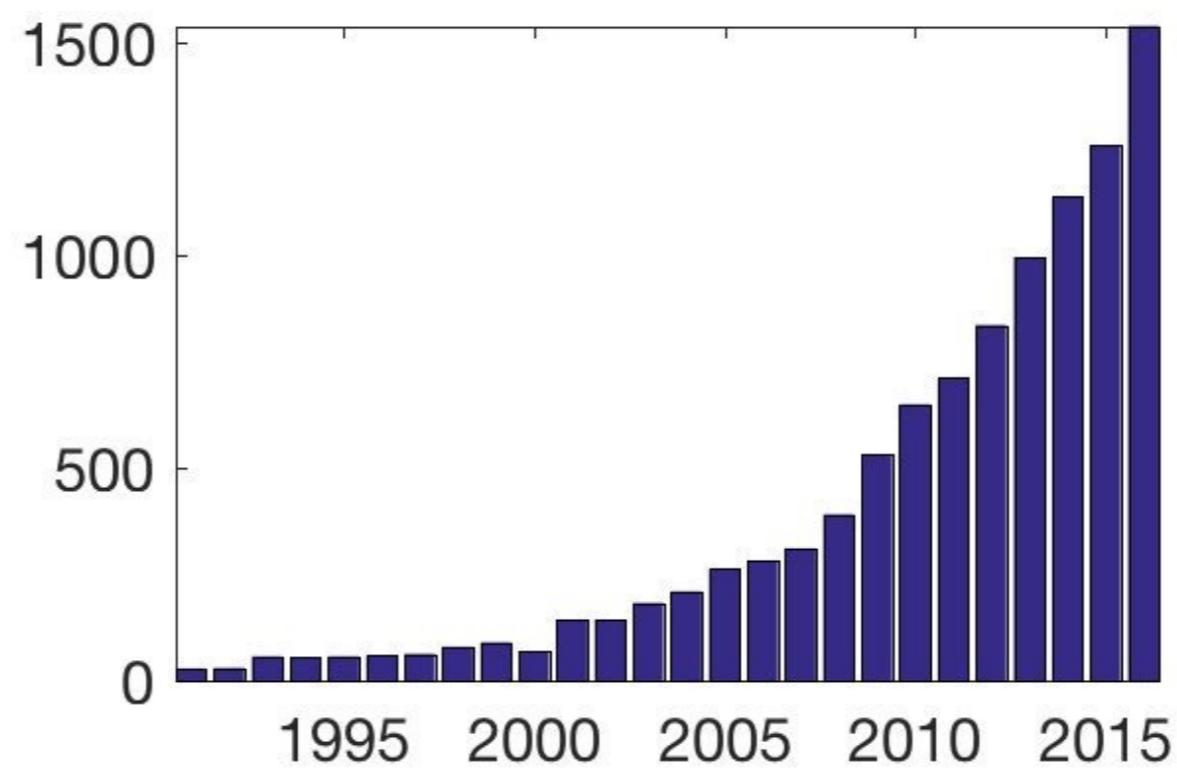
If  $(\mu, \nu)$  have densities bounded from below/above by  $0 < a < b < +\infty$ ,

$$b^{-1/2} \|\mu - \nu\|_{H^{-1}} \leq W_2(\mu, \nu) \leq a^{-1/2} \|\mu - \nu\|_{H^{-1}}$$

Wasserstein-2 and Energy Distance (Sobolev  $H^{-1}$ ) are equivalent for bounded densities. Sharp constant by R. Peyre <https://arxiv.org/abs/1104.4631v2>

$$W_c(\mu, \nu) = \sup_{f,g} \left\{ \int_X f d\mu + \int_Y g d\nu ; \forall (x,y), f(x) + g(y) \leq c(x,y) \right\}$$
$$W_c^*(f,g) = \inf_{x,y} c(x,y) - f(x) + g(y)$$

The Legendre transform of the Wasserstein distance has a nice expression! #ProbablyUseless



Google Scholar citations for "Optimal Transport". Soon trendier than Deep Learning!

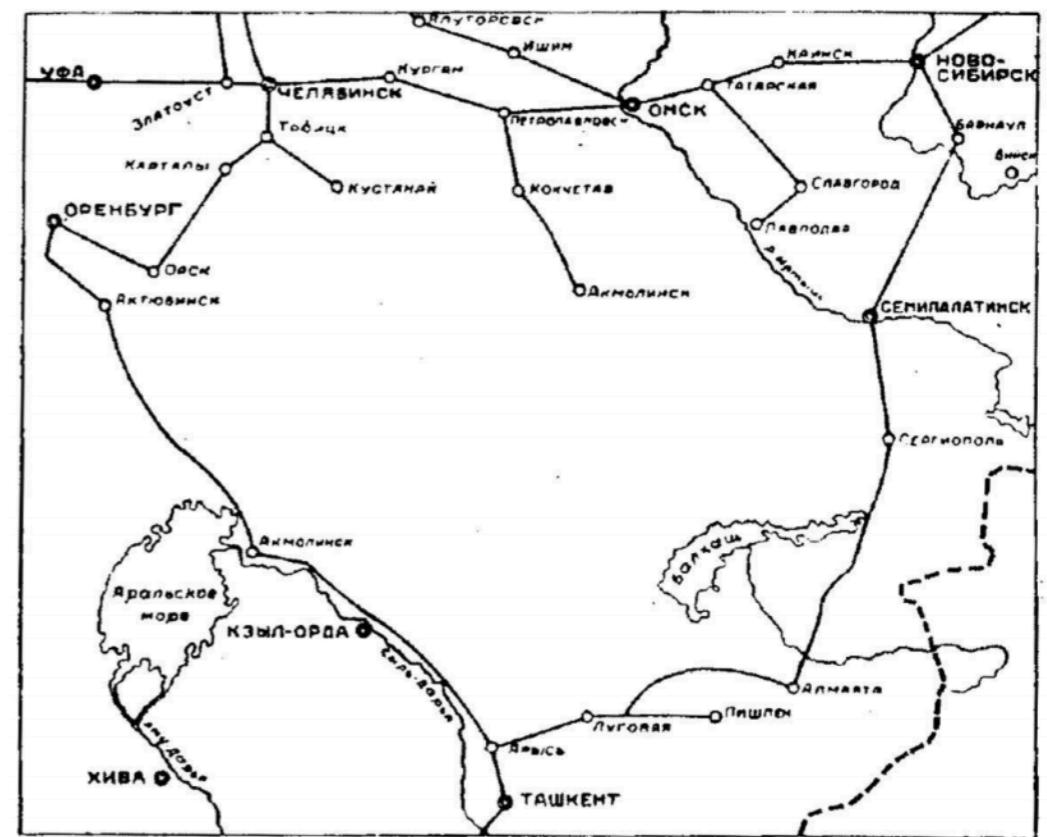
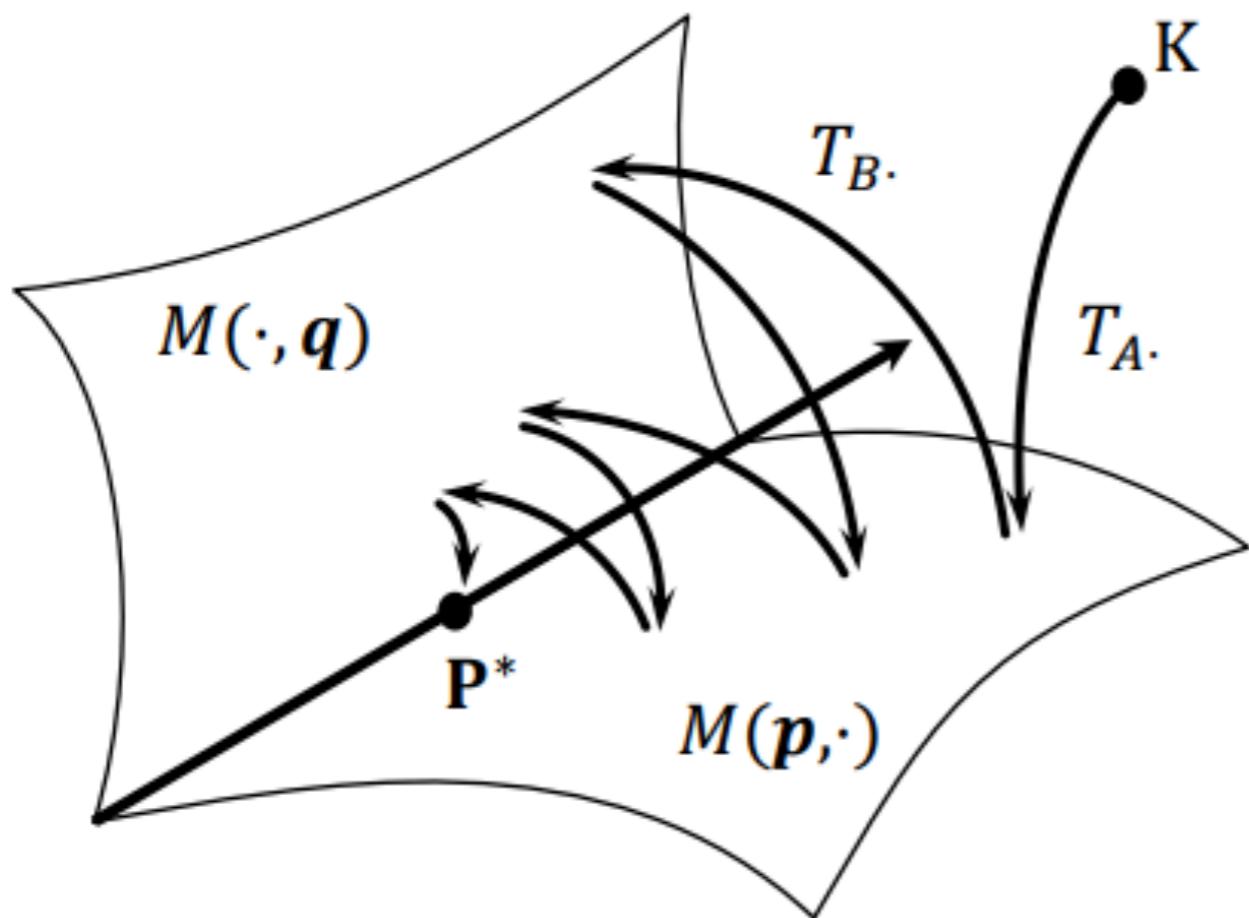
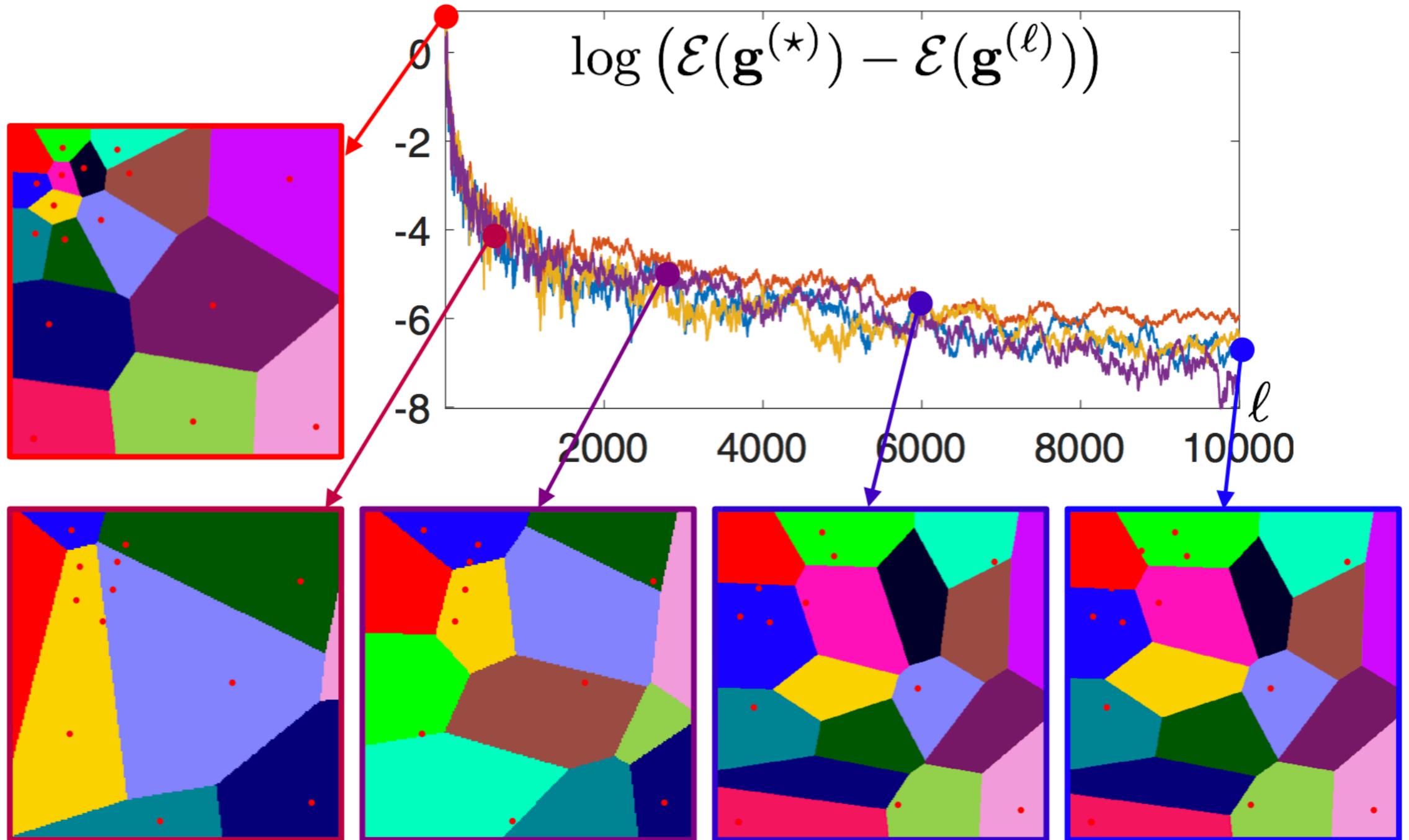
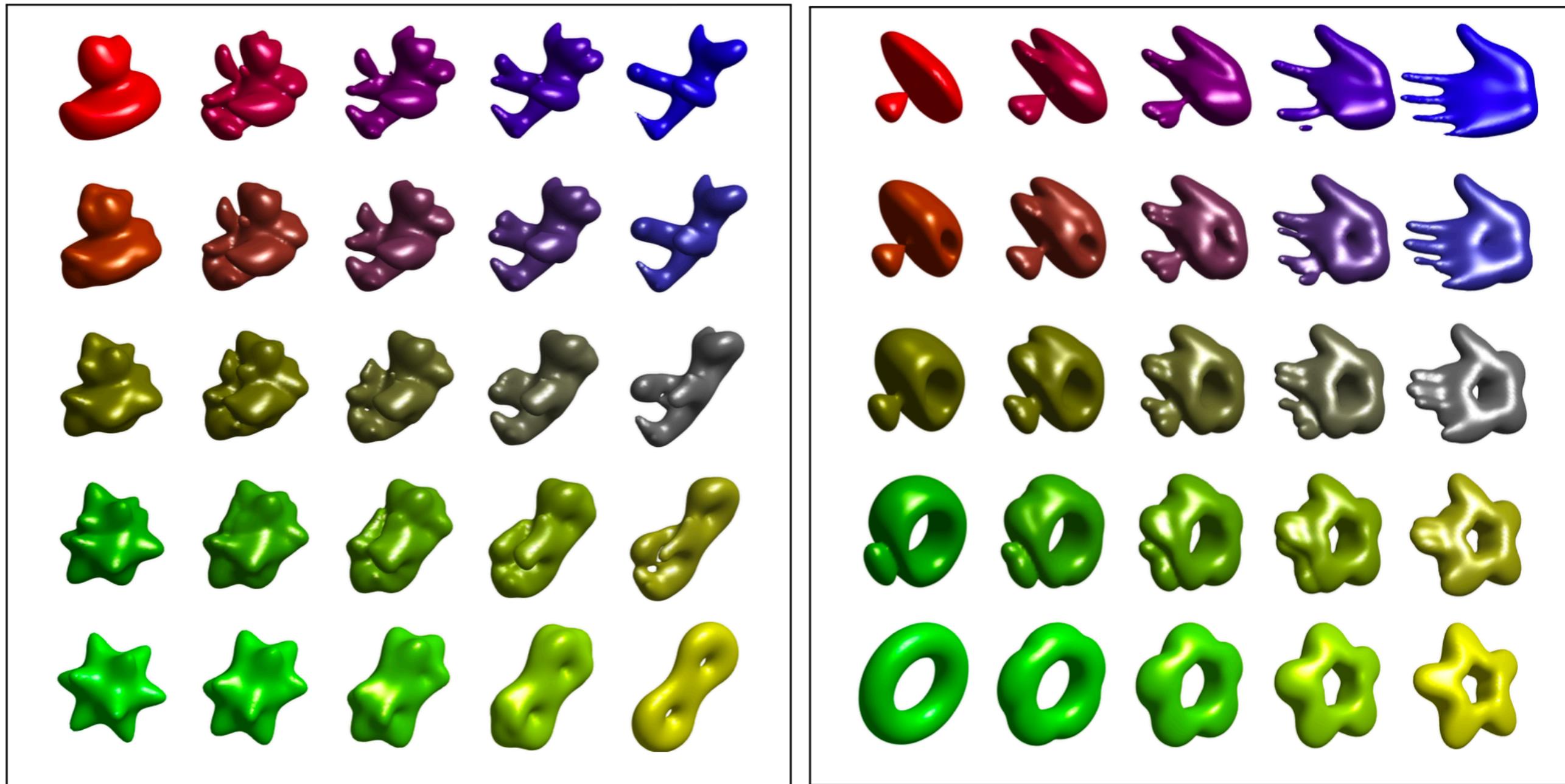


Figure 1: Figure from Tolstoi [1930] to illustrate a negative cycle

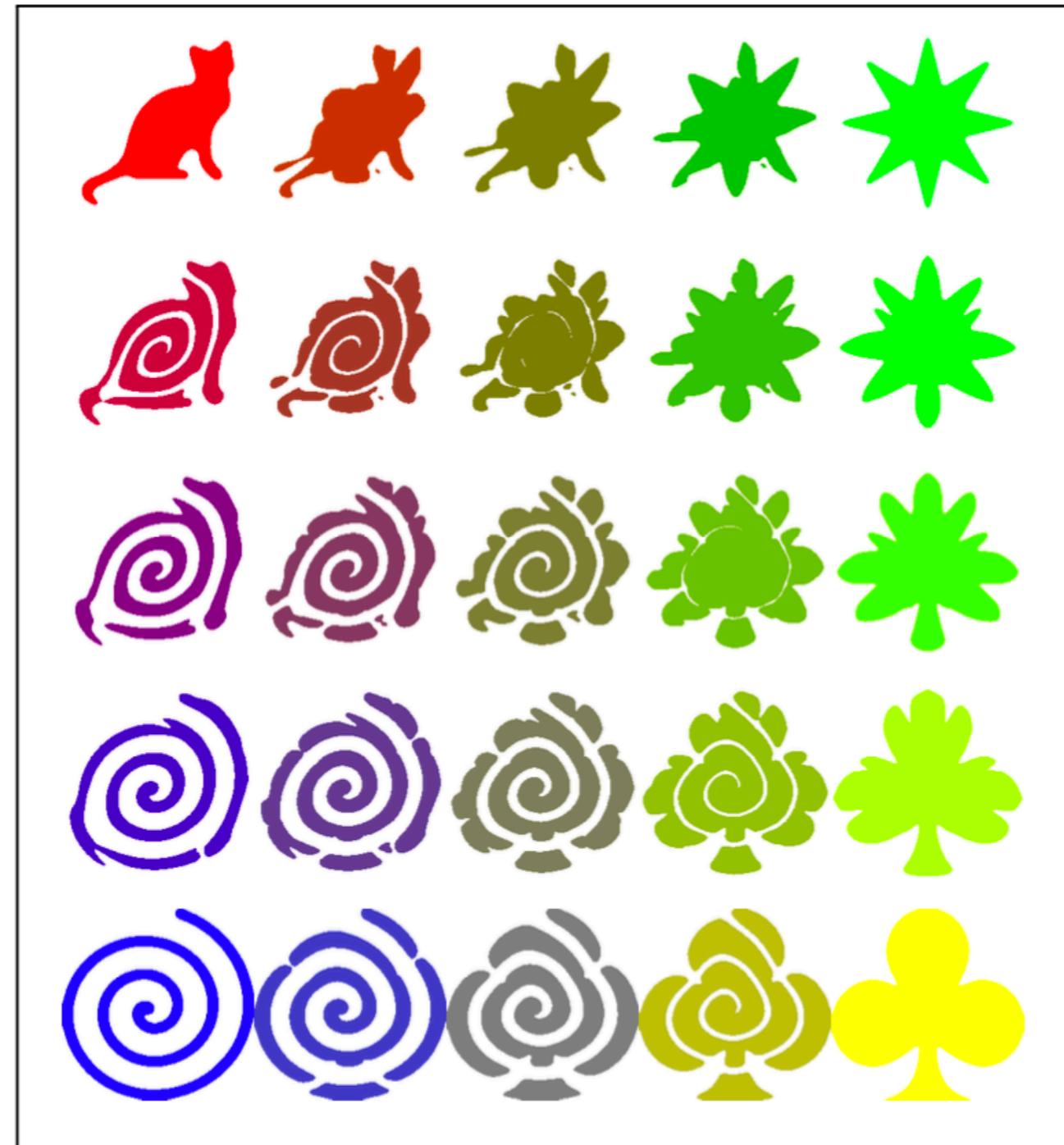
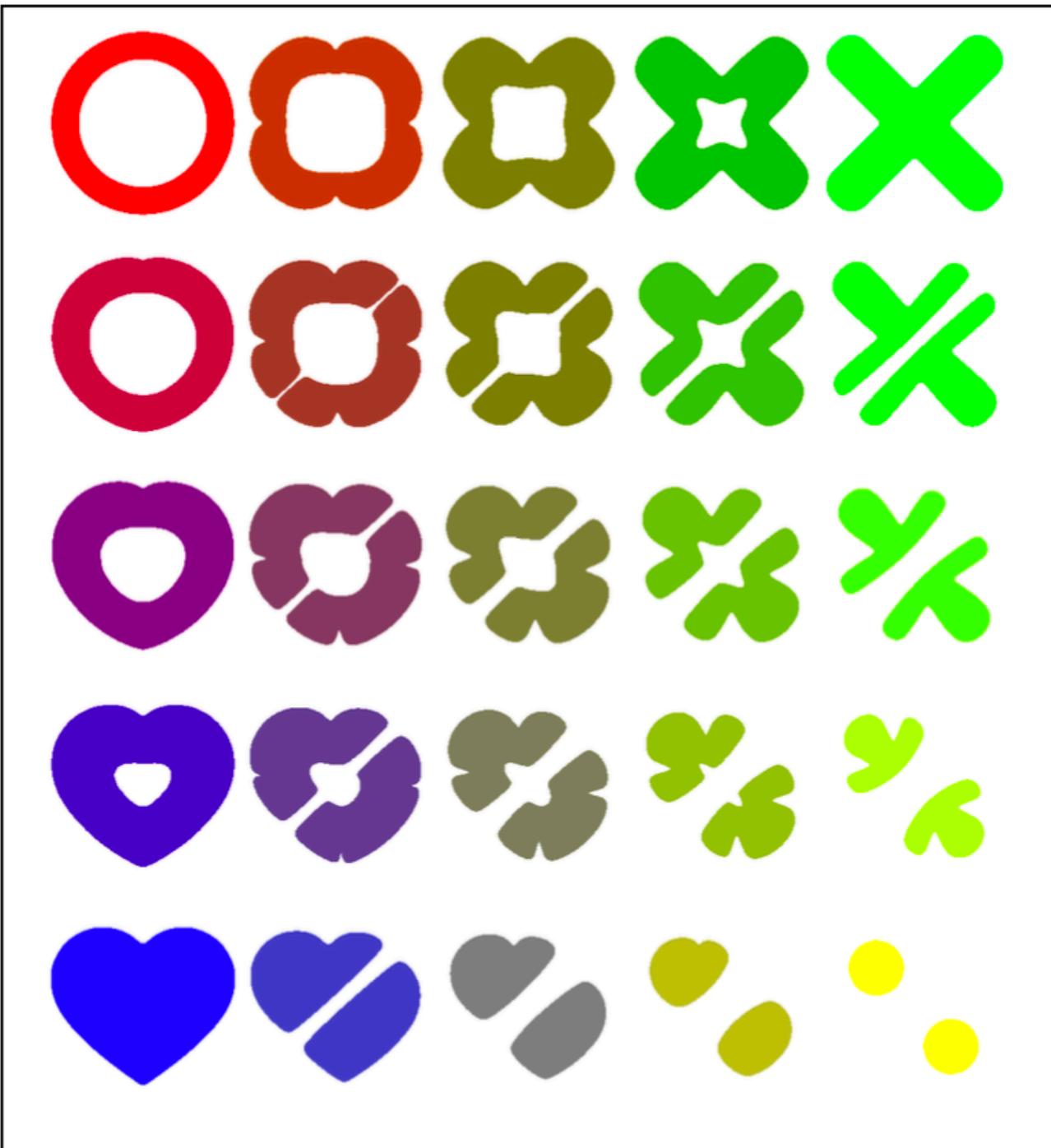
Optimal Transport was formulated in 1930 by A.N. Tolstoi, 12 years before Kantorovich. He even solved a "large scale"  $10 \times 68$  instance!



Stochastic gradient descent for the semi-discrete Optimal Transport, illustration of convergence and corresponding Laguerre cells. <https://arxiv.org/abs/1605.08527>

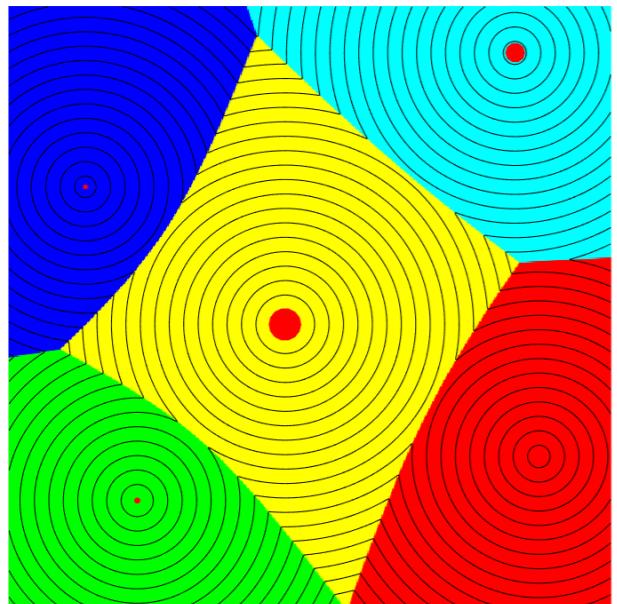


Optimal Transport barycenters. <https://hal.archives-ouvertes.fr/hal-01188953> <https://hal.archives-ouvertes.fr/hal-01096124>

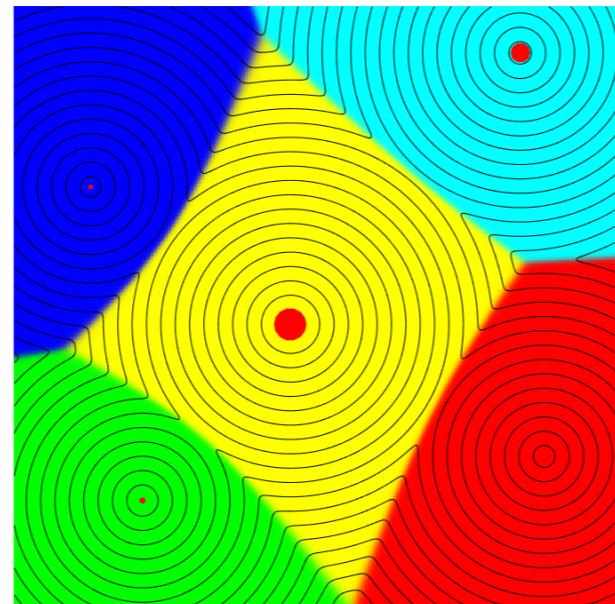


Heartbroken by Wasserstein barycenters!

Laguerre cells

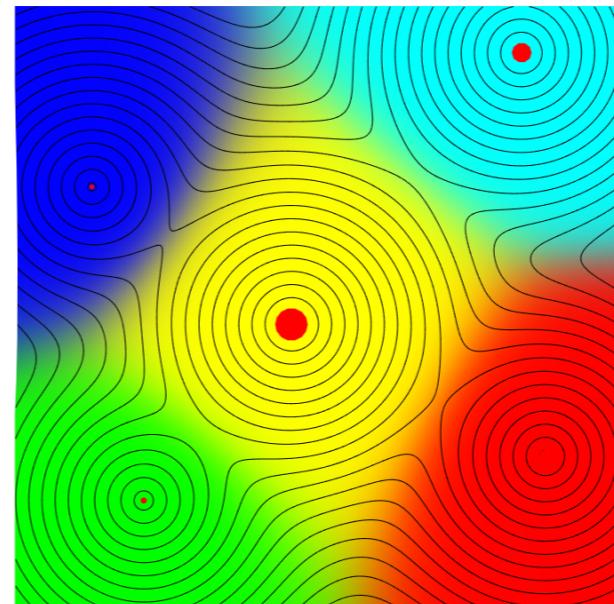


$$\varepsilon = 0$$

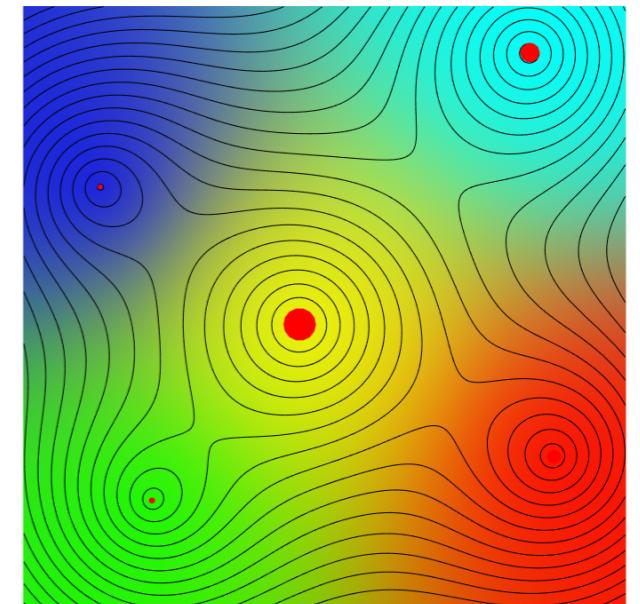


$$\varepsilon = 0.01$$

“Sinkhorn” Laguerre cells



$$\varepsilon = 0.1$$



$$\varepsilon = 0.3$$

Semi-discrete Optimal Transport is a multiclass SVM. Sinkhorn version is logistic regression. Decomposes space in entropic Laguerre diagrams

# Machine Learning & Inverse Problems

## Inverse Problems

$$y = Af + w$$



$$\begin{array}{lcl} A^\top y & = & (A^\top A)f + A^\top w \\ \stackrel{\text{def.}}{=} u & \stackrel{\text{def.}}{=} C & \stackrel{\text{def.}}{=} r \end{array}$$

Regularized inversion:

$$\min_f \frac{1}{2} \|Af - y\|^2 + \lambda \|f\|^2$$

$$f_\lambda = (\mathbf{C} + \lambda \text{Id}_p)^{-1} \mathbf{u}$$

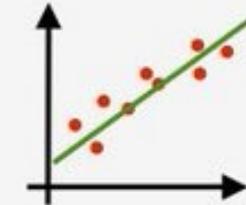
Exact covariance  $\mathbf{C}$

Deterministic bounded noise  $\mathbf{r}$

Noise level  $\varepsilon \stackrel{\text{def.}}{=} \|\mathbf{r}\|$

## Statistical Learning

$$y = Xf + \varepsilon$$



$$\begin{array}{lcl} \frac{1}{n} X^\top y & = & \frac{1}{n} (X^\top X)f + \frac{1}{n} X^\top \varepsilon \\ \stackrel{\text{def.}}{=} \mathbf{u}_n & \stackrel{\text{def.}}{=} \mathbf{C}_n & \stackrel{\text{def.}}{=} \mathbf{r}_n \\ \downarrow n \rightarrow +\infty & \downarrow (x_i, y_i)_i \text{ i.i.d.} & \\ \mathbf{u} = \mathbb{E}(yx) & \mathbf{C} = \mathbb{E}(xx^\top) & \end{array}$$

Empirical risk minimization:

$$\min_f \frac{1}{2n} \|Xf - y\|^2 + \lambda \|f\|^2$$

$$f_{\lambda,n} = (\mathbf{C}_n + \lambda \text{Id}_p)^{-1} \mathbf{u}_n$$

Exact covariance  $\mathbf{C}$

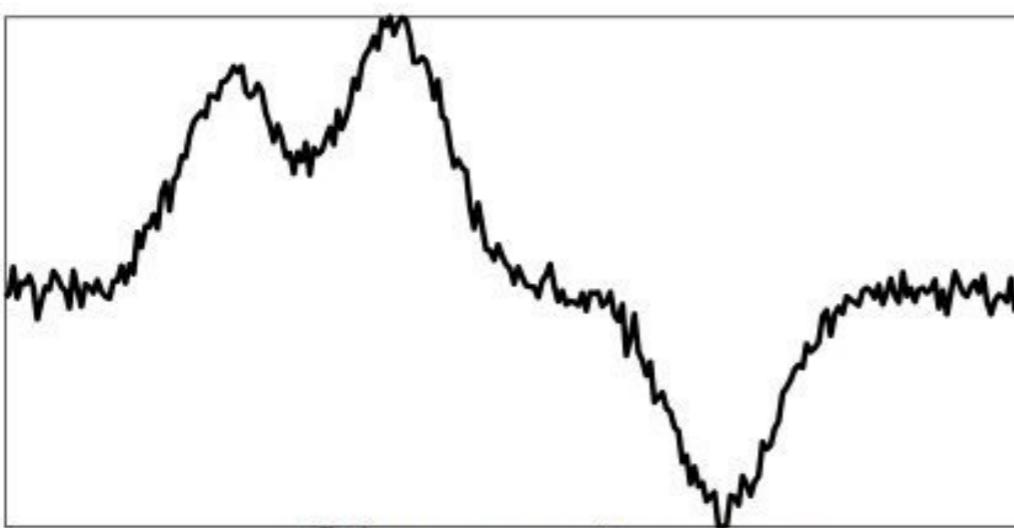
Noisy covariance  $\mathbf{C}_n$

Deterministic bounded noise  $\mathbf{r}$

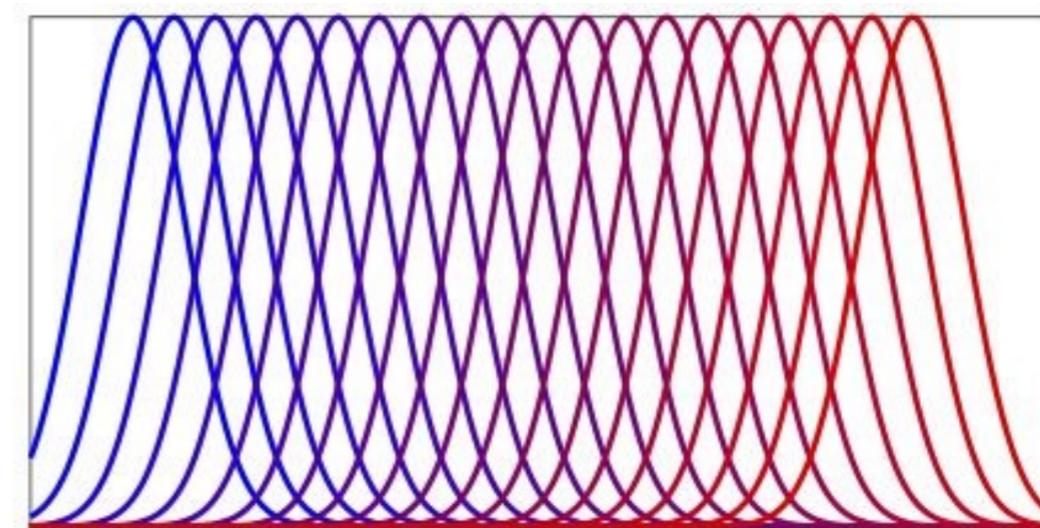
Random noise  $\mathbf{r}_n$

Noise level  $\varepsilon \stackrel{\text{def.}}{=} \|\mathbf{r}\|$

Noise level  $\|\mathbf{r}_n\| \sim \varepsilon = n^{-\frac{1}{2}}$

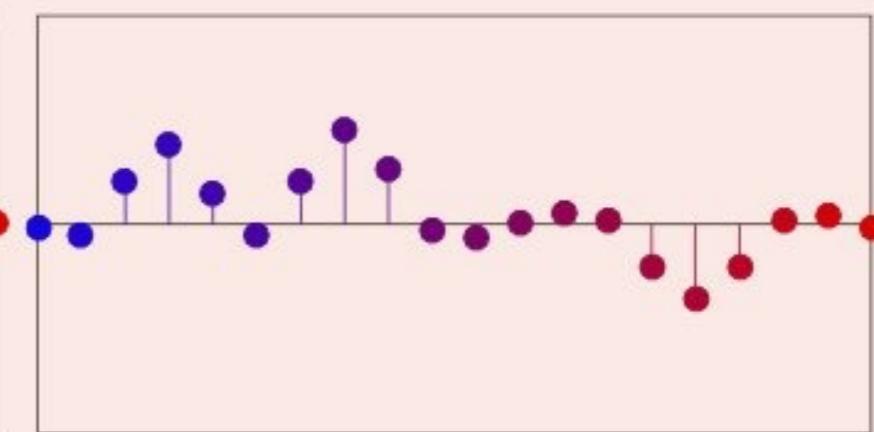
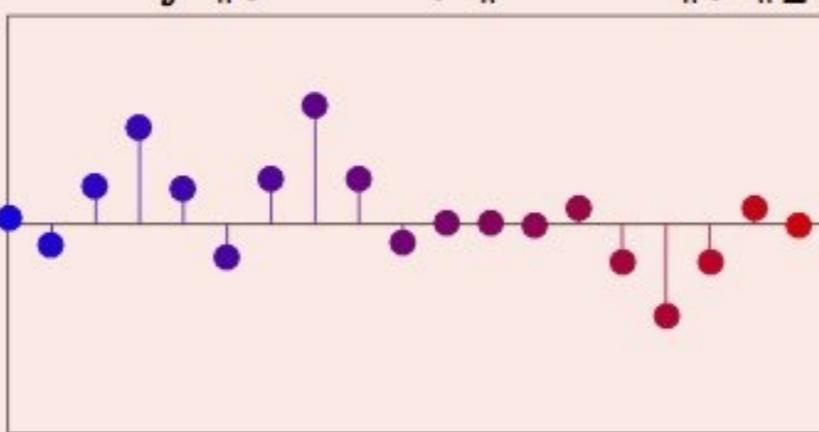
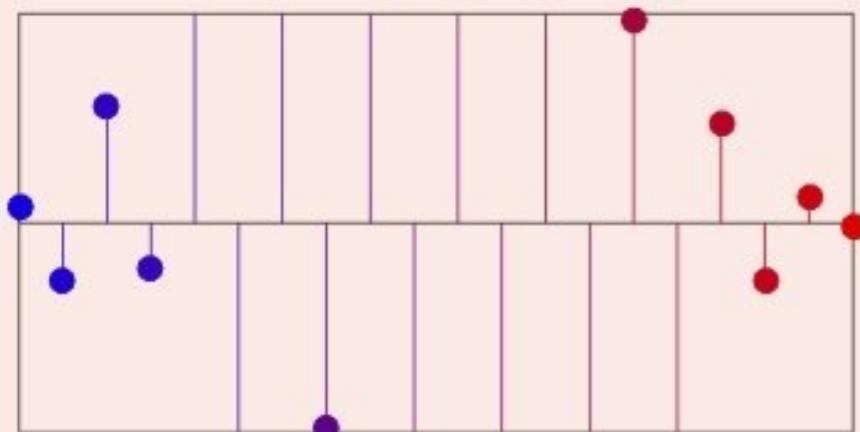


Observations  $y$

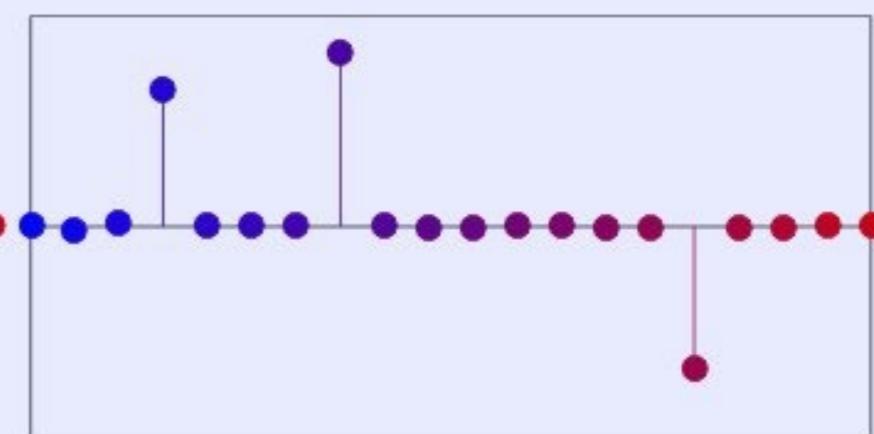
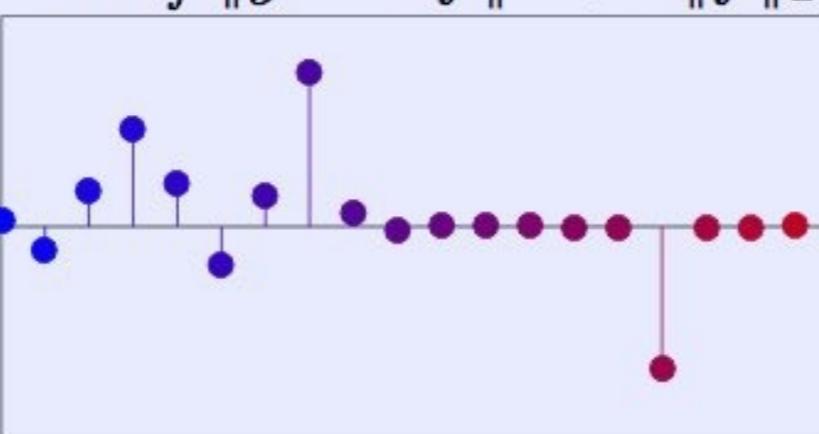
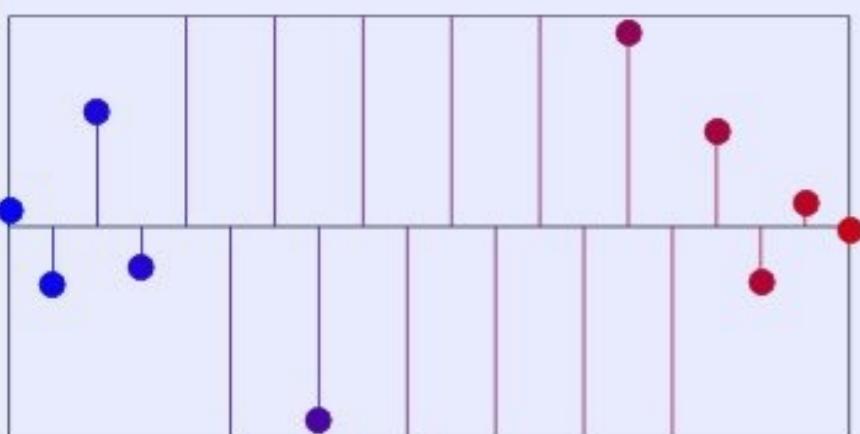


Columns of  $\Phi$

$$\min_f \|y - \Phi f\|^2 + \lambda \|f\|_2^2$$



$$\min_f \|y - \Phi f\|^2 + \lambda \|f\|_1$$



L2 vs L1 regularization for inverse problems / correlated designs.

Solving  $\textcolor{red}{y} \approx \textcolor{blue}{A} \textcolor{green}{x} \in \mathbb{R}^m \quad \textcolor{blue}{A} \in \mathbb{R}^{m \times n}$

Determined ( $m = n$ ):  $\textcolor{green}{x} = \textcolor{blue}{A}^{-1} \textcolor{red}{y}$

$$\begin{array}{c|c|c} \textcolor{red}{y} & = & \textcolor{blue}{A} \times \textcolor{green}{x} \end{array}$$

Over-determined ( $m > n$ ):  $\min_{\textcolor{green}{x}} \|\textcolor{blue}{A}\textcolor{green}{x} - \textcolor{red}{y}\|^2$

$$\textcolor{green}{x} = (\textcolor{blue}{A}^\top \textcolor{blue}{A})^{-1} \textcolor{blue}{A}^\top \textcolor{red}{y} \stackrel{\text{def.}}{=} \textcolor{blue}{A}^+ \textcolor{red}{y}$$

$$\begin{array}{c|c|c} \textcolor{red}{y} & \approx & \textcolor{blue}{A} \times \textcolor{green}{x} \end{array}$$

Under-determined ( $m < n$ ):  $\min_{\textcolor{green}{x}} \{\|\textcolor{green}{x}\| ; \textcolor{blue}{A}\textcolor{green}{x} = \textcolor{red}{y}\}$

$$\textcolor{green}{x} = \textcolor{blue}{A}^\top (\textcolor{blue}{A}\textcolor{blue}{A}^\top)^{-1} \textcolor{red}{y} \stackrel{\text{def.}}{=} \textcolor{blue}{A}^+ \textcolor{red}{y}$$

$$\begin{array}{c|c|c} \textcolor{red}{y} & = & \textcolor{blue}{A} \times \textcolor{green}{x} \end{array}$$

$A$  ill-posed and/or noise:  $\min_{\textcolor{green}{x}} \|\textcolor{blue}{A}\textcolor{green}{x} - \textcolor{red}{y}\|^2 + \lambda \|\textcolor{green}{x}\|^2$

$$\textcolor{green}{x} = (\textcolor{blue}{A}^\top \textcolor{blue}{A} + \lambda \text{Id}_n)^{-1} \textcolor{blue}{A}^\top \textcolor{red}{y} \xrightarrow{\lambda \rightarrow 0} \textcolor{blue}{A}^+ \textcolor{red}{y}$$

$$= \textcolor{blue}{A}^\top (\textcolor{blue}{A}\textcolor{blue}{A}^\top + \lambda \text{Id}_m)^{-1} \textcolor{red}{y} \quad (\text{Woodbury identity})$$

$$\begin{array}{c|c|c} \textcolor{red}{y} & \approx & \textcolor{blue}{A} \times \textcolor{green}{x} \end{array}$$

# Automatic Differentiation

How to compute  $\nabla \ell_{x,y}(\theta)$ ?  $\ell_{x,y}(\theta) \stackrel{\text{def.}}{=} L(f(x,\theta), y)$

Chain rule:  $\nabla \ell_{x,y}(\theta) = [\partial f(x, \theta)]^\top (\nabla L(f(x, \theta), y))$

Linear  $f(x, \theta) = \theta \times x$ :  $\partial f(x, \theta) = \theta$ .

Non-linear  $f(x, \theta)$ : painful ... but  $\ell_{x,y}$  it is just a computer program.

Computer program  $\Leftrightarrow$  directed acyclic graph  $\Leftrightarrow$  linear ordering of nodes  $(\theta_r)_r$

```

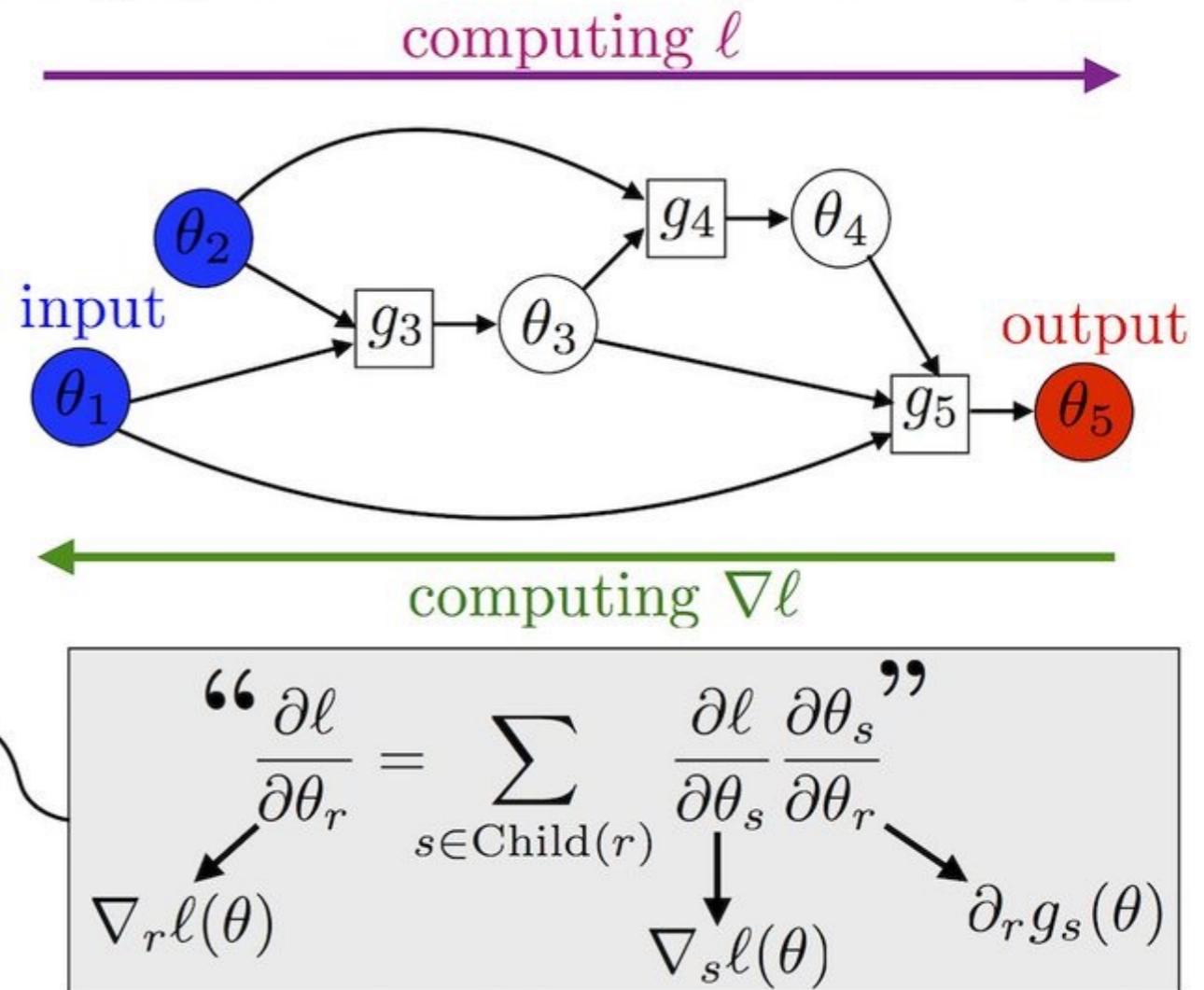
function  $\ell(\theta_1, \dots, \theta_M)$ 
  for  $r = M + 1, \dots, R$ 
    |  $\theta_r = g_r(\theta_{\text{Parents}(r)})$ 
  return  $\theta_R$ 

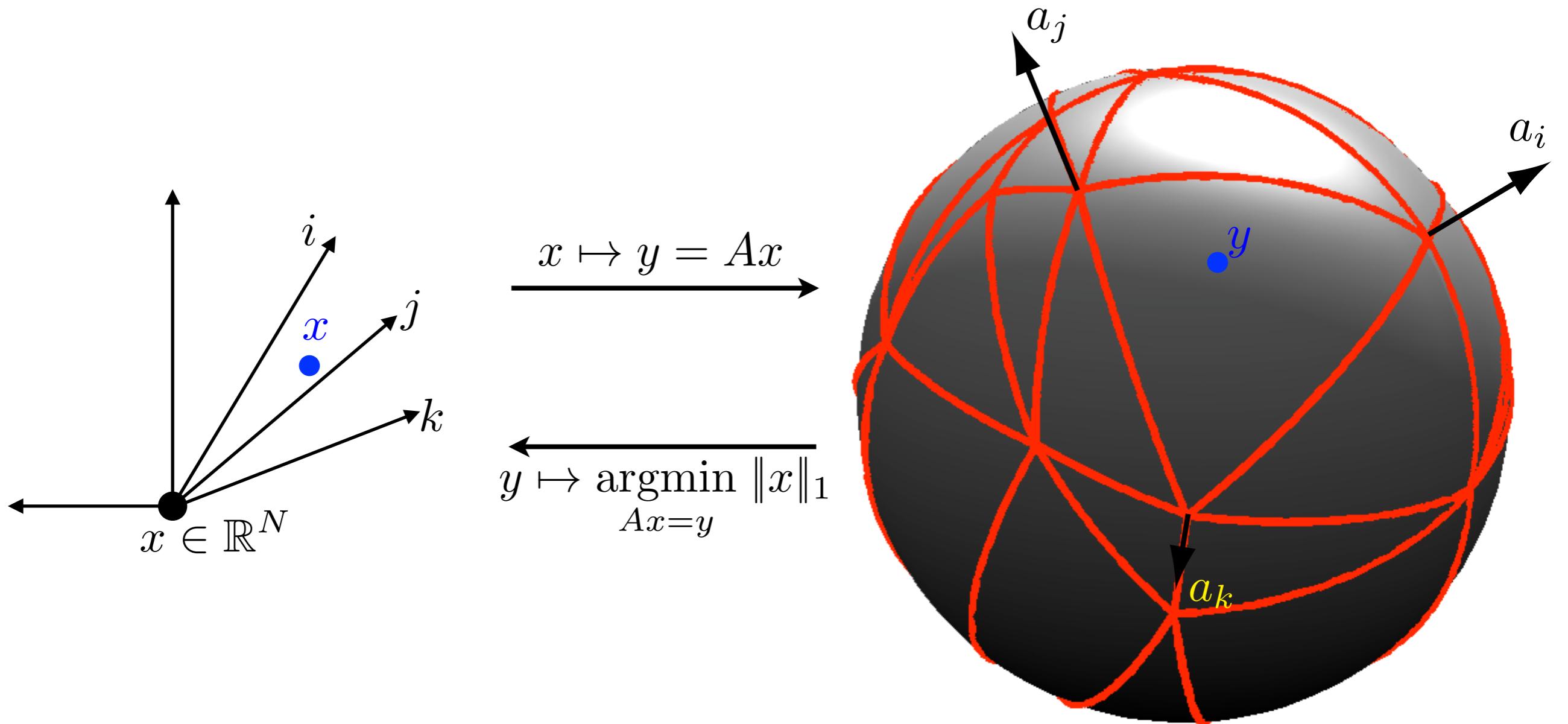
```

```

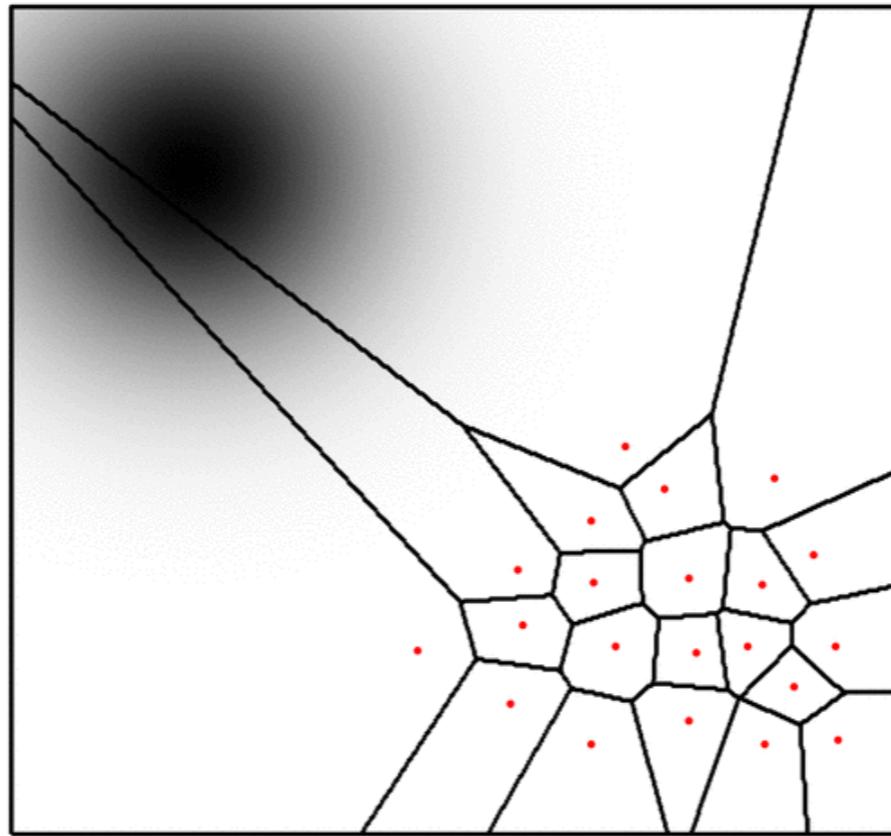
backward
function  $\nabla \ell(\theta_1, \dots, \theta_M)$ 
     $\nabla_R \ell = 1$ 
    for  $r = R - 1, \dots, 1$ 
         $\nabla_r \ell = \sum_{s \in \text{Child}(r)} \partial_r g_s(\theta) \nabla_s \ell$ 
    return  $(\nabla_1 \ell, \dots, \nabla_M \ell)$ 

```





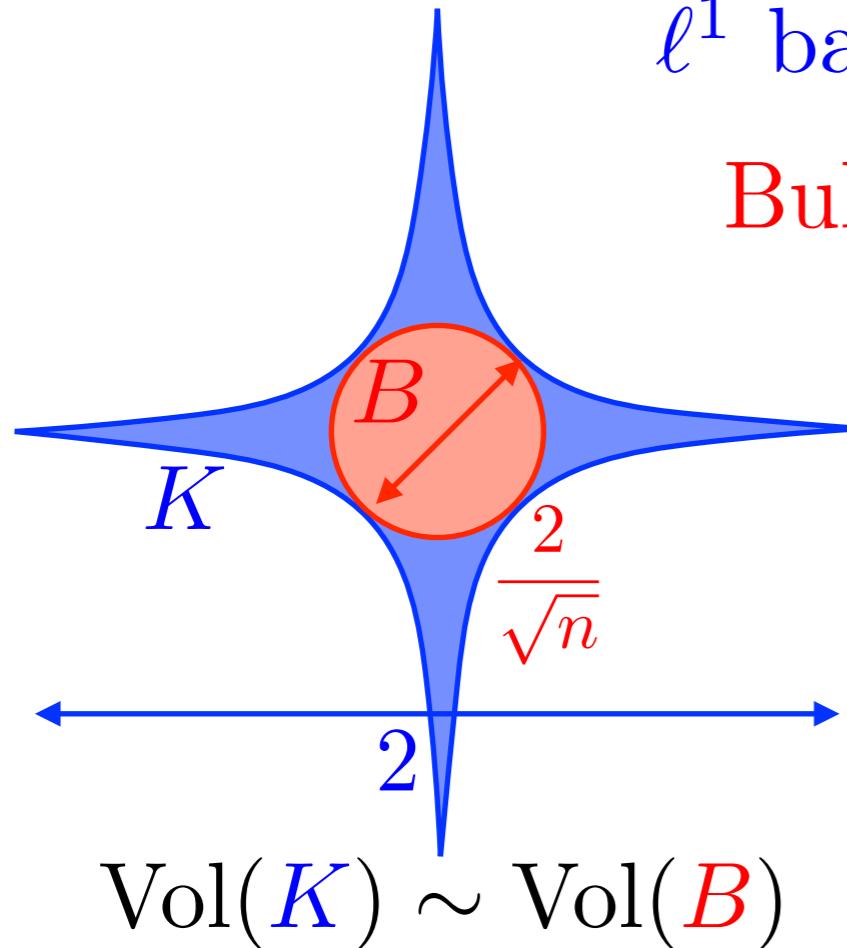
Signs of sparse vectors recovered by  $\ell_1$  minimization index a spherical Delaunay triangulation of the observation space. More on this in my lecture notes. <https://mathematical-tours.github.io/book-sources/chapters-pdf/sparse-theory.pdf> ...



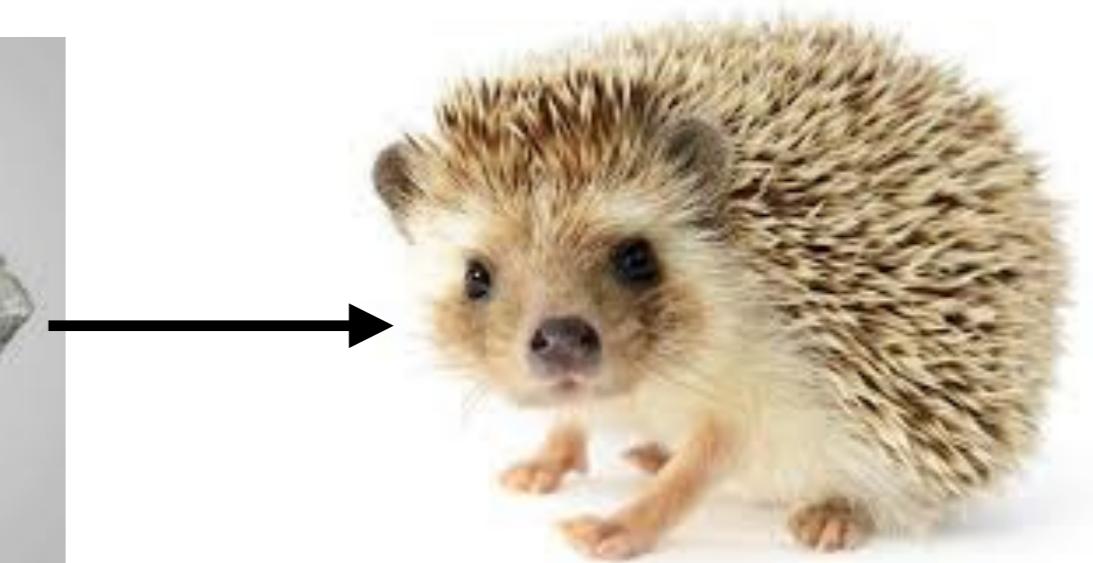
Lloyd algorithm (k-means) with non-uniform sampling density. [https://en.wikipedia.org/wiki/Lloyd%27s\\_algorithm ...](https://en.wikipedia.org/wiki/Lloyd%27s_algorithm)

$\ell^1$  ball:  $K \stackrel{\text{def.}}{=} \{x \in \mathbb{R}^n ; \sum_{i=1}^n |x_i| \leq 1\}$

Bulk:  $B \stackrel{\text{def.}}{=} \{x \in \mathbb{R}^n ; \sum_{i=1}^n |x_i|^2 \leq n^{-1}\}$



small  $n$



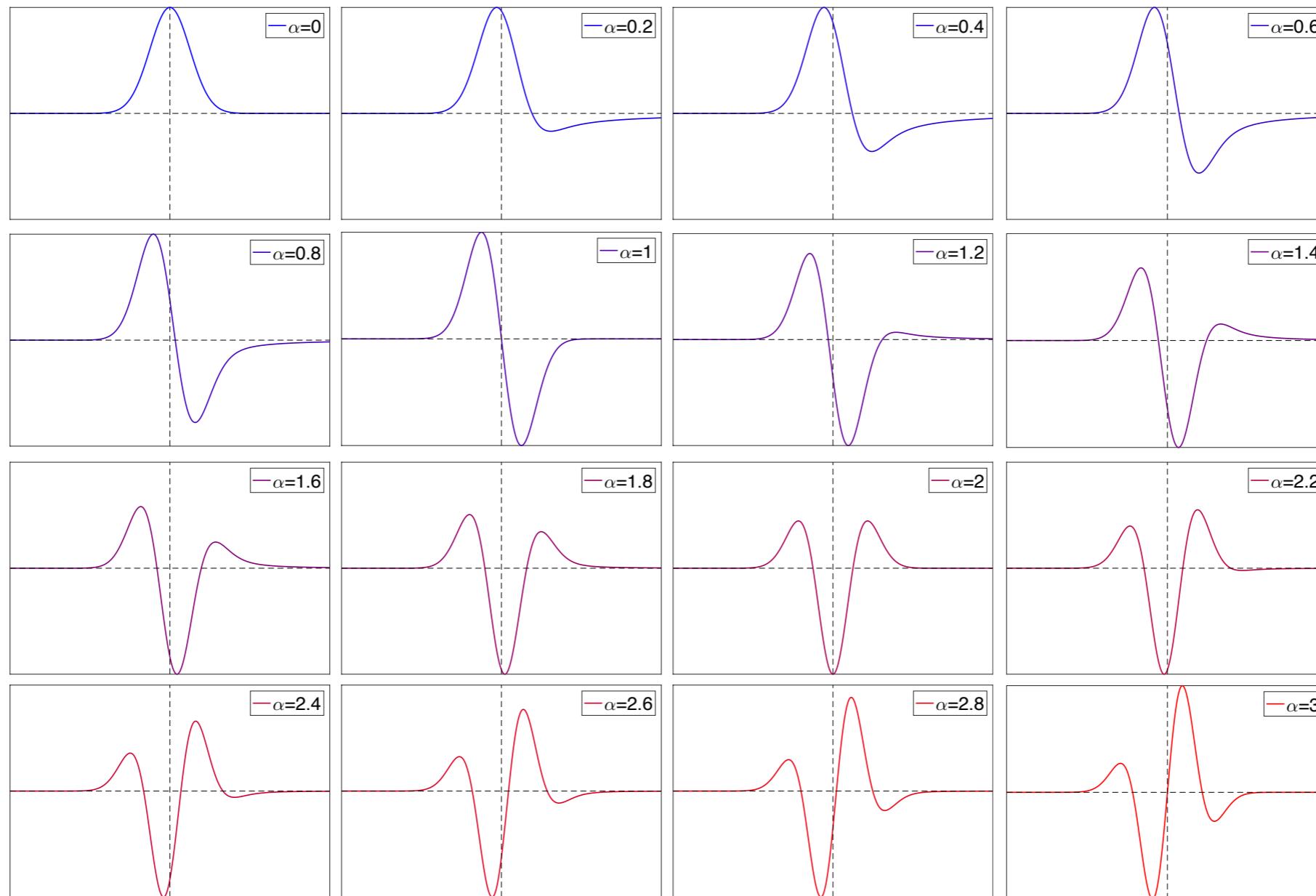
large  $n$

Vitali Milman's insight: convex sets in high dimension look like hedgehogs <http://www.math.tau.ac.il/~milman/files/survey4.ps>

# **Signal and Image Processing**

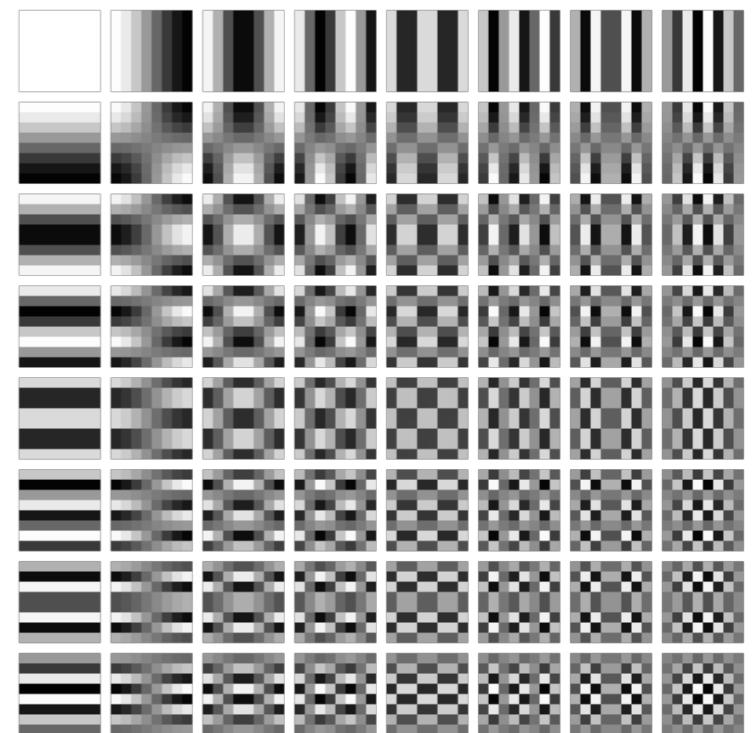
Fourier transform:  $\mathcal{F}(f)(\omega) \stackrel{\text{def.}}{=} \int_{\mathbb{R}} f(x)e^{-i\omega x}dx$

Fractional derivative:  $\mathcal{F}(f^{(\alpha)}) \stackrel{\text{def.}}{=} (i\omega)^{\alpha} \mathcal{F}(f)(\omega)$

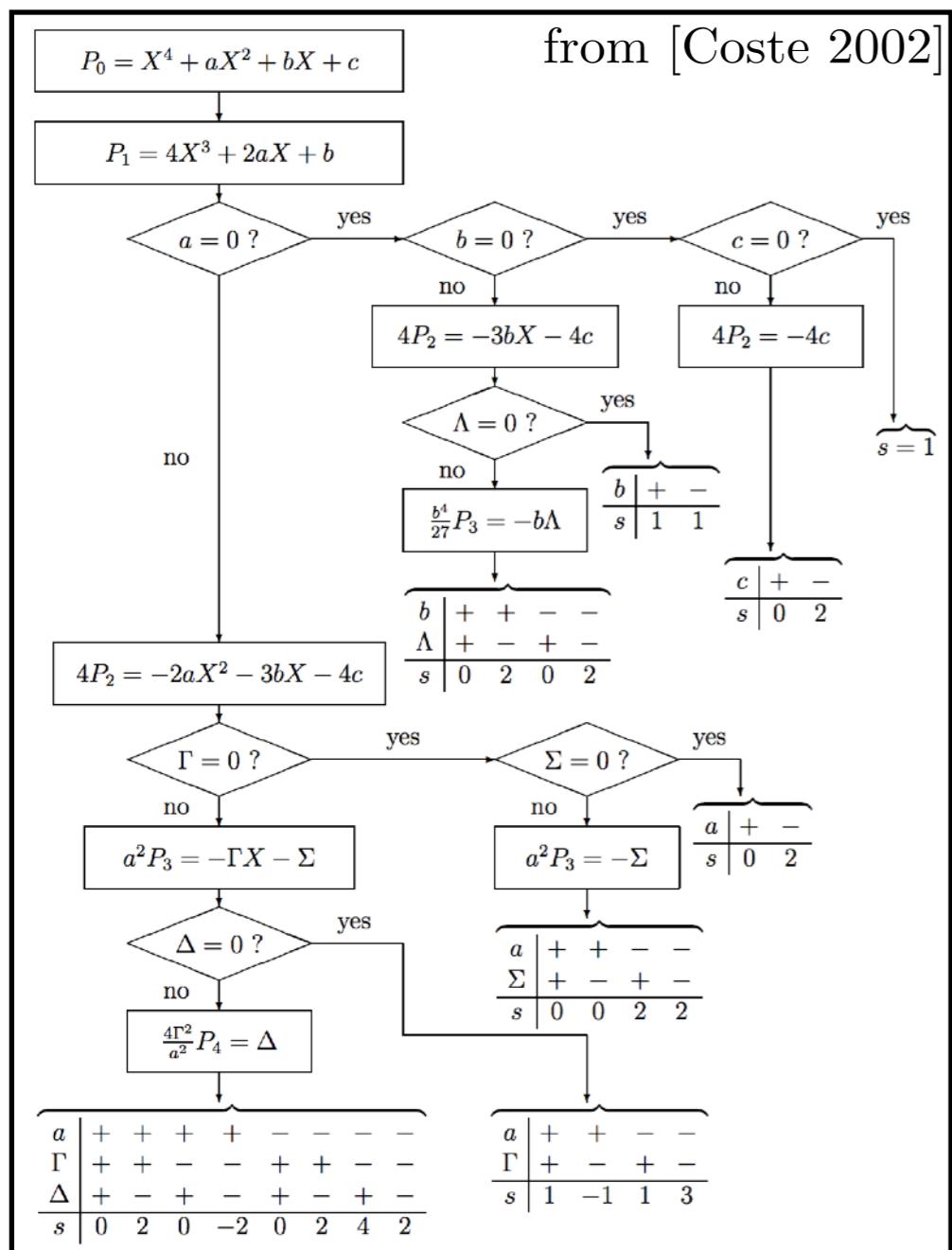


DCT-1:	$\cos jk \frac{\pi}{N-1}$	(divide by $\sqrt{2}$ when $j$ or $k$ is 0 or $N - 1$ )
DCT-2:	$\cos \left(j + \frac{1}{2}\right) k \frac{\pi}{N}$	(divide by $\sqrt{2}$ when $k = 0$ )
DCT-3:	$\cos j \left(k + \frac{1}{2}\right) \frac{\pi}{N}$	(divide by $\sqrt{2}$ when $j = 0$ )
DCT-4:	$\cos \left(j + \frac{1}{2}\right) \left(k + \frac{1}{2}\right) \frac{\pi}{N}$	

**The discrete case has a new level of variety and complexity, often appearing in the boundary conditions [G. Strang - SIAM review, 1999]**



# Ongoing



algebraic set

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

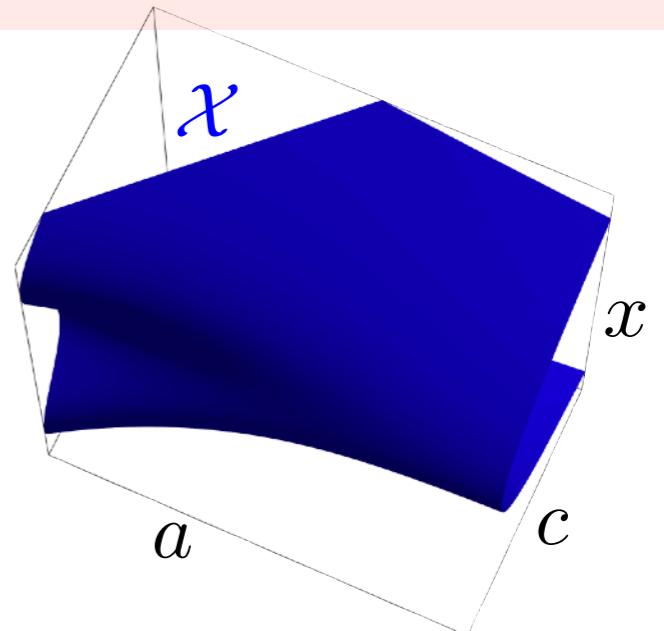
projection

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

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

semi-algebraic set

$$b = 0$$

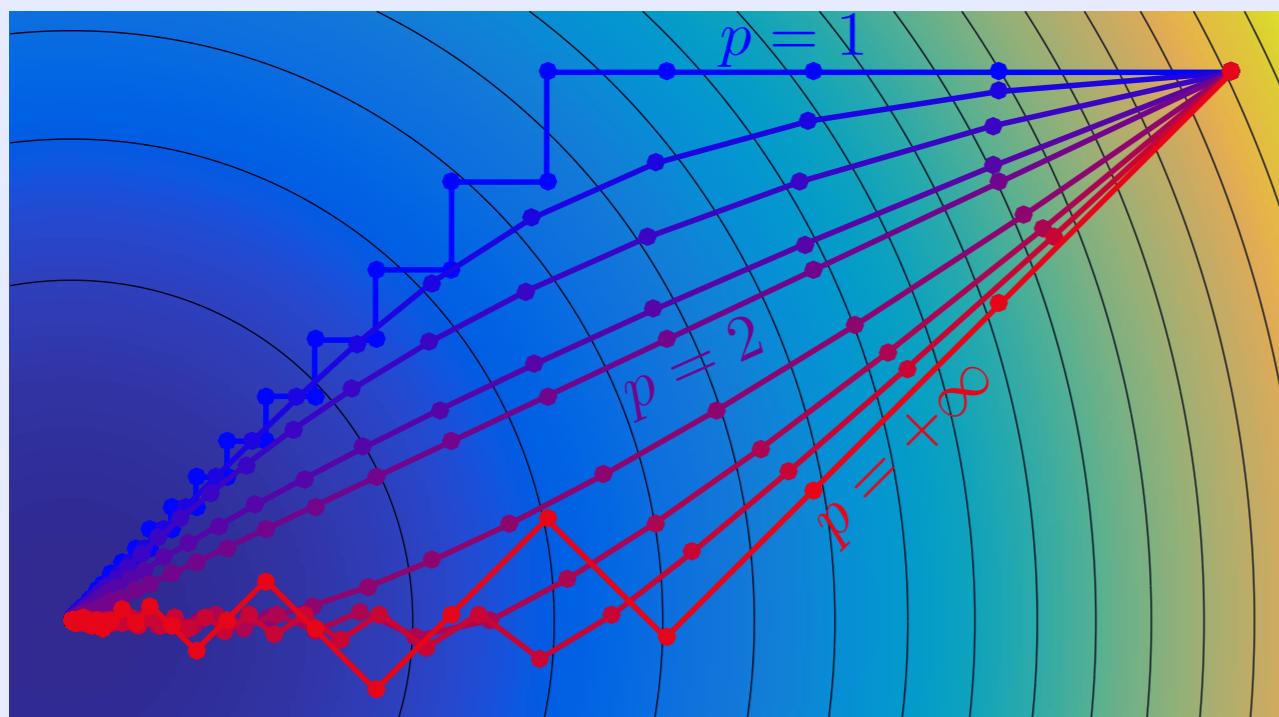


Tarski–Seidenberg in action: the set of polynomials having a real root is semi-algebraic. Coste provides the block diagram for degree 4 in his notes. Amazing.  
<http://gcomte.perso.math.cnrs.fr/M2/CosteIntroToSemialGeo.pdf>

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

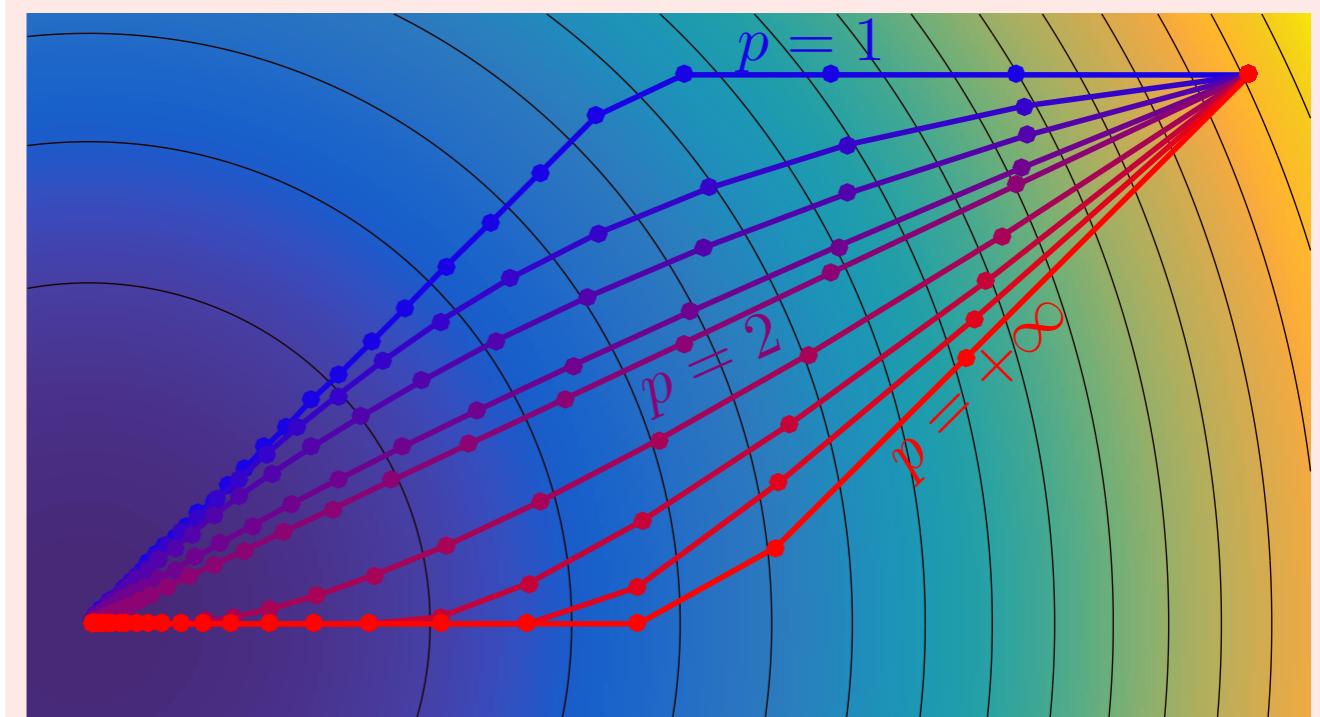
### Explicit

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



### Implicit

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



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

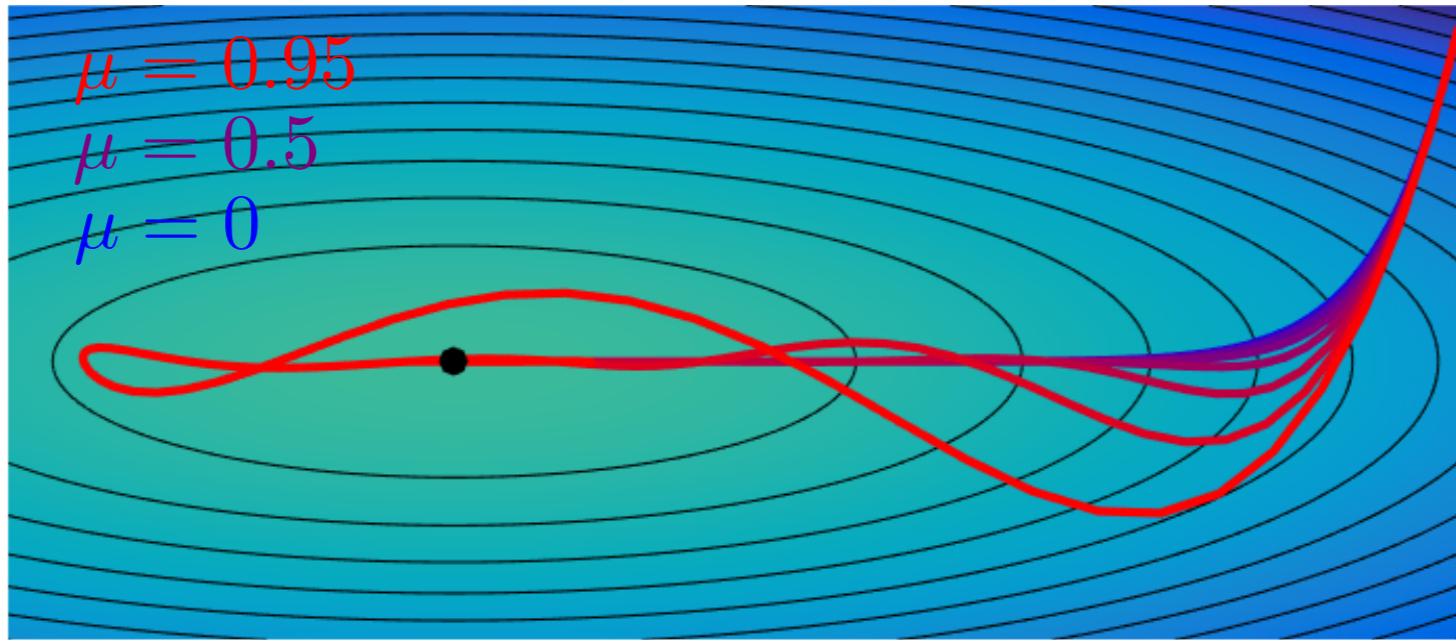
Comparison of explicit stepping (unstable) vs implicit stepping (stable) for gradient flow on metric space. Note how L1 flow is equivalent to a coordinate descent, and L2 is the usual gradient descent.

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

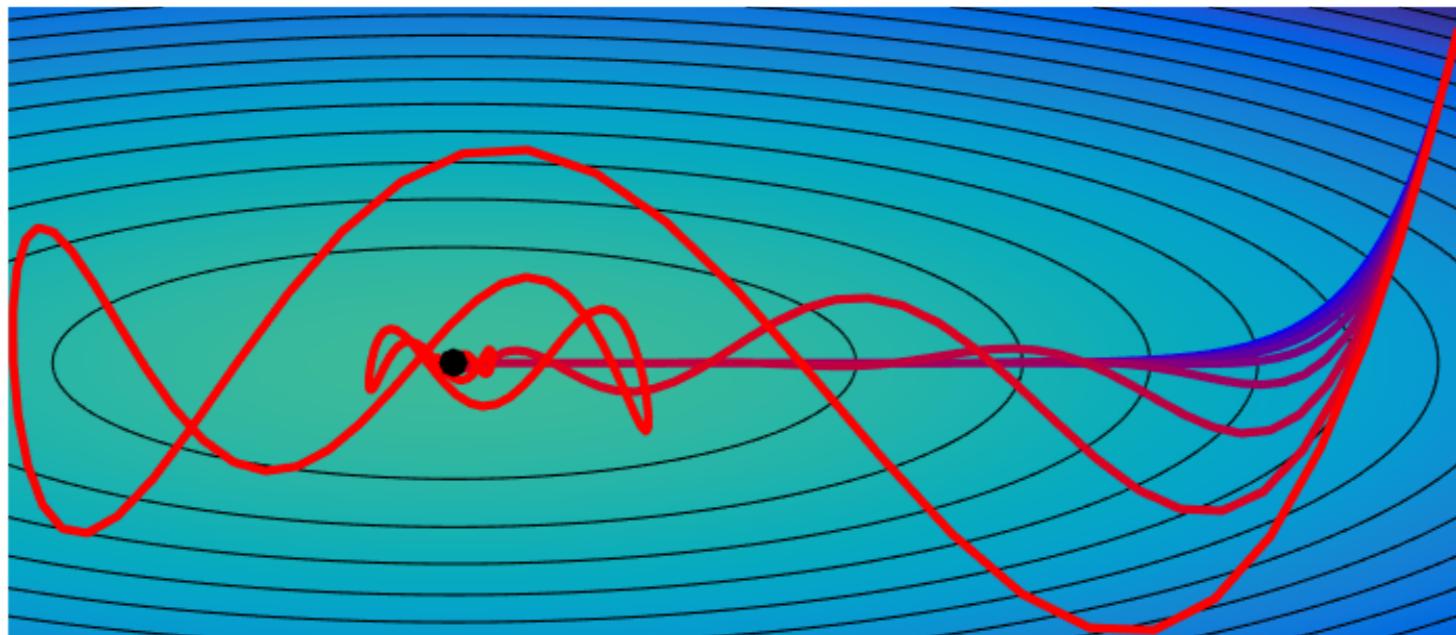
$$p_{k+1} = \mu p_k - \tau \begin{cases} \nabla f(x_k) & \text{Polyak} \\ \nabla f(x_k + \mu p_k) & \text{Nesterov} \end{cases}$$



Yurii  
Nesterov



Boris  
Polyak



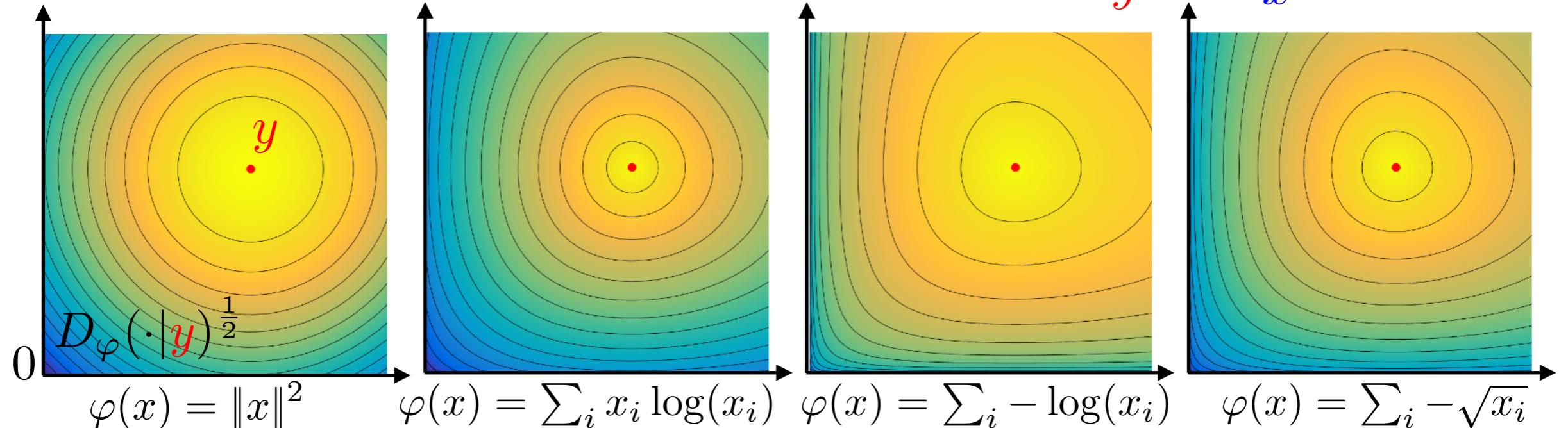
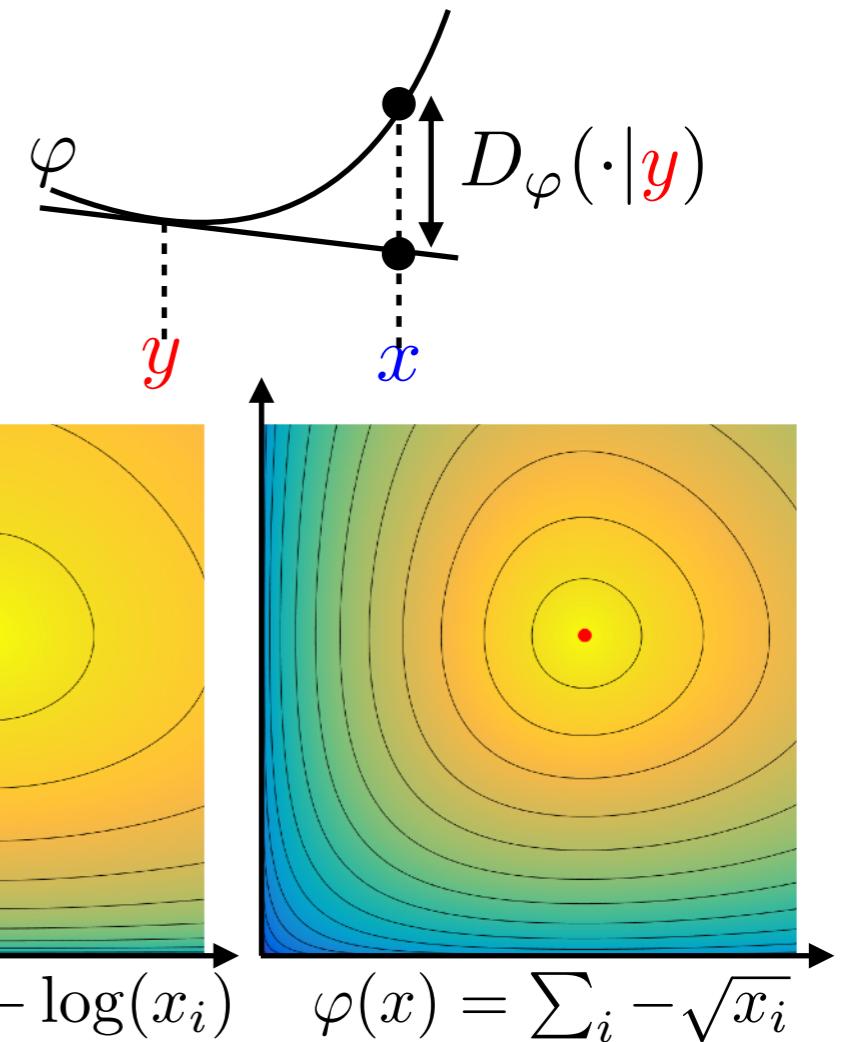
Boris Polyak's heavy ball and Yurii Nesterov's descent use momentum to accelerate but oscillate.

<http://lab7.ipu.ru/eng/people/polyak.html>

[https://fr.wikipedia.org/wiki/Yurii\\_Nesterov](https://fr.wikipedia.org/wiki/Yurii_Nesterov)

Bregman divergence:

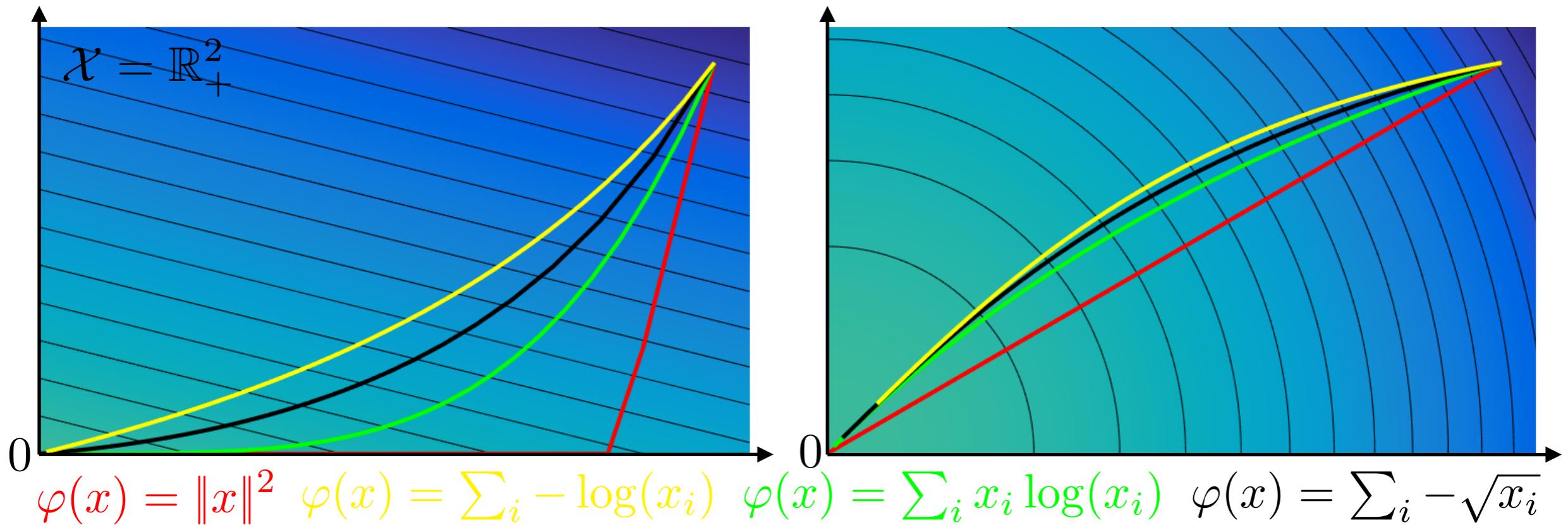
$$D_\varphi(\mathbf{x}|\mathbf{y}) \stackrel{\text{def.}}{=} \varphi(\mathbf{x}) - \varphi(\mathbf{y}) - \langle \mathbf{x} - \mathbf{y}, \nabla \varphi(\mathbf{y}) \rangle$$



Bregman divergences are convex distance-like functionals which are locally Euclidean. Any reasonable algorithm handling Euclidean distances generalizes to Bregman geometries.

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

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



Mirror descent generalizes gradient descent using Bregman geometries.

<https://blogs.princeton.edu/imabandit/2013/04/16/orf523-mirror-descent-part-iii/>