

Optimization flows landing on the Stiefel manifold

Simon Vary*

joint work with Bin Gao, Pierre Ablin***, Pierre-Antoine Absil***

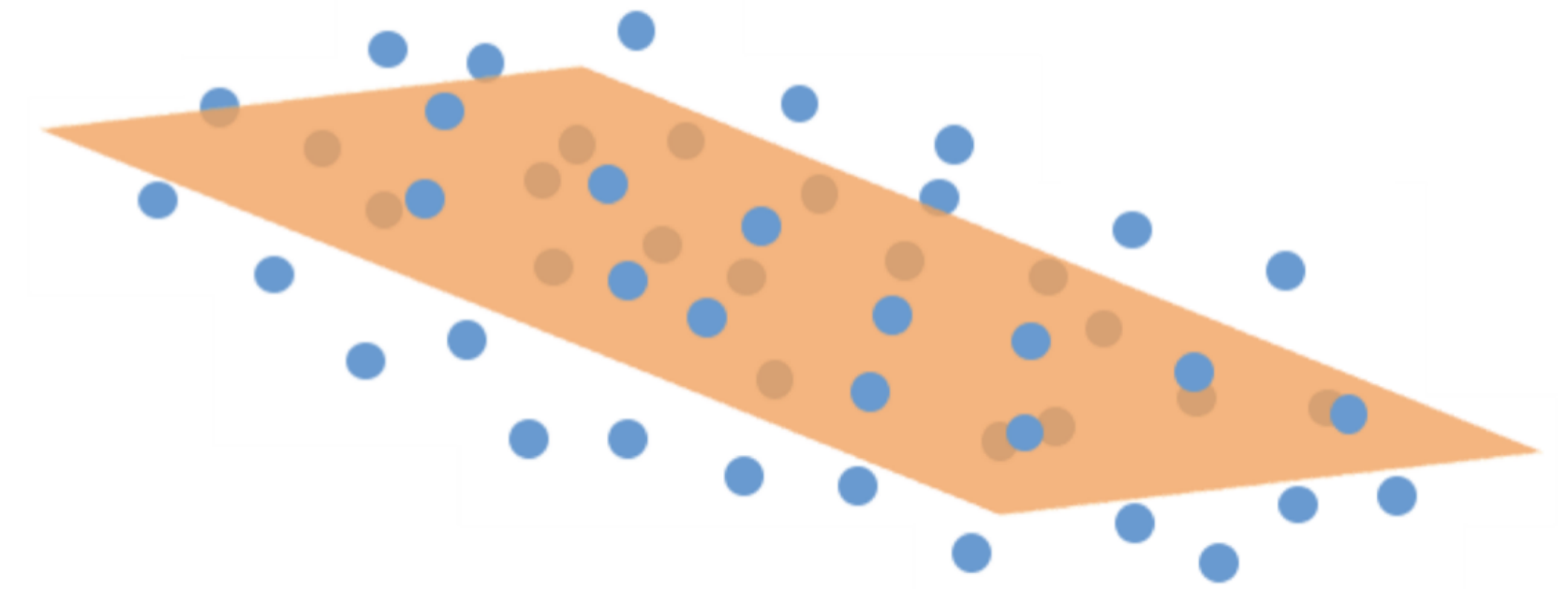
** ICTEAM, Université catholique de Louvain*

*** Institute for Applied Mathematics, University of Münster*

**** CNRS, Université Paris-Dauphine, PSL University
Machine Learning Group, Apple, Paris*

Optimization over the Stiefel manifold

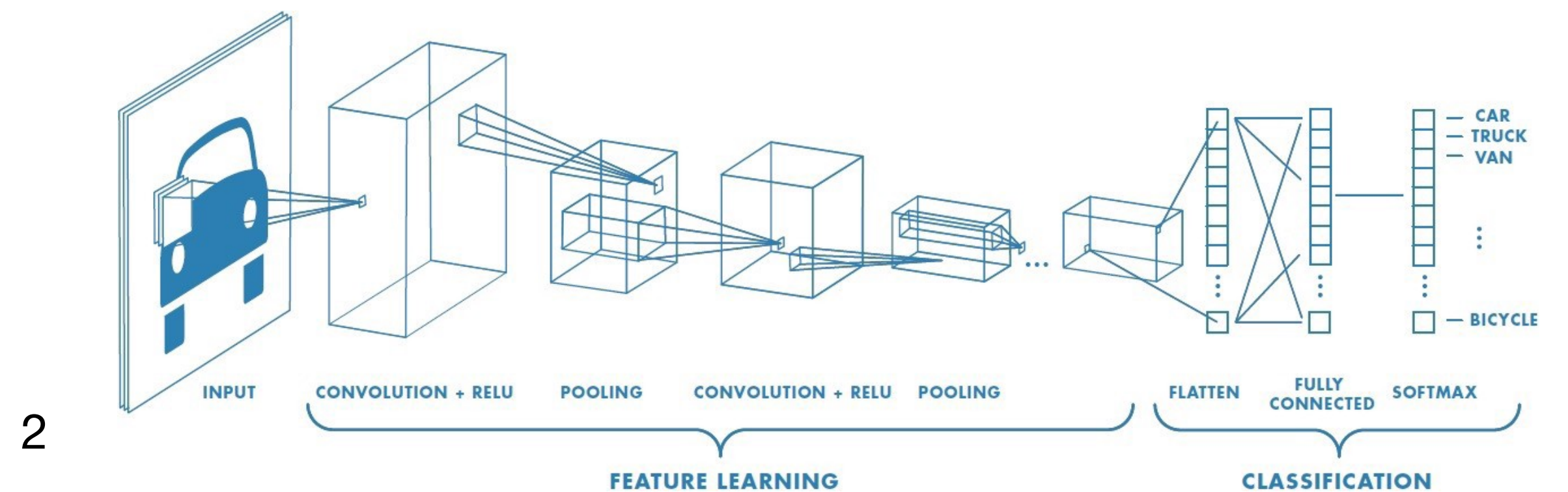
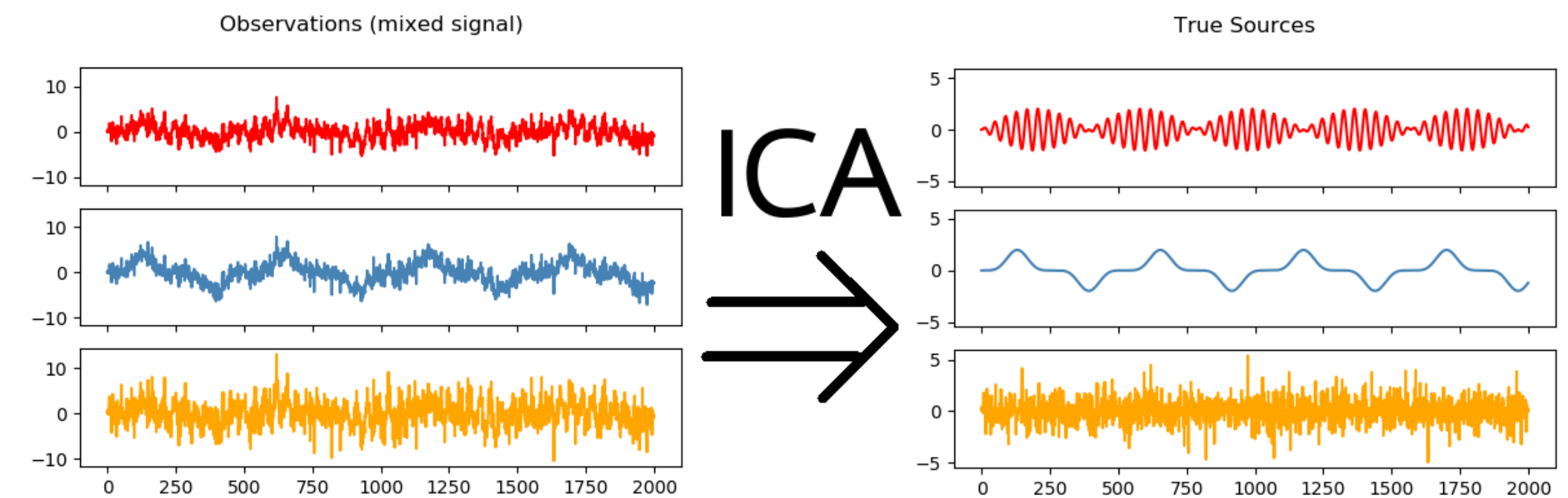
$$\min_{X \in \mathbb{R}^{n \times p}} f(x) \quad \text{s.t.} \quad \text{St}(p, n) := \left\{ X \in \mathbb{R}^{n \times p} : X^T X = I_p \right\}$$



Applications such as

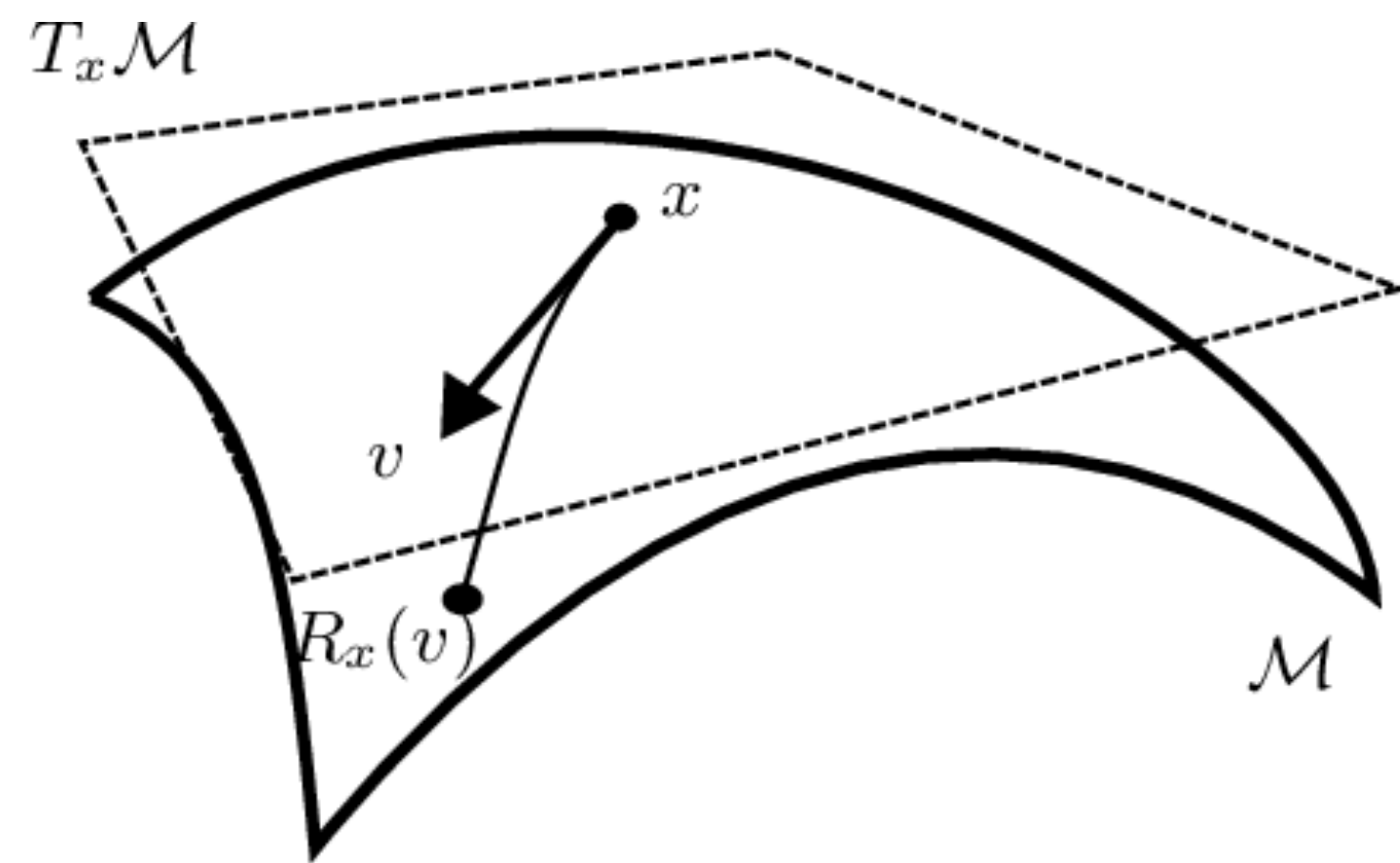
- Principal component analysis
- Independent component analysis
- Orthogonal weights in deep learning

$$\max_{X \in \text{St}(p, n)} \frac{1}{2} \|AX\|_F^2$$



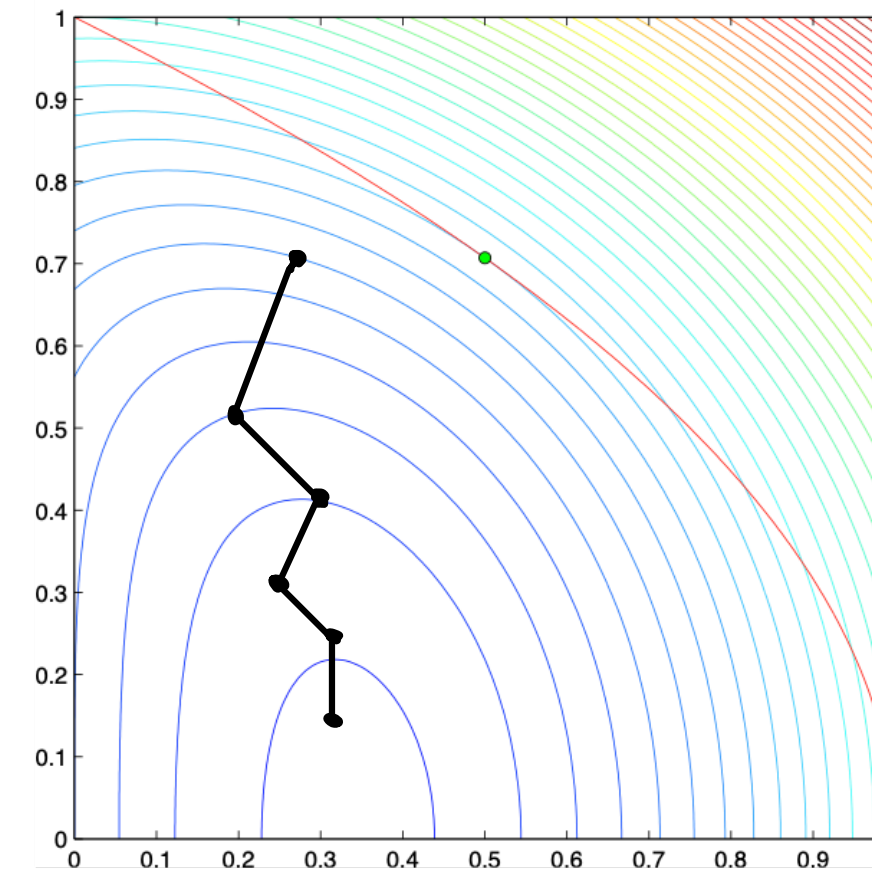
Optimization approaches

Riemannian optimization

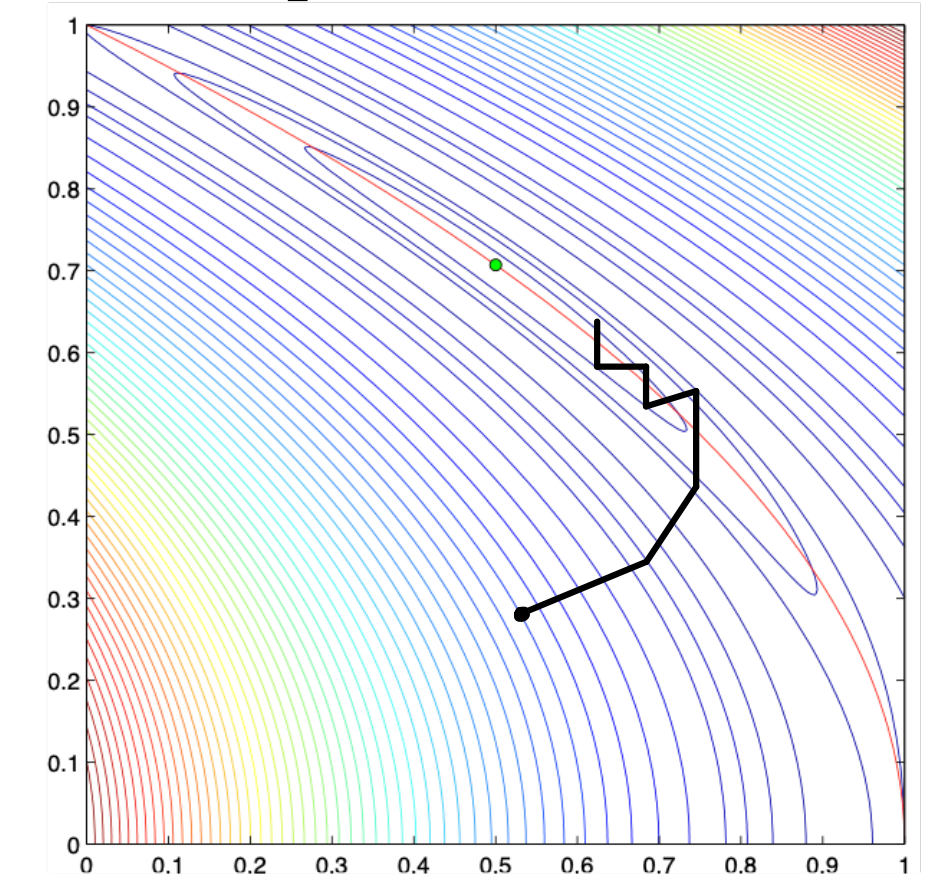


- $X^{t+1} = R_{X^t}(-\text{grad } f(X^t))$
- Retraction: QR $\mathcal{O}(np^2)$, Cayley $\mathcal{O}(p^3)$
- $\text{grad } f(X) = \text{skew}(\nabla f(X)X^\top)X$ ($\mathcal{O}(np^2)$)

Infeasible constrained optimization



$\sigma = 1$



$\sigma = 0.01$

from lecture slides of Coralia Cartis

- Penalty methods: $+\frac{1}{4\sigma}\|X^\top X - I_p\|_F^2$
- Augmented Lagrangian Method
- Adaptively choosing parameters can be tricky
- Gradient of the penalty: $X(X^\top X - I_p)$ (mat. mult. $\mathcal{O}(np^2)$)

Landing field

- Consider the following flow

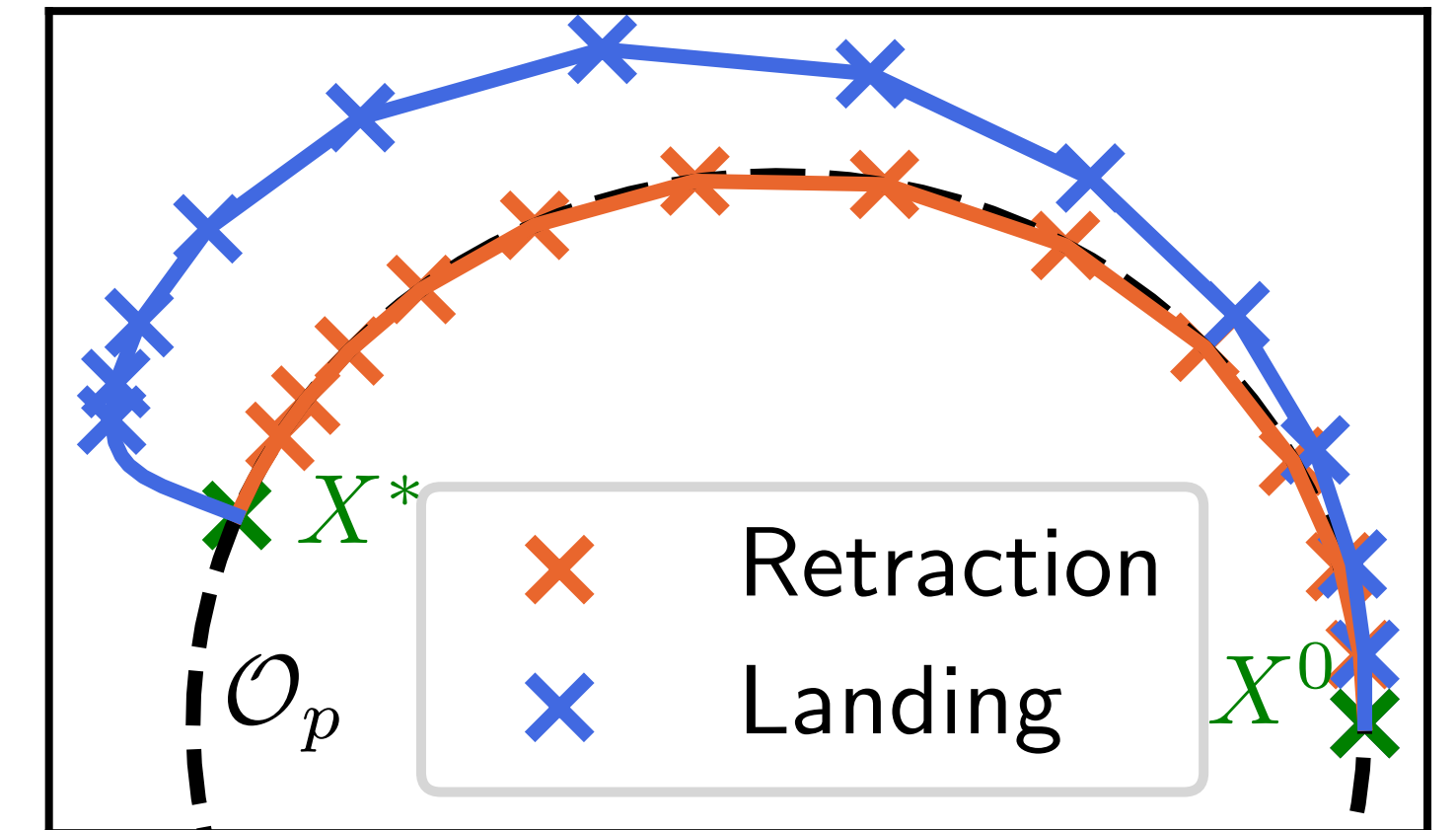
$$\dot{X}(t) = -\Lambda(X(t)),$$

where

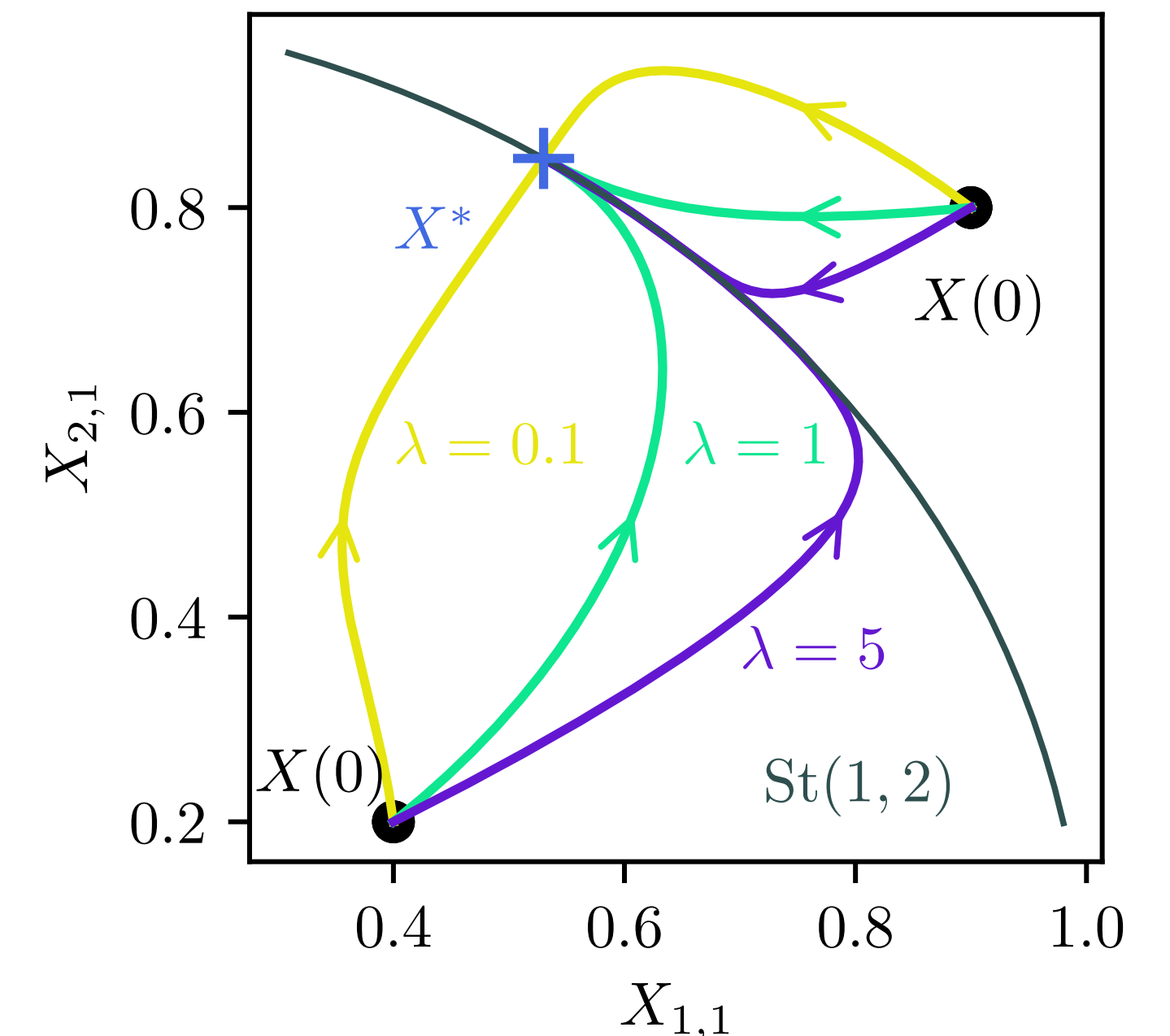
$$\Lambda(X) := \psi(X)X + \lambda \nabla \mathcal{N}(X),$$

and

- $\nabla \mathcal{N}(X) = X(X^\top X - I_p)$, with $\mathcal{N}(X) := \frac{1}{4} \left\| X^\top X - I_p \right\|_F^2$,
- $\psi(X) = 2\text{skew}(\nabla f(X)X^\top)$, with $\text{skew}(A) = \frac{1}{2}(A - A^\top)$



Fast and accurate optimization on the orthogonal manifold.
Ablin & Peyré, AISTATS 2022



Interpretation of the landing field (on the Stiefel manifold)

$$\begin{aligned}\Lambda(X) &:= \psi(X)X + \lambda \nabla \mathcal{N}(X) \\ &= 2\text{skew}(\nabla f(X)X^\top)X + \lambda X(X^\top X - I_p)\end{aligned}$$

- Orthogonality

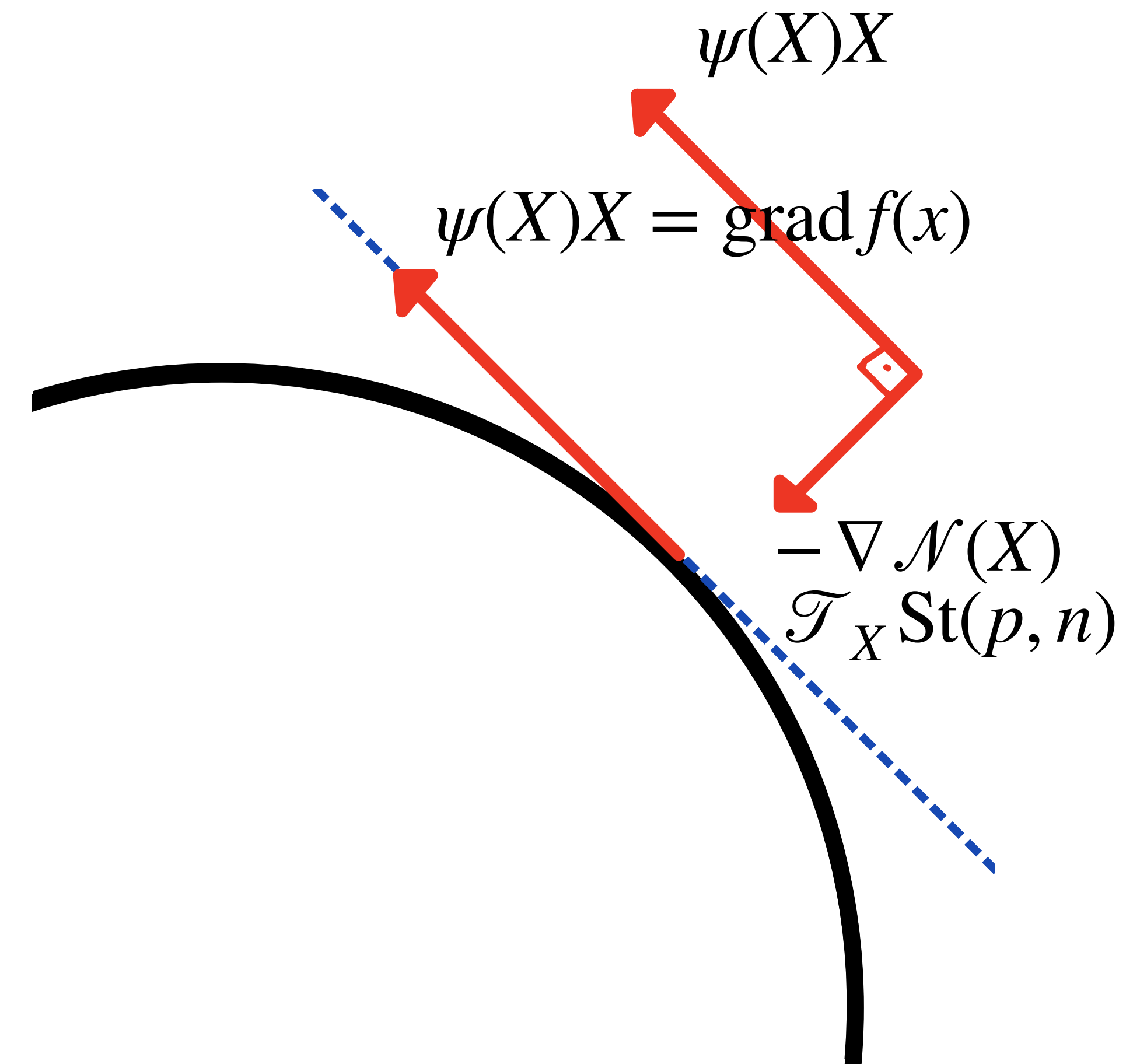
$$\langle \psi(X)X, X(X^\top X - I_p) \rangle = 0$$

- Interpretation of $\psi(X)X$ in the canonical metric $g_X^c(\cdot, \cdot)^*$

$$\psi(X)X = \text{Proj}_{\mathcal{T}_X \text{St}(p,n)} \nabla f(X) = \text{grad} f(X)$$

$$g_X^c(\xi, \zeta) := \langle \xi, (I_n - \frac{1}{2}XX^\top)\zeta \rangle \quad \text{for all } \xi, \zeta \in \mathbb{R}^{n \times p},$$

- What is the interpretation when away from the Stiefel manifold?



* The Geometry of Algorithms with Orthogonality Constraints.
Edelman et. al, SIAM. J. Matrix Anal. & Appl. 1998

Interpretation of the landing field (away from the Stiefel manifold)

- Consider a generalization of the Stiefel manifold and a map

$$\forall M > 0 : \text{St}_M(p, n) := \{Y \in \mathbb{R}^{n \times p} : Y^\top Y = M\} \quad \text{and} \quad \Phi_M : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{n \times p} : X \mapsto Y = XM^{\frac{1}{2}}.$$

- We can define a metric

$$g_Y(\xi, \zeta) := g_{\Phi_M^{-1}(Y)}^c(\Phi_M^{-1}(\xi), \Phi_M^{-1}(\zeta)) \quad \text{where} \quad g_X^c(\xi, \zeta) := \langle \xi, (I_n - \frac{1}{2}XX^\top)\zeta \rangle$$

$$(\text{St}(p, n), g^c) \quad \overset{\Phi_M \text{ isometry}}{\longleftrightarrow} \quad (\text{St}_M(p, n), g)$$

Interpretation of the landing field

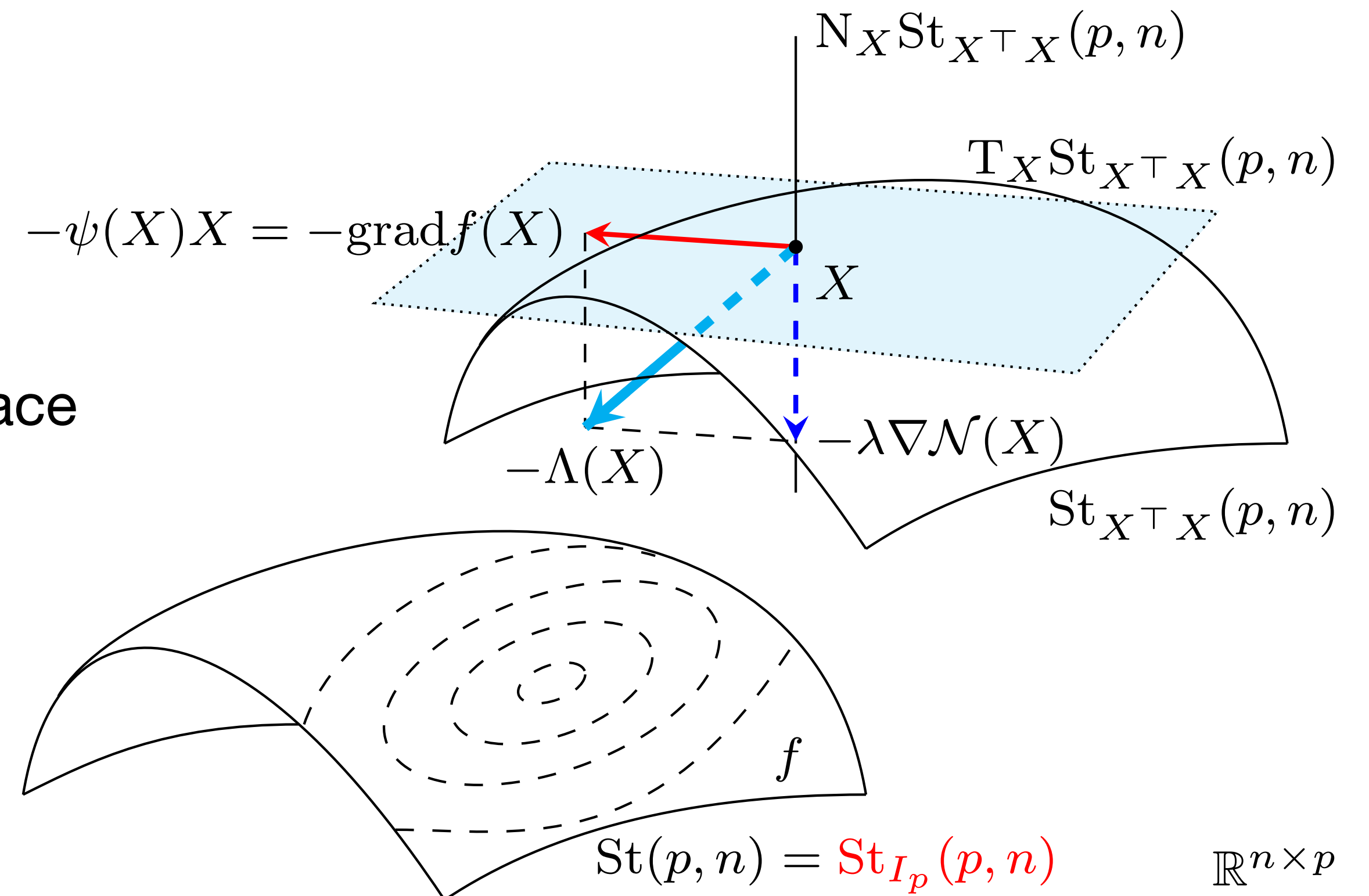
(on the generalized Stiefel manifold $\text{St}_{X^\top X}(p, n)$)

- Riemannian gradient of f on $(\text{St}_{X^\top X}(p, n), g_{X^\top X})$

$$\text{grad}_{\text{St}_{X^\top X}(p, n)} f(X) = \psi(X)X$$

- The normal component belongs to the normal space

$$\nabla \mathcal{N}(X) \in \text{N}_X \text{St}_{X^\top X}(p, n)$$



Convergence of the landing flow

(existence and uniqueness)

Differentiate $\mathcal{N}(X(t))$

$$\frac{d}{dt} \mathcal{N}(X(t)) = \langle \dot{X}(t), \nabla \mathcal{N}(X(t)) \rangle \quad \text{by orthogonality}$$

$$= - \langle \psi(X(t))X(t) + \lambda \nabla \mathcal{N}(X(t)), \nabla \mathcal{N}(X(t)) \rangle$$

$$= - \lambda \|\nabla \mathcal{N}(X(t))\|_F^2 \leq 0, \quad \text{for any } \lambda > 0$$

- hence $\mathcal{N}(X(t))$ is nonincreasing and remains in a closed compact set,
- for $\nabla f(X)$ locally Lipschitz, $\Lambda(X)$ is also locally Lipschitz and,
- by Picard–Lindelöf theorem, there is a unique solution $t \mapsto \varphi_t(X_0)$ such that $\varphi_0(X_0) = X_0$

Convergence of the landing flow (convergence to the Stiefel manifold)

- Let $\chi(t) = X(t)^\top X(t)$

$$\begin{aligned} \dot{\chi}(t) &= \dot{X}(t)^\top X(t) + X(t)^\top \dot{X}(t), & \text{where } \dot{X}(t) &= -\Lambda(X(t)) = \psi(X(t))X(t) + \lambda \nabla \mathcal{N}(X(t)) \\ &= -2\lambda \chi(t) \left(\chi(t) - I_p \right), \end{aligned}$$

- polynomial of a symmetric matrix $\Rightarrow \chi(t)$ has constant eigenvectors
- The eigenvalues of $\chi(t)$ evolve as:

$$\chi_i(t) = \frac{\chi_i(0)e^{2\lambda t}}{\chi_i(0)(e^{2\lambda t} - 1) + 1} \longrightarrow 1, \quad \text{for } \lambda > 0, \quad \text{thus } \lim_{t \rightarrow \infty} \mathcal{N}(\phi_t(X_0)) = 0.$$

Convergence of the landing flow (convergence to the local minima)

- Let $\mathcal{C} \subseteq \text{St}(p, n)$ be the set of critical points on the manifold, we have

$X^* \in \mathcal{C}$ if and only if $\Lambda(X^*) = 0$, by $\psi(X)X$ being a Riemannian gradient and

by orthogonality of $\psi(X)X$ and $\nabla \mathcal{N}(X^*)$

- For all $X_0 \in \mathbb{R}_*^{n \times p}$, the ω -limit points of $\varphi_t(X_0)$ belong to \mathcal{C} .
Therefore, the landing system $\dot{X}(t) = -\Lambda(X(t))$, converges to the set of critical points of f relative to $\text{St}(p, n)$.
- For all $X_0 \in \mathbb{R}_*^{n \times p}$, if X^* is a local minimum and isolated critical point of f relative to $\text{St}(p, n)$, and if X^* is an ω -limit point of $\varphi_t(X_0)$, then $\lim_{t \rightarrow \infty} \varphi_t(X_0) = X^*$.

Landing algorithm

- Discretize the flow
 - $X^{t+1} = X^t - \eta_t \Lambda (X^t)$, where $\psi(X)X + \lambda \nabla \mathcal{N}(X)$,
- Numerical experiments with Stochastic gradient descent (SGD) to compare
 - Riemannian SGD (Retraction -QR)
 - ℓ_2 -penalty (regularization with $+\frac{\lambda}{4}\|X^\top X - I_p\|_F^2$)
- Fixed step-size $\eta = 0.1$, landing parameter $\lambda = 1$

Online PCA

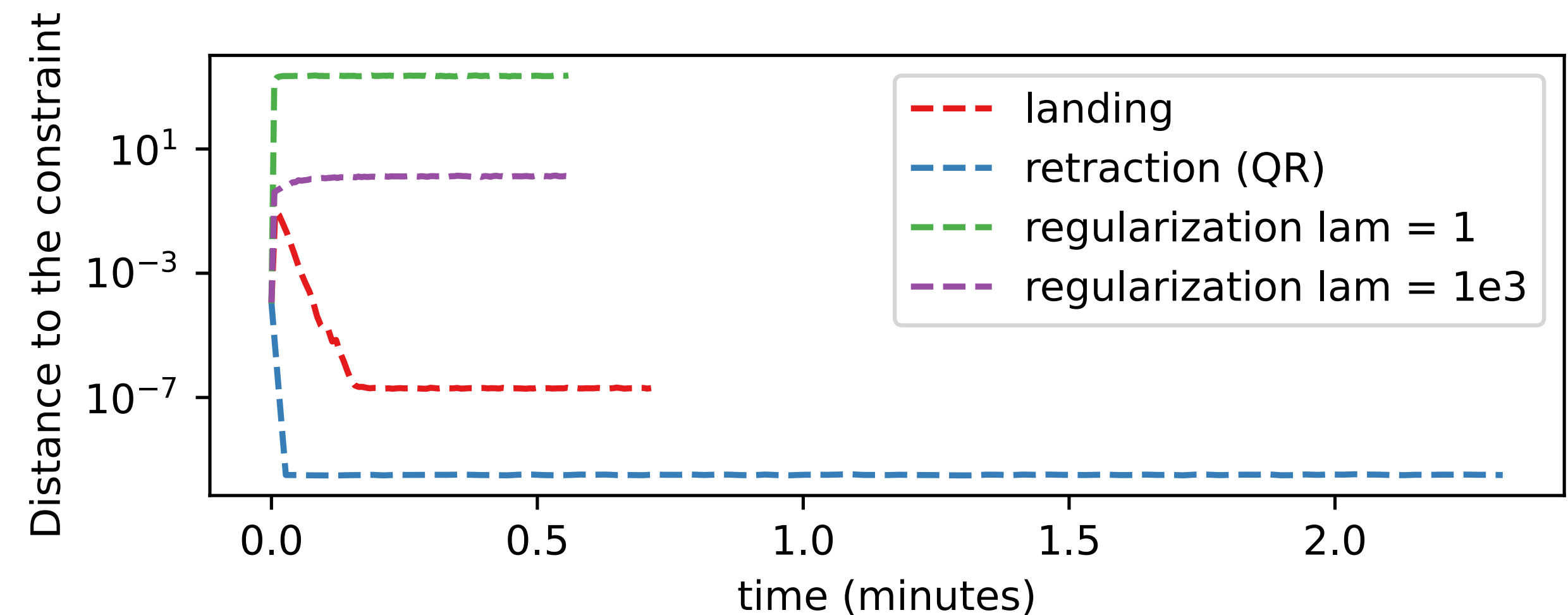
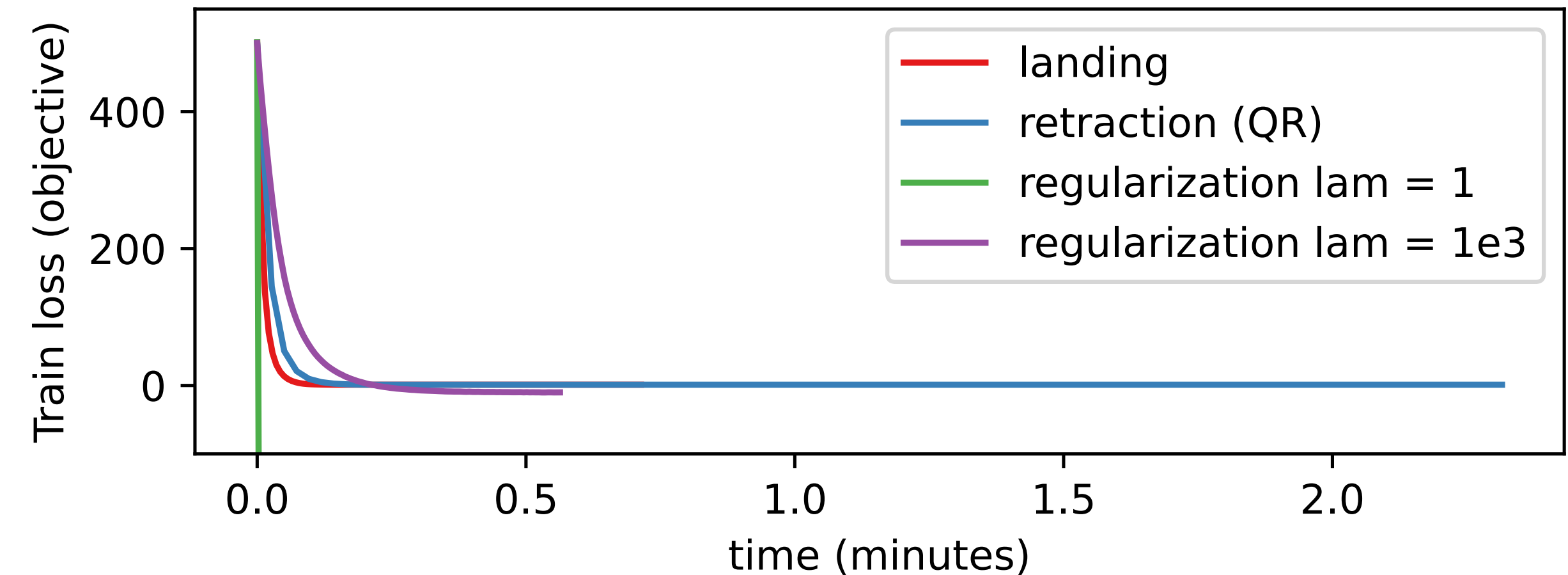
Consider

$$\min -\frac{1}{2} \|AX\|_F^2 \quad \text{s.t.} \quad X \in \text{St}(p, n),$$

where

- $A \in \mathbb{R}^{m \times n}$, with $m = 10\,000$, $n = 2\,000$
- $p = 1\,000$

with Stochastic gradient descent with batch size of 128 rows.



Orthogonal convolutional neural network*

Consider

$$\min \sum_i^N \ell(f_{\Theta}(x_i), y_i) \quad \text{s.t.} \quad \forall \theta \in \Theta_{\text{orth}} : \theta_i \in \text{St}(p, n)$$

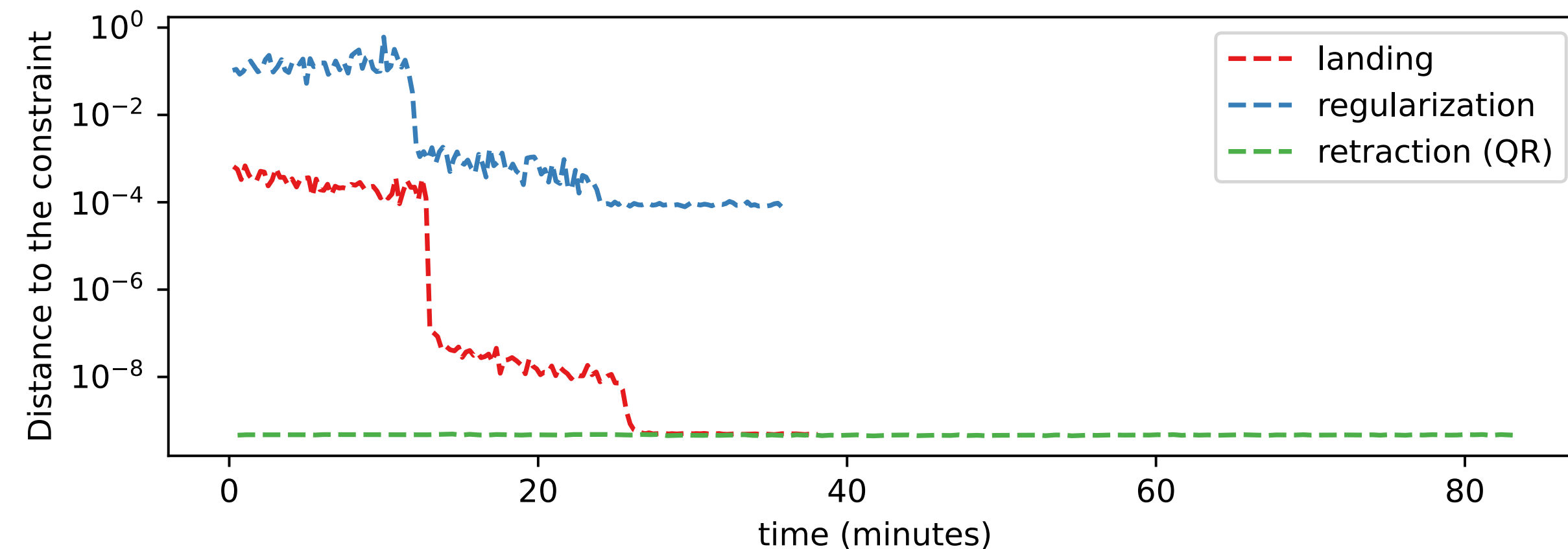
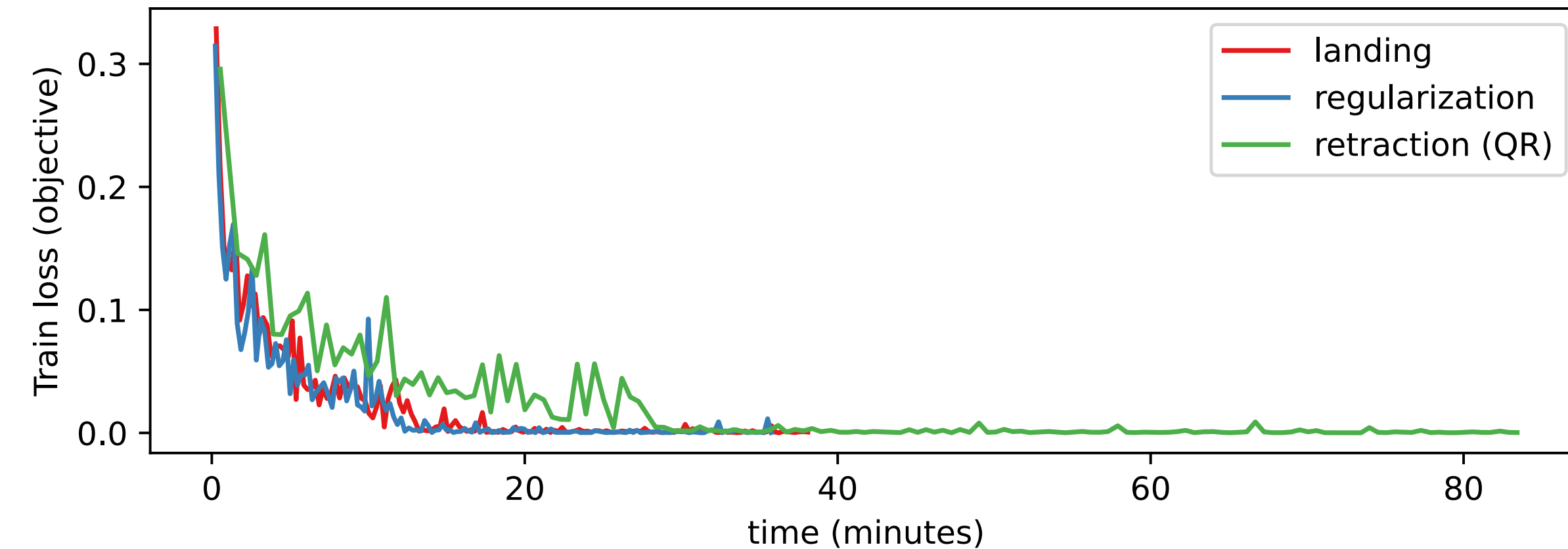
where

- $f_{\Theta}(\cdot)$ is a VGG16** convolutional neural network
- Θ_{orth} includes 13 matrices with size $\approx (1\,000)^2$
- trained on CIFAR-10

with a batch size of 128 samples, fixed step-size

* Orthogonal convolutional neural networks
Wang et al., CVPR 2020

** Very Deep Convolutional Net. for Large-Scale Image Recognition
Simonyan & Zisserman, ICLR 2015



Conclusions

- We propose a landing flow/algorithm and
 - provide a geometric interpretation of the field
 - analyze the continuous gradient flow and show its convergence to local minima
 - demonstrate with numerical experiments its efficiency
- Future work:
 - analysis of the discrete case
 - clever rules for the step-sizes in terms of λ and η
 - possible extensions, for example higher-order and acceleration