

# DS-607 - Optimization for Data Science

- focus on optimization algorithms over non-Euclidean geometry
- first half will be covering optimization over smooth manifolds
- At some point in second half, we will start Computational Optimal transport.
- Aim: Understand the concepts, read papers, implement algorithms, propose novel solutions for various applications

- Pre-requisites include some understanding of linear algebra, optimization, probability & statistics, and ML.
- References (provided in course logistics page)
- Attendance is compulsory.
- Grading: some combination of assignments, quizzes, midsem, endsem, course project.
- Scribing?

Consider the following problem:

$$\min_{x \in S} f(x) \quad f: S \rightarrow \mathbb{R}$$

What information will influence our algorithms

Say  $S = \mathbb{R}^n$ , i.e., a linear space.

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

$$\nabla f(x) < \cdot, \cdot \rangle$$

$$\frac{1}{2} \nabla^2 f(x) \geq \langle \cdot, \cdot \rangle$$

Say  $f(\cdot)$  is differentiable and we apply the standard gradient descent

What implicit choice are we making?

The choice of inner product  $\langle u, v \rangle = u^T v$

We choose to turn  $\mathbb{R}^n$  into a Euclidean space endowed with the standard inner product.

Say  $S = M$  is a *smooth manifold*.

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

We choose to turn  $M$  into a Riemannian manifold.

If  $f(\cdot)$  is smooth, we have a Riemannian gradient and can use it to develop Riemannian gradient descent algorithm.

Example of Optimization problem on manifolds.

e.g.)  $A \in \mathbb{R}^{n \times n}$ ,  $A = A^T$ . Obtain largest eigenvalue of  $A$ ?

If  $A = A^T$ ,  $A = \sum_{i=1}^n \lambda_i v_i v_i^T$   $\lambda_1 \geq \dots \geq \lambda_n \in \mathbb{R}$   
 $v_1, \dots, v_n$ ! orthogonal!

$$\|v_i\|=1, v_i^T v_j = 0$$

$$f(x) = x^T A x$$

$$\max_{x \in \mathbb{R}^n} x^T A x = \lambda_1$$

s.t.  $\|x\|=1$



$$S^{n-1} = \{x \in \mathbb{R}^n | x^T x = 1\}$$

is a manifold.



Eg:  $A \in \mathbb{R}^{m \times n}$ . Obtain largest singular value of A

$$A = \sum_{i=1}^r \sigma_i u_i v_i^T \quad (\text{SVD})$$

$\sigma_1 \geq \dots \geq \sigma_r$   
 $u_1, \dots, u_r$  are orthonormal in  $\mathbb{R}^m$   
 $v_1, \dots, v_r$  are " "  $\mathbb{R}^n$

$$f(x, y) = x^T A y$$

$$\max_{\substack{x \in \mathbb{R}^m, y \in \mathbb{R}^n \\ \|x\|=1, \|y\|=1}} x^T A y = \sigma_i$$

$S^{m-1} \times S^{n-1}$  is manifold

e.g.:  $A = A^T$ ,  $A = A^T$ , want top-k eigenspace of  $A$ .

Consider the problem

$$f(x) = \text{Tr}(x^T A x) \quad \text{where} \quad X \in \text{St}(n, k) = \left\{ X \in \mathbb{R}^{n \times k} \mid \begin{array}{l} \text{Stiefel manifold.} \\ x^T x = I_k \end{array} \right\}$$

$$\max_{\substack{x \in \mathbb{R}^{n \times k} \\ x^T x = I_k}} f(x) = \lambda_1 + \lambda_2 + \dots + \lambda_k.$$

Let  $\mathbb{O}(k) = \{ Q \in \mathbb{R}^{k \times k} \mid Q^T Q = Q Q^T = I_k \}$

Then we note that

$$\begin{aligned} f(x) &= \text{Tr}(x^T A x), \quad x \in \mathbb{S}^n(n, k) \\ &= \text{Tr}((xQ)^T A (xQ)), \quad x \in \mathbb{S}^n(n, k), Q \in \mathbb{O}(k) \\ &= \text{Tr}(Q^T x^T A x Q) \\ &= \text{Tr}(x^T A x) \end{aligned}$$

$\mathbb{O}(k)$  is a manifold.

Define:  $x_i^T x_i = I_k$ ,  $x_j^T x_j = I_k$

for  $x_1, x_2 \in St(n, k)$ ,  $x_1 \sim x_2 \Leftrightarrow x_1 = x_2 Q$  for some  $Q \in O(k)$

What is common between  $x_1$  and  $x_2$  s.t.  $x_1 \sim x_2$

Other than  $x_1 \in St(n, k)$  and  $x_2 \in St(n, k)$  ?

$$\text{span}(x_1) = \text{span}(x_2)$$

Let the equivalence class corresponding to  $x \in St(n, k)$  be denoted by  $[x]$

Then,  $[x] \in St(n, k)/\sim$  (Stiefel manifold modulo equivalence relations)

Visualize  $S\Gamma(n, k)$  as all possible ways of choosing a co-ordinate system for a  $k$ -dimensional subspace in  $\mathbb{R}^n$ .

Then  $S\Gamma(n, k)/\mathbb{R}$  is just the set of all the planes ignoring the co-ordinate system

$S\Gamma(n, k)/\mathbb{R}$  is called the quotient space and is a manifold.

## Applications

### 1. Low-rank matrix learning?

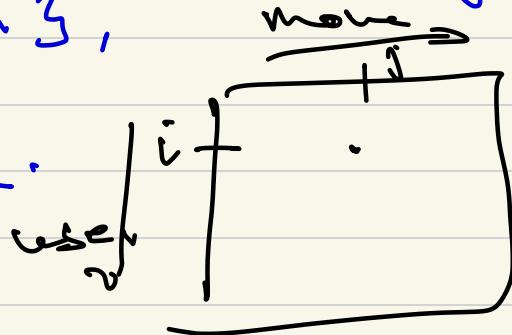
The set  $\mathbb{R}_r^{\text{man}} = \{ X \in \mathbb{R}^{m \times n} : \text{rank}(X) = r \}$  is a manifold.



- Netflix Challenge: Given 'm' users and 'n' movies,

$x_{ij}$  representing the score (1-5) that a user  $i$  gives to movie  $j$ , and  $S_2 = \{(i,j) \mid \text{score given by user } i \text{ to movie } j \text{ is known}\}$ ,

Complete  $X$  given  $X_{S_2}$ .



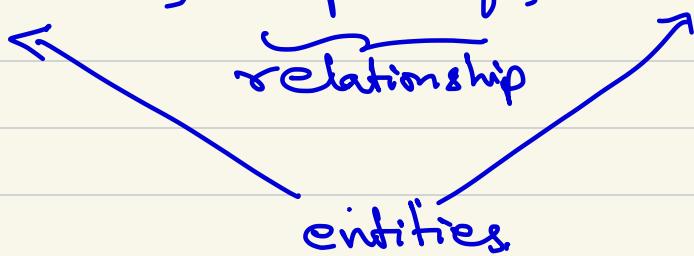
- Image or video completion - few pixels or frames are corrupted.

- Link predictions in knowledge graphs /social network!

new

Discover facts of form Subject-predicate-object.  
given a few such facts

e.g. { Mumbai, capital-of, Maharashtra }



These are modeled using low-rank tensors

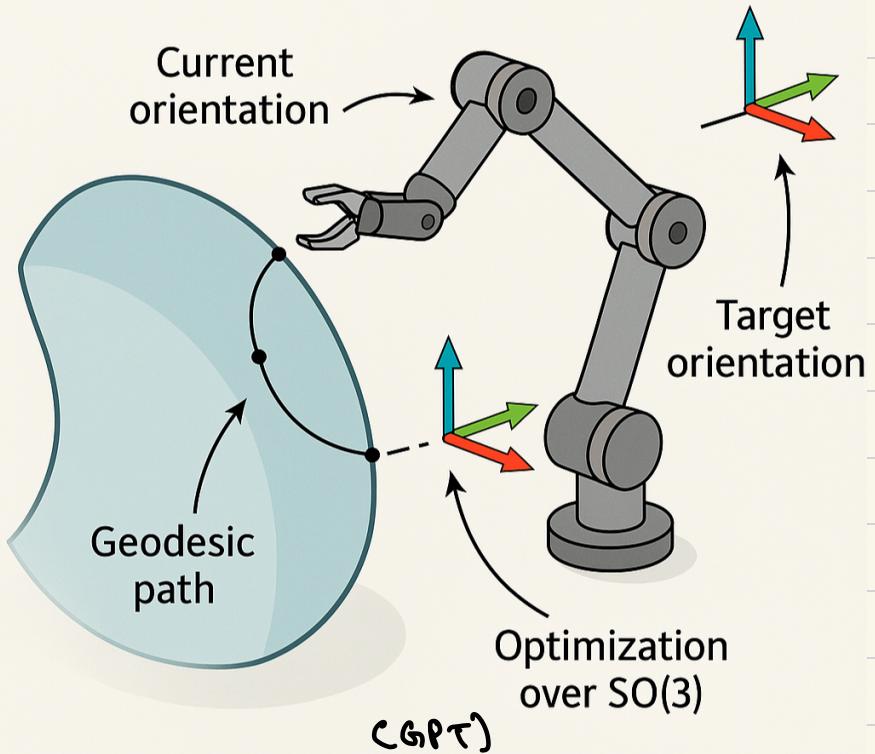
## 2. Symmetric Positive definite matrices:

$$\{ S \in \mathbb{R}^{d \times d} \mid S \succeq 0 \}$$

Used for modeling covariance matrices

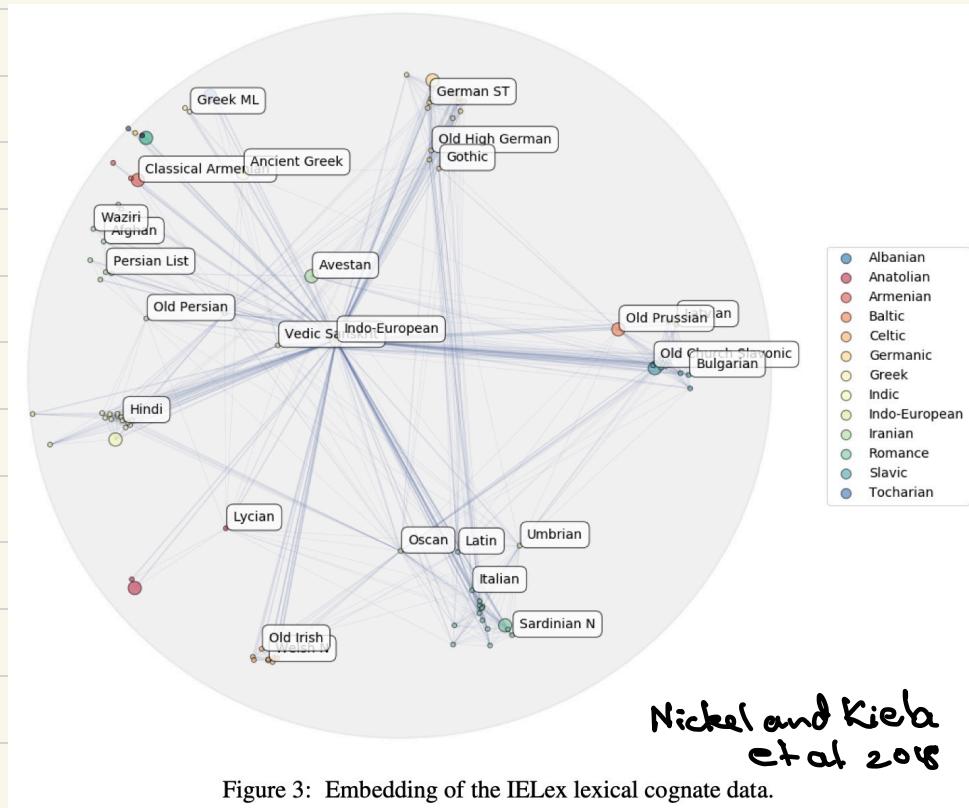
- Modeling mixture of Gaussians
- Classifying radar data, EEG recordings, or human action (using, e.g., SPDNet)
  - Raw motion data (from video or sensors)
  - features include joint positions, velocities over time
  - Covariance matrix  $\Sigma$  is computed.
  - Classify  $\Phi(x) (\succeq 0)$  obtained via SPD-Net

(3) Orientation learning via  $SO(3)$   
(set of rotation matrices)



(4)

## Hyperbolic spaces - representing hierarchies



## Hide & Seek: Transformer Symmetries Obscure Sharpness & Riemannian Geometry Finds It

Marvin E. da Silva<sup>1,2</sup> Felix Dangel<sup>2</sup> Sageev Oore<sup>1,2</sup>

## Efficient Optimization with Orthogonality Constraint: a Randomized Riemannian Submanifold Method

Andi Han<sup>1,2</sup> Pierre-Louis Poirion<sup>1</sup> Akiko Takeda<sup>1,3</sup>

## Score-based Pullback Riemannian Geometry: Extracting the Data Manifold Geometry using Anisotropic Flows

Willem Diepeveen<sup>\*1</sup> Georgios Batzolis<sup>\*2</sup> Zakhar Shumaylov<sup>3</sup> Carola-Bibiane Schönlieb<sup>3</sup>

## Fast, Accurate Manifold Denoising by Tunneling Riemannian Optimization

Shiyu Wang<sup>1,2</sup> Mariam Avagyan<sup>1,2</sup> Yihan Shen<sup>3,2</sup> Arnaud Lamy<sup>1,2</sup> Tingran Wang<sup>1,2</sup> Szabolcs Márka<sup>4,2</sup>  
Zsuzsa Márka<sup>5,2</sup> John Wright<sup>1,6,2</sup>

## Preconditioned Riemannian Gradient Descent Algorithm for Low-Multilinear-Rank Tensor Completion

Yuanwei Zhang<sup>1</sup> Fengmiao Bian<sup>2</sup> Xiaoqun Zhang<sup>1,3</sup> Jian-Feng Cai<sup>2</sup>

Some recent work in ICMl'25 using Riemannian optimization

Lets start with unstrained convex minimization problems

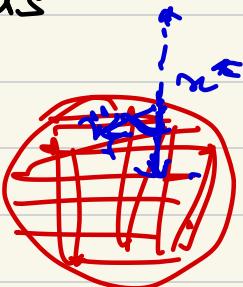
$$\min_{x \in \mathbb{R}^n} f(x)$$

Assume  $f \in C^0(\mathbb{R}^n)$

first order necessary conditions  
of optimality?

Now, we consider convex minimization problem over convex sets

$$\min_{x \in \mathcal{S}} f(x)$$



$\mathcal{S} \subset \mathbb{R}^d$  non-empty,  
closed,  
convex

first order necessary conditions  
for optimality?

constrained.

Assume  $x^*$  is a local sol<sup>n</sup>.  
Then  
 $\langle \nabla f(x^*), x - x^* \rangle \geq 0$   
 $\forall x \in \mathcal{S}$ .

Equivalent way:

$$\langle \nabla f(x^*), p \rangle \geq 0$$

$\forall p \in T_{\mathcal{S}}(x^*)$ , where

$$T_{\mathcal{S}}(x) = \text{closure} \{ h(y-x) \mid y \in \mathcal{S}, h \in \Sigma \}$$

Let us have a more concrete problem:

$$\begin{array}{ll} \min_{x \in \mathbb{R}^d} & f(x) \text{ s.t } h(x) = 0 \\ & (\text{e.g. } x_1 + x_2) \\ & (\text{e.g. } x_1^2 + x_2^2 - 1 = 0) \end{array}$$



K.K.T conditions of optimality (necessary)  
(for  $x^*$  to be a local optimizer)

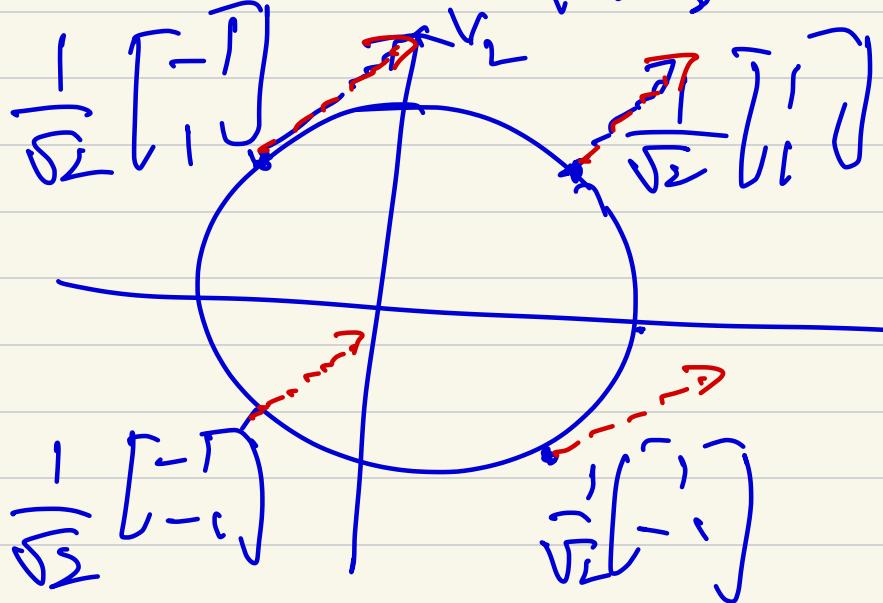
①  $\underline{\nabla f(x)} - \lambda \underline{\nabla h(x)} = 0$

②  $h(x) = 0$

for some  $(x, \lambda) = (x^*, \lambda^*)$

$$f_{\text{Car}} = x_1 + x_2$$

$$h_{\text{Car}} = x_1^2 + x_2^2 - 1$$



$$\nabla f_{\text{Car}} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\nabla h_{\text{Car}} = \begin{bmatrix} 2x_1 \\ 2x_2 \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ 1 \end{bmatrix} - \lambda \begin{bmatrix} 2x_1 \\ 2x_2 \end{bmatrix} = 0$$

$$x_1 \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \frac{1}{2\lambda} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

## first order analysis

$$\cdot \nabla h^T s = 0$$

$$h(x+s) \approx h(x) + \underbrace{\nabla h(x)^T s}_{\nabla h(x)^T s = 0} = 0$$

$$f(x+s) \approx f(x) + \nabla f(x)^T s$$

If  $f(x)$  is locally optimal

$$\Rightarrow \nabla f(x)^T s \geq 0$$

$$\nabla h(\omega)^T s = 0$$

$$\nabla f(\omega)^T s < 0$$