

# Optimization On Manifolds

Pierre-Antoine Absil  
Robert Mahony  
Rodolphe Sepulchre

Based on ‘‘Optimization Algorithms on Matrix Manifolds’’, Princeton  
University Press, January 2008

Compiled on August 21, 2008

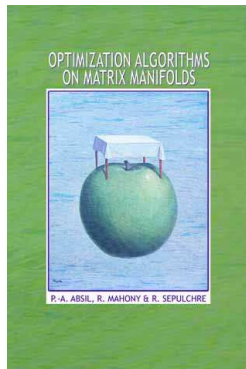
## Collaboration

Chris Baker  
(Florida State University and Sandia National Laboratories)

Kyle Gallivan  
(Florida State University)

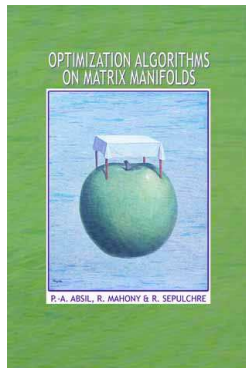
Paul Van Dooren  
(Université catholique de Louvain)

## Reference



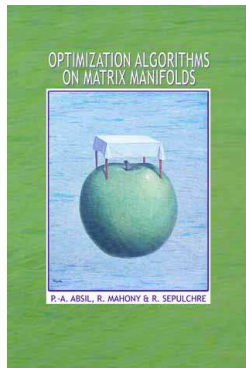
*Optimization Algorithms on Matrix Manifolds*  
P.-A. Absil, R. Mahony, R. Sepulchre  
Princeton University Press, January 2008

## About the reference



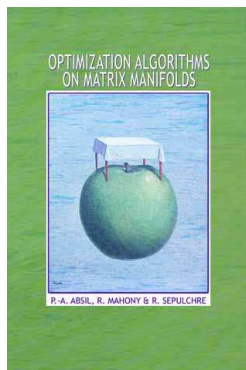
- ▶ The publisher, Princeton University Press, has been a non-profit company since 1910.
- ▶ Official publication date is January 2008.
- ▶ Copies are already shipping.

## Reference: contents



1. Introduction
2. Motivation and applications
3. Matrix manifolds: first-order geometry
4. Line-search algorithms
5. Matrix manifolds: second-order geometry
6. Newton's method
7. Trust-region methods
8. A constellation of superlinear algorithms

# Matrix Manifolds: first-order geometry



## Chap 3: Matrix Manifolds: first-order geometry

1. Charts, atlases, manifolds
2. Differentiable functions
3. Embedded submanifolds
4. Quotient manifolds
5. Tangent vectors and differential maps
6. Riemannian metric, distance, gradient

## Smooth optimization in $\mathbb{R}^n$

General unconstrained optimization problem in  $\mathbb{R}^n$ :

Let

$$f : \mathbb{R}^n \rightarrow \mathbb{R},$$

The real-valued function  $f$  is termed the *cost function* or *objective function*.

Problem: find  $x_* \in \mathbb{R}^n$  such that there exists  $\epsilon > 0$  for which

$$f(x) \geq f(x_*) \text{ whenever } \|x - x_*\| < \epsilon.$$

Such a point  $x_*$  is called a *local minimizer* of  $f$ .

## Smooth optimization in $\mathbb{R}^n$

General unconstrained optimization problem in  $\mathbb{R}^n$ :

Let

$$f : \mathbb{R}^n \rightarrow \mathbb{R},$$

The real-valued function  $f$  is termed the *cost function* or *objective function*.

Problem: find  $x_* \in \mathbb{R}^n$  such that there exists a **neighborhood**  $\mathcal{N}$  of  $x_*$  such that

$$f(x) \geq f(x_*) \text{ whenever } x \in \mathcal{N}.$$

Such a point  $x_*$  is called a *local minimizer* of  $f$ .



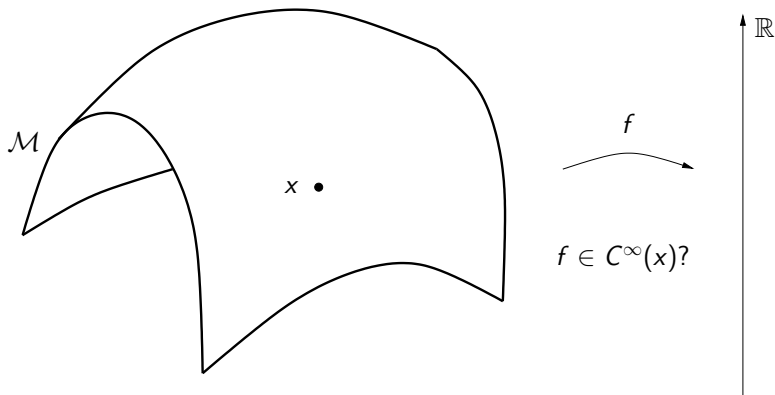
## Smooth optimization *beyond* $\mathbb{R}^n$

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

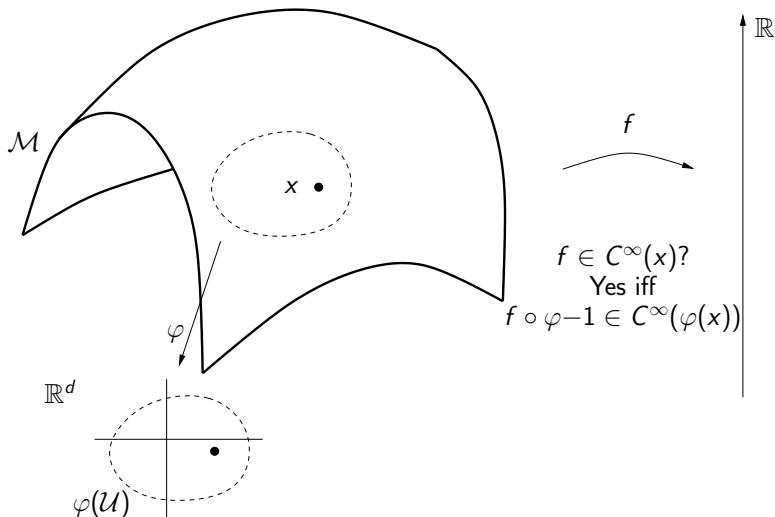
- ▶ Several optimization techniques require the cost function to be differentiable to some degree:
  - ▶ Steepest-descent at  $x$  requires  $Df(x)$ .
  - ▶ Newton's method at  $x$  requires  $D^2f(x)$ .
- ▶ Can we go beyond  $\mathbb{R}^n$  without losing the concept of differentiability?

$$\arg \min_{x \in \mathbb{R}^n} f(x) \quad \rightsquigarrow \quad \arg \min_{x \in \mathcal{M}} f(x)$$

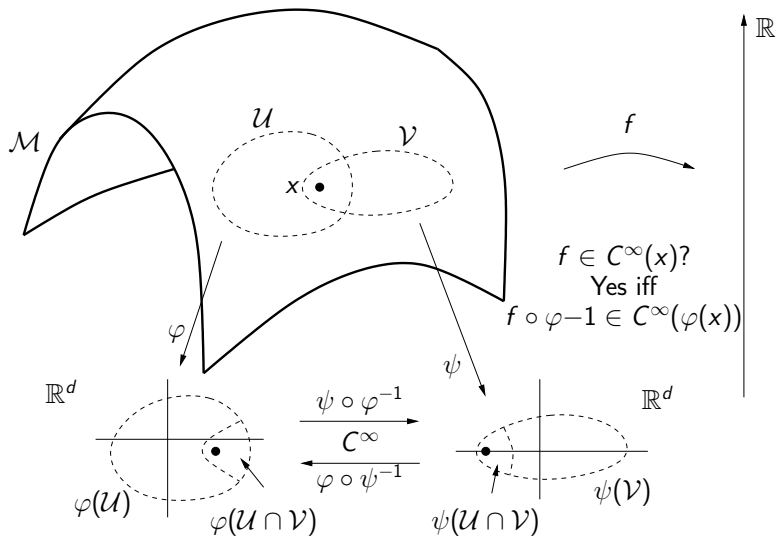
## Smooth optimization on a manifold: what “smooth” means



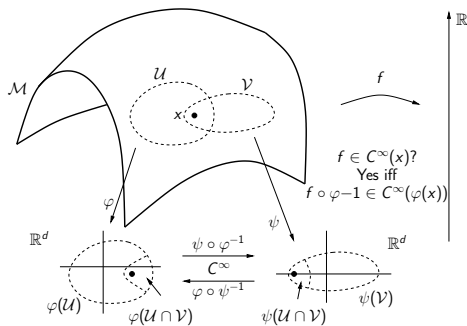
## Smooth optimization on a manifold: what “smooth” means



## Smooth optimization on a manifold: what “smooth” means



## Smooth optimization on a manifold: what “smooth” means

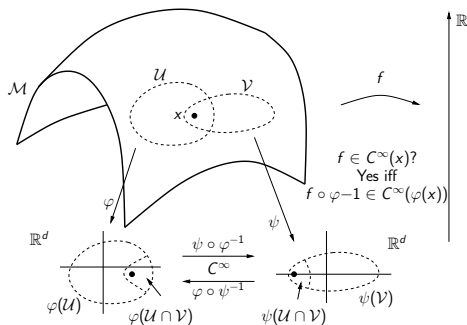


**Chart:**  $U \xrightarrow[\text{bij.}]{\varphi} \varphi(U)$

**Atlas:** Collection of “compatible chars” that cover  $\mathcal{M}$

**Manifold:** Set with an atlas

# Optimization on manifolds in its most abstract formulation

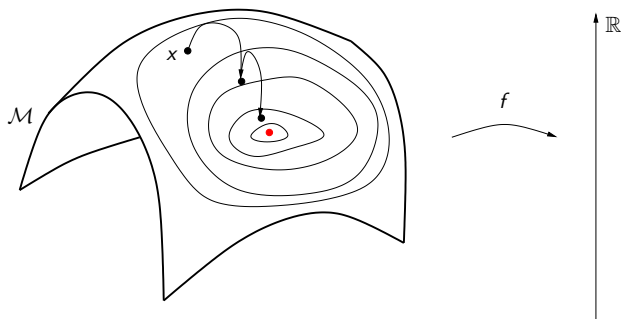


Given:

- ▶ A set  $\mathcal{M}$  endowed (explicitly or implicitly) with a manifold structure (i.e., a collection of compatible charts).
- ▶ A function  $f : \mathcal{M} \rightarrow \mathbb{R}$ , smooth in the sense of the manifold structure.

Task: Compute a local minimizer of  $f$ .

## Optimization on manifolds: algorithms

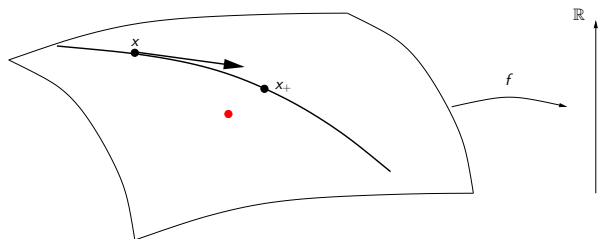


Given:

- ▶ A set  $\mathcal{M}$  endowed (explicitly or implicitly) with a manifold structure (i.e., a collection of compatible charts).
- ▶ A function  $f : \mathcal{M} \rightarrow \mathbb{R}$ , smooth in the sense of the manifold structure.

Task: Compute a local minimizer of  $f$ .

## Previous work on Optimization On Manifolds



Luenberger (1973), *Introduction to linear and nonlinear programming*. Luenberger mentions the idea of performing line search along geodesics, “which we would use if it were computationally feasible (which it definitely is not)”.



## The purely Riemannian era

**Gabay (1982)**, *Minimizing a differentiable function over a differential manifold*. Stepest descent along geodesics; Newton's method along geodesics; Quasi-Newton methods along geodesics.

**Smith (1994)**, *Optimization techniques on Riemannian manifolds*.

Levi-Civita connection  $\nabla$ ; Riemannian exponential; parallel translation.  
But Remark 4.9: If Algorithm 4.7 (Newton's iteration on the sphere for the Rayleigh quotient) is simplified by replacing the exponential update with the update

$$x_{k+1} = \frac{x_k + \eta_k}{\|x_k + \eta_k\|}$$

then we obtain the Rayleigh quotient iteration.

## The pragmatic era

Manton (2002), *Optimization algorithms exploiting unitary constraints*  
“The present paper breaks with tradition by not moving along geodesics”. The geodesic update  $\text{Exp}_x \eta$  is replaced by a projective update  $\pi(x + \eta)$ , the *projection* of the point  $x + \eta$  onto the manifold.

Adler, Dedieu, Shub, et al. (2002), *Newton's method on Riemannian manifolds and a geometric model for the human spine*. The exponential update is relaxed to the general notion of *retraction*. The geodesic can be replaced by any (smoothly prescribed) curve tangent to the search direction.

## Looking ahead: Newton on abstract manifolds

Required: Riemannian manifold  $\mathcal{M}$ ; retraction  $R$  on  $\mathcal{M}$ ; affine connection  $\nabla$  on  $\mathcal{M}$ ; real-valued function  $f$  on  $\mathcal{M}$ .

Iteration  $x_k \in \mathcal{M} \mapsto x_{k+1} \in \mathcal{M}$  defined by

1. Solve the Newton equation

$$\text{Hess } f(x_k)\eta_k = -\text{grad } f(x_k)$$

for the unknown  $\eta_k \in T_{x_k}\mathcal{M}$ , where

$$\text{Hess } f(x_k)\eta_k := \nabla_{\eta_k} \text{grad } f.$$

2. Set

$$x_{k+1} := R_{x_k}(\eta_k).$$

## Looking ahead: Newton on submanifolds of $\mathbb{R}^n$

Required: Riemannian submanifold  $\mathcal{M}$  of  $\mathbb{R}^n$ ; retraction  $R$  on  $\mathcal{M}$ ; real-valued function  $f$  on  $\mathcal{M}$ .

Iteration  $x_k \in \mathcal{M} \mapsto x_{k+1} \in \mathcal{M}$  defined by

1. Solve the Newton equation

$$\text{Hess } f(x_k)\eta_k = -\text{grad } f(x_k)$$

for the unknown  $\eta_k \in T_{x_k}\mathcal{M}$ , where

$$\text{Hess } f(x_k)\eta_k := P_{T_{x_k}\mathcal{M}}\text{grad } f(x_k).$$

2. Set

$$x_{k+1} := R_{x_k}(\eta_k).$$

Looking ahead: Newton on the unit sphere  $S^{n-1}$ 

Required: real-valued function  $f$  on  $S^{n-1}$ .

Iteration  $x_k \in \mathcal{M} \mapsto x_{k+1} \in S^{n-1}$  defined by

1. Solve the Newton equation

$$\begin{cases} P_{x_k} D(\text{grad } f)(x_k)[\eta_k] = -\text{grad } f(x_k) \\ x^T \eta_k = 0, \end{cases}$$

for the unknown  $\eta_k \in \mathbb{R}^n$ , where

$$P_{x_k} = (I - x_k x_k^T).$$

2. Set

$$x_{k+1} := \frac{x_k + \eta_k}{\|x_k + \eta_k\|}.$$

## Looking ahead: Newton for Rayleigh quotient optimization on unit sphere

Iteration  $x_k \in S^{n-1} \mapsto x_{k+1} \in S^{n-1}$  defined by

1. Solve the Newton equation

$$\begin{cases} P_{x_k} A P_{x_k} \eta_k - \eta_k x_k^T A x_k = -P_{x_k} A x_k, \\ x_k^T \eta_k = 0, \end{cases}$$

for the unknown  $\eta_k \in \mathbb{R}^n$ , where

$$P_{x_k} = (I - x_k x_k^T).$$

2. Set

$$x_{k+1} := \frac{x_k + \eta_k}{\|x_k + \eta_k\|}.$$

# Programme

- ▶ Provide background in differential geometry instrumental for algorithmic development
- ▶ Present manifold versions of some classical optimization algorithms: steepest-descent, Newton, conjugate gradients, trust-region methods
- ▶ Show how to turn these abstract geometric algorithms into practical implementations
- ▶ Illustrate several problems that can be rephrased as optimization problems on manifolds.

## Some important manifolds

- ▶ Stiefel manifold  $St(p, n)$ : set of all orthonormal  $n \times p$  matrices.
- ▶ Grassmann manifold  $Grass(p, n)$ : set of all  $p$ -dimensional subspaces of  $\mathbb{R}^n$
- ▶ Euclidean group  $SE(3)$ : set of all rotations-translations
- ▶ Flag manifold, shape manifold, oblique manifold...
- ▶ Several unnamed manifolds



# A manifold-based approach to the symmetric eigenvalue problem

OPT



EVP

OPT

Opt algorithms

for  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ 

⋮

EVP

Algorithms

for EVP

OPT

Opt algorithms  
for  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

$f \equiv$  Rayleigh quotient

EVP

Algorithms  
for EVP

## Rayleigh quotient

Rayleigh quotient of  $(A, B)$ :

$$f : \mathbb{R}_*^n \rightarrow \mathbb{R} : f(y) = \frac{y^T A y}{y^T B y}$$

Let  $A, B$  in  $\mathbb{R}^{n \times n}$ ,  $A = A^T$ ,  $B = B^T \succ 0$ ,

$$A v_i = \lambda_i B v_i$$

with  $\lambda_1 < \lambda_2 \leq \dots \leq \lambda_n$ .

Stationary points of  $f$ :  $\alpha v_i$ , for all  $\alpha \neq 0$ .

Local (and global) minimizers of  $f$ :  $\alpha v_1$ , for all  $\alpha \neq 0$ .

OPT

Opt algorithms  
for  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

$f \equiv$  Rayleigh quotient

EVP

Algorithms  
for EVP

## “Block” Rayleigh quotient

Let  $\mathbb{R}_*^{n \times p}$  denote the set of all full-column-rank  $n \times p$  matrices.  
Generalized (“block”) Rayleigh quotient:

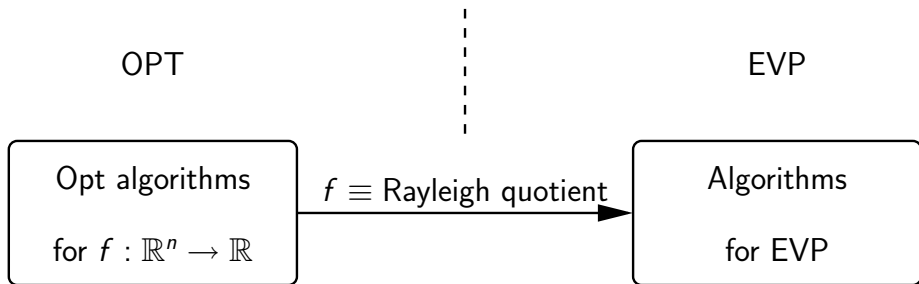
$$f : \mathbb{R}_*^{n \times p} \rightarrow \mathbb{R} : f(Y) = \text{trace} \left( (Y^T B Y)^{-1} Y^T A Y \right)$$

Stationary points of  $f$ :

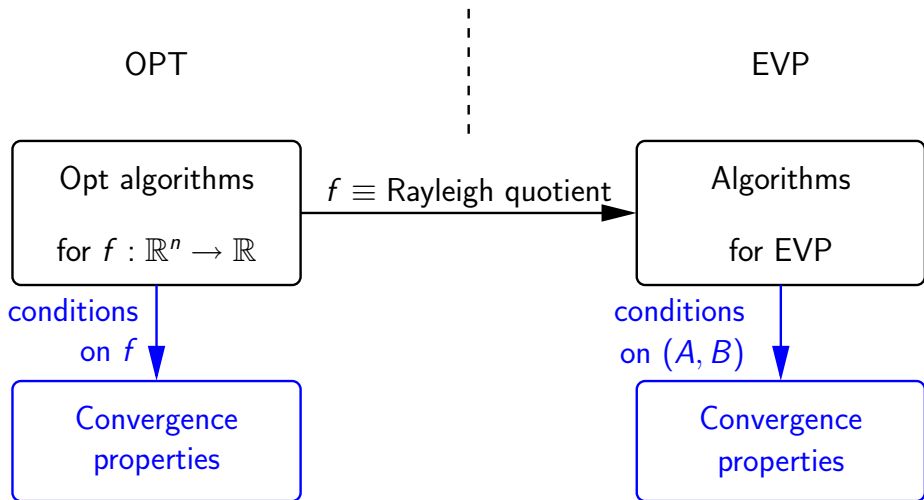
$$\begin{bmatrix} v_{i_1} & \dots & v_{i_p} \end{bmatrix} M, \quad \text{for all } M \in \mathbb{R}_*^{p \times p}.$$

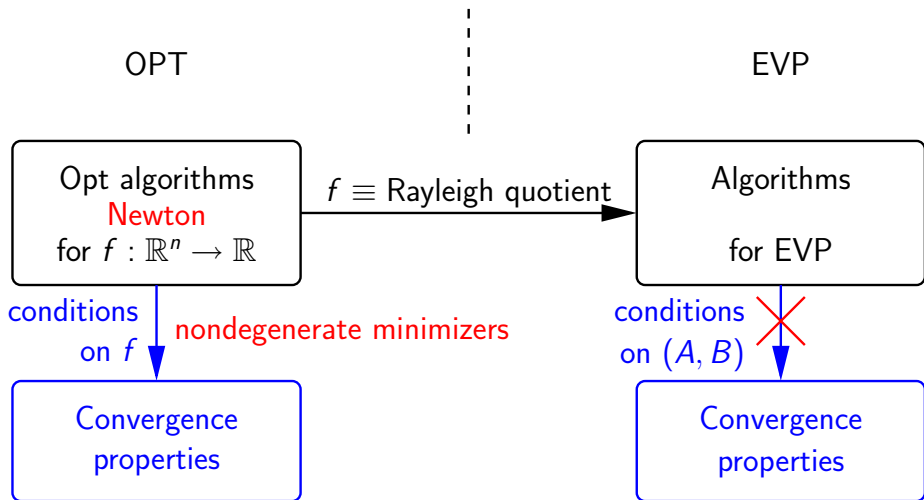
Minimizers of  $f$ :

$$\begin{bmatrix} v_1 & \dots & v_p \end{bmatrix} M, \quad \text{for all } M \in \mathbb{R}_*^{p \times p}.$$









Newton for Rayleigh quotient in  $\mathbb{R}_0^n$ 

Let  $f$  denote the Rayleigh quotient of  $(A, B)$ .

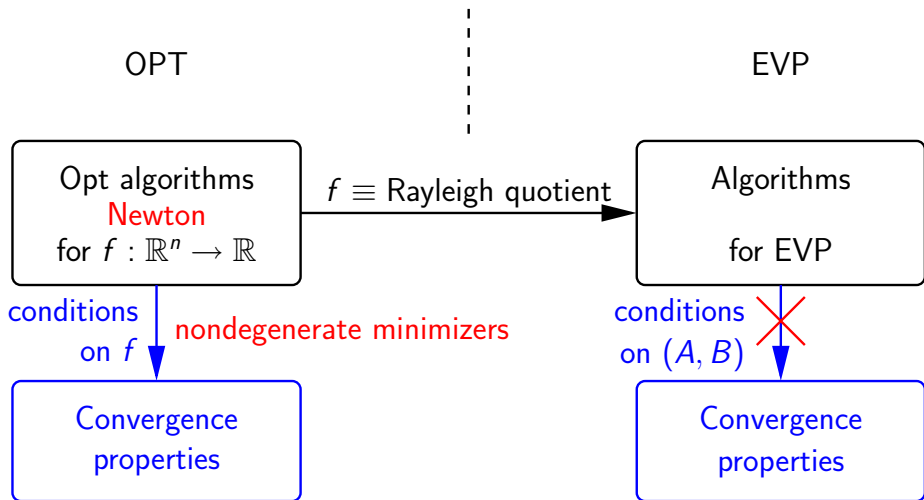
Let  $x \in \mathbb{R}_0^n$  be any point such that  $f(x) \notin \text{spec}(B^{-1}A)$ .

Then the Newton iteration

$$x \mapsto x - (D^2f(x))^{-1} \cdot \text{grad } f(x)$$

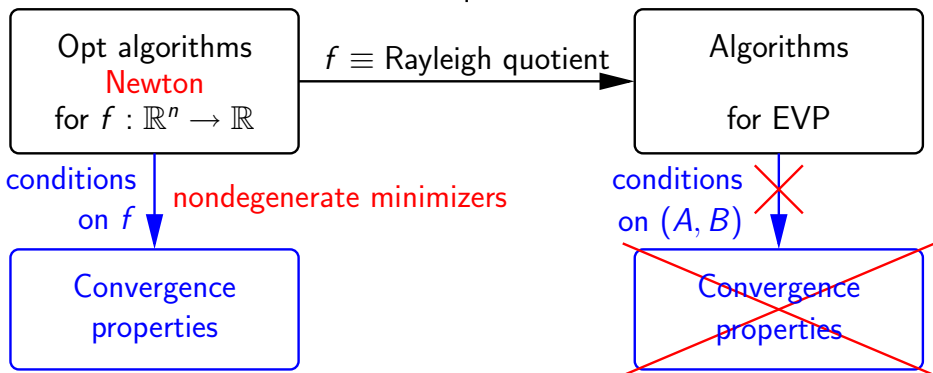
reduces to the iteration

$$x \mapsto 2x.$$



OPT

EVP



## Invariance properties of the Rayleigh quotient

*Rayleigh quotient* of  $(A, B)$ :

$$f : \mathbb{R}_*^n \rightarrow \mathbb{R} : f(y) = \frac{y^T A y}{y^T B y}$$

Invariance:  $f(\alpha y) = f(y)$  for all  $\alpha \in \mathbb{R}_0$ .

## Invariance properties of the Rayleigh quotient

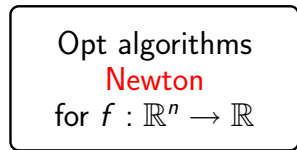
Generalized (“block”) Rayleigh quotient:

$$f : \mathbb{R}_*^{n \times p} \rightarrow \mathbb{R} : f(Y) = \text{trace} \left( (Y^T B Y)^{-1} Y^T A Y \right)$$

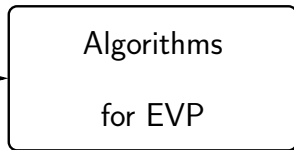
Invariance:  $f(YM) = f(Y)$  for all  $M \in \mathbb{R}_*^{p \times p}$ .

OPT

EVP



$f \equiv$  Rayleigh quotient



conditions  
 on  $f$

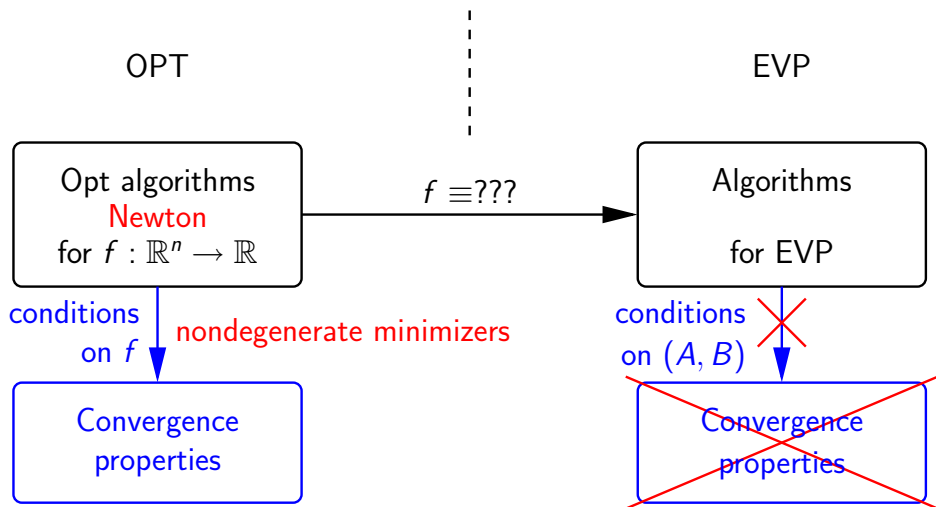
nondegenerate minimizers



conditions  
 on  $(A, B)$





Remedy 1: modify  $f$ 

## Remedy 1: modify $f$

Consider

$$P_A : \mathbb{R}^n \rightarrow \mathbb{R} : x \mapsto P_A(x) := (x^T x)^2 - 2x^T Ax.$$

### Theorem

(i)

$$\min_{x \in \mathbb{R}^n} P_A(x) = -\lambda_n^2$$

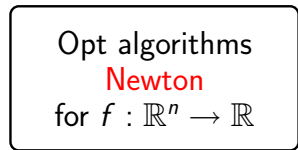
The minimum is attained at any  $\sqrt{\lambda_n} v_n$ , where  $v_n$  is a **unitary** eigenvector related to  $\lambda_n$ .

(ii) The set of critical points of  $P_A$  is  $\{0\} \cup \{\sqrt{\lambda_k} v_k\}$ .

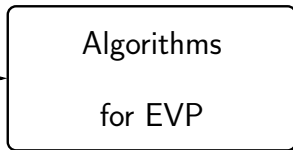
References: Auchmuty (1989), Mongeau and Torki (2004).

OPT

EVP



$f \equiv$  Rayleigh quotient



conditions  
 on  $f$

nondegenerate minimizers

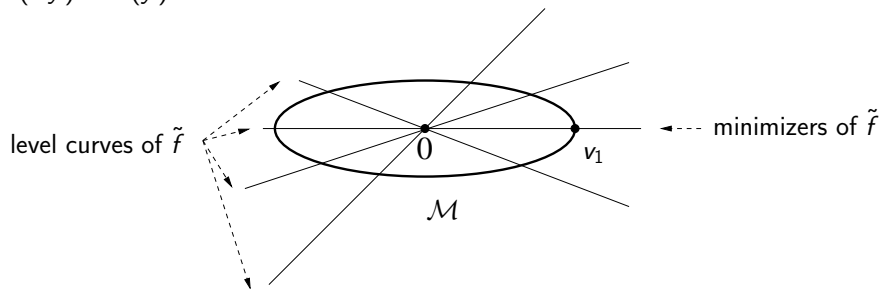


conditions  
 on  $(A, B)$



## EVP: optimization on ellipsoid

$$f(\alpha y) = f(y)$$



## Remedy 2: modify the search space

Instead of

$$f : \mathbb{R}_*^n \rightarrow \mathbb{R} : f(y) = \frac{y^T A y}{y^T B y},$$

minimize

$$f : \mathcal{M} \rightarrow \mathbb{R} : f(y) = \frac{y^T A y}{y^T B y},$$

where

$$\mathcal{M} = \{y \in \mathbb{R}^n : y^T B y = 1\}.$$

Stationary points of  $f$ :  $\pm v_i$ .

Local (and global) minimizers of  $f$ :  $\pm v_1$ .

## Remedy 2: modify search space: block case

Instead of generalized (“block”) Rayleigh quotient:

$$f : \mathbb{R}_*^{n \times p} \rightarrow \mathbb{R} : f(Y) = \text{trace} \left( (Y^T B Y)^{-1} Y^T A Y \right),$$

minimize

$$f : \text{Grass}(p, n) \rightarrow \mathbb{R} : f(\text{col}(Y)) = \text{trace} \left( (Y^T B Y)^{-1} Y^T A Y \right),$$

where  $\text{Grass}(p, n)$  denotes the set of all  $p$ -dimensional subspaces of  $\mathbb{R}^n$ , called the *Grassmann manifold*.

Stationary points of  $f$ :  $\text{col}([v_{i_1} \ \dots \ v_{i_p}])$ .

Minimizer of  $f$ :  $\text{col}([v_1 \ \dots \ v_p])$ .

OPT

Opt algorithms  
**Newton**  
for  $f : \mathcal{M} \rightarrow \mathbb{R}$

conditions  
on  $f$

**nondegenerate minimizers**

Convergence  
properties

$f \equiv$  Rayleigh quotient

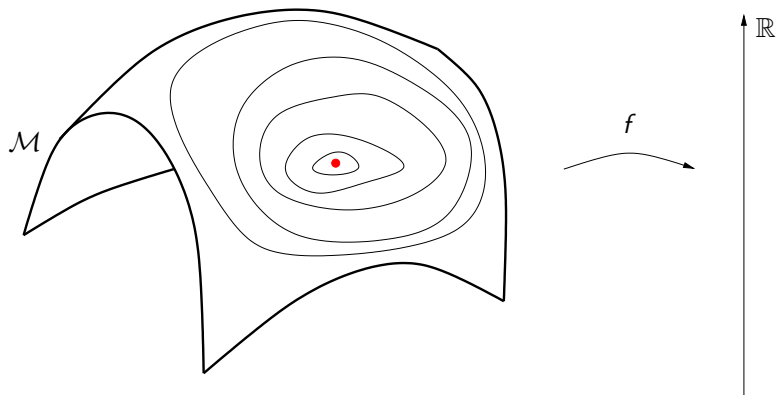
EVP

Algorithms  
for EVP

conditions  
on  $(A, B)$

Convergence  
properties

## Smooth optimization on a manifold: big picture





## Smooth optimization on a manifold: tools

	Purely Riemannian way	Pragmatic way
Search direction	Tangent vector	Tangent vector
Steepest descent dir.	$-\text{grad } f(x)$	$-\text{grad } f(x)$
Derivative of vector field	Levi-Civita connection $\frac{g}{\nabla}$	<b>Any</b> connection $\nabla$
Update	Search along the geodesic tangent to the search direction	Search along <b>any</b> curve tangent to the search direction (scribed by a <i>retraction</i> )
Displacement of tgt vectors	Parallel translation induced by $\frac{g}{\nabla}$	<b>Vector Transport</b>

OPT

Opt algorithms  
**Newton**  
for  $f : \mathcal{M} \rightarrow \mathbb{R}$

conditions  
on  $f$

**nondegenerate minimizers**

Convergence  
properties

$f \equiv$  Rayleigh quotient

EVP

Algorithms  
for EVP

conditions  
on  $(A, B)$

Convergence  
properties

## Newton's method on abstract manifolds

Required: Riemannian manifold  $\mathcal{M}$ ; retraction  $R$  on  $\mathcal{M}$ ; affine connection  $\nabla$  on  $\mathcal{M}$ ; real-valued function  $f$  on  $\mathcal{M}$ .

Iteration  $x_k \in \mathcal{M} \mapsto x_{k+1} \in \mathcal{M}$  defined by

1. Solve the Newton equation

$$\text{Hess } f(x_k)\eta_k = -\text{grad } f(x_k)$$

for the unknown  $\eta_k \in T_{x_k}\mathcal{M}$ , where  $\text{Hess } f(x_k)\eta_k := \nabla_{\eta_k}\text{grad } f$ .

2. Set

$$x_{k+1} := R_{x_k}(\eta_k).$$

OPT

Opt algorithms  
**Newton**  
 for  $f : \mathcal{M} \rightarrow \mathbb{R}$

conditions  
 on  $f$

**nondegenerate minimizers**

Convergence  
 properties

$f \equiv$  Rayleigh quotient

EVP

Algorithms  
 for EVP

conditions  
 on  $(A, B)$

Convergence  
 properties

## Convergence of Newton's method on abstract manifolds

### Theorem

Let  $x_* \in \mathcal{M}$  be a **nongenerate critical point** of  $f$ , i.e.,  $\text{grad } f(x_*) = 0$  and  $\text{Hess } f(x_*)$  invertible.

Then there exists a neighborhood  $\mathcal{U}$  of  $x_*$  in  $\mathcal{M}$  such that, for all  $x_0 \in \mathcal{U}$ , Newton's method generates an infinite sequence  $(x_k)_{k=0,1,\dots}$  **converging superlinearly** (at least quadratically) to  $x_*$ .

OPT

Opt algorithms  
**Newton**  
for  $f : \mathcal{M} \rightarrow \mathbb{R}$

conditions  
on  $f$

**nondegenerate minimizers**

Convergence  
properties

$f \equiv$  Rayleigh quotient

EVP

Algorithms  
for EVP

conditions  
on  $(A, B)$

Convergence  
properties

## Geometric Newton for Rayleigh quotient optimization

Iteration  $x_k \in S^{n-1} \mapsto x_{k+1} \in S^{n-1}$  defined by

1. Solve the Newton equation

$$\begin{cases} P_{x_k} A P_{x_k} \eta_k - \eta_k x_k^T A x_k = -P_{x_k} A x_k, \\ x_k^T \eta_k = 0, \end{cases}$$

for the unknown  $\eta_k \in \mathbb{R}^n$ , where

$$P_{x_k} = (I - x_k x_k^T).$$

2. Set

$$x_{k+1} := \frac{x_k + \eta_k}{\|x_k + \eta_k\|}.$$

## Geometric Newton for Rayleigh quotient optimization: block case

Iteration  $\text{col}(Y_k) \in \text{Grass}(p, n) \mapsto \text{col}(Y_{k+1}) \in \text{Grass}(p, n)$  defined by

1. Solve the linear system

$$\begin{cases} P_{Y_k}^h (AZ_k - Z_k(Y_k^T Y_k)^{-1} Y_k^T A Y_k) = -P_{Y_k}^h (A Y_k) \\ Y_k^T Z_k = 0 \end{cases}$$

for the unknown  $Z_k \in \mathbb{R}^{n \times p}$ , where

$$P_{Y_k}^h = (I - Y_k(Y_k^T Y_k)^{-1} Y_k^T).$$

2. Set

$$Y_{k+1} = (Y_k + Z_k)N_k$$

where  $N_k$  is a nonsingular  $p \times p$  matrix chosen for normalization.



OPT

Opt algorithms  
**Newton**  
for  $f : \mathcal{M} \rightarrow \mathbb{R}$

conditions  
on  $f$

**nondegenerate minimizers**

Convergence  
properties

$f \equiv$  Rayleigh quotient

EVP

Algorithms  
for EVP

conditions  
on  $(A, B)$

Convergence  
properties

## Convergence of the EVP algorithm

### Theorem

Let  $Y_* \in \mathbb{R}^{n \times p}$  be such that  $\text{col}(Y_*)$  is a **spectral** invariant subspace of  $B^{-1}A$ . Then there exists a neighborhood  $\mathcal{U}$  of  $\text{col}(Y_*)$  in  $\text{Grass}(p, n)$  such that, for all  $Y_0 \in \mathbb{R}^{n \times p}$  with  $\text{col}(Y_0) \in \mathcal{U}$ , Newton's method generates an infinite sequence  $(Y_k)_{k=0,1,\dots}$  such that  $(\text{col}(Y_k))_{k=0,1,\dots}$  **converges superlinearly** (at least quadratically) to  $\text{col}(Y_*)$  on  $\text{Grass}(p, n)$ .

OPT

Opt algorithms  
**Newton**  
for  $f : \mathcal{M} \rightarrow \mathbb{R}$

conditions  
on  $f$

**nondegenerate minimizers**

Convergence  
properties

$f \equiv$  Rayleigh quotient

⋮

EVP

Algorithms  
for EVP

conditions  
on  $(A, B)$

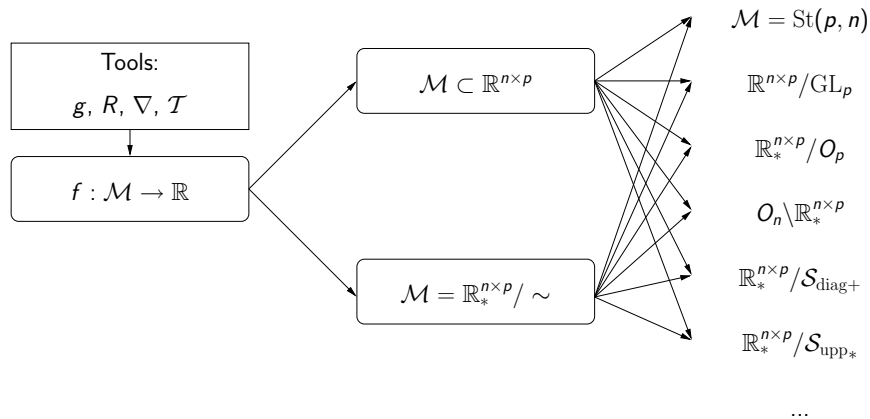
Convergence  
properties

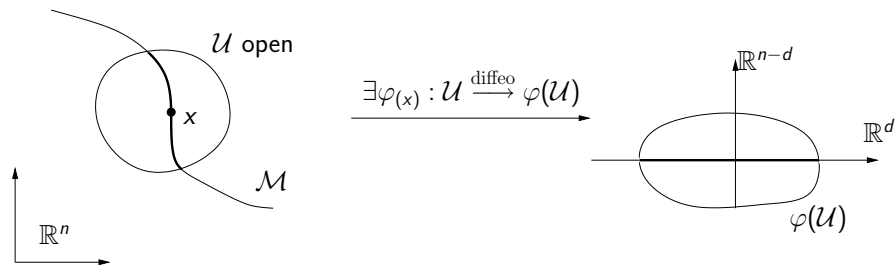
## Other optimization methods

- ▶ Trust-region methods: PAA, C. G. Baker, K. A. Gallivan, *Trust-region methods on Riemannian manifolds*, Foundations of Computational Mathematics, 2007.
- ▶ “Implicit” trust-region methods: PAA, C. G. Baker, K. A. Gallivan, submitted.

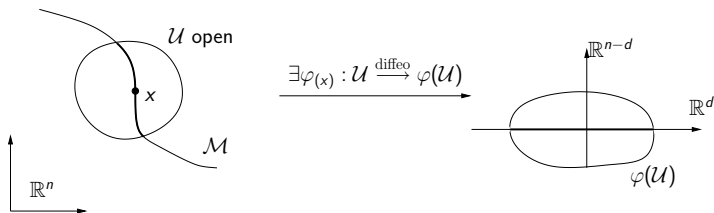
# Manifolds

## Manifolds, submanifolds, quotient manifolds



Submanifolds of  $\mathbb{R}^n$ 

The set  $\mathcal{M} \subset \mathbb{R}^n$  is termed a *submanifold* of  $\mathbb{R}^n$  if the situation described above holds for all  $x \in \mathcal{M}$ .

Submanifolds of  $\mathbb{R}^n$ 

The manifold structure on  $\mathcal{M}$  is defined in a unique way as the manifold

structure generated by the atlas  $\left\{ \left[ \begin{array}{c} e_1^T \\ \vdots \\ e_d^T \end{array} \right] \varphi(x)|_{\mathcal{M}} : x \in \mathcal{M} \right\}$ .



Back to the basics: partial derivatives in  $\mathbb{R}^n$ 

Let  $F : \mathbb{R}^n \rightarrow \mathbb{R}^q$ .

Define  $\partial_i F : \mathbb{R}^n \rightarrow \mathbb{R}^q$  by

$$\partial_i F(x) = \lim_{t \rightarrow 0} \frac{F(x + te_i) - F(x)}{t}.$$

If  $\partial_i F$  is defined and continuous on  $\mathbb{R}^n$ , then  $F$  is termed *continuously differentiable*, denoted by  $F \in C^1$ .

Back to the basics: (Fréchet) derivative in  $\mathbb{R}^n$ 

If  $F \in C^1$ , then

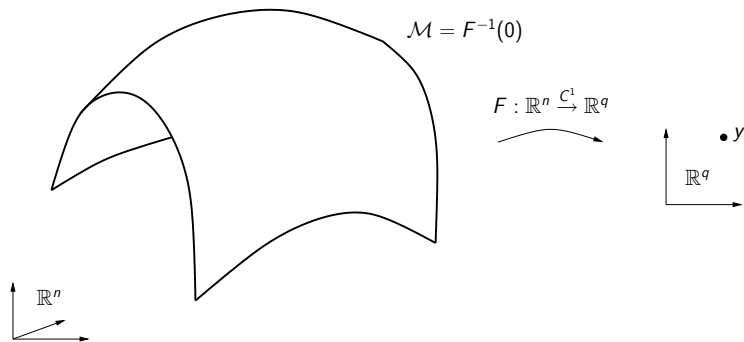
$$DF(x) : \mathbb{R}^n \xrightarrow{\text{lin}} \mathbb{R}^q : z \mapsto DF(x)[z] := \lim_{t \rightarrow 0} \frac{F(x + tz) - F(x)}{t}$$

is the *derivative* (or *differential*) of  $F$  at  $x$ .

We have  $DF(x)[z] = J_F(x)z$ , where the matrix

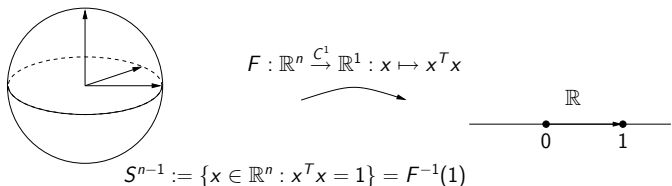
$$J_F(x) = \begin{bmatrix} \partial_1(e_1^T F)(x) & \cdots & \partial_n(e_1^T F)(x) \\ \vdots & \ddots & \vdots \\ \partial_1(e_q^T F)(x) & \cdots & \partial_n(e_q^T F)(x) \end{bmatrix}$$

is the *Jacobian matrix* of  $F$  at  $x$ .

Submanifolds of  $\mathbb{R}^n$ : sufficient condition

$y \in \mathbb{R}^q$  is a *regular value* of  $F$  if, for all  $x \in F^{-1}(y)$ ,  $DF(x)$  is an onto function (*surjection*).

**Theorem (submersion theorem):** If  $y \in \mathbb{R}^q$  is a regular value of  $F$ , then  $F^{-1}(y)$  is a submanifold of  $\mathbb{R}^n$ .

Submanifolds of  $\mathbb{R}^n$ : sufficient condition: application

The unit sphere

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

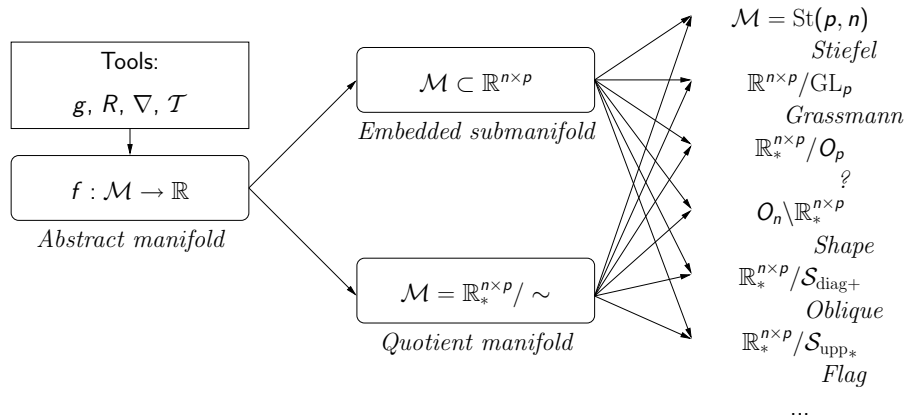
is a submanifold of  $\mathbb{R}^n$ .

Indeed, for all  $x \in S^{n-1}$ , we have that

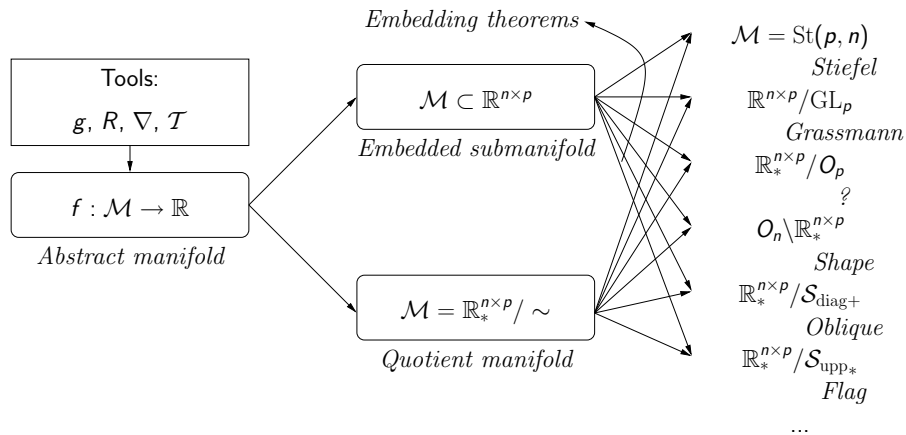
$$DF(x) : \mathbb{R}^n \rightarrow \mathbb{R} : z \mapsto DF(x)[z] = x^T z + z^T x$$

is an onto function.

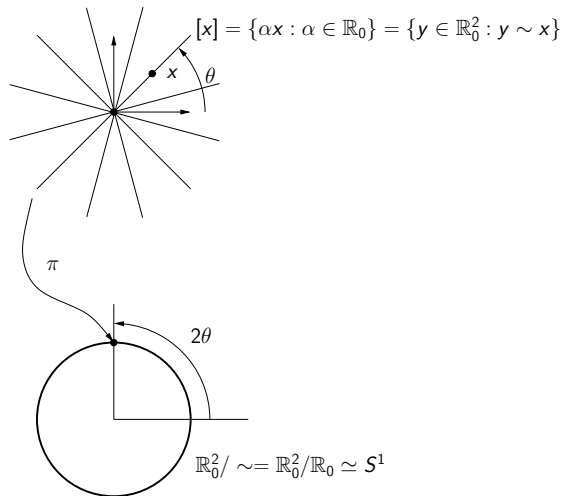
## Manifolds, submanifolds, quotient manifolds



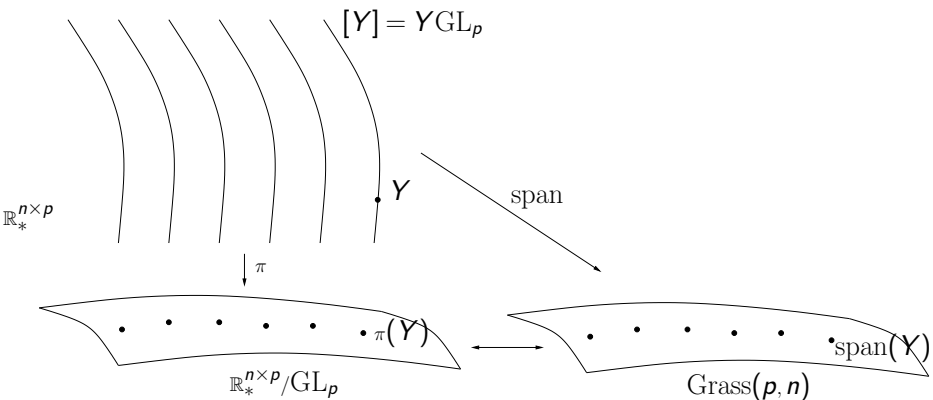
## Manifolds, submanifolds, quotient manifolds



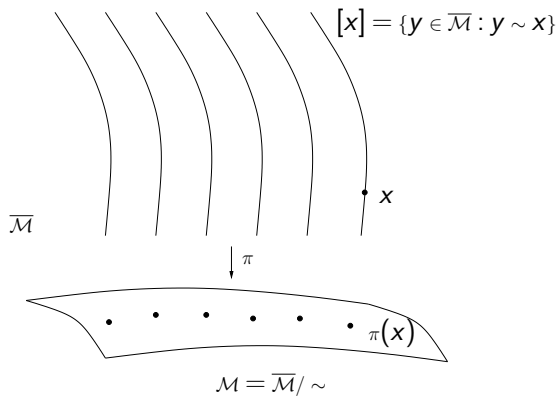
## A simple quotient set: the projective space

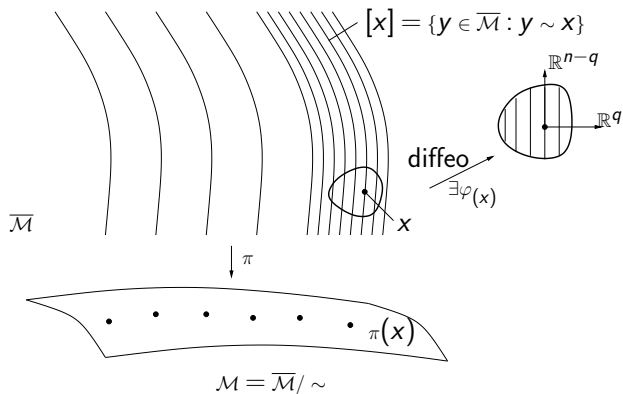


A slightly less simple quotient set:  $\mathbb{R}_*^{n \times p} / \text{GL}_p$

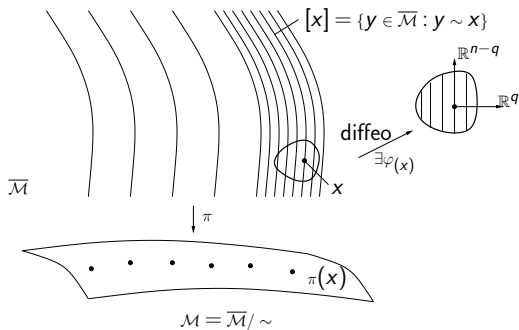




Abstract quotient set  $\overline{\mathcal{M}}/\sim$ 

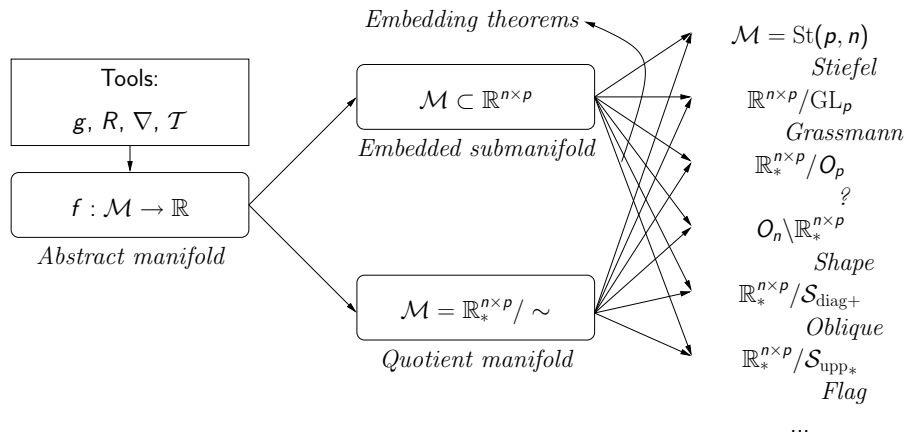
Abstract quotient manifold  $\overline{\mathcal{M}}/\sim$ 

The set  $\overline{\mathcal{M}}/\sim$  is termed a *quotient manifold* if the situation described above holds for all  $x \in \overline{\mathcal{M}}$ .

Abstract quotient manifold  $\overline{\mathcal{M}}/\sim$ 

The manifold structure on  $\overline{\mathcal{M}}/\sim$  is defined in a unique way as the manifold structure generated by the atlas  $\left\{ \left[ \begin{array}{c} e_1^T \\ \vdots \\ e_q^T \end{array} \right] \varphi(x) \circ \pi^{-1} : x \in \overline{\mathcal{M}} \right\}$ .

## Manifolds, submanifolds, quotient manifolds



## Manifolds, and where they appear

- ▶ Stiefel manifold  $\text{St}(p, n)$  and orthogonal group  $O_p = \text{St}(n, n)$

$$\text{St}(p, n) = \{X \in \mathbb{R}^{n \times p} : X^T X = I_p\}$$

Applications: computer vision; principal component analysis; independent component analysis...

- ▶ Grassmann manifold  $\text{Grass}(p, n)$

Set of all  $p$ -dimensional subspaces of  $\mathbb{R}^n$

Applications: various dimension reduction problems...

- ▶  $\mathbb{R}_*^{n \times p} / O_p$

$$X \sim Y \Leftrightarrow \exists Q \in O_p : Y = XQ$$

Applications: Low-rank approximation of symmetric matrices; low-rank approximation of tensors...

## Manifolds, and where they appear

- ▶ Shape manifold  $O_n/\mathbb{R}_*^{n \times p}$

$$Y \sim Y \Leftrightarrow \exists U \in O_n : Y = UX$$

Applications: shape analysis

- ▶ Oblique manifold  $\mathbb{R}_*^{n \times p}/\mathcal{S}_{\text{diag}+}$

$$\mathbb{R}_*^{n \times p}/\mathcal{S}_{\text{diag}+} \simeq \{Y \in \mathbb{R}_*^{n \times p} : \text{diag}(Y^T Y) = I_p\}$$

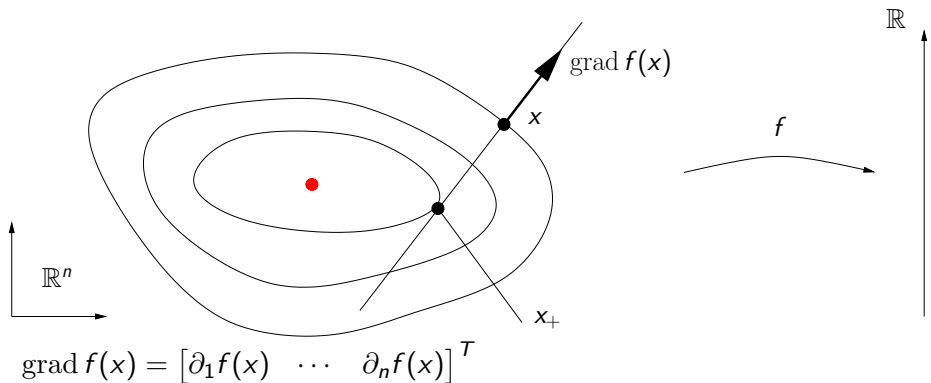
Applications: independent component analysis; factor analysis (oblique Procrustes problem)...

- ▶ Flag manifold  $\mathbb{R}_*^{n \times p}/\mathcal{S}_{\text{upp}+}$

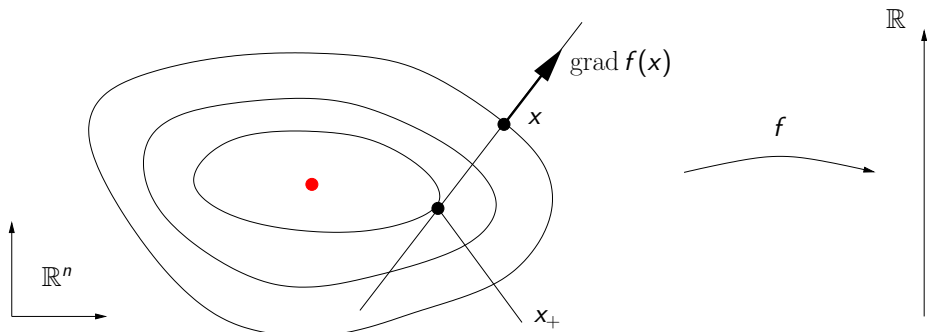
Elements of the flag manifold can be viewed as a  $p$ -tuple of linear subspaces  $(\mathcal{V}_1, \dots, \mathcal{V}_p)$  such that  $\dim(\mathcal{V}_i) = i$  and  $\mathcal{V}_i \subset \mathcal{V}_{i+1}$ .

Applications: analysis of QR algorithm...

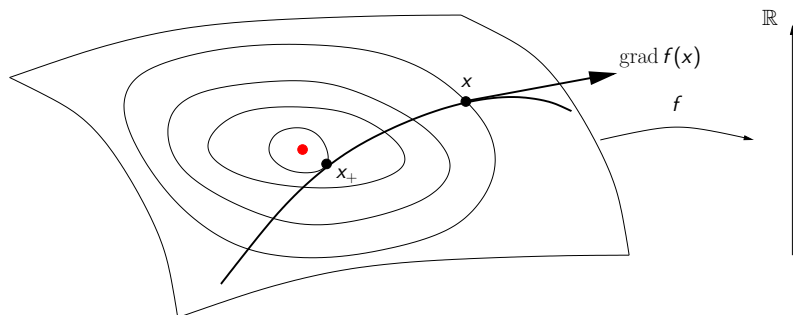
# Steepest-descent methods on manifolds

Steepest-descent in  $\mathbb{R}^n$ 



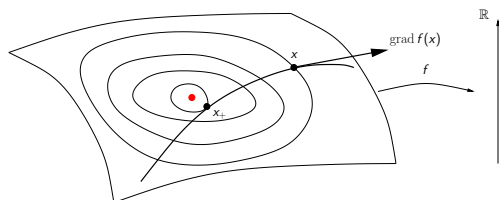
Steepest-descent: from  $\mathbb{R}^n$  to manifolds

	$\mathbb{R}^n$	Manifold
Search direction	Vector at $x$	Tangent vector at $x$
Steepest-desc. dir.	$-\text{grad } f(x)$	$-\text{grad } f(x)$
Curve	$\gamma : t \mapsto x - t \text{grad } f(x)$	$\gamma$ s.t. $\gamma(0) = x$ and $\dot{\gamma}(0) = -\text{grad } f(x)$

Steepest-descent: from  $\mathbb{R}^n$  to manifolds

	$\mathbb{R}^n$	Manifold
Search direction	Vector at $x$	Tangent vector at $x$
Steepest-desc. dir.	$-\text{grad } f(x)$	$-\text{grad } f(x)$
Curve	$\gamma : t \mapsto x - t \text{grad } f(x)$	$\gamma$ s.t. $\gamma(0) = x$ and $\dot{\gamma}(0) = -\text{grad } f(x)$

## Update directions: tangent vectors



Let  $\gamma$  be a curve in the manifold  $\mathcal{M}$  with  $\gamma(0) = x$ .

For an abstract manifold, the definition  $\dot{\gamma}(0) = \frac{d\gamma}{dt}(0) = \lim_{t \rightarrow 0} \frac{\gamma(t) - \gamma(0)}{t}$  is meaningless.

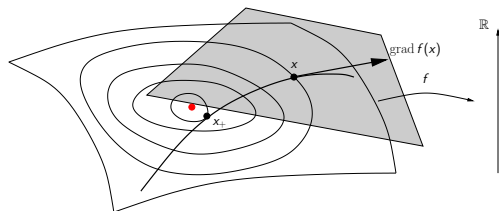
Instead, define:  $Df(x)[\dot{\gamma}(0)] := \left. \frac{d}{dt} f(\gamma(t)) \right|_{t=0}$

If  $\mathcal{M} \subset \mathbb{R}^n$  and  $f = \bar{f}|_{\mathcal{M}}$ , then

$$Df(x)[\dot{\gamma}(0)] = D\bar{f}(x) \left[ \frac{d\gamma}{dt}(0) \right].$$

The application  $\dot{\gamma}(0) : f \mapsto Df(x)[\dot{\gamma}(0)]$  is a *tangent vector* at  $x$ .

## Update directions: tangent spaces



The set

$$T_x \mathcal{M} = \{ \dot{\gamma}(0) : \gamma \text{ curve in } \mathcal{M} \text{ through } x \text{ at } t = 0 \}$$

is the *tangent space* to  $\mathcal{M}$  at  $x$ .

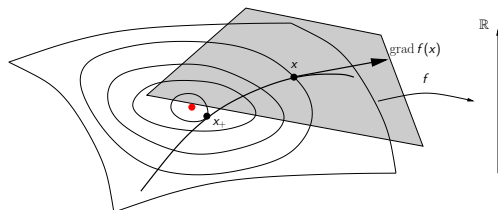
With the definition

$$\alpha \dot{\gamma}_1(0) + \beta \dot{\gamma}_2(0) : f \mapsto \alpha Df(x)[\dot{\gamma}_1(0)] + \beta Df(x)[\dot{\gamma}_2(0)],$$

the tangent space  $T_x \mathcal{M}$  becomes a *linear space*.

The *tangent bundle*  $T\mathcal{M}$  is the set of all tangent vectors to  $\mathcal{M}$ .

## Tangent vectors: submanifolds of Euclidean spaces

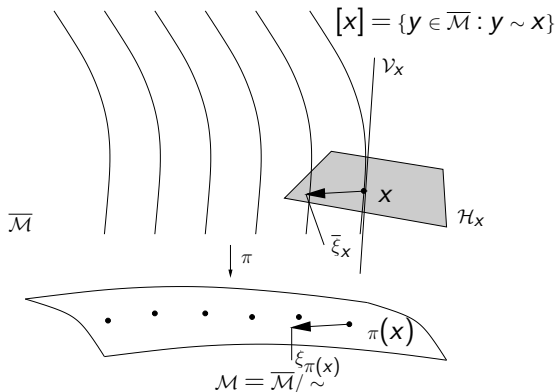


If  $\mathcal{M}$  is a submanifold of  $\mathbb{R}^n$  and  $f = \bar{f}|_{\mathcal{M}}$ , then

$$Df(x)[\dot{\gamma}(0)] = D\bar{f}(x) \left[ \frac{d\gamma}{dt}(0) \right].$$

Proof: The left-hand side is equal to  $\left. \frac{d}{dt} f(\gamma(t)) \right|_{t=0}$ . This is equal to  $\left. \frac{d}{dt} \bar{f}(\gamma(t)) \right|_{t=0}$  because  $\gamma(t) \in \mathcal{M}$  for all  $t$ . The classical chain rule yields the right-hand side.

## Tangent vectors: quotient manifolds



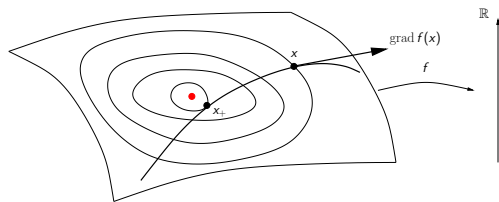
Let  $\overline{\mathcal{M}}/\sim$  be a quotient manifold. Then  $[x]$  is a submanifold of  $\overline{\mathcal{M}}$ . The tangent space  $T_x[x]$  is the *vertical space*  $\mathcal{V}_x$ . A *horizontal space* is a subspace of  $T_x \overline{\mathcal{M}}$  complementary to  $\mathcal{V}_x$ .

Let  $\xi_{\pi(x)}$  be a tangent vector to  $\overline{\mathcal{M}}/\sim$  at  $\pi(x)$ .

Theorem: In  $\mathcal{H}_x$  there is one and only one  $\bar{\xi}_x$  such that

$$D\pi(x)[\bar{\xi}_x] = \xi_{\pi(x)}.$$

## Steepest-descent: norm of tangent vectors



The steepest ascent direction is along

$$\arg \max_{\substack{\xi \in T_x \mathcal{M} \\ \|\xi\|=1}} Df(x)[\xi].$$

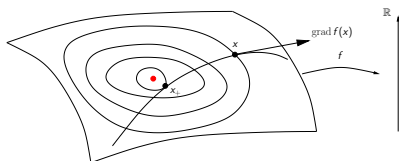
To this end, we need a norm on  $T_x \mathcal{M}$ .

For all  $x \in \mathcal{M}$ , let  $g_x$  denote an inner product in  $T_x \mathcal{M}$ , and define

$$\|\xi_x\| := \sqrt{g_x(\xi_x, \xi_x)}.$$

When  $g_x$  “smoothly” depends on  $x$ , we say that  $(\mathcal{M}, g)$  is a *Riemannian manifold*.

## Steepest-descent: gradient



There is a unique  $\text{grad } f(x)$ , called the *gradient* of  $f$  at  $x$ , such that

$$\begin{cases} \text{grad } f(x) \in T_x \mathcal{M} \\ \mathbf{g}_x(\text{grad } f(x), \xi_x) = Df(x)[\xi_x], \quad \forall \xi_x \in T_x \mathcal{M}. \end{cases}$$

We have

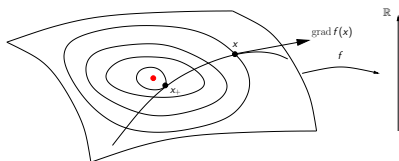
$$\frac{\text{grad } f(x)}{\|\text{grad } f(x)\|} = \arg \max_{\substack{\xi \in T_x \mathcal{M} \\ \|\xi\|=1}} Df(x)[\xi]$$

and

$$\|\text{grad } f(x)\| = Df(x) \left[ \frac{\text{grad } f(x)}{\|\text{grad } f(x)\|} \right].$$



## Steepest-descent: Riemannian submanifolds



Let  $(\overline{\mathcal{M}}, \overline{g})$  be a Riemannian manifold and  $\mathcal{M}$  be a submanifold of  $\overline{\mathcal{M}}$ . Then

$$g_x(\xi_x, \zeta_x) := \overline{g}_x(\xi_x, \eta_x), \quad \forall \xi_x, \zeta_x \in T_x \mathcal{M}$$

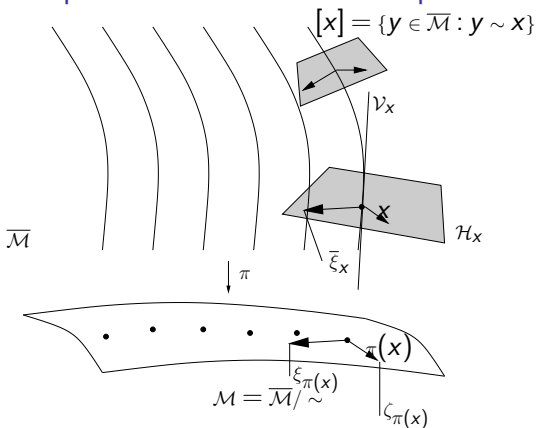
defines a Riemannian metric  $g$  on  $\mathcal{M}$ . With this Riemannian metric,  $\mathcal{M}$  is a *Riemannian submanifold* of  $\overline{\mathcal{M}}$ .

Every  $z \in T_x \overline{\mathcal{M}}$  admits a decomposition  $z = \underbrace{P_x z}_{\in T_x \mathcal{M}} + \underbrace{P_x^\perp z}_{\in T_x^\perp \mathcal{M}}$ .

If  $\overline{f} : \overline{\mathcal{M}} \rightarrow \mathbb{R}$  and  $f = \overline{f}|_{\mathcal{M}}$ , then

$$\text{grad } f(x) = P_x \text{grad } \overline{f}(x).$$

## Steepest-descent: Riemannian quotient manifolds

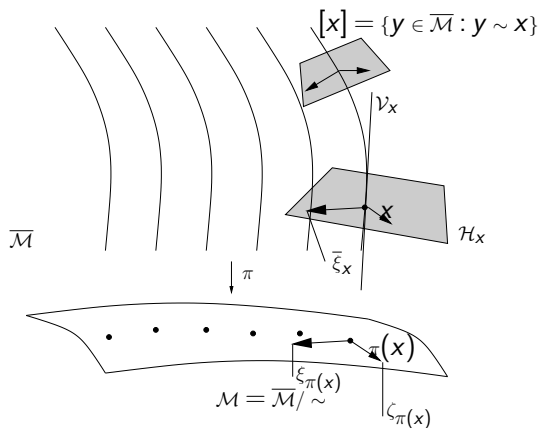


Let  $\tilde{g}$  be a Riemannian metric on  $\overline{\mathcal{M}}$ .

Suppose that, for all  $\xi_{\pi(x)}$  and  $\zeta_{\pi(x)}$  in  $T_{\pi(x)}\overline{\mathcal{M}} / \sim$ , and all  $\tilde{x} \in \pi^{-1}(\pi(x))$ , we have

$$\overline{g}_{\tilde{x}}(\overline{\xi}_{\tilde{x}}, \overline{\zeta}_{\tilde{x}}) = \overline{g}_x(\overline{\xi}_x, \overline{\zeta}_x).$$

## Steepest-descent: Riemannian quotient manifolds

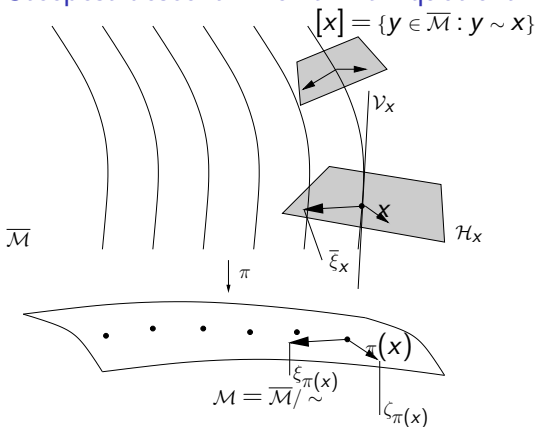


Then

$$g_{\pi(x)}(\xi_{\pi(x)}, \zeta_{\pi(x)}) := \overline{g}_x(\overline{\xi}_x, \overline{\zeta}_x).$$

defines a Riemannian metric on  $\overline{\mathcal{M}} / \sim$ . This turns  $\overline{\mathcal{M}} / \sim$  into a Riemannian quotient manifold.

## Steepest-descent: Riemannian quotient manifolds



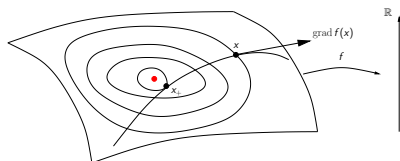
Let  $f : \overline{M} / \sim \rightarrow \mathbb{R}$ . Let  $P_x^{h, \overline{g}}$  denote the orthogonal projection onto  $\mathcal{H}_x$ .

$$\overline{\text{grad}} f_x = P_x^{h, \overline{g}} \text{grad}(f \circ \pi)(x).$$

If  $\mathcal{H}_x$  is the orthogonal complement of  $\nu_x$  in the sense of  $\overline{g}$  ( $\pi$  is a *Riemannian submersion*), then  $\text{grad}(f \circ \pi)(x)$  is already in  $\mathcal{H}_x$ , and thus

$$\overline{\text{grad}} f_x = \text{grad}(f \circ \pi)(x).$$

## Steepest-descent: choosing the search curve



It remains to choose a curve  $\gamma$  through  $x$  at  $t = 0$  such that

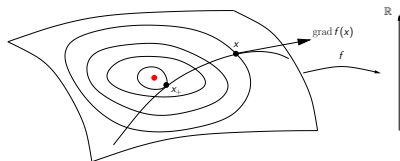
$$\dot{\gamma}(0) = -\text{grad } f(x).$$

Let  $R : T\mathcal{M} \rightarrow \mathcal{M}$  be a *retraction* on  $\mathcal{M}$ , that is

1.  $R(0_x) = x$ , where  $0_x$  denotes the origin of  $T_x\mathcal{M}$ ;
2.  $\frac{d}{dt}R(t\xi_x) = \xi_x$ .

Then choose  $\gamma : t \mapsto R(-t\text{grad } f(x))$ .

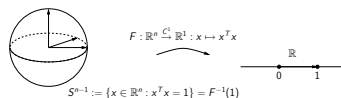
## Steepest-descent: line-search procedure



Find  $t$  such that  $f(\gamma(t))$  is “sufficiently smaller” than  $f(\gamma(0))$ . Since  $t \mapsto f(\gamma(t))$  is just a function from  $\mathbb{R}$  to  $\mathbb{R}$ , we can use the step selection techniques that are available for classical line-search methods.

For example: exact minimization, Armijo backtracking,...

## Steepest-descent: Rayleigh quotient on unit sphere



Let the manifold be the unit sphere

$$S^{n-1} = \{x \in \mathbb{R}^n : x^T x = 1\} = F^{-1}(1),$$

where  $F: \mathbb{R}^n \rightarrow \mathbb{R} : x \mapsto x^T x$ .

Let  $A = A^T \in \mathbb{R}^{n \times n}$  and let the cost function be the Rayleigh quotient

$$f: S^{n-1} \rightarrow \mathbb{R} : x \mapsto x^T A x.$$

The **tangent space** to  $S^{n-1}$  at  $x$  is

$$T_x S^{n-1} = \ker(DF(x)) = \{z \in \mathbb{R}^n : x^T z = 0\}.$$

## Derivation formulas

If  $F$  is linear, then

$$DF(x)[z] = F(z).$$

Chain rule: If  $\text{range}(F) \subseteq \text{dom}(G)$ , then

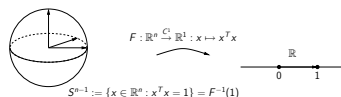
$$D(G \circ F)(x)[z] = DG(F(x))[DF(x)[z]].$$

Product rule: If the ranges of  $F$  and  $G$  are in matrix spaces of compatible dimension, then

$$D(FG)(x)[z] = DF(x)[z]G(x) + F(x)DG(x)[z].$$



## Steepest-descent: Rayleigh quotient on unit sphere



Rayleigh quotient:

$$f: S^{n-1} \rightarrow \mathbb{R} : x \mapsto x^T A x.$$

The tangent space to  $S^{n-1}$  at  $x$  is

$$T_x S^{n-1} = \ker(DF(x)) = \{z \in \mathbb{R}^n : x^T z = 0\}.$$

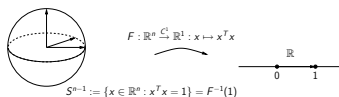
Product rule:

$$D(FG)(x)[z] = DF(x)[z]G(x) + F(x)DG(x)[z].$$

**Differential** of  $f$  at  $x \in S^{n-1}$ :

$$Df(x)[z] = x^T A z + z^T A x = 2z^T A x, \quad z \in T_x S^{n-1}.$$

## Steepest-descent: Rayleigh quotient on unit sphere



“Natural” Riemannian metric on  $S^{n-1}$ :

$$g_x(z_1, z_2) = z_1^T z_2, \quad z_1, z_2 \in T_x S^{n-1}.$$

Differential of  $f$  at  $x \in S^{n-1}$ :

$$Df(x)[z] = 2z^T Ax = 2g_x(z, Ax), \quad z \in T_x S^{n-1}.$$

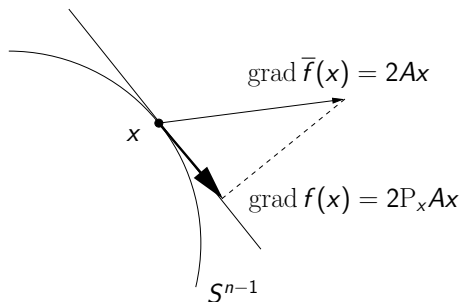
Gradient:

$$\text{grad } f(x) = 2P_x Ax = 2(I - xx^T)Ax.$$

Check:

$$\begin{cases} \text{grad } f(x) \in T_x S^{n-1} \\ Df(x)[z] = g_x(\text{grad } f(x), z), \quad \forall z \in T_x S^{n-1}. \end{cases}$$

## Steepest-descent: Rayleigh quotient on unit sphere



$$f : S^{n-1} \rightarrow \mathbb{R} : x \mapsto x^T Ax$$

$$\bar{f} : \mathbb{R}^n \rightarrow \mathbb{R} : x \mapsto x^T Ax$$

$$\text{grad } \bar{f}(x) = 2Ax$$

$$\text{grad } f(x) = 2P_x Ax = 2(I - xx^T)Ax.$$

# Newton's method on manifolds

Newton in  $\mathbb{R}^n$ 

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ .

Recall  $\text{grad } f(x) = [\partial_1 f(x) \ \cdots \ \partial_n f(x)]^T$ .

Newton's iteration:

1. Solve, for the unknown  $z \in \mathbb{R}^n$ ,

$$D(\text{grad } f)(x)[z] = -\text{grad } f(x).$$

2. Set

$$x_+ = x + z.$$

Newton in  $\mathbb{R}^n$ : how it may fail

Let  $f : \mathbb{R}_0^n \rightarrow \mathbb{R} : x \mapsto \frac{x^T A x}{x^T x}$ .

Newton's iteration:

1. Solve, for the unknown  $z \in \mathbb{R}^n$ ,

$$D(\text{grad } f)(x)[z] = -\text{grad } f(x).$$

2. Set

$$x_+ = x + z.$$

Proposition: For all  $x$  such that  $f(x)$  is not an eigenvalue of  $A$ , we have

$$x_+ = 2x.$$

## Newton: how to make it work for RQ

Let  $f : S^{n-1} \rightarrow \mathbb{R} : x \mapsto \frac{x^T A x}{x^T x}$ .

Newton's iteration:

1. Solve, for the unknown  $z \in \mathbb{R}^n \rightsquigarrow \eta_x \in T_x S^{n-1}$

$$D(\text{grad } f)(x)[z] = -\text{grad } f(x) \rightsquigarrow \boxed{?}(\text{grad } f)(x)[\eta_x] = -\text{grad } f(x)$$

2. Set

$$x_+ = x + z \rightsquigarrow x_+ = R(\eta_x)$$

## Newton's equation on an abstract manifold

Let  $\mathcal{M}$  be a manifold and let  $f : \mathcal{M} \rightarrow \mathbb{R}$  be a cost function.  
 The mapping  $x \in \mathcal{M} \mapsto \text{grad } f(x) \in T_x\mathcal{M}$  is a *vector field*.

$$D(\text{grad } f)(x)[z] = -\text{grad } f(x) \quad \rightsquigarrow \quad \boxed{?}(\text{grad } f)(x)[\eta_x] = -\text{grad } f(x)$$

The new object has to be such that

- ▶ In  $\mathbb{R}^n$ ,  $\boxed{?}$  reduces to the classical derivative
- ▶  $\boxed{?}(\text{grad } f)(x)[\eta_x]$  belongs to  $T_x\mathcal{M}$
- ▶  $\boxed{?}$  has the same linearity properties and multiplication rule as the classical derivative.

Differential geometry offers a concept that matches these conditions: the concept of an *affine connection*.



## Newton: affine connections

Let  $\mathfrak{X}(\mathcal{M})$  denote the set of smooth vector fields on  $\mathcal{M}$  and  $\mathfrak{F}(\mathcal{M})$  the set of real-valued functions on  $\mathcal{M}$ .

An *affine connection*  $\nabla$  on a manifold  $\mathcal{M}$  is a mapping

$$\nabla : \mathfrak{X}(\mathcal{M}) \times \mathfrak{X}(\mathcal{M}) \rightarrow \mathfrak{X}(\mathcal{M}),$$

which is denoted by  $(\eta, \xi) \xrightarrow{\nabla} \nabla_{\eta}\xi$  and satisfies the following properties:

- i)  $\mathfrak{F}(\mathcal{M})$ -linearity in  $\eta$ :  $\nabla_{f\eta+g\chi}\xi = f\nabla_{\eta}\xi + g\nabla_{\chi}\xi$ ,
- ii)  $\mathbb{R}$ -linearity in  $\xi$ :  $\nabla_{\eta}(a\xi + b\zeta) = a\nabla_{\eta}\xi + b\nabla_{\eta}\zeta$ ,
- iii) Product rule (Leibniz' law):  $\nabla_{\eta}(f\xi) = (\eta f)\xi + f\nabla_{\eta}\xi$ ,

in which  $\eta, \chi, \xi, \zeta \in \mathfrak{X}(\mathcal{M})$ ,  $f, g \in \mathfrak{F}(\mathcal{M})$ , and  $a, b \in \mathbb{R}$ .

## Newton's method on abstract manifolds

Cost function:  $f : \mathbb{R}^n \rightarrow \mathbb{R} \rightsquigarrow f : \mathcal{M} \rightarrow \mathbb{R}$ .

Newton's iteration:

1. Solve, for the unknown  $z \in \mathbb{R}^n \rightsquigarrow \eta_x \in T_x \mathcal{M}$

$$D(\text{grad } f)(x)[z] = -\text{grad } f(x) \rightsquigarrow \nabla(\text{grad } f)(x)[\eta_x] = -\text{grad } f(x)$$

2. Set

$$x_+ = x + z \rightsquigarrow x_+ = R(\eta_x)$$

In the algorithm above,  $\nabla$  is an affine connection on  $\mathcal{M}$  and  $R$  is a retraction on  $\mathcal{M}$ .

Newton's method on  $S^{n-1}$ 

If  $\mathcal{M}$  is a Riemannian submanifold of  $\mathbb{R}^n$ , then  $\nabla$  defined by

$$\nabla_{\eta_x} \xi = P_x D\xi(x)[\eta_x], \quad \eta_x \in T_x \mathcal{M}, \quad \xi \in \mathfrak{X}(\mathcal{M})$$

is a particular affine connection, called *Riemannian connection*.

For the unit sphere  $S^{n-1}$ , this yields

$$\nabla_{\eta_x} \xi = (I - xx^T)D\xi(x)[\eta_x], \quad x^T \eta_x = 0.$$

# Newton's method for Rayleigh quotient on $S^{n-1}$

$$\text{Let } f : \begin{cases} \mathbb{R}^n \\ \mathcal{M} \\ S^{n-1} \end{cases} \rightarrow \mathbb{R} : x \mapsto \begin{cases} f(x) \\ f(x) \\ \frac{x^T A x}{x^T x} \end{cases} .$$

Newton's iteration:

1. Solve, for the unknown  $z \in \mathbb{R}^n \rightsquigarrow \eta_x \in T_x \mathcal{M} \rightsquigarrow x^T \eta_x = 0$

$$\begin{aligned} D(\text{grad } f)(x)[z] &= -\text{grad } f(x) \\ &\rightsquigarrow \nabla(\text{grad } f)(x)[\eta_x] = -\text{grad } f(x) \\ &\rightsquigarrow (I - xx^T)(A - f(x)I)\eta_x = -(I - xx^T)Ax \end{aligned}$$

2. Set

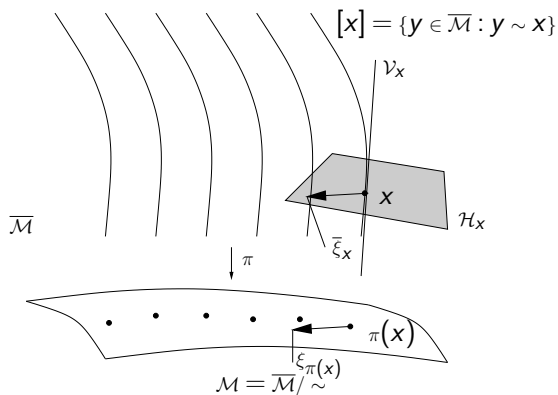
$$x_+ = x + z \rightsquigarrow x_+ = R(\eta_x) \rightsquigarrow x_+ = \frac{x + \eta_x}{\|x + \eta_x\|}$$

Newton for RQ on  $S^{n-1}$ : a closer look

$$\begin{aligned}(I - xx^T)(A - f(x)I)\eta_x &= -(I - xx^T)Ax \\ \Rightarrow (I - xx^T)(A - f(x)I)(x + \eta_x) &= 0 \\ \Rightarrow (A - f(x)I)(x + \eta_x) &= \alpha x\end{aligned}$$

Therefore,  $x_+$  is collinear with  $(A - f(x)I)^{-1}x$ , which is the vector computed by the Rayleigh quotient iteration.

## Newton method on quotient manifolds

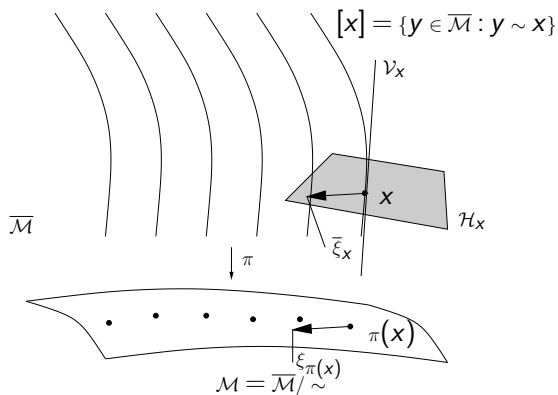


Affine connection: choose  $\nabla$  defined by

$$\overline{\nabla}_{\eta} \bar{\xi}_x = P_x^h \overline{\nabla}_{\bar{\eta}_x} \bar{\xi},$$

provided that this really defines a horizontal lift. This requires special choices of  $\bar{\nabla}$ .

## Newton method on quotient manifolds



If  $\pi : \overline{\mathcal{M}} \rightarrow \overline{\mathcal{M}} / \sim$  is a Riemannian submersion, then the Riemannian connection on  $\overline{\mathcal{M}} / \sim$  is given by

$$\overline{\nabla}_{\eta} \bar{\xi}_x = P_x^h \overline{\nabla}_{\bar{\eta}_x} \bar{\xi},$$

where  $\overline{\nabla}$  denotes the Riemannian connection on  $\overline{\mathcal{M}}$ .

## A detailed exercise

Newton's method for the Rayleigh  
quotient on the Grassmann  
manifold



## Manifold: Grassmann

The manifold is the Grassmann manifold of  $p$ -planes in  $\mathbb{R}^n$ :

$$\text{Grass}(p, n) \simeq \text{ST}(p, n)/\text{GL}_p.$$

The one-to-one correspondence is

$$\text{Grass}(p, n) \ni \mathcal{Y} \leftrightarrow Y \text{GL}_p \in \text{ST}(p, n)/\text{GL}_p$$

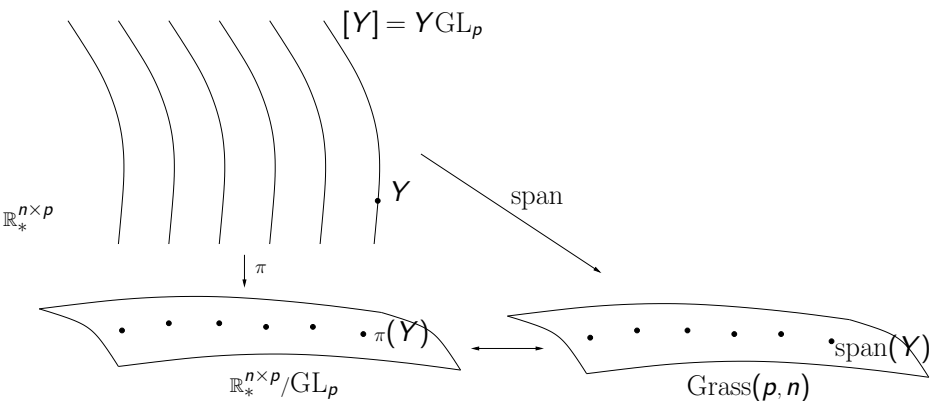
such that  $\mathcal{Y}$  is the column space of  $Y$ .

The quotient map

$$\pi : \text{ST}(p, n) \rightarrow \text{Grass}(p, n)$$

is the “column space” or “span” operation.

## Grassmann and its quotient representation



## Total space: the noncompact Stiefel manifold

The total space of the quotient is

$$\text{ST}(\rho, n) = \{Y \in \mathbb{R}^{n \times \rho} : \text{rank}(Y) = \rho\}.$$

This is an open submanifold of the Euclidean space  $\mathbb{R}^{n \times \rho}$ .

Tangent spaces:  $T_Y \text{ST}(\rho, n) \simeq \mathbb{R}^{n \times \rho}$ .

## Riemannian metric on the total space

Define a Riemannian metric  $\bar{g}$  on  $ST(p, n)$  by

$$\bar{g}_Y(Z_1, Z_2) = \text{trace} \left( (Y^T Y)^{-1} Z_1^T Z_2 \right).$$

This is not the canonical Riemannian metric, but it will allow us to turn the quotient map  $\pi : ST(p, n) \rightarrow \text{Grass}(p, n)$  into a Riemannian submersion.

## Vertical and horizontal spaces

The vertical spaces are the tangent spaces to the equivalence classes:

$$\mathcal{V}_Y := T_Y(YGL_p) = Y T_Y GL_p = Y \mathbb{R}^{p \times p}.$$

Choice of horizontal space:

$$\begin{aligned} \mathcal{H}_Y &:= (\mathcal{V}_Y)^\perp \\ &= \{Z \in T_Y ST(p, n) : \bar{g}_Y(Z, V) = 0, \forall V \in \mathcal{V}_Y\} \\ &= \{Z \in \mathbb{R}^{n \times p} : Y^T Z = 0\}. \end{aligned}$$

Horizontal projection:

$$P_Y^h = (I - Y(Y^T Y)^{-1} Y^T).$$

## Compatibility equation for horizontal lifts

Given  $\xi \in T_\pi(Y)\text{Grass}(p, n)$ , we have

$$\bar{\xi}_{YM} = \bar{\xi}_Y M.$$

To see this, observe that  $\bar{\xi}_Y M$  is in  $\mathcal{H}_{YM}$ ; moreover, since  $YM + t\bar{\xi}_Y M$  and  $Y + t\bar{\xi}_Y$  have the same column space for all  $t$ , one has

$$D\pi(YM)[\bar{\xi}_Y M] = D\pi(Y)[\bar{\xi}_Y] = \xi_{\pi(Y)}.$$

Thus  $\bar{\xi}_Y M$  satisfies the conditions to be  $\bar{\xi}_{YM}$ .

## Riemannian metric on the quotient

On  $\text{Grass}(p, n) \simeq \text{ST}(p, n)/\text{GL}_p$ , define the Riemannian metric  $g$  by

$$g_{\pi(Y)}(\xi_{\pi(Y)}, \zeta_{\pi(Y)}) = \bar{g}_Y(\bar{\xi}_Y, \bar{\zeta}_Y).$$

This is well defined, because for all  $\tilde{Y} \in \pi^{-1}(\pi(Y)) = Y\text{GL}_p$ , we have  $\tilde{Y} = YM$  for some invertible  $M$ , and

$$\bar{g}_{YM}(\bar{\xi}_{YM}, \bar{\zeta}_{YM}) = \bar{g}_Y(\bar{\xi}_Y, \bar{\zeta}_Y).$$

This definition of  $g$  turns

$$\pi : (\text{ST}(p, n), \bar{g}) \rightarrow (\text{Grass}(p, n), g)$$

into a Riemannian submersion.

## Cost function: Rayleigh quotient

Consider the cost function

$$f : \text{Grass}(p, n) \rightarrow \mathbb{R} : \text{span}(Y) \mapsto \text{trace} \left( (Y^T Y)^{-1} Y^T A Y \right).$$

This is the *projection* of

$$\bar{f} : \text{ST}(p, n) \rightarrow \mathbb{R} : Y \mapsto \text{trace} \left( (Y^T Y)^{-1} Y^T A Y \right).$$

That is,  $\bar{f} = f \circ \pi$ .



## Gradient of the cost function

For all  $Z \in \mathbb{R}^{n \times p}$ ,

$$D\bar{f}(Y)[Z] = 2 \operatorname{trace} \left( (Y^T Y)^{-1} Z^T (AY - Y(Y^T Y)^{-1} Y^T AY) \right).$$

Hence

$$\operatorname{grad} \bar{f}(Y) = 2 \left( AY - Y(Y^T Y)^{-1} Y^T AY \right),$$

and

$$\overline{\operatorname{grad} f}_Y = 2 \left( AY - Y(Y^T Y)^{-1} Y^T AY \right).$$

## Riemannian connection

The quotient map is a Riemannian submersion. Therefore

$$\overline{\nabla_{\eta} \xi} = P_Y^h (\overline{\nabla_{\bar{\eta}_Y} \bar{\xi}})$$

It turns out that

$$\overline{\nabla_{\eta} \xi} = P_Y^h (D\bar{\xi}(Y) [\bar{\eta}_Y]).$$

(This is because the Riemannian metric  $\bar{g}$  is “horizontally invariant”.)

For the Rayleigh quotient  $f$ , this yields

$$\begin{aligned} \overline{\nabla_{\eta} \text{grad } f} &= P_Y^h (D\overline{\text{grad } f}(Y) [\bar{\eta}_Y]) \\ &= 2 P_Y^h \left( A\bar{\eta}_Y - \bar{\eta}_Y (Y^T Y)^{-1} Y^T A Y \right). \end{aligned}$$

## Newton's equation

Newton's equation at  $\pi(Y)$  is

$$\nabla_{\eta_{\pi(Y)}} \text{grad } f = -\text{grad } f(\pi(Y))$$

for the unknown  $\eta_{\pi(Y)} \in T_{\pi(Y)} \text{Grass}(p, n)$ .

To turn this equation into a matrix equation, we take its horizontal lift.

This yields

$$P_Y^h \left( A\bar{\eta}_Y - \bar{\eta}_Y(Y^T Y)^{-1} Y^T A Y \right) = -P_Y^h A Y, \quad \bar{\eta}_Y \in \mathcal{H}_Y,$$

whose solution  $\bar{\eta}_Y$  in the horizontal space  $\mathcal{H}_Y$  is the horizontal lift of the solution  $\eta$  of the Newton equation.

## Retraction

Newton's method sends  $\pi(Y)$  to  $\mathcal{Y}_+$  according to

$$\begin{aligned}\nabla_{\eta_{\pi(Y)}} \text{grad } f &= -\text{grad } f(\pi(Y)) \\ \mathcal{Y}_+ &= R_{\pi(Y)}(\eta_{\pi(Y)}).\end{aligned}$$

It remains to pick the retraction  $R$ .

Choice:  $R$  defined by

$$R_{\pi(Y)}\xi_{\pi(Y)} = \pi(Y + \bar{\xi}_Y).$$

(This is a well-defined retraction.)

## Newton's iteration for RQ on Grassmann

**Require:** Symmetric matrix  $A$ .

**Input:** Initial iterate  $Y_0 \in \text{ST}(p, n)$ .

**Output:** Sequence of iterates  $\{Y_k\}$  in  $\text{ST}(p, n)$ .

- 1: **for**  $k = 0, 1, 2, \dots$  **do**
- 2:     Solve the linear system

$$\begin{cases} P_{Y_k}^h (AZ_k - Z_k(Y_k^T Y_k)^{-1} Y_k^T A Y_k) = -P_{Y_k}^h (A Y_k) \\ Y_k^T Z_k = 0 \end{cases}$$

for the unknown  $Z_k$ , where  $P_Y^h$  is the orthogonal projector onto  $\mathcal{H}_Y$ . (The condition  $Y_k^T Z_k = 0$  expresses that  $Z_k$  belongs to the horizontal space  $\mathcal{H}_{Y_k}$ .)

- 3:     Set

$$Y_{k+1} = (Y_k + Z_k)N_k$$

where  $N_k$  is a nonsingular  $p \times p$  matrix chosen for normalization purposes.

- 4: **end for**

A new tool for Optimization On  
Manifolds:  
Vector Transport

## Filling a gap

	Purely Riemannian way	Pragmatic way
Update	Search along the geodesic tangent to the search direction	Search along <b>any</b> curve tangent to the search direction (described by a <i>retraction</i> )
Displacement of tgt vectors	Parallel translation induced by $\frac{g}{\nabla}$	??

## Where do we use parallel translation?

In CG. Quoting (approximately) Smith (1994):

1. Select  $x_0 \in \mathcal{M}$ , compute  $H_0 = -\text{grad } f(x_0)$ , and set  $k = 0$
2. Compute  $t_k$  such that  $f(\text{Exp}_{x_k}(t_k H_k)) \leq f(\text{Exp}_{x_k}(t H_k))$  for all  $t \geq 0$ .
3. Set  $x_{k+1} = \text{Exp}_{x_k}(t_k H_k)$ .
4. Set  $H_{k+1} = -\text{grad } f(x_{k+1}) + \beta_k \tau H_k$ , where  $\tau$  is the **parallel translation** along the geodesic from  $x_k$  to  $x_{k+1}$ .



## Where do we use parallel translation?

In **BFGS**. Quoting (approximately) Gabay (1982):

$x_{k+1} = \text{Exp}_{x_k}(t_k \xi_k)$  (update along geodesic)

$\text{grad } f(x_{k+1}) - \tau_0^{t_k} \text{grad } f(x_k) = B_{k+1} \tau_0^{t_k}(t_k \xi_k)$  (requirement on approximate Jacobian  $B$ )

This leads to the a *generalized BFGS update formula* involving parallel translation.

## Where else could we use parallel translation?

In **finite-difference quasi-Newton**.

Let  $\xi$  be a vector field on a Riemannian manifold  $\mathcal{M}$ . Exact Jacobian of  $\xi$  at  $x \in \mathcal{M}$ :  $J_\xi(x)[\eta] = \nabla_\eta \xi$ .

Finite difference approximation to  $J_\xi$ : choose a basis  $(E_1, \dots, E_d)$  of  $T_x \mathcal{M}$  and define  $\tilde{J}(x)$  as the linear operator that satisfies

$$\tilde{J}(x)[E_i] = \frac{\tau_h^0 \xi_{\text{Exp}_x(hE_i)} - \xi_x}{h}.$$

## Filling a gap

	Purely Riemannian way	Pragmatic way
Update	Search along the geodesic tangent to the search direction	Search along <b>any</b> pres-curve tangent to the search direction
Displacement of tgt vectors	Parallel translation induced by $\frac{g}{\nabla}$	??

## Parallel translation can be tough

Edelman et al (1998): We are unaware of any closed form expression for the parallel translation on the Stiefel manifold (defined with respect to the Riemannian connection induced by the embedding in  $\mathbb{R}^{n \times p}$ ).

Parallel transport along geodesics on Grassmannians:

$$\overline{\xi(t)}_{Y(t)} = -Y_0 V \sin(\Sigma t) U^T \overline{\xi(0)}_{Y_0} + U \cos(\Sigma t) U^T \overline{\xi(0)}_{Y_0} + (I - UU^T) \overline{\xi(0)}_{Y_0}.$$

where  $\overline{\dot{Y}(0)}_{Y_0} = U \Sigma V^T$  is a thin SVD.

## Alternatives found in the literature

Edelman et al (1998): “extrinsic” CG algorithm. “Tangency of the search direction at the new point is imposed via the projection  $I - YY^T$ ” (instead of via parallel translation).

Brace & Manton (2006), *An improved BFGS-on-manifold algorithm for computing weighted low rank approximation*. “The second change is that parallel translation is not defined with respect to the Levi-Civita connection, but rather is all but ignored.”

## Filling a gap

	Purely Riemannian way	Pragmatic way
Update	Search along the geodesic tangent to the search direction	Search along <b>any</b> curve tangent to the search direction (described by a <i>retraction</i> )
Displacement of tgt vectors	Parallel translation induced by $\frac{g}{\nabla}$	??

## Filling a gap: Vector Transport

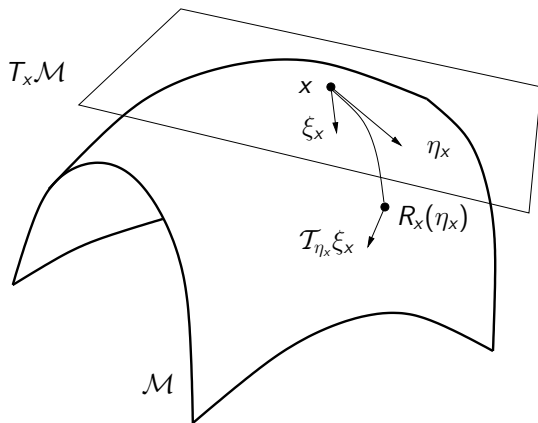
	Purely Riemannian way	Pragmatic way
Update	Search along the geodesic tangent to the search direction	Search along <b>any</b> curve tangent to the search direction (described by a <i>retraction</i> )
Displacement of tgt vectors	Parallel translation induced by $\frac{g}{\nabla}$	<b>Vector Transport</b>

## Still to come

- ▶ Vector transport in one picture
- ▶ Formal definition
- ▶ Particular vector transports
- ▶ Applications: finite-difference Newton, BFGS, CG.



## The concept of vector transport



## Retraction

A *retraction* on a manifold  $\mathcal{M}$  is a smooth mapping

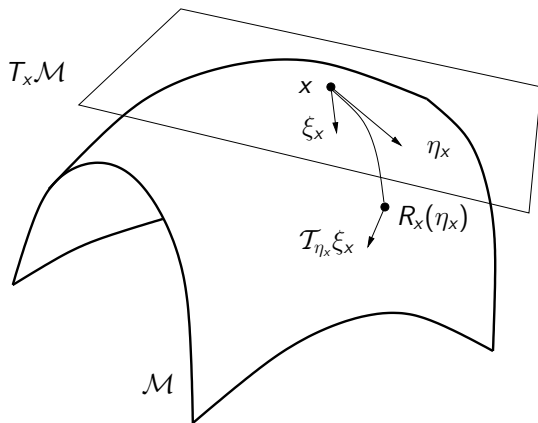
$$R : T\mathcal{M} \rightarrow \mathcal{M}$$

such that

1.  $R(0_x) = x$  for all  $x \in \mathcal{M}$ , where  $0_x$  denotes the origin of  $T_x\mathcal{M}$ ;
2.  $\frac{d}{dt}R(t\xi_x)|_{t=0} = \xi_x$  for all  $\xi_x \in T_x\mathcal{M}$ .

Consequently, the curve  $t \mapsto R(t\xi_x)$  is a curve on  $\mathcal{M}$  tangent to  $\xi_x$ .

## The concept of vector transport – Whitney sum



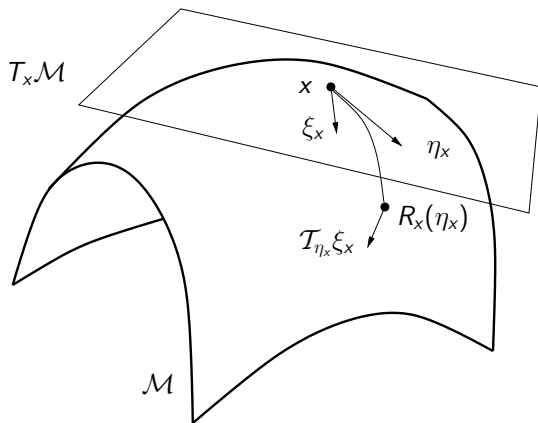
## Whitney sum

Let  $T\mathcal{M} \oplus T\mathcal{M}$  denote the set

$$T\mathcal{M} \oplus T\mathcal{M} = \{(\eta_x, \xi_x) : \eta_x, \xi_x \in T_x\mathcal{M}, x \in \mathcal{M}\}.$$

This set admits a natural manifold structure.

## The concept of vector transport – definition



## Vector transport: definition

A *vector transport* on a manifold  $\mathcal{M}$  on top of a retraction  $R$  is a smooth map

$$T\mathcal{M} \oplus T\mathcal{M} \rightarrow T\mathcal{M} : (\eta_x, \xi_x) \mapsto \mathcal{T}_{\eta_x}(\xi_x) \in T\mathcal{M}$$

satisfying the following properties for all  $x \in \mathcal{M}$ :

1. (Underlying retraction)  $\mathcal{T}_{\eta_x}\xi_x$  belongs to  $T_{R_x(\eta_x)}\mathcal{M}$ .
2. (Consistency)  $\mathcal{T}_{0_x}\xi_x = \xi_x$  for all  $\xi_x \in T_x\mathcal{M}$ ;
3. (Linearity)  $\mathcal{T}_{\eta_x}(a\xi_x + b\zeta_x) = a\mathcal{T}_{\eta_x}(\xi_x) + b\mathcal{T}_{\eta_x}(\zeta_x)$ .

## Inverse vector transport

When it exists,  $(\mathcal{T}_{\eta_x})^{-1}(\xi_{R_x(\eta_x)})$  belongs to  $T_x\mathcal{M}$ . If  $\eta$  and  $\xi$  are two vector fields on  $\mathcal{M}$ , then  $(\mathcal{T}_\eta)^{-1}\xi$  is naturally defined as the vector field satisfying

$$((\mathcal{T}_\eta)^{-1}\xi)_x = (\mathcal{T}_{\eta_x})^{-1}(\xi_{R_x(\eta_x)}).$$

## Still to come

- ▶ Vector transport in one picture
- ▶ Formal definition
- ▶ Particular vector transports
- ▶ Applications: finite-difference Newton, BFGS, CG.



## Parallel translation is a vector transport

### Proposition

If  $\nabla$  is an affine connection and  $R$  is a retraction on a manifold  $\mathcal{M}$ , then

$$\mathcal{T}_{\eta_x}(\xi_x) := P_{\gamma}^{1 \leftarrow 0} \xi_x \quad (1)$$

is a vector transport with associated retraction  $R$ , where  $P_{\gamma}$  denotes the parallel translation induced by  $\nabla$  along the curve  $t \mapsto \gamma(t) = R_x(t\eta_x)$ .

## Vector transport on Riemannian submanifolds

If  $\mathcal{M}$  is an embedded submanifold of a Euclidean space  $\mathcal{E}$  and  $\mathcal{M}$  is endowed with a retraction  $R$ , then we can rely on the natural inclusion  $T_y\mathcal{M} \subset \mathcal{E}$  for all  $y \in \mathcal{N}$  to simply define the vector transport by

$$\mathcal{T}_{\eta_x}\xi_x := P_{R_x(\eta_x)}\xi_x, \quad (2)$$

where  $P_x$  denotes the orthogonal projector onto  $T_x\mathcal{N}$ .

## Still to come

- ▶ Vector transport in one picture
- ▶ Formal definition
- ▶ Particular vector transports
- ▶ Applications: finite-difference Newton, BFGS, CG.

## Vector transport in finite differences

Let  $\mathcal{M}$  be a manifold endowed with a vector transport  $\mathcal{T}$  on top of a retraction  $R$ . Let  $x \in \mathcal{M}$  and let  $(E_1, \dots, E_d)$  be a basis of  $T_x\mathcal{M}$ . Given a smooth vector field  $\xi$  and a real constant  $h > 0$ , let

$\tilde{J}_\xi(x) : T_x\mathcal{M} \rightarrow T_x\mathcal{M}$  be the linear operator that satisfies, for  $i = 1, \dots, d$ ,

$$\tilde{J}_\xi(x)[E_i] = \frac{(\mathcal{T}_{hE_i})^{-1}\xi_{R(hE_i)} - \xi_x}{h}. \quad (3)$$

### Lemma (finite differences)

Let  $x_*$  be a nondegenerate zero of  $\xi$ . Then there is  $c > 0$  such that, for all  $x$  sufficiently close to  $x_*$  and all  $h$  sufficiently small, it holds that

$$\|\tilde{J}_\xi(x)[E_i] - J(x)[E_i]\| \leq c(h + \|\xi_x\|). \quad (4)$$

## Convergence of Newton's method with finite differences

### Proposition

*Consider the geometric Newton method where the exact Jacobian  $J(x_k)$  is replaced by the operator  $\tilde{J}_\xi(x_k)$  with  $h := h_k$ . If*

$$\lim_{k \rightarrow \infty} h_k = 0,$$

*then the convergence to nondegenerate zeros of  $\xi$  is superlinear. If, moreover, there exists some constant  $c$  such that*

$$h_k \leq c \|\xi_{x_k}\|$$

*for all  $k$ , then the convergence is (at least) quadratic.*

## Vector transport in BFGS

With the notation

$$s_k := \mathcal{T}_{\eta_k} \eta_k \in T_{x_{k+1}} \mathcal{M},$$

$$y_k := \text{grad } f(x_{k+1}) - \mathcal{T}_{\eta_k}(\text{grad } f(x_k)) \in T_{x_{k+1}} \mathcal{M},$$

we define the operator  $A_{k+1} : T_{x_{k+1}} \mathcal{M} \mapsto T_{x_{k+1}} \mathcal{M}$  by

$$A_{k+1} \eta = \tilde{A}_k \eta - \frac{\langle s_k, \tilde{A}_k \eta \rangle}{\langle s_k, \tilde{A}_k s_k \rangle} \tilde{A}_k s_k + \frac{\langle y_k, \eta \rangle}{\langle y_k, s_k \rangle} y_k \quad \text{for all } \eta \in T_{x_{k+1}} \mathcal{M},$$

with

$$\tilde{A}_k = \mathcal{T}_{\eta_k} \circ A_k \circ (\mathcal{T}_{\eta_k})^{-1}.$$

## Vector transport in CG

Compute a step size  $\alpha_k$  and set

$$x_{k+1} = R_{x_k}(\alpha_k \eta_k). \quad (5)$$

Compute  $\beta_{k+1}$  and set

$$\eta_{k+1} = -\text{grad } f(x_{k+1}) + \beta_{k+1} \mathcal{T}_{\alpha_k \eta_k}(\eta_k). \quad (6)$$

## Filling a gap: Vector Transport

	Purely Riemannian way	Pragmatic way
Update	Search along the geodesic tangent to the search direction	Search along <b>any</b> curve tangent to the search direction (described by a <i>retraction</i> )
Displacement of tgt vectors	Parallel translation induced by $\frac{g}{\nabla}$	<b>Vector Transport</b>



## Ongoing work

- ▶ Use vector transport wherever we can.
- ▶ Extend convergence analyses.
- ▶ Develop recipes for building efficient vector transports.

# Trust-region methods on Riemannian manifolds

## Motivating application: Mechanical vibrations

Mass matrix  $M$ , stiffness matrix  $K$ .

Equation of vibrations (for undamped discretized linear structures):

$$Kx = \omega^2 Mx$$

were

- ▶  $\omega$  is an angular frequency of vibration
- ▶  $x$  is the corresponding mode of vibration

Task: find lowest modes of vibration.

## Generalized eigenvalue problem

Given  $n \times n$  matrices  $A = A^T$  and  $B = B^T \succ 0$ , there exist  $v_1, \dots, v_n$  in  $\mathbb{R}^n$  and  $\lambda_1 \leq \dots \leq \lambda_n$  in  $\mathbb{R}$  such that

$$Av_i = \lambda_i Bv_i$$

$$v_i^T Bv_j = \delta_{ij}.$$

Task: find  $\lambda_1, \dots, \lambda_p$  and  $v_1, \dots, v_p$ .

We assume throughout that  $\lambda_p < \lambda_{p+1}$ .

Case  $p = 1$ : optimization in  $\mathbb{R}^n$ 

$$Av_i = \lambda_i Bv_i$$

Consider the Rayleigh quotient

$$\tilde{f} : \mathbb{R}_*^n \rightarrow \mathbb{R} : f(y) = \frac{y^T Ay}{y^T By}$$

Invariance:  $\tilde{f}(\alpha y) = \tilde{f}(y)$ .

Stationary points of  $\tilde{f}$ :  $\alpha v_i$ , for all  $\alpha \neq 0$ .

Minimizers of  $\tilde{f}$ :  $\alpha v_1$ , for all  $\alpha \neq 0$ .

Difficulty: the minimizers are not isolated.

Remedy: optimization on manifold.

Case  $p = 1$ : optimization on ellipsoid

$$\tilde{f} : \mathbb{R}_*^n \rightarrow \mathbb{R} : f(y) = \frac{y^T A y}{y^T B y}$$

Invariance:  $\tilde{f}(\alpha y) = \tilde{f}(y)$ .

Remedy 1:

- ▶  $\mathcal{M} := \{y \in \mathbb{R}^n : y^T B y = 1\}$ , *submanifold* of  $\mathbb{R}^n$ .
- ▶  $f : \mathcal{M} \rightarrow \mathbb{R} : f(y) = y^T A y$ .

Stationary points of  $f$ :  $\pm v_1, \dots, \pm v_n$ .

Minimizers of  $f$ :  $\pm v_1$ .

Case  $p = 1$ : optimization on projective space

$$\tilde{f} : \mathbb{R}_*^n \rightarrow \mathbb{R} : f(y) = \frac{y^T A y}{y^T B y}$$

Invariance:  $\tilde{f}(\alpha y) = \tilde{f}(y)$ .

Remedy 2:

- ▶  $[y] := y\mathbb{R} := \{y\alpha : \alpha \in \mathbb{R}\}$
- ▶  $\mathcal{M} := \mathbb{R}_*^n / \mathbb{R} = \{[y]\}$
- ▶  $f : \mathcal{M} \rightarrow \mathbb{R} : f([y]) := \tilde{f}(y)$

Stationary points of  $f$ :  $[v_1], \dots, [v_n]$ .

Minimizer of  $f$ :  $[v_1]$ .

Case  $p \geq 1$ : optimization on the Grassmann manifold

$$\tilde{f} : \mathbb{R}_*^{n \times p} \rightarrow \mathbb{R} : \tilde{f}(Y) = \text{trace} \left( (Y^T B Y)^{-1} Y^T A Y \right)$$

Invariance:  $\tilde{f}(YR) = \tilde{f}(Y)$ .

Define:

- ▶  $[Y] := \{YR : R \in \mathbb{R}_*^{p \times p}\}, \quad Y \in \mathbb{R}_*^{n \times p}$
- ▶  $\mathcal{M} := \text{Grass}(p, n) := \{[Y]\}$
- ▶  $f : \mathcal{M} \rightarrow \mathbb{R} : f([Y]) := \tilde{f}(Y)$

Stationary points of  $f$ :  $\text{span}\{v_{i_1}, \dots, v_{i_p}\}$ .

Minimizer of  $f$ :  $[Y] = \text{span}\{v_1, \dots, v_p\}$ .



## Optimization on Manifolds

- ▶ Luenberger [Lue73], Gabay [Gab82]: optimization on submanifolds of  $\mathbb{R}^n$ .
- ▶ Smith [Smi93, Smi94] and Udriște [Udr94]: optimization on general Riemannian manifolds (steepest descent, Newton, CG).
- ▶ ...
- ▶ PAA, Baker and Gallivan [ABG07]: trust-region methods on Riemannian manifolds.
- ▶ PAA, Mahony, Sepulchre [AMS08]: *Optimization Algorithms on Matrix Manifolds*, textbook.

## The Problem : Leftmost Eigenpairs of Matrix Pencil

Given  $n \times n$  matrix pencil  $(A, B)$ ,  $A = A^T$ ,  $B = B^T \succ 0$  with (unknown) eigen-decomposition

$$A [v_1 | \dots | v_n] = B [v_1 | \dots | v_n] \text{diag}(\lambda_1, \dots, \lambda_n)$$

$$[v_1 | \dots | v_n]^T B [v_1 | \dots | v_n] = I, \quad \lambda_1 < \lambda_2 \leq \dots \leq \lambda_n.$$

The problem is to **compute the minor eigenvector  $\pm v_1$** .

## The ideal algorithm

Given  $(A, B)$ ,  $A = A^T$ ,  $B = B^T \succ 0$  with (unknown) eigenvalues  $0 < \lambda_1 \leq \dots \leq \lambda_n$  and associated eigenvectors  $v_1, \dots, v_n$ .

1. **Global convergence:**

- ▶ Convergence to some eigenvector for **all** initial conditions.
- ▶ **Stable** convergence to the “leftmost” eigenvector  $\pm v_1$  **only**.

2. **Superlinear** (cubic) local convergence to  $\pm v_1$ .

3. **“Matrix-free”** (no factorization of  $A, B$ )  
but possible use of **preconditioner**.

4. **Minimal storage** space required.

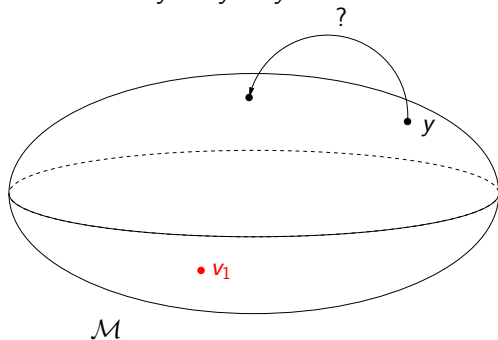
# Strategy

- ▶ Rewrite computation of leftmost eigenpair as an **optimization problem (on a manifold)**.
- ▶ Use a **model-trust-region** scheme to solve the problem.  
     $\rightsquigarrow$  **Global convergence**.
- ▶ Take the **exact quadratic model** (at least, close to the solution).  
     $\rightsquigarrow$  **Superlinear convergence**.
- ▶ Solve the trust-region subproblems using the **(Steihaug-Toint) truncated CG (tCG)** algorithm.  
     $\rightsquigarrow$  **“Matrix-free”**, preconditioned iteration.  
     $\rightsquigarrow$  **Minimal storage** of iteration vectors.

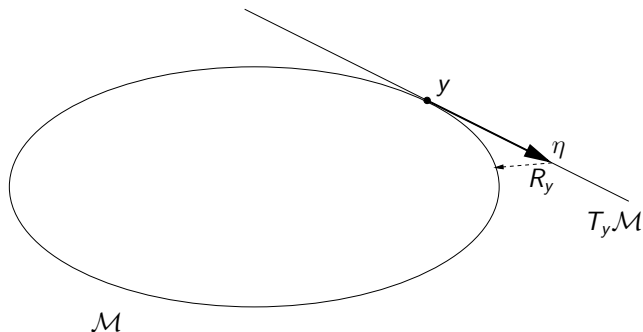
## Iteration on the manifold

Manifold: ellipsoid  $\mathcal{M} = \{y \in \mathbb{R}^n : y^T B y = 1\}$ .

Cost function:  $f : \mathcal{M} \rightarrow \mathbb{R} : y \mapsto y^T A y$



## Tangent space and retraction (2D picture)



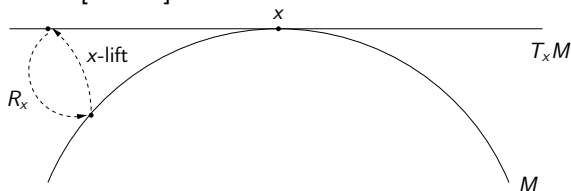
Tangent space:  $T_y \mathcal{M} := \{\eta \in \mathbb{R}^n : y^T B \eta = 0\}$ .

Retraction:  $R_y \eta := (y + \eta) / \|y + \eta\|_B$ .

Lifted cost function:  $\hat{f}_y(\eta) := f(R_y \eta) = \frac{(y+\eta)^T A (y+\eta)}{(y+\eta)^T B (y+\eta)}$ .

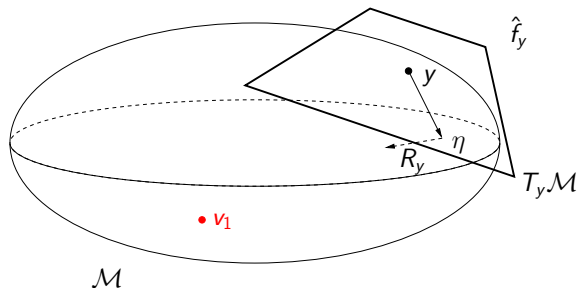
## Concept of retraction

Introduced by Shub [Shu86].



1.  $R_x$  is defined and one-to-one in a neighbourhood of  $0_x$  in  $T_x M$ .
2.  $R_x(0_x) = x$ .
3.  $DR_x(0_x) = \text{id}_{T_x M}$ , the identity mapping on  $T_x M$ , with the canonical identification  $T_{0_x} T_x M \simeq T_x M$ .

## Tangent space and retraction



Tangent space:  $T_y \mathcal{M} := \{\eta \in \mathbb{R}^n : y^T B \eta = 0\}$ .

Retraction:  $R_y \eta := (y + \eta) / \|y + \eta\|_B$ .

Lifted cost function:  $\hat{f}_y(\eta) := f(R_y \eta) = \frac{(y + \eta)^T A (y + \eta)}{(y + \eta)^T B (y + \eta)}$ .



## Quadratic model

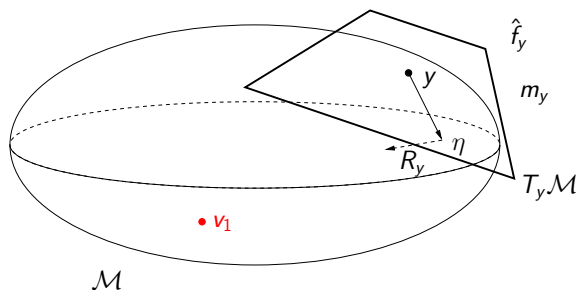
$$\begin{aligned}\hat{f}_y(\eta) &= \frac{y^T A y}{y^T B y} + 2 \frac{y^T A \eta}{y^T B y} + \frac{1}{y^T B y} \left( \eta^T A \eta - \frac{y^T A y}{y^T B y} \eta^T B \eta \right) + \dots \\ &= f(y) + 2 \langle P A y, \eta \rangle + \frac{1}{2} \langle 2P(A - f(y)B)P \eta, \eta \rangle + \dots\end{aligned}$$

where  $\langle u, v \rangle = u^T v$  and  $P = I - B y (y^T B^2 y)^{-1} y^T B$ .

Model:

$$m_y(\eta) = f(y) + 2 \langle P A y, \eta \rangle + \frac{1}{2} \langle P(A - f(y)B)P \eta, \eta \rangle, \quad y^T B \eta = 0.$$

## Quadratic model



$$m_y(\eta) = f(y) + 2\langle PAy, \eta \rangle + \frac{1}{2}\langle P(A - f(y)B)P\eta, \eta \rangle, \quad y^T B \eta = 0.$$

## Newton vs Trust-Region

Model:

$$m_y(\eta) = f(y) + 2\langle PAy, \eta \rangle + \frac{1}{2}\langle P(A - f(y)B)P\eta, \eta \rangle, \quad y^T B \eta = 0. \quad (7)$$

**Newton method:** Compute the **stationary point** of the model, i.e., solve

$$P(A - f(y)B)P\eta = -PAy.$$

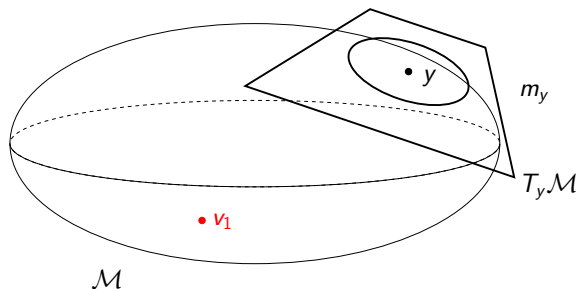
Instead, compute (approximately) the **minimizer** of  $m_y$  within a **trust-region**

$$\{\eta \in T_x \mathcal{M} : \eta^T \eta \leq \Delta^2\}.$$

## Trust-region subproblem

Minimize

$$m_y(\eta) = f(y) + 2\langle PAy, \eta \rangle + \frac{1}{2}\langle P(A - f(y)B)P\eta, \eta \rangle, \quad y^T B \eta = 0.$$

subject to  $\eta^T \eta \leq \Delta^2$ .

## Truncated CG method for the TR subproblem (1)

Let  $\langle \cdot, \cdot \rangle$  denote the standard inner product and let  $\mathcal{H}_{x_k} := P(A - f(x_k)B)P$  denote the Hessian operator.

**Initializations:**

Set  $\eta_0 = 0$ ,  $r_0 = P_{x_k} A x_k = A x_k - B x_k (x_k^T B^2 x_k)^{-1} x_k^T B A x_k$ ,  $\delta_0 = -r_0$ ;

Then repeat the following loop on  $j$ :

**Check for negative curvature**

**if**  $\langle \delta_j, \mathcal{H}_{x_k} \delta_j \rangle \leq 0$

    Compute  $\tau$  such that  $\eta = \eta_j + \tau \delta_j$  minimizes  $m(\eta)$  in (7) and satisfies  $\|\eta\| = \Delta$ ;

**return**  $\eta$ ;

## Truncated CG method for the TR subproblem (2)

**Generate next inner iterate**

Set  $\alpha_j = \langle r_j, r_j \rangle / \langle \delta_j, \mathcal{H}_{x_k} \delta_j \rangle$ ;

Set  $\eta_{j+1} = \eta_j + \alpha_j \delta_j$ ;

**Check trust-region**

**if**  $\|\eta_{j+1}\| \geq \Delta$

    Compute  $\tau \geq 0$  such that  $\eta = \eta_j + \tau \delta_j$  satisfies  $\|\eta\| = \Delta$ ;

**return**  $\eta$ ;

## Truncated CG method for the TR subproblem (3)

**Update residual and search direction**

Set  $r_{j+1} = r_j + \alpha_j \mathcal{H}_{x_k} \delta_j$ ;

Set  $\beta_{j+1} = \langle r_{j+1}, r_{j+1} \rangle / \langle r_j, r_j \rangle$ ;

Set  $\delta_{j+1} = -r_{j+1} + \beta_{j+1} \delta_j$ ;

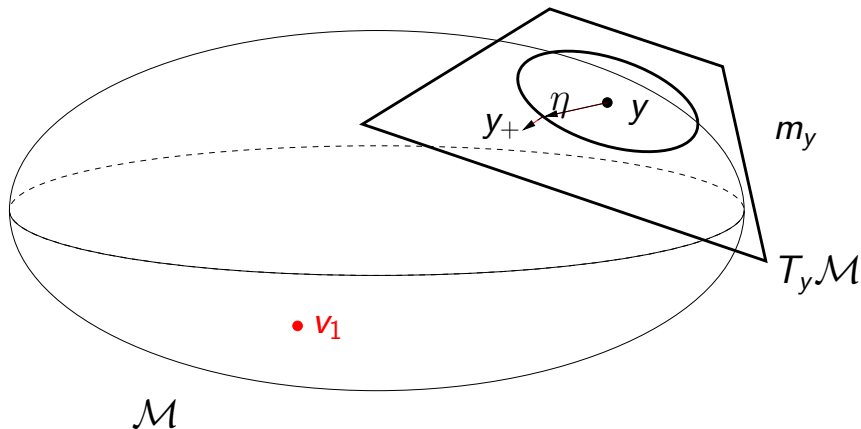
$j \leftarrow j + 1$ ;

**Check residual**

If  $\|r_j\| \leq \|r_0\| \min(\|r_0\|^\theta, \kappa)$  for some prescribed  $\theta$  and  $\kappa$

**return**  $\eta_j$ ;

## Overall iteration





## The outer iteration – manifold trust-region (1)

**Data:** symmetric  $n \times n$  matrices  $A$  and  $B$ , with  $B$  positive definite.

**Parameters:**  $\bar{\Delta} > 0$ ,  $\Delta_0 \in (0, \bar{\Delta})$ , and  $\rho' \in (0, \frac{1}{4})$ .

**Input:** initial iterate  $x_0 \in \{y : y^T B y = 1\}$ .

**Output:** sequence of iterates  $\{x_k\}$  in  $\{y : y^T B y = 1\}$ .

**Initialization:**  $k = 0$

Repeat the following:

## The outer iteration – manifold trust-region (2)

- ▶ Obtain  $\eta_k$  using the Steihaug-Toint truncated conjugate-gradient method to approximately solve the trust-region subproblem

$$\min_{x_k^T B \eta = 0} m_{x_k}(\eta) \quad \text{s.t.} \quad \|\eta\| \leq \Delta_k, \quad (8)$$

where  $m$  is defined in (7).

## The outer iteration – manifold trust-region (3)

- ▶ Evaluate

$$\rho_k = \frac{\hat{f}_{x_k}(0) - \hat{f}_{x_k}(\eta_k)}{m_{x_k}(0) - m_{x_k}(\eta_k)} \quad (9)$$

where  $\hat{f}_{x_k}(\eta) = \frac{(x_k + \eta)^T A(x_k + \eta)}{(x_k + \eta)^T B(x_k + \eta)}$ .

- ▶ Update the trust-region radius:

**if**  $\rho_k < \frac{1}{4}$

$$\Delta_{k+1} = \frac{1}{4} \Delta_k$$

**else if**  $\rho_k > \frac{3}{4}$  **and**  $\|\eta_k\| = \Delta_k$

$$\Delta_{k+1} = \min(2\Delta_k, \bar{\Delta})$$

**else**

$$\Delta_{k+1} = \Delta_k;$$

## The outer iteration – manifold trust-region (4)

- Update the iterate:

**if**  $\rho_k > \rho'$

$$x_{k+1} = (x_k + \eta_k) / \|x_k + \eta_k\|_B; \quad (10)$$

**else**

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

$$k \leftarrow k + 1$$

# Strategy

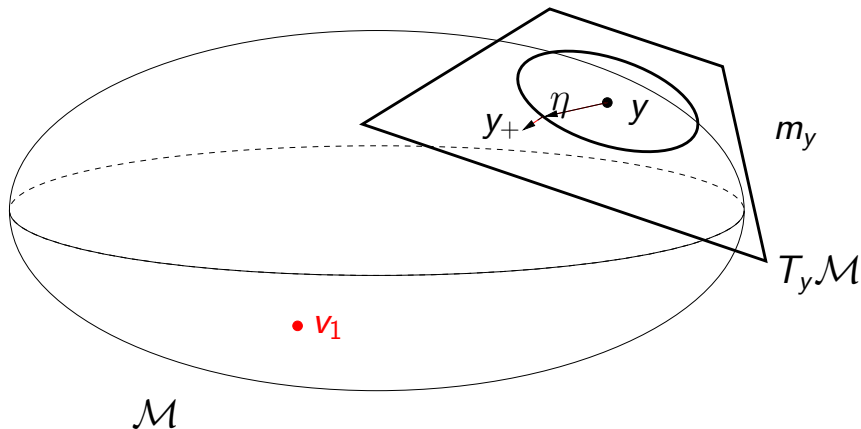
- ▶ Rewrite computation of leftmost eigenpair as an **optimization problem (on a manifold)**.
- ▶ Use a **model-trust-region** scheme to solve the problem.  
~> **Global convergence**.
- ▶ Take the **exact quadratic model** (at least, close to the solution).  
~> **Superlinear convergence**.
- ▶ Solve the trust-region subproblems using the **(Steihaug-Toint) truncated CG (tCG)** algorithm.  
~> **“Matrix-free”**, preconditioned iteration.  
~> **Minimal storage** of iteration vectors.

## Summary

We have obtained a trust-region algorithm for minimizing the Rayleigh quotient over an ellipsoid.

Generalization to trust-region algorithms for minimizing functions on manifolds: the **Riemannian Trust-Region (RTR)** method [ABG07].

## Convergence analysis



## Global convergence of Riemannian Trust-Region algorithms

Let  $\{x_k\}$  be a sequence of iterates generated by the RTR algorithm with  $\rho' \in (0, \frac{1}{4})$ . Suppose that  $f$  is  $C^2$  and bounded below on the level set  $\{x \in M : f(x) < f(x_0)\}$ . Suppose that  $\|\text{grad } f(x)\| \leq \beta_g$  and  $\|\text{Hess } f(x)\| \leq \beta_H$  for some constants  $\beta_g, \beta_H$ , and all  $x \in M$ . Moreover suppose that

$$\left\| \frac{D}{dt} \frac{d}{dt} Rt\xi \right\| \leq \beta_D \quad (11)$$

for some constant  $\beta_D$ , for all  $\xi \in TM$  with  $\|\xi\| = 1$  and all  $t < \delta_D$ , where  $\frac{D}{dt}$  denotes the covariant derivative along the curve  $t \mapsto Rt\xi$ . Further suppose that all approximate solutions  $\eta_k$  of the trust-region subproblems produce a decrease of the model that is at least a fixed fraction of the Cauchy decrease.



## Global convergence (cont'd)

It then follows that

$$\lim_{k \rightarrow \infty} \text{grad } f(x_k) = 0.$$

And only the local minima are stable (the saddle points and local maxima are unstable).

## Local convergence of Riemannian Trust-Region algorithms

Consider the RTR-tCG algorithm. Suppose that  $f$  is a  $C^2$  cost function on  $M$  and that

$$\|\mathcal{H}_k - \text{Hess } \hat{f}_{x_k}(0_k)\| \leq \beta_{\mathcal{H}} \|\text{grad } f(x_k)\|. \quad (12)$$

Let  $v \in M$  be a **nondegenerate local minimum** of  $f$ , (i.e.,  $\text{grad } f(v) = 0$  and  $\text{Hess } f(v)$  is positive definite). Further assume that  $\text{Hess } \hat{f}_{x_k}$  is Lipschitz-continuous at  $0_x$  uniformly in  $x$  in a neighborhood of  $v$ , i.e., there exist  $\beta_1 > 0$ ,  $\delta_1 > 0$  and  $\delta_2 > 0$  such that, for all  $x \in B_{\delta_1}(v)$  and all  $\xi \in B_{\delta_2}(0_x)$ , it holds

$$\|\text{Hess } \hat{f}_{x_k}(\xi) - \text{Hess } \hat{f}_{x_k}(0_{x_k})\| \leq \beta_{L2} \|\xi\|. \quad (13)$$

## Local convergence (cont'd)

Then there exists  $c > 0$  such that, for all sequences  $\{x_k\}$  generated by the RTR-tCG algorithm converging to  $v$ , there exists  $K > 0$  such that for all  $k > K$ ,

$$\text{dist}(x_{k+1}, v) \leq c (\text{dist}(x_k, v))^{\min\{\theta+1, 2\}}, \quad (14)$$

where  $\theta$  governs the stopping criterion of the tCG inner iteration.

## Convergence of trust-region-based eigensolver

### Theorem:

Let  $(A, B)$  be an  $n \times n$  symmetric/positive-definite matrix pencil with eigenvalues  $\lambda_1 < \lambda_2 \leq \dots \leq \lambda_{n-1} \leq \lambda_n$  and an associated  $B$ -orthonormal basis of eigenvectors  $(v_1, \dots, v_n)$ .

Let  $\mathcal{S}_i = \{y : Ay = \lambda_i By, y^T By = 1\}$  denote the intersection of the eigenspace of  $(A, B)$  associated to  $\lambda_i$  with the set  $\{y : y^T By = 1\}$ .

...

## Convergence (global)

- (i) Let  $\{x_k\}$  be a sequence of iterates generated by the Algorithm. Then  $\{x_k\}$  converges to the eigenspace of  $(A, B)$  associated to one of its eigenvalues. That is, there exists  $i$  such that  $\lim_{k \rightarrow \infty} \text{dist}(x_k, \mathcal{S}_i) = 0$ .
- (ii) Only the set  $\mathcal{S}_1 = \{\pm v_1\}$  is stable.

## Convergence (local)

(iii) There exists  $c > 0$  such that, for all sequences  $\{x_k\}$  generated by the Algorithm converging to  $\mathcal{S}_1$ , there exists  $K > 0$  such that for all  $k > K$ ,

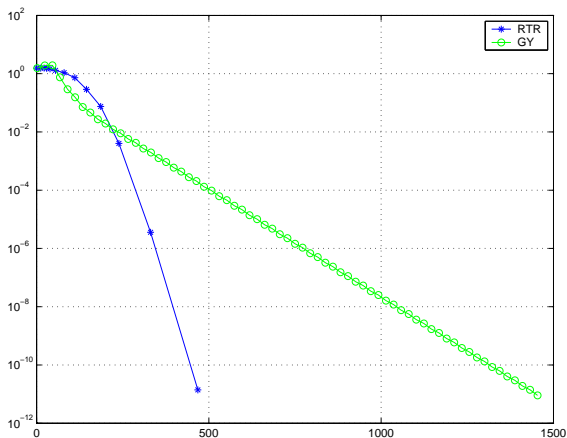
$$\text{dist}(x_{k+1}, \mathcal{S}_1) \leq c (\text{dist}(x_k, \mathcal{S}_1))^{\min\{\theta+1, 2\}} \quad (15)$$

with  $\theta > 0$ .

# Strategy

- ▶ Rewrite computation of leftmost eigenpair as an **optimization problem (on a manifold)**.
- ▶ Use a **model-trust-region** scheme to solve the problem.  
     $\rightsquigarrow$  **Global convergence**.
- ▶ Take the **exact quadratic model** (at least, close to the solution).  
     $\rightsquigarrow$  **Superlinear convergence**.
- ▶ Solve the trust-region subproblems using the **(Steihaug-Toint) truncated CG (tCG)** algorithm.  
     $\rightsquigarrow$  **“Matrix-free”**, preconditioned iteration.  
     $\rightsquigarrow$  **Minimal storage** of iteration vectors.

## Numerical experiments: RTR vs Krylov [GY02]



Distance to target versus matrix-vector multiplications.  
Symmetric/positive-definite generalized eigenvalue problem.



## Conclusion: A Three-Step Approach

- ▶ Formulation of the computational problem as a geometric optimization problem.
- ▶ Generalization of optimization algorithms on abstract manifolds.
- ▶ Exploit flexibility and additional structure to build numerically efficient algorithms.

## A few pointers

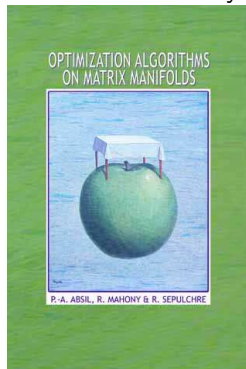
- ▶ Optimization on manifolds: Luenberger [Lue73], Gabay [Gab82], Smith [Smi93, Smi94], Udriște [Udr94], Manton [Man02], Mahony and Manton [MM02], PAA *et al.* [ABG04, ABG07]...
- ▶ Trust-region methods: Powell [Pow70], Moré and Sorensen [MS83], Moré [Mor83], Conn *et al.* [CGT00].
- ▶ Truncated CG: Steihaug [Ste83], Toint [Toi81], Conn *et al.* [CGT00]...
- ▶ Retractions: Shub [Shu86], Adler *et al.* [ADM<sup>+</sup>02]...

# THE END

## *Optimization Algorithms on Matrix Manifolds*





P.-A. Absil, R. Mahony, R. Sepulchre






Princeton University Press, January 2008












1. Introduction
2. Motivation and applications
3. Matrix manifolds: first-order geometry
4. Line-search algorithms
5. Matrix manifolds: second-order geometry
6. Newton's method
7. Trust-region methods
8. A constellation of superlinear algorithms

-  P.-A. Absil, C. G. Baker, and K. A. Gallivan, *Trust-region methods on Riemannian manifolds with applications in numerical linear algebra*, Proceedings of the 16th International Symposium on Mathematical Theory of Networks and Systems (MTNS2004), Leuven, Belgium, 5–9 July 2004, 2004.
-  \_\_\_\_\_, *Trust-region methods on Riemannian manifolds*, Found. Comput. Math. **7** (2007), no. 3, 303–330.
-  Roy L. Adler, Jean-Pierre Dedieu, Joseph Y. Margulies, Marco Martens, and Mike Shub, *Newton's method on Riemannian manifolds and a geometric model for the human spine*, IMA J. Numer. Anal. **22** (2002), no. 3, 359–390.
-  P.-A. Absil, R. Mahony, and R. Sepulchre, *Optimization algorithms on matrix manifolds*, Princeton University Press, Princeton, NJ, January 2008.

-  Andrew R. Conn, Nicholas I. M. Gould, and Philippe L. Toint, *Trust-region methods*, MPS/SIAM Series on Optimization, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2000. MR MR1774899 (2003e:90002)
-  D. Gabay, *Minimizing a differentiable function over a differential manifold*, J. Optim. Theory Appl. **37** (1982), no. 2, 177–219. MR MR663521 (84h:49071)
-  Gene H. Golub and Qiang Ye, *An inverse free preconditioned Krylov subspace method for symmetric generalized eigenvalue problems*, SIAM J. Sci. Comput. **24** (2002), no. 1, 312–334.
-  Magnus R. Hestenes and William Karush, *A method of gradients for the calculation of the characteristic roots and vectors of a real symmetric matrix*, J. Research Nat. Bur. Standards **47** (1951), 45–61.

-  Uwe Helmke and John B. Moore, *Optimization and dynamical systems*, Communications and Control Engineering Series, Springer-Verlag London Ltd., London, 1994, With a foreword by R. Brockett. MR MR1299725 (95j:49001)
-  David G. Luenberger, *Introduction to linear and nonlinear programming*, Addison-Wesley, Reading, MA, 1973.
-  Jonathan H. Manton, *Optimization algorithms exploiting unitary constraints*, IEEE Trans. Signal Process. **50** (2002), no. 3, 635–650. MR MR1895067 (2003i:90078)
-  Robert Mahony and Jonathan H. Manton, *The geometry of the Newton method on non-compact Lie groups*, J. Global Optim. **23** (2002), no. 3-4, 309–327, Nonconvex optimization in control. MR MR1923049 (2003g:90114)
-  J. J. Moré, *Recent developments in algorithms and software for trust region methods*, Mathematical programming: the state of the art (Bonn, 1982) (Berlin), Springer, 1983, pp. 258–287.

-  Jorge J. Moré and D. C. Sorensen, *Computing a trust region step*, SIAM J. Sci. Stat. Comput. **4** (1983), no. 3, 553–572. MR MR723110 (86b:65063)
-  M. Mongeau and M. Torki, *Computing eigenelements of real symmetric matrices via optimization*, Comput. Optim. Appl. **29** (2004), no. 3, 263–287. MR MR2101850 (2005h:65061)
-  M. J. D. Powell, *A new algorithm for unconstrained optimization*, Nonlinear Programming (Proc. Sympos., Univ. of Wisconsin, Madison, Wis., 1970), Academic Press, New York, 1970, pp. 31–65.
-  Michael Shub, *Some remarks on dynamical systems and numerical analysis*, Proc. VII ELAM. (L. Lara-Carrero and J. Lewowicz, eds.), Equinoccio, U. Simón Bolívar, Caracas, 1986, pp. 69–92.
-  Steven Thomas Smith, *Geometric optimization methods for adaptive filtering*, Ph.D. thesis, Division of Applied Sciences, Harvard University, Cambridge, MA, May 1993.

-  Steven T. Smith, *Optimization techniques on Riemannian manifolds*, Hamiltonian and gradient flows, algorithms and control, Fields Inst. Commun., vol. 3, Amer. Math. Soc., Providence, RI, 1994, pp. 113–136. MR MR1297990 (95g:58062)
-  Trond Steihaug, *The conjugate gradient method and trust regions in large scale optimization*, SIAM J. Numer. Anal. **20** (1983), no. 3, 626–637. MR MR701102 (84g:49047)
-  Ph. L. Toint, *Towards an efficient sparsity exploiting Newton method for minimization*, Sparse Matrices and Their Uses (I. S. Duff, ed.), Academic Press, London, 1981, pp. 57–88.
-  Constantin Udriște, *Convex functions and optimization methods on Riemannian manifolds*, Mathematics and its Applications, vol. 297, Kluwer Academic Publishers Group, Dordrecht, 1994. MR MR1326607 (97a:49038)