

CS468 – Lecture V

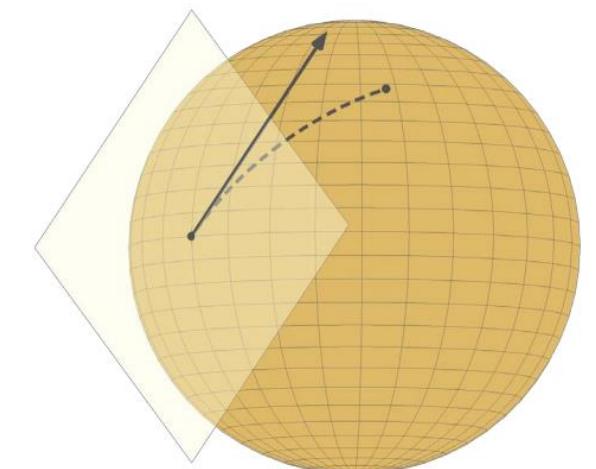
Non-Euclidean Methods in Machine Learning

OPTIMIZATION ON RIEMANNIAN MANIFOLDS

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

TOLGA BIRDAL
INES CHAMI
LEONIDAS GUIBAS

CS468.STANFORD.EDU



Optimization on Riemannian Manifolds

Partial Slide Credits: Nicolas Boumal

An introduction to optimization on smooth manifolds ↗

<https://web.math.princeton.edu/~nboumal/book/index.html>

Constrained vs. Riemannian Optimization

$$\min_x f(x) \text{ subject to } x \in \mathcal{M}$$

Constraints on Parameters

Linear spaces

Unconstrained; linear equality constraints

Low rank (matrices, tensors)

Recommender systems, large-scale Lyapunov equations, ...

Orthonormality (Grassmann, Stiefel, rotations)

Dictionary learning, SfM, SLAM, PCA, ICA, SBM, Electr. Struct. Comp....

Positivity (positive definiteness, positive orthant)

Metric learning, Gaussian mixtures, diffusion tensor imaging, ...

Symmetry (quotient manifolds)

Invariance under group actions

Constrained vs. Riemannian Optimization

$$\min_x f(x) \text{ subject to } x \in \mathcal{M}$$



Constraints on Parameters

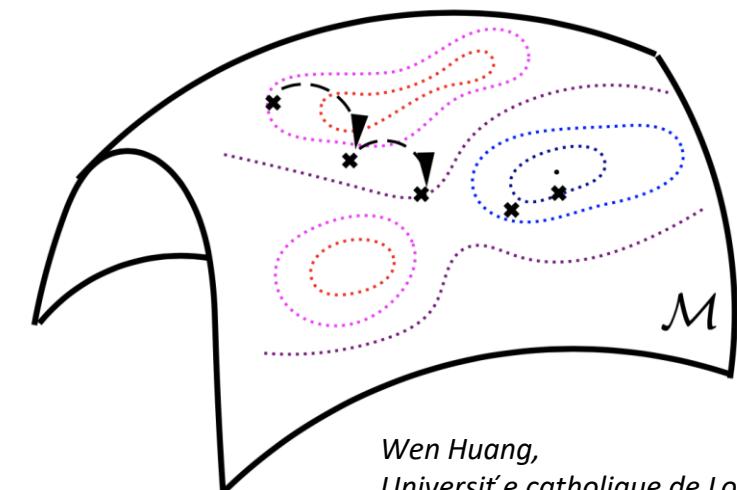


Unconstrained
Optimization

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

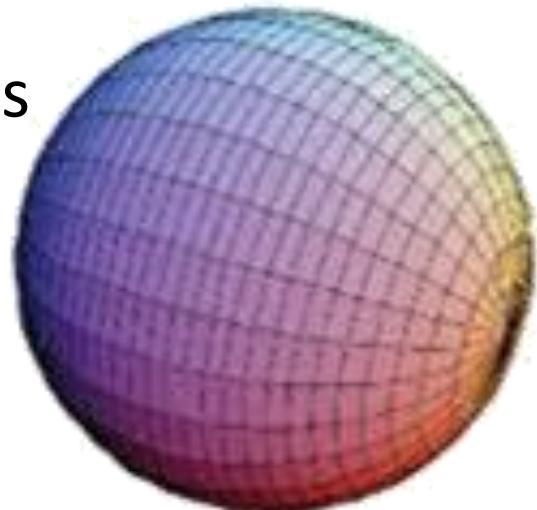
Constrained vs. Riemannian Optimization

1. We can use unconstrained optimization tools (gradient descent, newton etc.).
2. No need to consider Lagrange multipliers or penalty functions.
3. Theoretical guarantees usually transfer from Euclidean space to Riemannian manifolds.
4. Can be cheaper in terms of resource use depending upon the application.

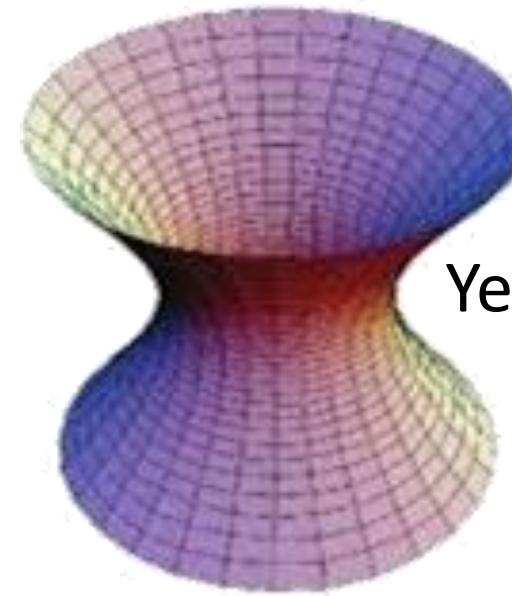


Which Manifolds?

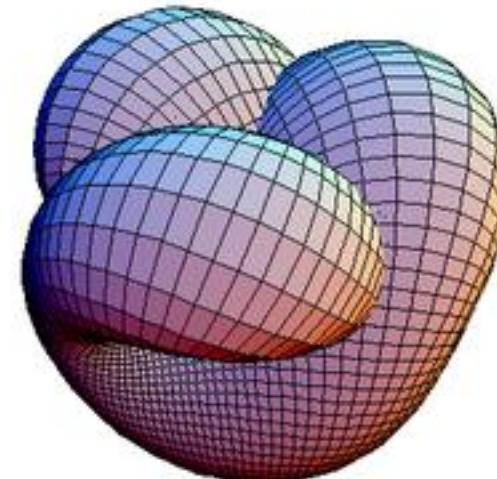
Yes



Yes



No

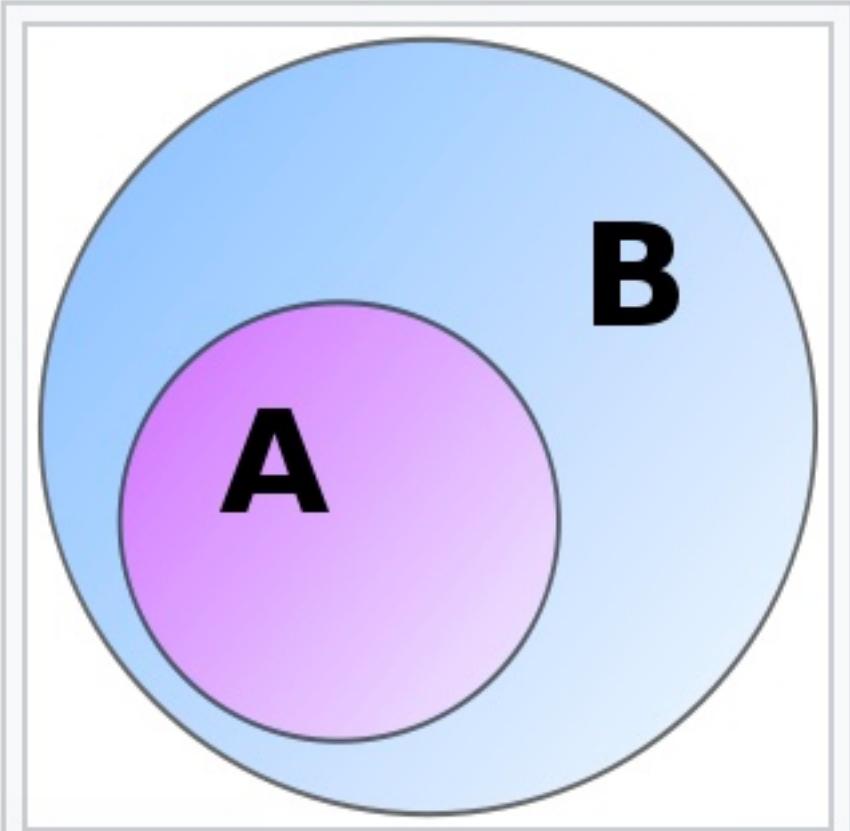


We restrict ourselves to *embedded submanifolds* of **linear spaces**.

Let N be a smooth manifold. One possible definition (I believe) for an embedded submanifold of N is some $M \subset N$ that is a (smooth) manifold such that the inclusion $i : M \hookrightarrow N$ is an embedding.

In [mathematics](#), if A is a [subset](#) of B , then the **inclusion map** (also [inclusion function](#), [insertion](#),^[1] or [canonical injection](#)) is the [function](#) ι that sends each element x of A to x , treated as an element of B :

$$\iota : A \rightarrow B, \quad \iota(x) = x.$$



A is a subset of B , and \square
 B is a superset of A .

For smoothness of \mathcal{M} , our model space is the unit sphere in \mathbb{R}^n :

$$S^{n-1} = \{x \in \mathbb{R}^n : \|x\| = 1\}, \quad (3.2)$$

where $\|x\| = \sqrt{x^\top x}$ is the Euclidean norm on \mathbb{R}^n . Intuitively, we think of S^{n-1} as a smooth nonlinear space in \mathbb{R}^n . Our definitions below are compatible with this intuition, and we call S^{n-1} an *embedded submanifold* of \mathbb{R}^n .

An important element in these definitions is to capture the idea that S^{n-1} can be locally approximated by a linear space around any point x : we call these *tangent spaces*, denoted by $T_x S^{n-1}$. This is as opposed to a cube for which no good linearization exists at the edges. More specifically for our example, S^{n-1} is defined by the constraint $x^\top x = 1$, and we expect that differentiating this constraint should yield a suitable linearization:

$$T_x S^{n-1} = \{v \in \mathbb{R}^n : v^\top x + x^\top v = 0\} = \{v \in \mathbb{R}^n : x^\top v = 0\}. \quad (3.3)$$

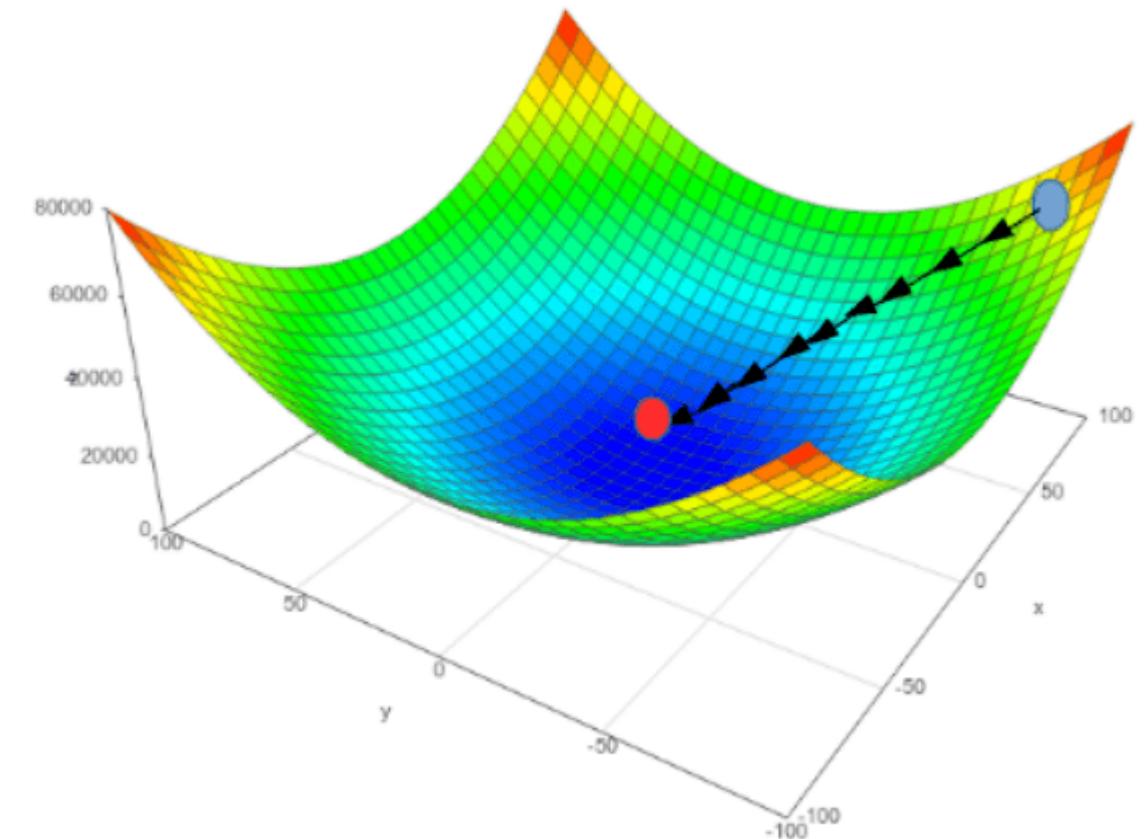
In the same spirit, it stands to reason that linear spaces, or open subsets of linear spaces, should also be considered smooth.

Optimization on a Euclidean Space

$$x_{k+1} = x_k + \Delta x_k = x_k + \tau_k d_k$$

↓
previous iterate

movement
↓
step



Optimization on a Euclidean Space

- Steepest (Gradient) Descent:

$$x_{k+1} = x_k + \Delta x_k = x_k + \tau_k d_k$$

↓
movement
 $x_{k+1} = x_k - \tau_k \nabla f(x_k)$
update step size
↓
step
previous iterate

- Newton's Method:

$$x_{k+1} = x_k - [\nabla^2 f(x_k)]^{-1} \nabla f(x_k)$$

Hessian Gradient

- Other methods to compute the direction and the step size.

A Naïve Manifold Optimizer: Projected Gradient Descent

Loop:

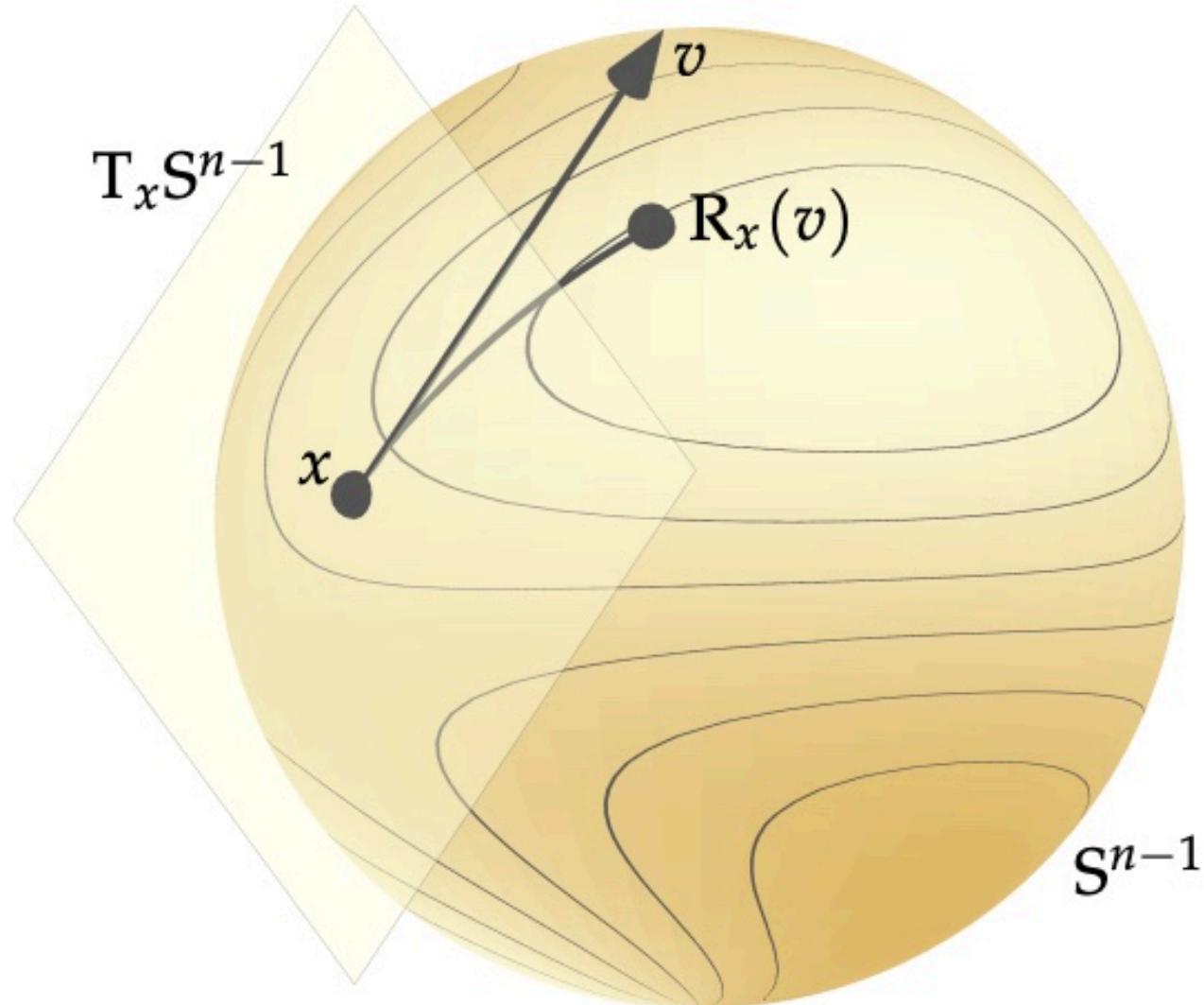
1. Perform the gradient update:

$$\hat{x} = x_k - \tau_k \nabla f(x_k)$$

2. Project on the manifold:

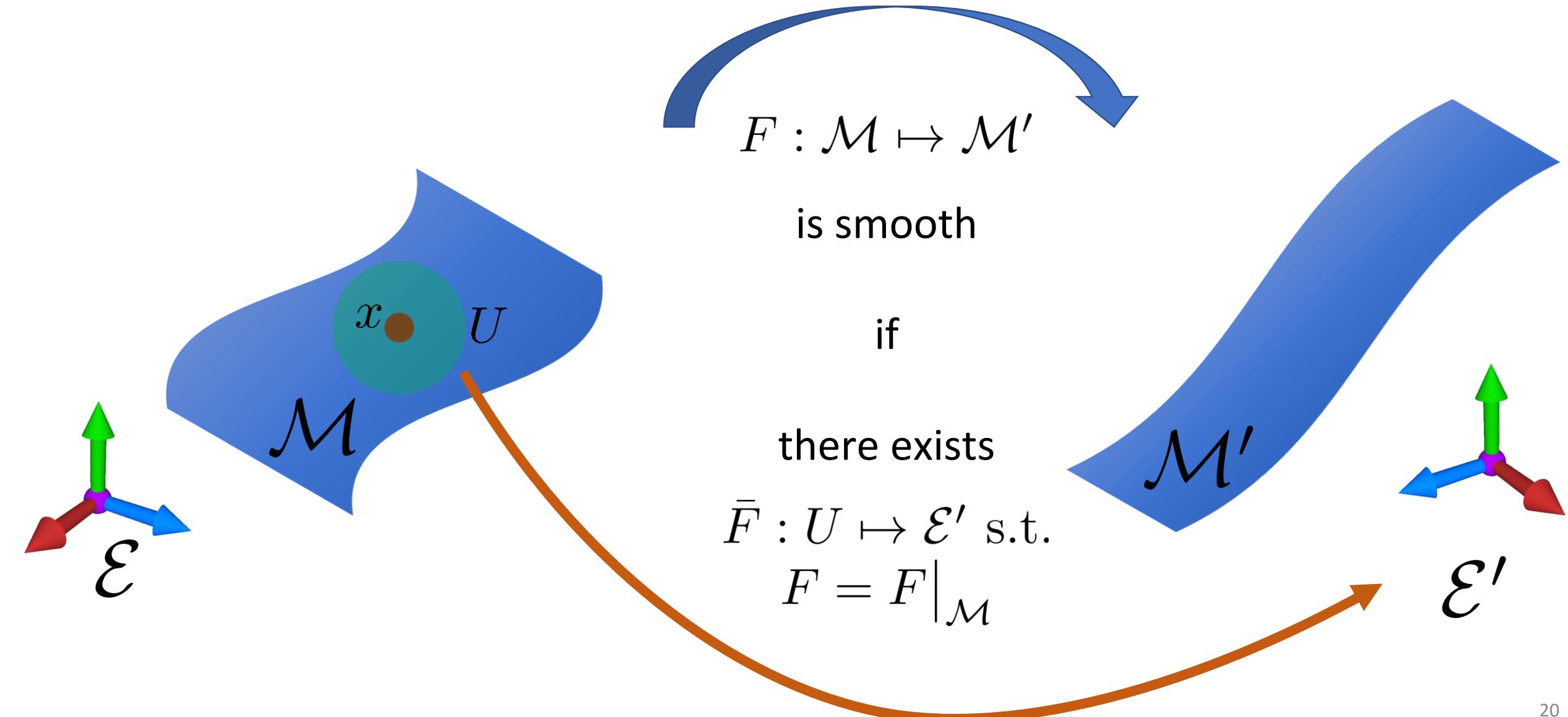
$$x_{k+1} = \Pi_{\mathcal{M}}(\hat{x})$$

Inaccurate & Slow Convergence



Retraction $R_x(v) = \frac{x+v}{\|x+v\|}$ on the sphere.

Smooth Map



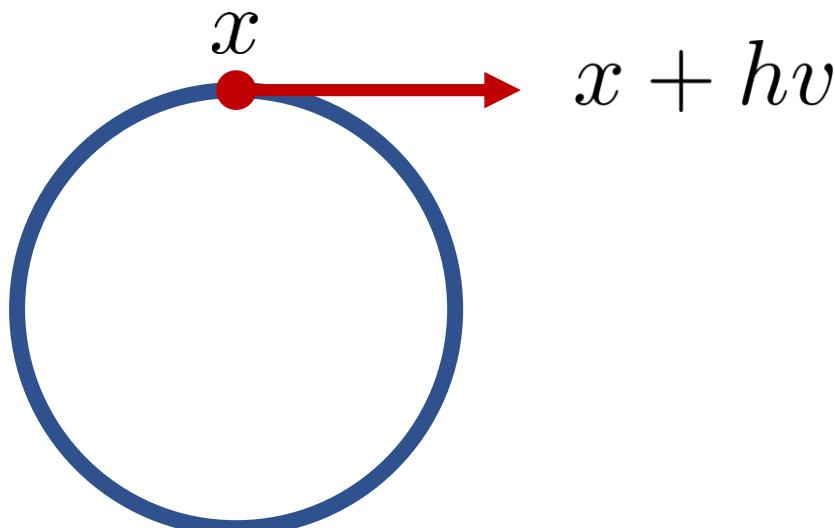
The Differential of a Smooth Map

Directional derivative of a smooth function $f : \mathbb{R}^n \mapsto \mathbb{R}^n$ along a vector v :

$$D_v f = Df(x) \cdot v = \lim_{h \rightarrow 0} \frac{f(x + hv) - f(x)}{h}$$

Cannot apply to a smooth manifold \mathcal{M} as $(x + hv)$ might not be belong to \mathcal{M} .

For \bar{F} , this
is always defined,
but not for F .



The Differential of a Smooth Map

For $\bar{F} : \mathcal{E} \mapsto \mathcal{E}'$ we have:

$$D_v \bar{F} = D\bar{F}(x) \cdot v = \lim_{h \rightarrow 0} \frac{\bar{F}(x + hv) - \bar{F}(x)}{h}$$

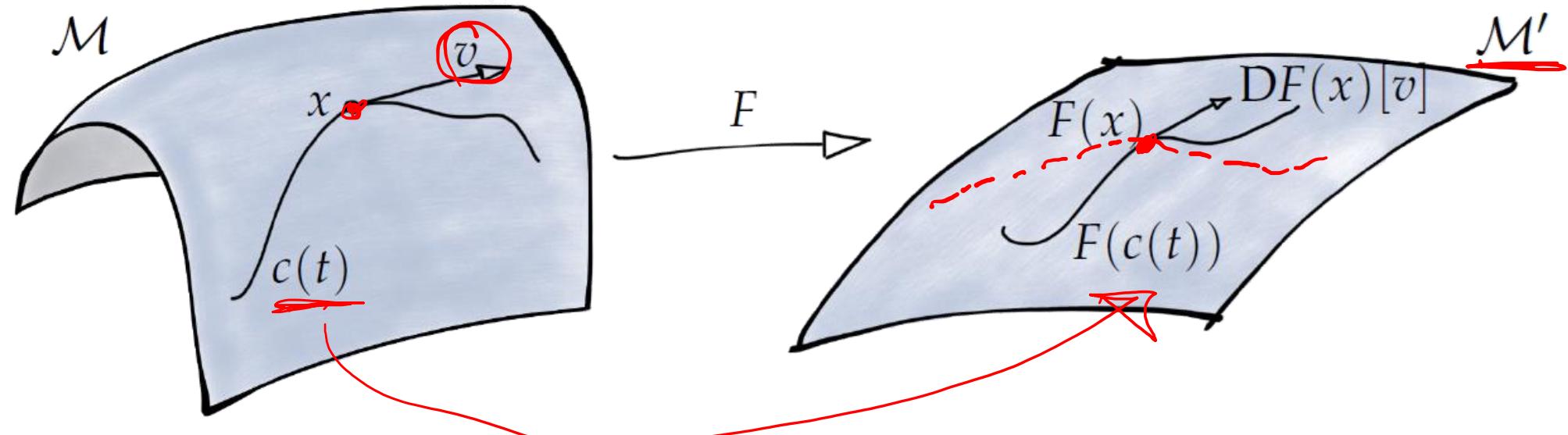
For $F : \mathcal{M} \mapsto \mathcal{M}'$ and $v \in T_x \mathcal{M}$, we write:

$$D_v F = D_v \bar{F}$$

Does not depend on the choice
of the smooth extension!

The Differential of a Smooth Map

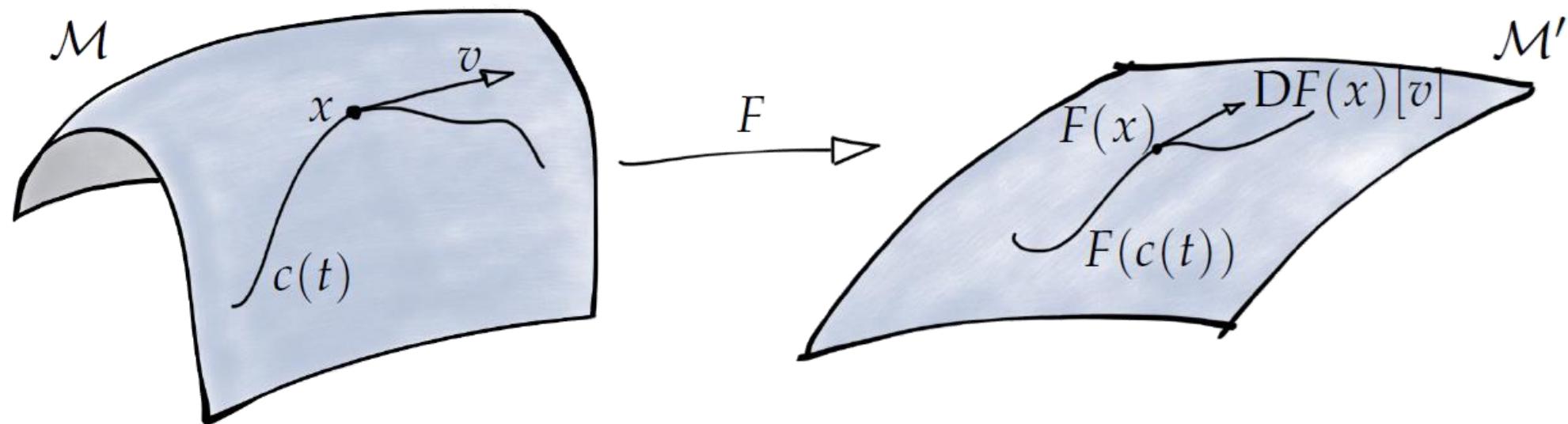
Equivalently, consider a smooth curve $c(t) \in \mathcal{M}$ with $F(c(t)) \in \mathcal{M}'$ being a curve in \mathcal{M}' passing through $F(x)$ with velocity $DF(x) \cdot v$.



$$DF(x)[v] : \mathcal{T}_x \mathcal{M} \mapsto \mathcal{T}_{F(x)} \mathcal{M}' = \frac{d}{dt} F(c(t)) \Big|_{t=0}$$

Does not depend on the choice of the curve!

The Differential of a Smooth Map



$$DF(x) \cdot v = DF(x)[v] = \lim_{h \rightarrow 0} \frac{F(c(t)) - F(x)}{h} = (F \circ c)'(0)$$

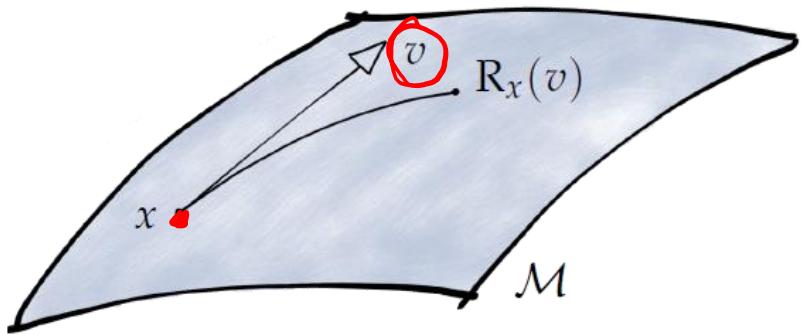
Does not depend on the
choice of the curve!

$$DF(x) \cdot v = DF(x)[v] = \lim_{h \rightarrow 0} \frac{F(c(t+h)) - F(c(t))}{h}|_{t=0} = (F \circ c)'(0)$$

$$DF(x) \cdot v = DF(x)[v] = \lim_{h \rightarrow 0} \frac{F(c(t+h)) - F(c(t))}{h} = (F \circ c)'(t)$$

How to choose these curves == Retractions

- Smooth choice of curves over the tangent bundle.
- Maps tangent vectors back to the manifold.
- Defines curves in a given direction.



$$c(t) = R(x, tv) = R_x(tv)$$

~~$c(0) = x$ and $c'(0) = v$~~

A **Retraction** map $R : \mathcal{T}_x \mathcal{M} \mapsto \mathcal{M}$ satisfies:

1. R is continuously differentiable.
2. $R_x(0) = \underline{x}$ (centering)
3. $DR_x(0)[v] = \underline{v}$ (local rigidity)
Identity

In differential geometry, the **tangent bundle** of a differentiable manifold M is a manifold TM which assembles all the tangent vectors in M . As a set, it is given by the disjoint union^[note 1] of the tangent spaces of M . That is,

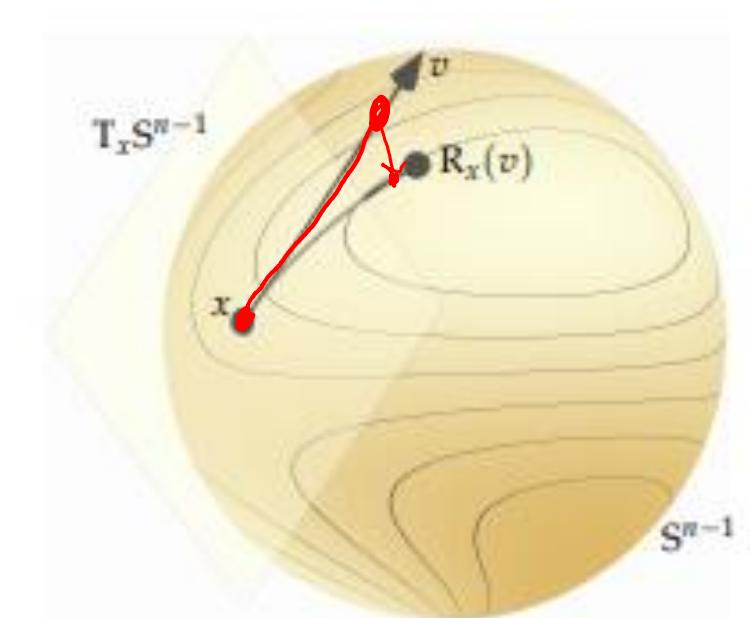
$$\begin{aligned} TM &= \bigsqcup_{x \in M} T_x M \\ &= \bigcup_{x \in M} \{x\} \times T_x M \\ &= \bigcup_{x \in M} \{(x, y) \mid y \in T_x M\} \\ &= \{(x, y) \mid x \in M, y \in T_x M\} \end{aligned}$$

Projection as a Retraction

One way to define a retraction on the sphere \mathbb{S}^{d-1} :

$$R_x(v) = \frac{x + v}{\|x + v\|}$$

Exercise: Verify that this is indeed a valid retraction.



Retraction $R_x(v) = \frac{x+v}{\|x+v\|}$ on the sphere.

Example 3.40. Let x be a point on the sphere S^{n-1} and let v be tangent at x , that is, $x^\top v = 0$. To move away from x along v while remaining on the sphere, one way is to take the step in \mathbb{R}^n then to project back to the sphere:

$$R_x(v) \triangleq \frac{x + v}{\|x + v\|} = \frac{x + v}{\sqrt{1 + \|v\|^2}}. \quad (3.31)$$

Consider the curve $c: \mathbb{R} \rightarrow S^{n-1}$ defined by:

$$c(t) = R_x(tv) = \frac{x + tv}{\sqrt{1 + t^2 \|v\|^2}}.$$

Evidently, $c(0) = x$. This holds because $R_x(0) = x$. Furthermore, one can compute $c'(0) = v$, that is: locally around x , up to first order, the retraction curve moves along v . Another way to state this is via the chain rule:

$$v = c'(0) = DR_x(0)[v],$$

where $DR_x(0)$ is understood as per Definition 3.27 for $R_x: T_x S^{n-1} \rightarrow S^{n-1}$. In other words: $DR_x(0): T_x S^{n-1} \rightarrow T_x S^{n-1}$ is the identity map.

Needed for Optimization

1. A representation for points $x \in \mathcal{M}$, for tangent spaces $\mathcal{T}_x\mathcal{M}$ and Riemannian metrics $g_x(\cdot, \cdot)$,
2. A map from the tangent space to the manifold: $R_x : \mathcal{T}_x\mathcal{M} \mapsto \mathcal{M}$
3. Expressions for $f(x)$, $\text{grad } f(x)$ and $\text{Hess } f(x)$
4. Notion of *vector transport* for second order methods

Do we need the true geodesics to move?

NO!

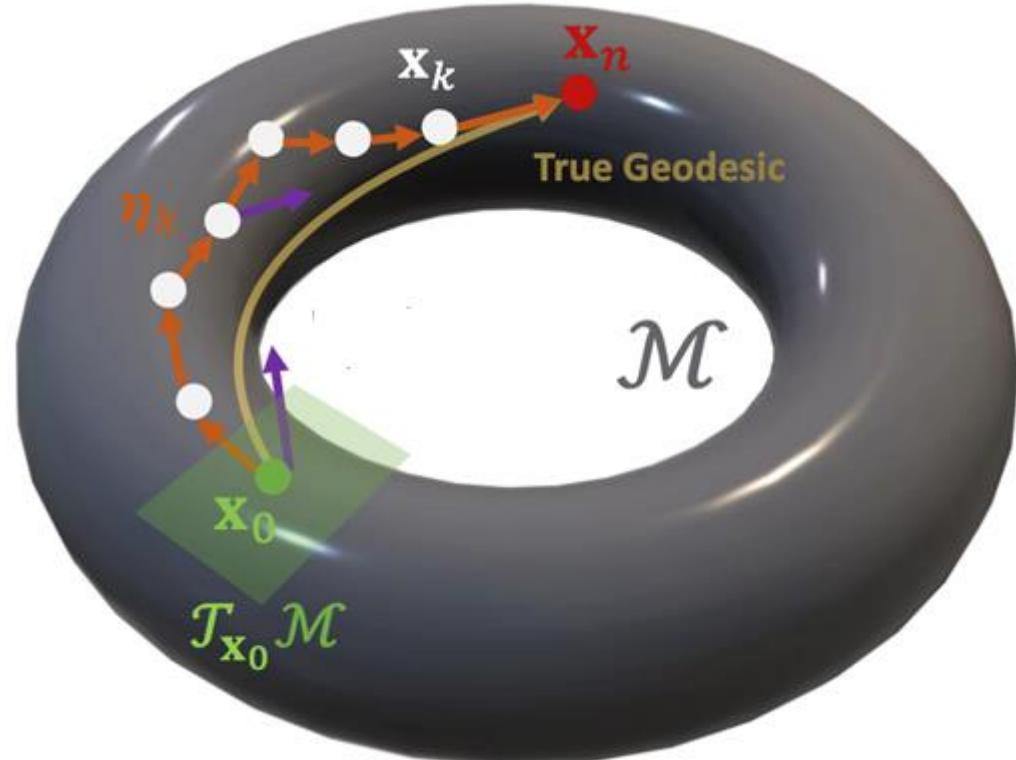
A Less Naïve Manifold Optimizer

Euclidean

$$x_{k+1} = x_k + \tau_k d_k$$

Riemannian

$$x_{k+1} = R_{x_k}(\tau_k \eta_k), \quad \eta_k \in T_{x_k} \mathcal{M}$$



Riemannian Gradient & Hessian

Riemannian gradient of $f(x)$ at x is the unique **tangent vector** in $\mathcal{T}_x\mathcal{M}$ satisfying $\forall v \in \mathcal{T}_x\mathcal{M}$:

$$DF(x)[v] = \langle \text{grad}f(x), v \rangle_x$$

Direction of
steepest ascent

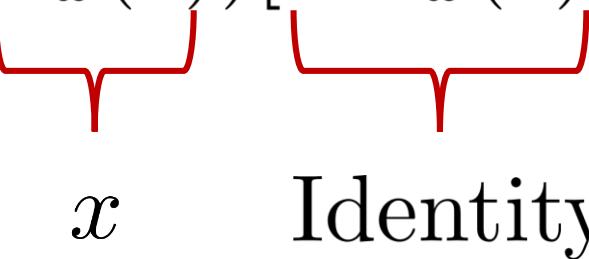
If x is a local optimum of f , then $\text{grad}f(x) = 0$.

Riemannian Hessian of $f(x)$ at x is a linear operator from $\mathcal{T}_x\mathcal{M} \mapsto \mathcal{T}_x\mathcal{M}$:

$$\text{Hess}f(x) : \mathcal{T}_x\mathcal{M} \mapsto \mathcal{T}_x\mathcal{M} : v \mapsto \nabla_v \text{grad}f$$

More on that
later.

Riemannian Gradient

1. $\text{D}(f \circ R_x)(0)[v] = \text{D}f(R_x(0))[\text{D}R_x(0)[v]] = \text{D}f(x)[v]$


x Identity
2. Hence, $\langle \text{grad}(f \circ R_x)(0), v \rangle_x = \langle \text{grad}f(x), v \rangle_x$
3. $\text{D}f(x)[v] = \text{D}\bar{f}(x)[v] = \langle v, \text{grad}\bar{f}(x) \rangle$ (Euclidean)

Riemannian Gradient

Observe that $\mathcal{T}_x\mathcal{M} \subset \mathcal{E}$ and $\text{grad}\bar{f}(x) \in \mathcal{E}$. Hence:

$$\text{grad}\bar{f}(x) = \text{grad}\bar{f}(x)_{||} + \text{grad}\bar{f}(x)_{\perp}$$

Tangential Orthogonal

Since $\forall v \in \mathcal{T}'_x$ $\langle v, \text{grad}\bar{f}(x)_{\perp} \rangle = 0$:

$$\begin{aligned}\langle v, \text{grad}f(x) \rangle_x &= \langle v, \text{grad}\bar{f}(x) \rangle_x \\ &= \langle v, \text{grad}\bar{f}(x)_{||} + \text{grad}\bar{f}(x)_{\perp} \rangle \\ &= \langle v, \text{grad}\bar{f}(x)_{||} \rangle\end{aligned}$$

Riemannian Gradient

$$\text{grad } f(x) = \text{grad } \bar{f}(x)_{||}$$

To compute the **Riemannian gradient**:

1. Obtain an expression for the classical gradient.
2. Orthogonally project to the tangent space
(cancel the normal component)

Riemannian Gradient

$$\text{grad}f(x) = \text{Proj}_x(\text{grad}\bar{f}(x)) = \Pi_x(\text{grad}\bar{f}(x))$$

where $\Pi_x \triangleq \text{Proj}_x$ is a **projector**:

$$\langle u - \Pi_x(u), v \rangle = 0 \quad \forall v \in \mathcal{T}_x \mathcal{M} \text{ and } u \in \mathcal{E}$$

$$\Pi_x \circ \Pi_x = \Pi_x$$

Consider the Sphere

$$T_x \mathcal{S}^{d-1} = \{v \in \mathbb{R}^n : x^\top v = 0\} = \{v \in \mathbb{R}^n : \langle x, v \rangle = 0\}$$

$$\Pi_x(u) = u - (x^\top u)x = (I - xx^\top)u$$

A Less Naïve Manifold Optimizer

Loop:

1. Perform the gradient update:

$$\hat{x} = x_k - \tau_k \text{grad} f(x_k)$$

2. Project on the manifold:

$$x_{k+1} = \Pi_{\mathcal{M}}(\hat{x})$$

Inaccurate & Slow Convergence

Riemannian Gradient Descent

Loop:

$$x_{k+1} = x_k + \tau_k \text{ grad } f(x_k)$$

Provable convergence to local minimum.

Riemannian Gradient Descent – Taylor Perspective

The composition $f \circ R_x: T_x \mathcal{M} \rightarrow \mathbf{R}$ is on a linear space, hence we may Taylor expand it:

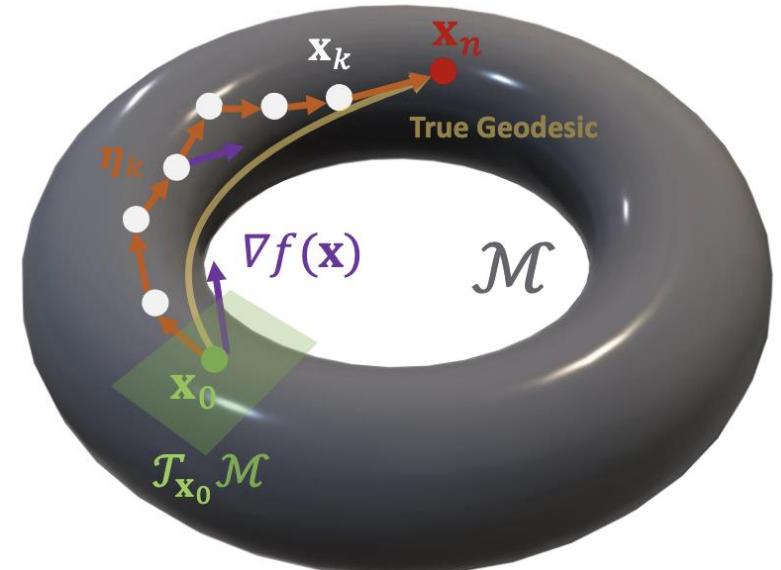
$$\begin{aligned} f(R_x(v)) &= f(R_x(0)) + \langle \text{grad}(f \circ R_x)(0), v \rangle_x + O(\|v\|_x^2) \\ &= f(x) + \langle \text{grad}f(x), v \rangle_x + O(\|v\|_x^2) \end{aligned}$$

This will take us to 2nd order methods.

A Generic Riemannian Optimizer

Sufficient to describe a large family of optimization methods.

- 1 **input:** A Riemannian manifold \mathcal{M} , a retraction operator R
- 2 **while** \mathbf{x}_k does not sufficiently minimize f **do**
- 3 Pick a gradient related descent direction $\eta_k \in \mathcal{T}_{\mathbf{x}_k} \mathcal{M}$.
- 4 Choose a retraction $R_{\mathbf{x}_k} : \mathcal{T}_{\mathbf{x}_k} \mathcal{M} \rightarrow \mathcal{M}$.
- 5 Choose a step length $\tau_k \in \mathbb{R}$.
- 6 Set $\mathbf{x}_{k+1} \leftarrow R_{\mathbf{x}_k}(\tau_k \eta_k)$.
- 7 $k \leftarrow k + 1$.



Riemannian Steepest Descent with Armijo Line Search

- Accelerated scheme with adaptive step size
- Determines the step size via Armijo Condition

1 **input:** A Riemannian manifold \mathcal{M} , a retraction operator R , the projection operator onto the tangent space $\Pi_{\mathbf{x}_k} : \mathbb{R}^n \rightarrow T_{\mathbf{x}_k} \mathcal{M}$, a real-valued, differentiable potential energy f , initial iterate $\mathbf{x}_0 \in \mathcal{M}$ and the Armijo line search scalars including c .

2 **while** \mathbf{x}_k does not sufficiently minimize f **do**

// Euclidean gradient to Riemannian direction
3 $\eta_k \leftarrow -\text{grad } f(\mathbf{x}_k) \triangleq \Pi_{\mathbf{x}_k}(-\nabla f(\mathbf{x}_k))$.

4 Select \mathbf{x}_{k+1} such that:

5 where τ_k is the Armijo step size.

6 $k \leftarrow k + 1$.

$$x_{k+1} = x_k + \tau_k \eta_k \quad (22)$$

Armijo Condition

Determining the Step Size

Options:

1. Fix $\tau_k = \tau \forall k$
2. Optimal τ_k that minimizes $f(R_{x_k}(-t\text{grad}f(x_k)))$ can be found cheaply in rare cases.
3. **Backtracking:** Start with an initial guess and iteratively reduce until the step size is deemed acceptable. (conditions vary)