

# Optimization on Manifolds for ML

---

**Reshad Hosseini**

**Department of Machine Intelligence and Robotics  
University of Tehran**

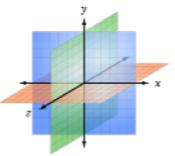
**Invited Talk at Amirkabir AI Summer Symposium**



## Joint Work with Suvrit Sra



## ► Vector spaces



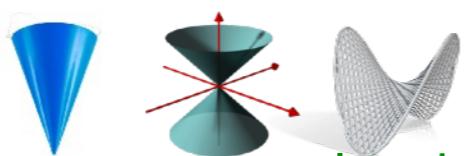
## ► Manifolds

(hypersphere, orthogonal matrices, complicated surfaces)



## ► Convex sets

(probability simplex, semidefinite cone, polyhedra)



## ► Metric spaces

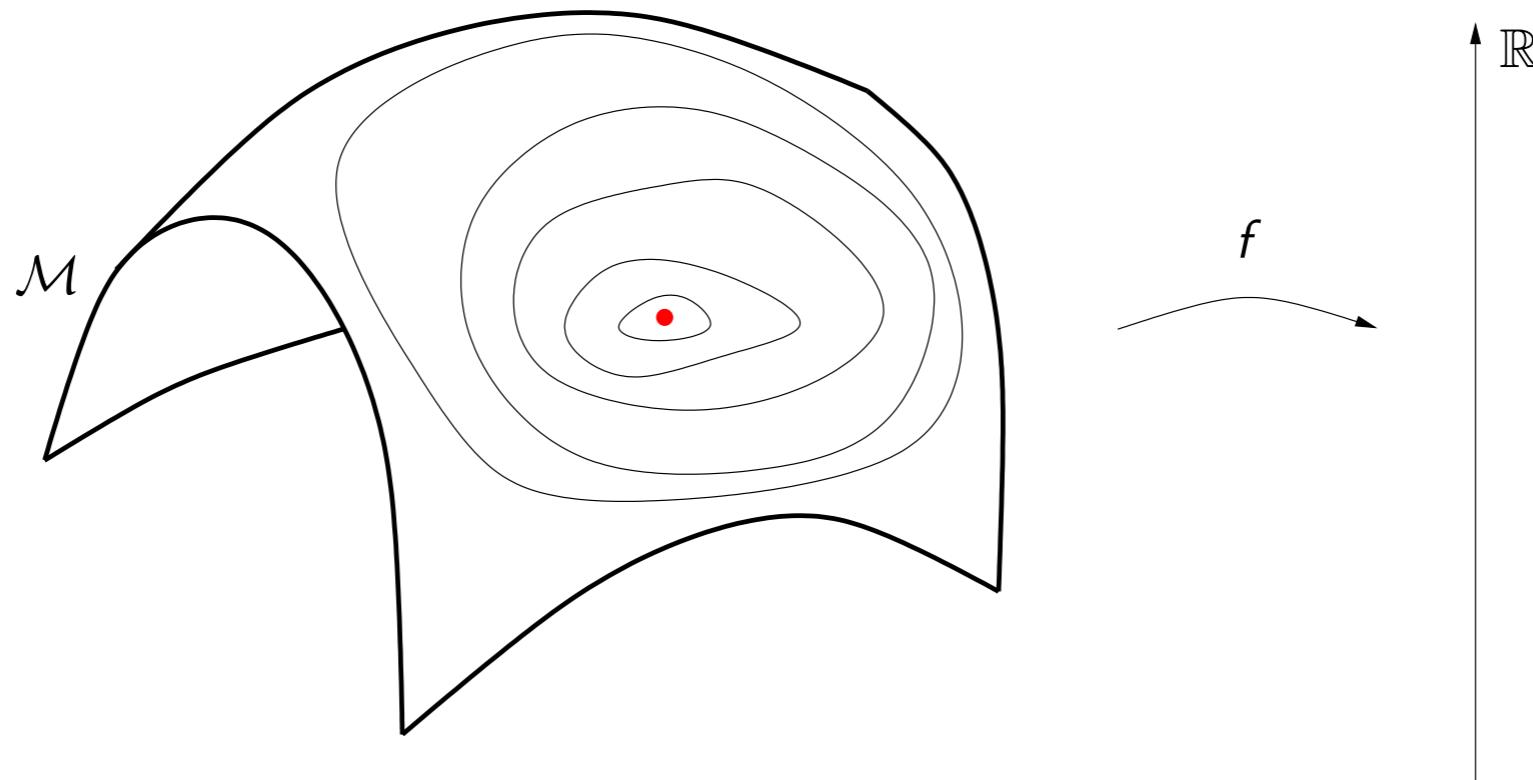
(tree space, Wasserstein spaces, CAT(0), space-of-spaces)



Machine Learning  
Graphics  
Robotics  
Vision  
BCI  
NLP  
Statistics

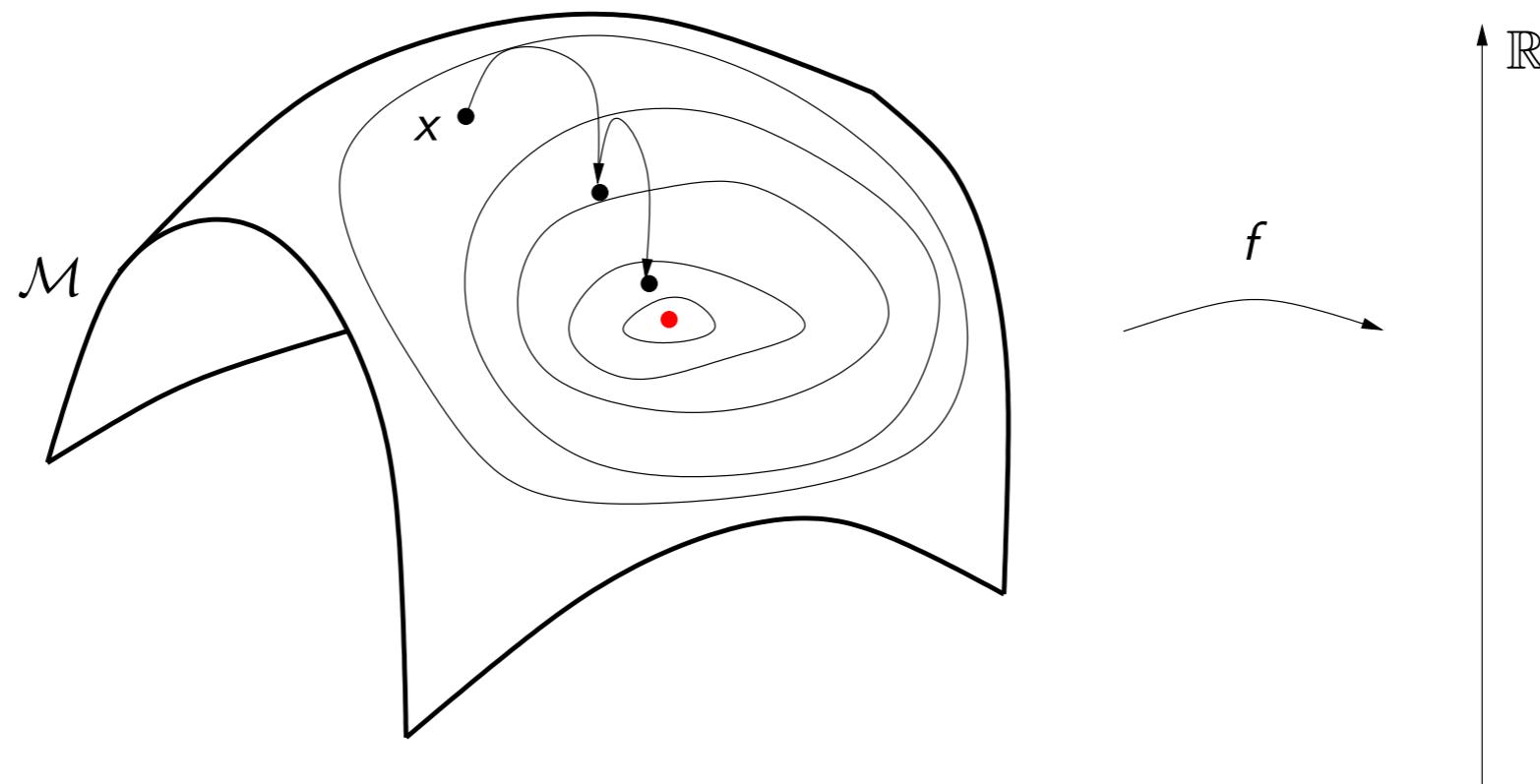
# Optimization on Manifolds

# Optimization on manifolds in one picture



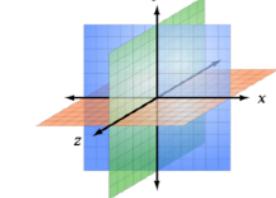
Slide Credit: Pierre-Antoine Absil

# Optimization on manifolds in one picture



Slide Credit: Pierre-Antoine Absil

# Example: Riemannian optimization



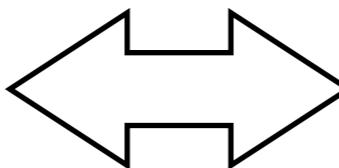
**Vector space  
optimization**

Orthogonality  
constraint

Fixed-rank  
constraint

Positive  
semi-definite  
constraint

....



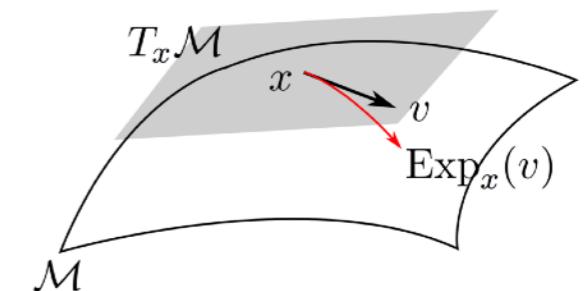
Stiefel  
manifold

Grassmann  
manifold

PSD  
manifold

....

**Riemannian  
optimization**

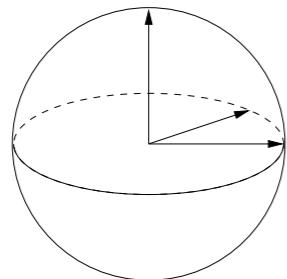


[Udriste, 1994; Absil et al., 2009]

# Example: Riemannian optimization

Given  $A = A^T \in \mathbb{R}^{n \times n}$   
and  $N = \text{diag}(1, \dots, p)$ ,

$$\begin{aligned}\min f(X) &= \text{trace}(X^T A X N) \\ \text{subj. to } X &\in \mathbb{R}^{n \times p} : X^T X = I\end{aligned}$$

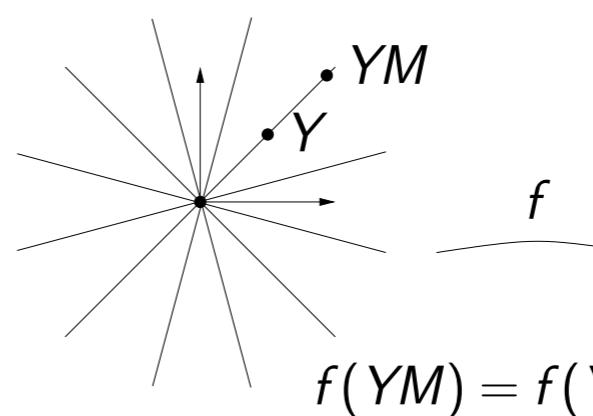


Domain:  $\text{St}(p, n)$   
 $= \{X \in \mathbb{R}^{n \times p} : X^T X = I\}$

Embedded submanifold

Given  $A = A^T \in \mathbb{R}^{n \times n}$ ,

$$\begin{aligned}\min f(Y) &= \text{trace}((Y^T Y)^{-1} (Y^T A Y)) \\ \text{subj. to } Y &\in \mathbb{R}_*^{n \times p} \text{ (i.e., } Y \text{ full rank)}\end{aligned}$$

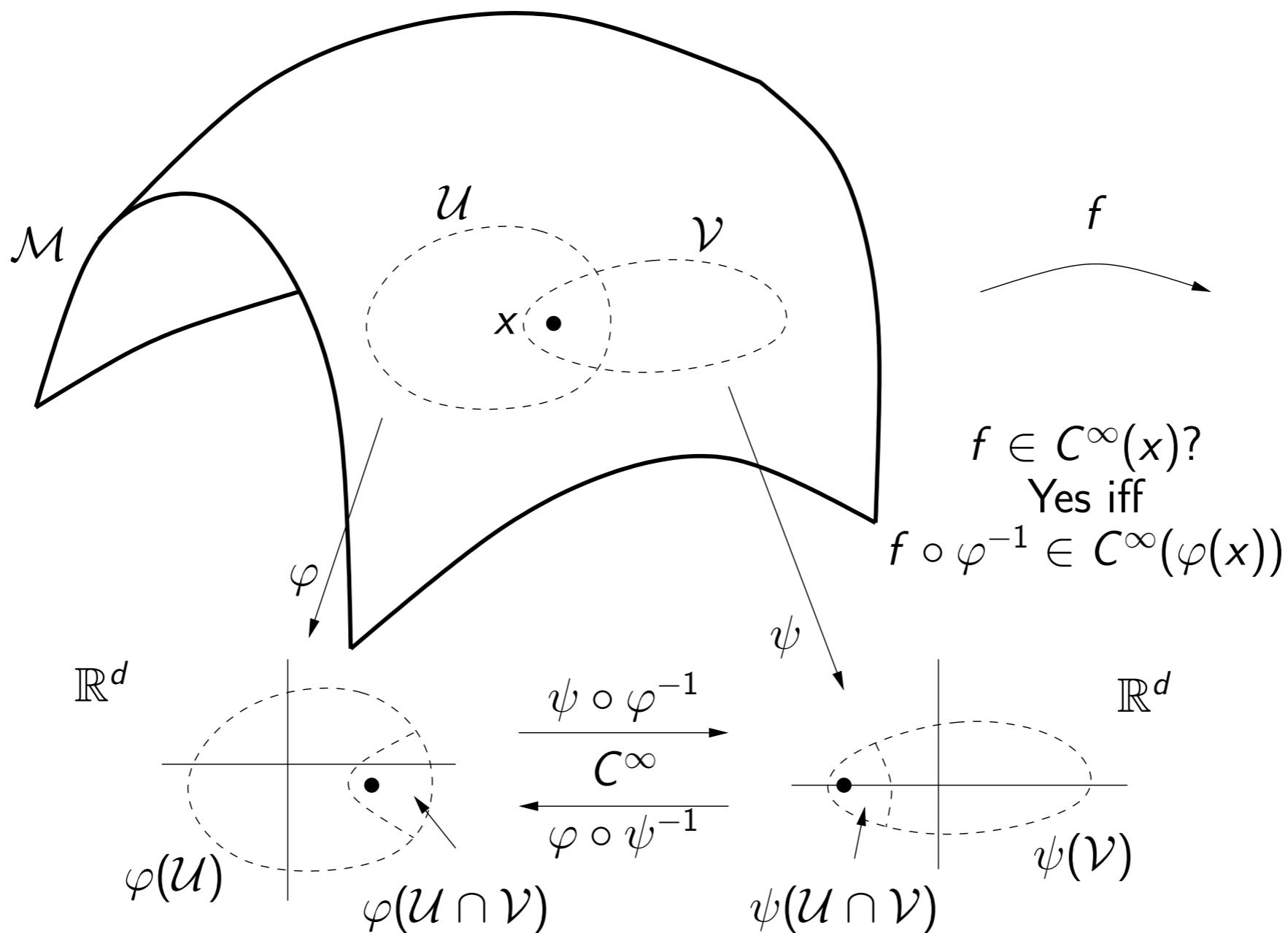


Domain:  $\text{Gr}(p, n)$   
 $= \left\{ \{YM : M \in \mathbb{R}_*^{p \times p}\} : Y \in \mathbb{R}_*^{n \times p} \right\}$

Quotient manifold

Slide Credit: Pierre-Antoine Absil

# Smooth manifolds and smooth functions



Slide Credit: Pierre-Antoine Absil

# Algorithm on Riemannian manifolds

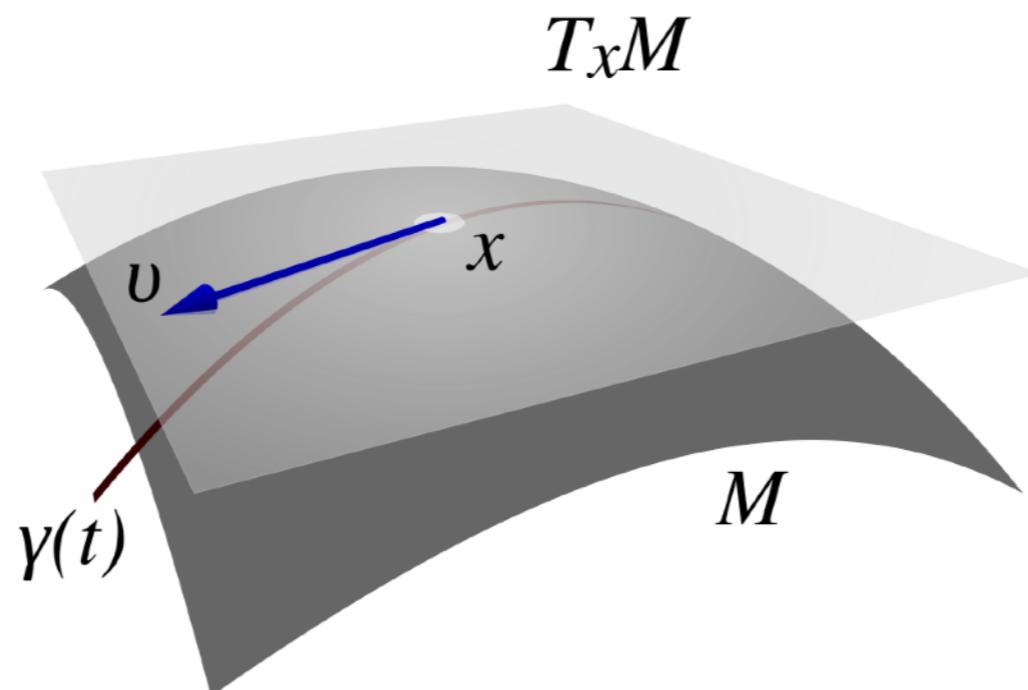
---

- **Steepest Descent** – Needs Retraction
- **Newton** – Needs Affine Connection and Retraction
- **Conjugate Gradient** – Needs Retraction and Vector Transport
- **BFGS** – Needs Retraction and Vector Transport
- **Trust Region** – Needs Retraction

# Tangent space

---

- Doing Calculus needs **differential structure**
- This can be obtained by defining a tangent bundle, that is, the set of tangent spaces for all points



# Riemannian metric

---

- Define an inner product called **Riemannian metric**

$$g_x(\xi, \eta) : T_x M \times T_x M \rightarrow \mathbb{R}$$

# Riemannian metric

---

- Define an inner product called **Riemannian metric**

$$g_x(\xi, \eta) : T_x M \times T_x M \rightarrow \mathbb{R}$$

- This metric induces a norm

# Riemannian metric

---

- Define an inner product called **Riemannian metric**

$$g_x(\xi, \eta) : T_x M \times T_x M \rightarrow \mathbb{R}$$

- This metric induces a norm

$$\|\xi\|_x = \sqrt{g_x(\xi, \xi)}$$

# Riemannian metric

---

- Define an inner product called **Riemannian metric**

$$g_x(\xi, \eta) : T_x M \times T_x M \rightarrow \mathbb{R}$$

- This metric induces a norm

$$\|\xi\|_x = \sqrt{g_x(\xi, \xi)}$$

- This can be used to measure the **length** of a curve

# Riemannian metric

---

- Define an inner product called **Riemannian metric**

$$g_x(\xi, \eta) : T_x M \times T_x M \rightarrow \mathbb{R}$$

- This metric induces a norm

$$\|\xi\|_x = \sqrt{g_x(\xi, \xi)}$$

- This can be used to measure the **length** of a curve

$$L(x(t)) = \int_a^b \sqrt{\|\dot{x}(t)\|}$$

# Riemannian metric

---

- Define an inner product called **Riemannian metric**

$$g_x(\xi, \eta) : T_x M \times T_x M \rightarrow \mathbb{R}$$

- This metric induces a norm

$$\|\xi\|_x = \sqrt{g_x(\xi, \xi)}$$

- This can be used to measure the **length** of a curve

$$L(x(t)) = \int_a^b \sqrt{\|\dot{x}(t)\|}$$

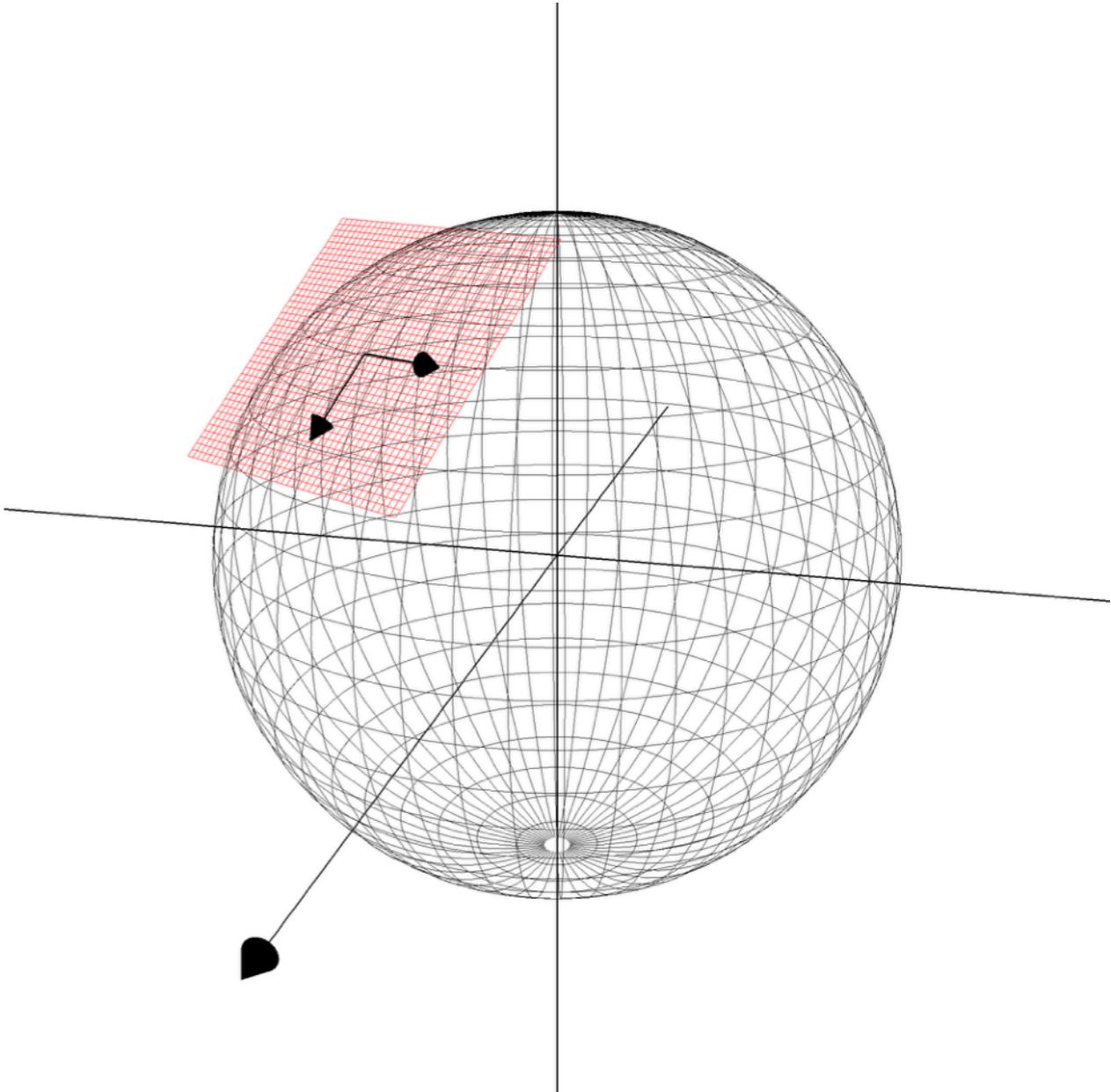
- **Geodesics** are locally minimizing curves between two points

# Example: n-sphere

---

- The tangent space is given by

$$T_x M = \{v \in \mathbb{R}^n | v^T x = 0\}$$

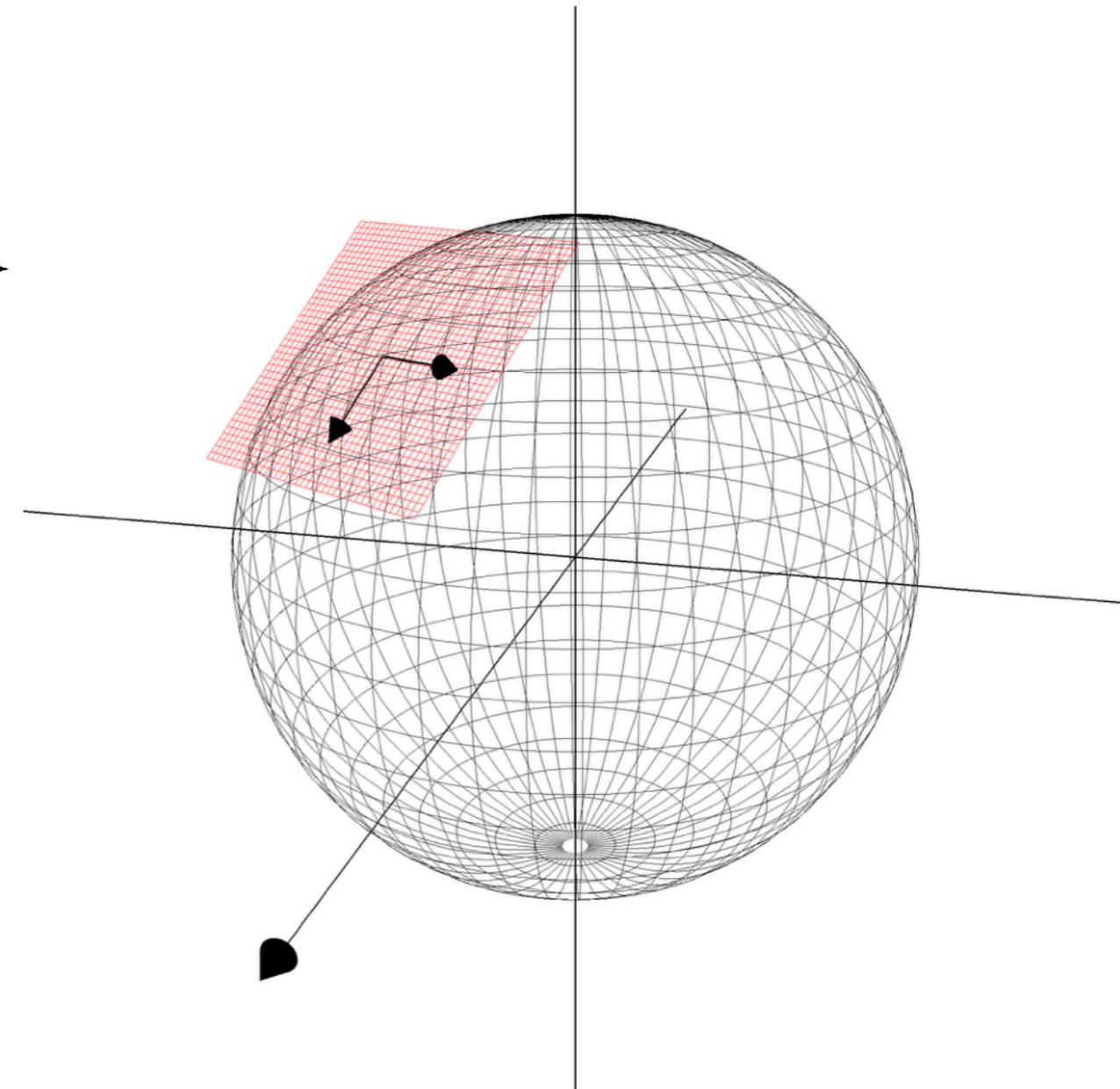


# Example: n-sphere

- The tangent space is given by

$$T_x M = \{v \in \mathbb{R}^n | v^T x = 0\}$$

- The inner product is inherited by the embedding Euclidean space and is **standard inner product**

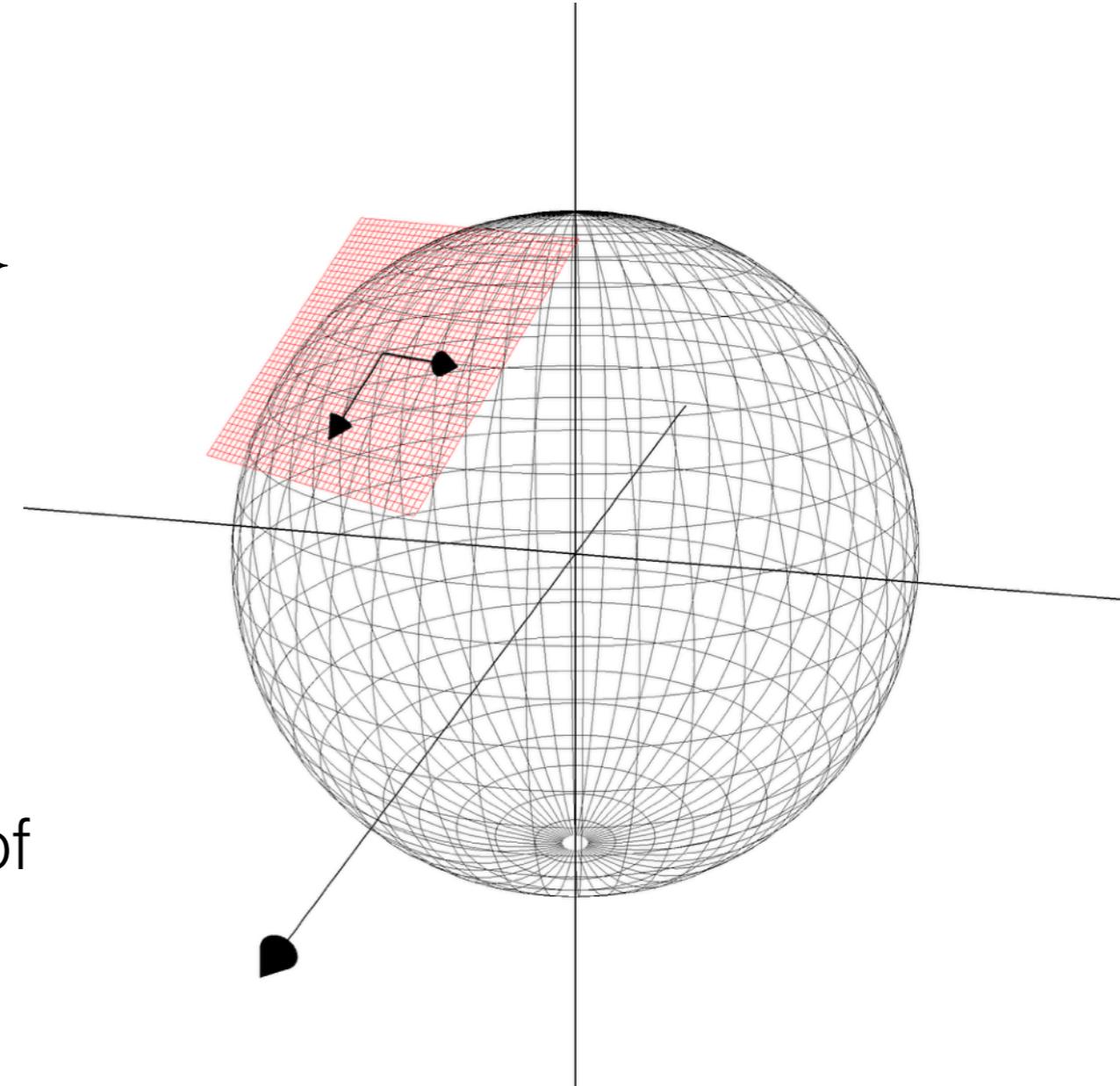


# Example: n-sphere

- The tangent space is given by

$$T_x M = \{v \in \mathbb{R}^n | v^T x = 0\}$$

- The inner product is inherited by the embedding Euclidean space and is **standard inner product**
- Geodesics are the **great circles** of the sphere



# Riemannian gradient

---

- **Riemannian gradient** is the direction of the steepest descent

# Riemannian gradient

---

- **Riemannian gradient** is the direction of the steepest descent

$$g_x(\nabla f, \xi) = D_\xi f(x)$$

# Riemannian gradient

---

- **Riemannian gradient** is the direction of the steepest descent

$$g_x(\nabla f, \xi) = D_\xi f(x)$$

- Where

# Riemannian gradient

---

- **Riemannian gradient** is the direction of the steepest descent

$$g_x(\nabla f, \xi) = D_\xi f(x)$$

- Where

$$D_\xi f(x) = \frac{d}{dt} f(\delta(t)), \quad \dot{\delta}(t) = \xi$$

# Riemannian gradient descent

---

- Consider Euclidean gradient descent

$$x_{t+1} = x_t - \lambda \nabla f(x_t)$$

# Riemannian gradient descent

---

- Consider Euclidean gradient descent

$$x_{t+1} = x_t - \lambda \nabla f(x_t)$$

- For arbitrary manifold this does not make sense

# Riemannian gradient descent

---

- Consider Euclidean gradient descent

$$x_{t+1} = x_t - \lambda \nabla f(x_t)$$

- For arbitrary manifold this does not make sense
- Moving along straight line is replaced by moving along geodesics that is denoted by **Exponential map**

# Riemannian gradient descent

---

- Consider Euclidean gradient descent

$$x_{t+1} = x_t - \lambda \nabla f(x_t)$$

- For arbitrary manifold this does not make sense
- Moving along straight line is replaced by moving along geodesics that is denoted by **Exponential map**

$$x_{t+1} = \text{Exp}_{x_t}(-\lambda \nabla f(x_t))$$

# Riemannian gradient descent

---

- Consider Euclidean gradient descent

$$x_{t+1} = x_t - \lambda \nabla f(x_t)$$

- For arbitrary manifold this does not make sense
- Moving along straight line is replaced by moving along geodesics that is denoted by **Exponential map**

$$x_{t+1} = \text{Exp}_{x_t}(-\lambda \nabla f(x_t))$$

- Exponential map is moving along geodesic  $\gamma(t)$

# Riemannian gradient descent

---

- Consider Euclidean gradient descent

$$x_{t+1} = x_t - \lambda \nabla f(x_t)$$

- For arbitrary manifold this does not make sense
- Moving along straight line is replaced by moving along geodesics that is denoted by **Exponential map**

$$x_{t+1} = \text{Exp}_{x_t}(-\lambda \nabla f(x_t))$$

- Exponential map is moving along geodesic  $\gamma(t)$

$$\gamma(0) = x_t, \quad \dot{\gamma}(0) = \nabla f(x_t)$$

# Retraction map

---

- **Exponential maps** can be **hard to compute**. We can replace that with **retraction**

# Retraction map

---

- **Exponential maps** can be **hard to compute**. We can replace that with **retraction**
- Retraction is a mapping from the tangent vector to a point in manifold.

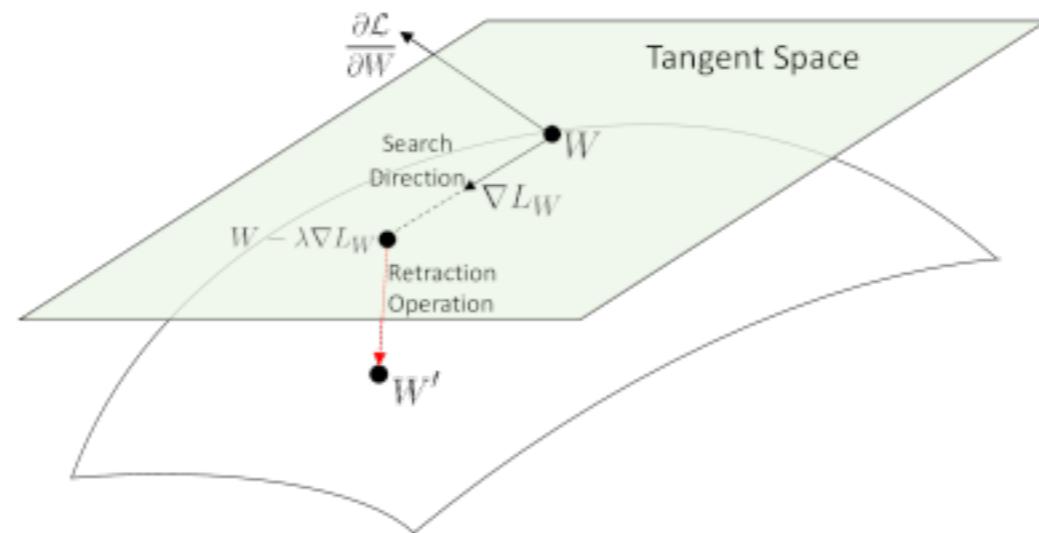
# Retraction map

---

- **Exponential maps** can be **hard to compute**. We can replace that with **retraction**
- Retraction is a mapping from the tangent vector to a point in manifold.
- It is like moving along a curve with velocity equal to input tangent vector

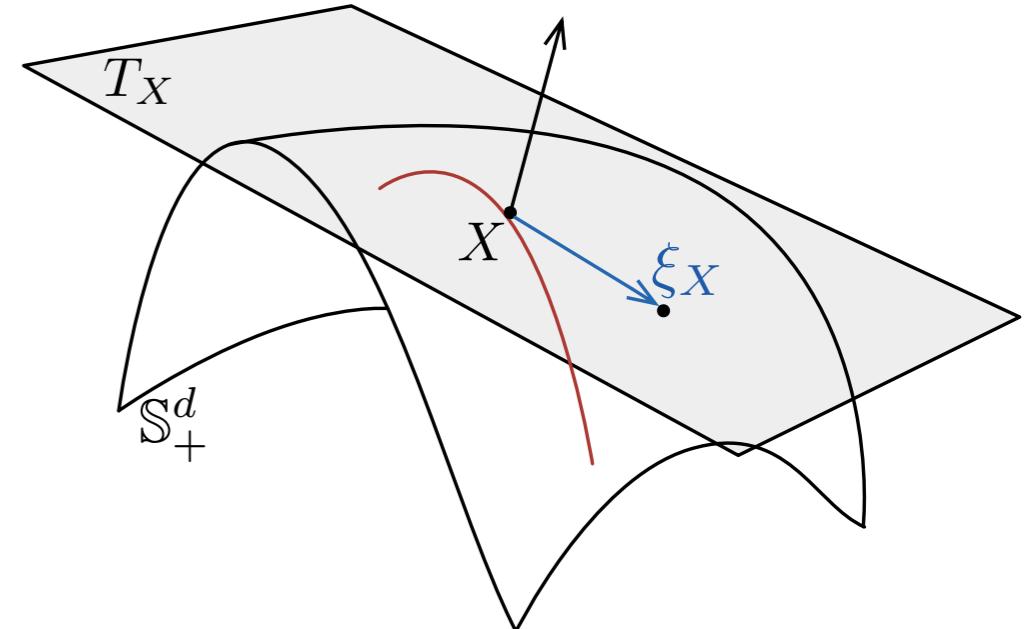
# Retraction map

- **Exponential maps** can be **hard to compute**. We can replace that with **retraction**
- Retraction is a mapping from the tangent vector to a point in manifold.
- It is like moving along a curve with velocity equal to input tangent vector



# Riemannian L-BFGS

(own version)



Retraction

## Strong-Wolfe line-search

initial  
stepsize

$$f(R_{X_k}(\alpha \xi_k)) \leq f(X_k) + c_1 \alpha Df(X_k) \xi_k,$$
$$|Df(X_{k+1}) \xi_{k+1}| \leq c_2 Df(X_k) \xi_k$$

$$\alpha^* = 2 \frac{f(X_k) - f(X_{k-1})}{Df(X_{k-1}) \xi_{k-1}}.$$



[github.com/utvisionlab/mixest](https://github.com/utvisionlab/mixest)

# Riemannian Newton Method

---

- For finding Newton direction solve

$$Hf(x)\xi = -\nabla f(x)$$

# Riemannian Newton Method

---

- For finding Newton direction solve

$$Hf(x)\xi = -\nabla f(x)$$

- Where

# Riemannian Newton Method

---

- For finding Newton direction solve

$$Hf(x)\xi = -\nabla f(x)$$

- Where

$$Hf(x)\xi = \partial_\xi \nabla f(x), \partial \text{ is affine connection}$$

# Riemannian Newton Method

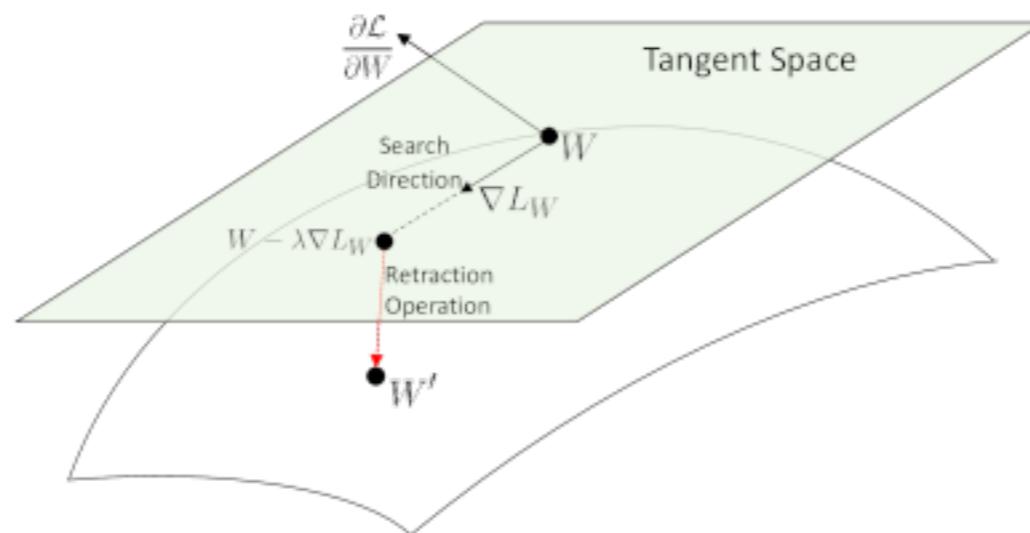
---

- For finding Newton direction solve

$$Hf(x)\xi = -\nabla f(x)$$

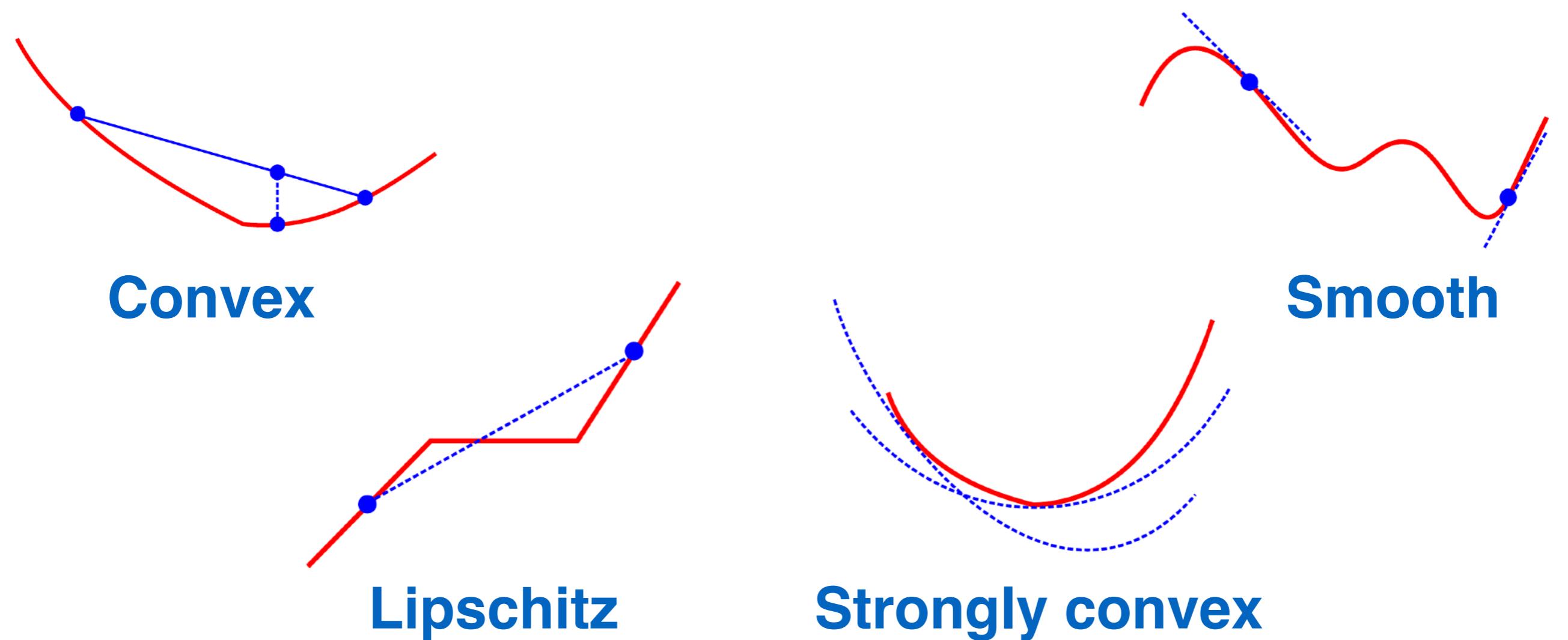
- Where

$$Hf(x)\xi = \partial_\xi \nabla f(x), \partial \text{ is affine connection}$$

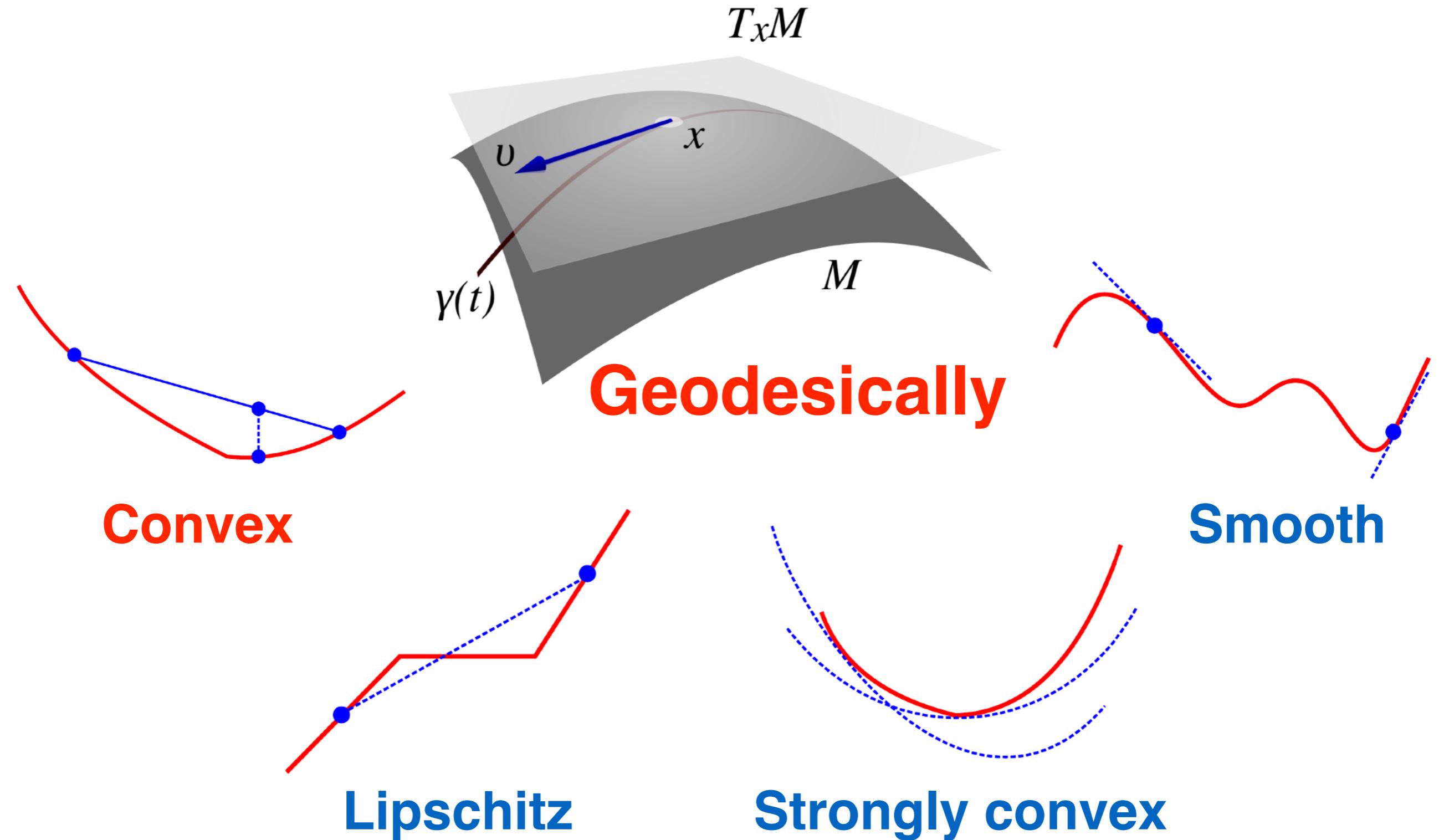


# Classes of function in optimization

---

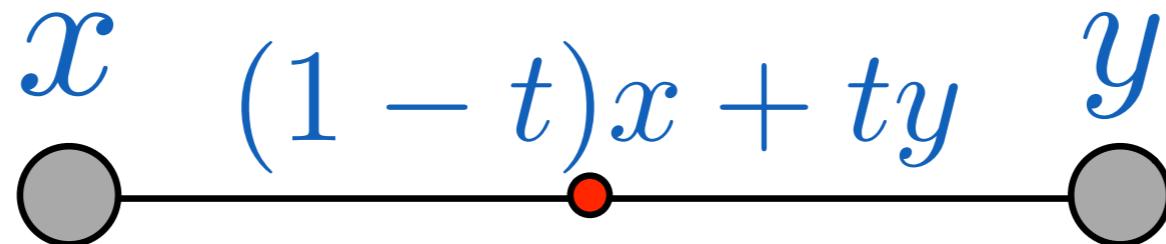


# Classes of function in optimization

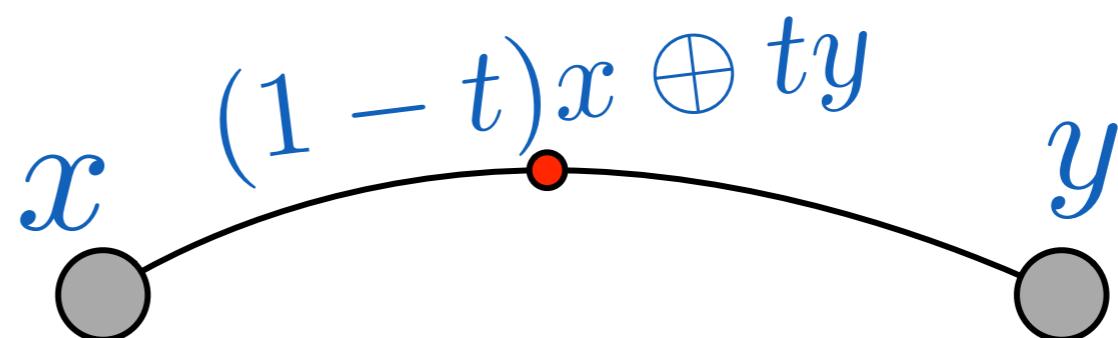


# What is geodesic convexity?

**Convexity**



**Geodesic convexity**



$$f((1 - t)x \oplus ty) \leq (1 - t)f(x) + tf(y)$$

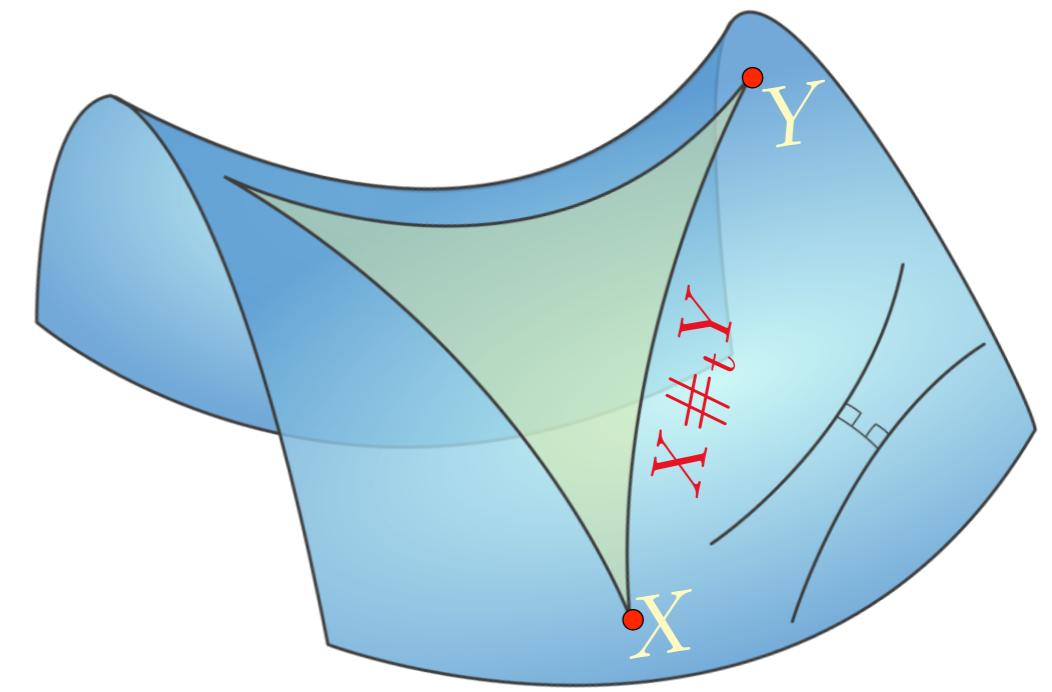
*on a Riemannian manifold*  $f(y) \geq f(x) + \langle g_x, \text{Exp}_x^{-1}(y) \rangle_x$

Metric spaces & curvature: [Menger; Alexandrov; Busemann; Bridson, Häflinger; Gromov; Perelman]

# Example: Positive definite matrices

*Geodesic*

$$\begin{aligned} X \#_t Y &:= X^{\frac{1}{2}} (X^{-\frac{1}{2}} Y X^{-\frac{1}{2}})^t X^{\frac{1}{2}} \\ &= (1 - t)X \oplus tY \end{aligned}$$



*Examples*

$$f(X) = \begin{cases} \log \det(X), & \log \text{tr}(X), \\ \text{tr}(X^\alpha), & \|X^\alpha\|. \quad (\alpha > 0) \end{cases}$$

*Exercise: verify for above examples*

$$f(X \#_t Y) \leq (1 - t)f(X) + tf(Y)$$

# PSD manifold and g-convexity

## Recognizing, constructing, and optimizing g-convex functions



[Sra, Hosseini (2013,2015)]

- [Wiesel 2012]
- [Rápcsák 1984]
- [Udriste 1994]

## Corollaries

$$X \mapsto \log \det(B + \sum_i A_i^* X A_i)$$

$$X \mapsto \log \text{per}(B + \sum_i A_i^* X A_i)$$

$$\delta_R^2(X, Y), \quad \delta_S^2(X, Y)$$

(jointly g-convex)

Many more theorems and corollaries

One-D version known as: **Geometric Programming**

[www.stanford.edu/~boyd/papers/gp\\_tutorial.html](http://www.stanford.edu/~boyd/papers/gp_tutorial.html)  
[Boyd, Kim, Vandenberghe, Hassibi (2007). 61pp.]

# PSD manifold and g-convexity

## Recognizing, constructing, and optimizing g-convex functions



[Sra, Hosseini (2013,2015)]

- [Wiesel 2012]
- [Rápcsák 1984]
- [Udriste 1994]

## Corollaries

$$X \mapsto \log \det(B + \sum_i A_i^* X A_i)$$

$$X \mapsto \log \text{per}(B + \sum_i A_i^* X A_i)$$

$$\delta_R^2(X, Y), \quad \delta_S^2(X, Y)$$

(jointly g-convex)

Many more theorems and corollaries

One-D version known as: **Geometric Programming**

[www.stanford.edu/~boyd/papers/gp\\_tutorial.html](http://www.stanford.edu/~boyd/papers/gp_tutorial.html)  
[Boyd, Kim, Vandenberghe, Hassibi (2007). 61pp.]

# PSD manifold and g-convexity

## Recognizing, constructing, and optimizing g-convex functions



[Sra, Hosseini (2013,2015)]

- [Wiesel 2012]
- [Rápcsák 1984]
- [Udriste 1994]

## Corollaries

$$X \mapsto \log \det(B + \sum_i A_i^* X A_i)$$

$$X \mapsto \log \text{per}(B + \sum_i A_i^* X A_i)$$

$$\delta_R^2(X, Y), \quad \delta_S^2(X, Y)$$

(jointly g-convex)

Many more theorems and corollaries

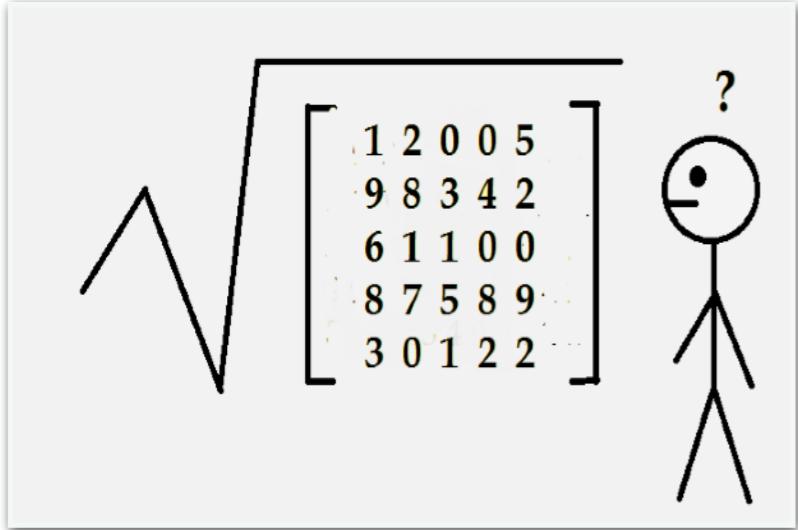
One-D version known as: **Geometric Programming**

[www.stanford.edu/~boyd/papers/gp\\_tutorial.html](http://www.stanford.edu/~boyd/papers/gp_tutorial.html)  
[Boyd, Kim, Vandenberghe, Hassibi (2007). 61pp.]

# Examples

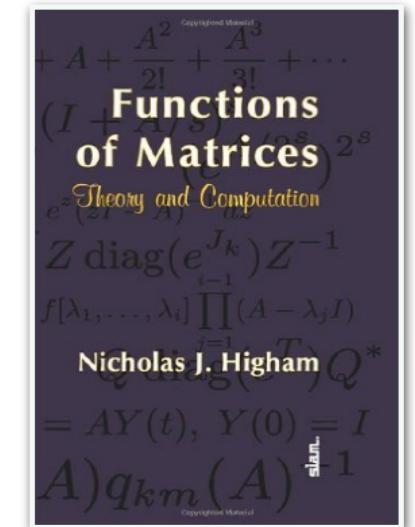
$$X \succ 0$$

# Matrix square root



Broadly applicable

Key to ‘logm’



# Matrix square root

---



Nonconvex optimization through the Euclidean lens

*[Jain, Jin, Kakade, Netrapalli; Jul 2015]*

$$\min_{X \in \mathbb{R}^{n \times n}} \|M - X^2\|_F^2$$

# Matrix square root

---



Nonconvex optimization through the Euclidean lens

[Jain, Jin, Kakade, Netrapalli; Jul 2015]

$$\min_{X \in \mathbb{R}^{n \times n}} \|M - X^2\|_F^2$$

## Gradient descent

$$X_{t+1} \leftarrow X_t - \eta(X_t^2 - M)X_t - \eta X_t(X_t^2 - M)$$

Simple algorithm; linear convergence; **nontrivial** analysis

# Matrix square root

---

*Geodesic*

$$X \#_t Y := X^{\frac{1}{2}} (X^{-\frac{1}{2}} Y X^{-\frac{1}{2}})^t X^{\frac{1}{2}}$$

*Midpoint*

$$A^{\frac{1}{2}} = A \#_{\frac{1}{2}} I$$

# Matrix square root

---



Nonconvex optimization through **non-Euclidean** lens

[Sra; Jul 2015]

$$\min_{X \succ 0} \quad \delta_S^2(X, A) + \delta_S^2(X, I)$$

$$\delta_S^2(X, Y) := \frac{1}{2} \log \det \left( \frac{X+Y}{2} \right) - \frac{1}{2} \log \det(XY)$$

(Curiously, this is the Jensen-Shannon divergence between multivariate Gaussians!)

# Matrix square root



Nonconvex optimization through **non-Euclidean** lens

[Sra; Jul 2015]

$$\min_{X \succ 0} \quad \delta_S^2(X, A) + \delta_S^2(X, I)$$

## Fixed-point iteration

$$X_{k+1} \leftarrow [(X_k + A)^{-1} + (X_k + I)^{-1}]^{-1}$$

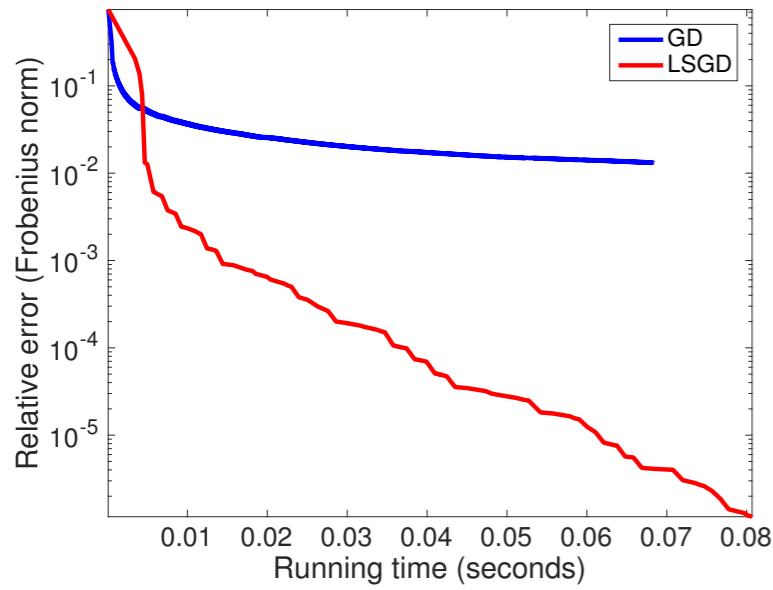
Simple method; linear convergence; 1/2 page analysis!

**Global optimality thanks to geodesic convexity**

$$\delta_S^2(X, Y) := \frac{1}{2} \log \det \left( \frac{X+Y}{2} \right) - \frac{1}{2} \log \det(XY)$$

(Curiously, this is the Jensen-Shannon divergence between multivariate Gaussians!)

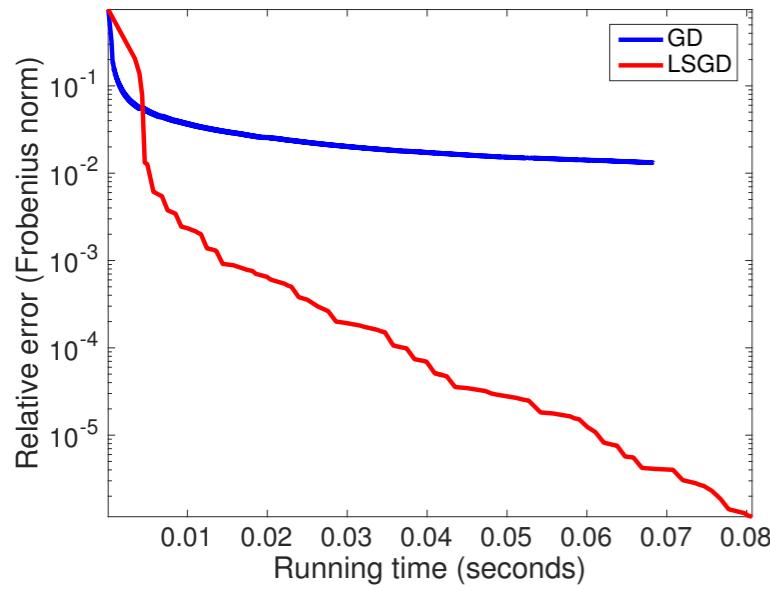
# Matrix square root



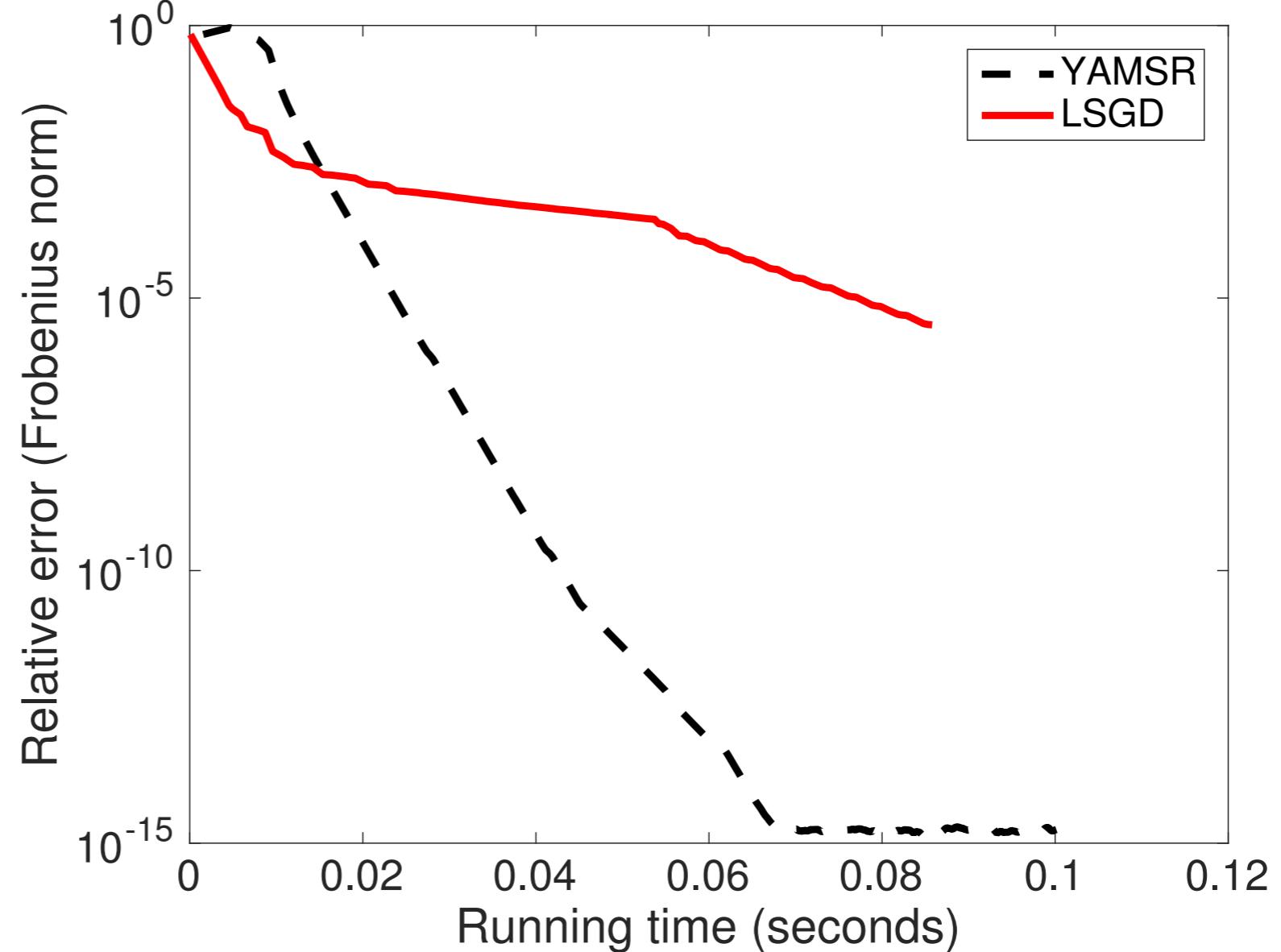
$50 \times 50$  matrix  $I + \beta UU^T$

$$\kappa \approx 64$$

# Matrix square root



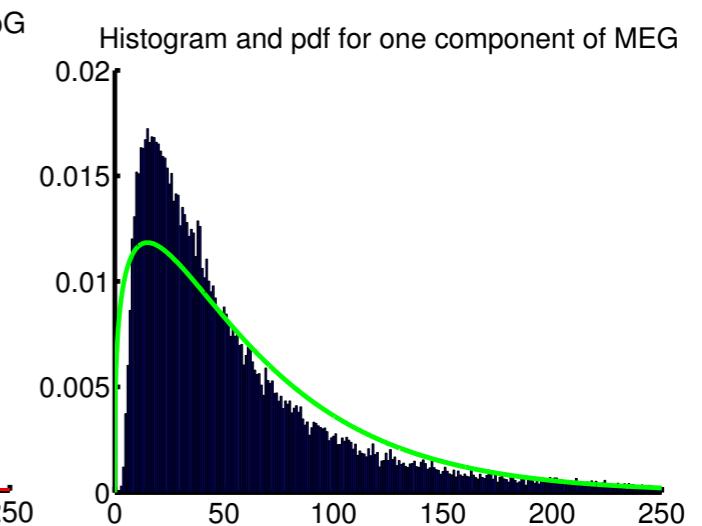
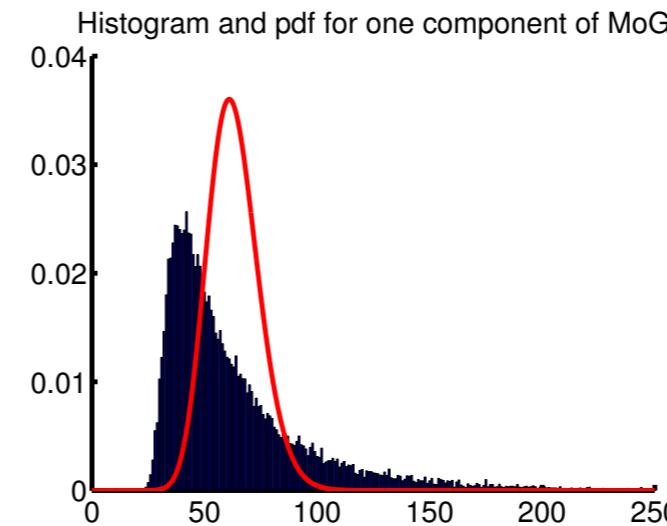
$50 \times 50$  matrix  $I + \beta UU^T$   
 $\kappa \approx 64$



# M-estimators: Elliptical Distributions

## Natural Image Statistics

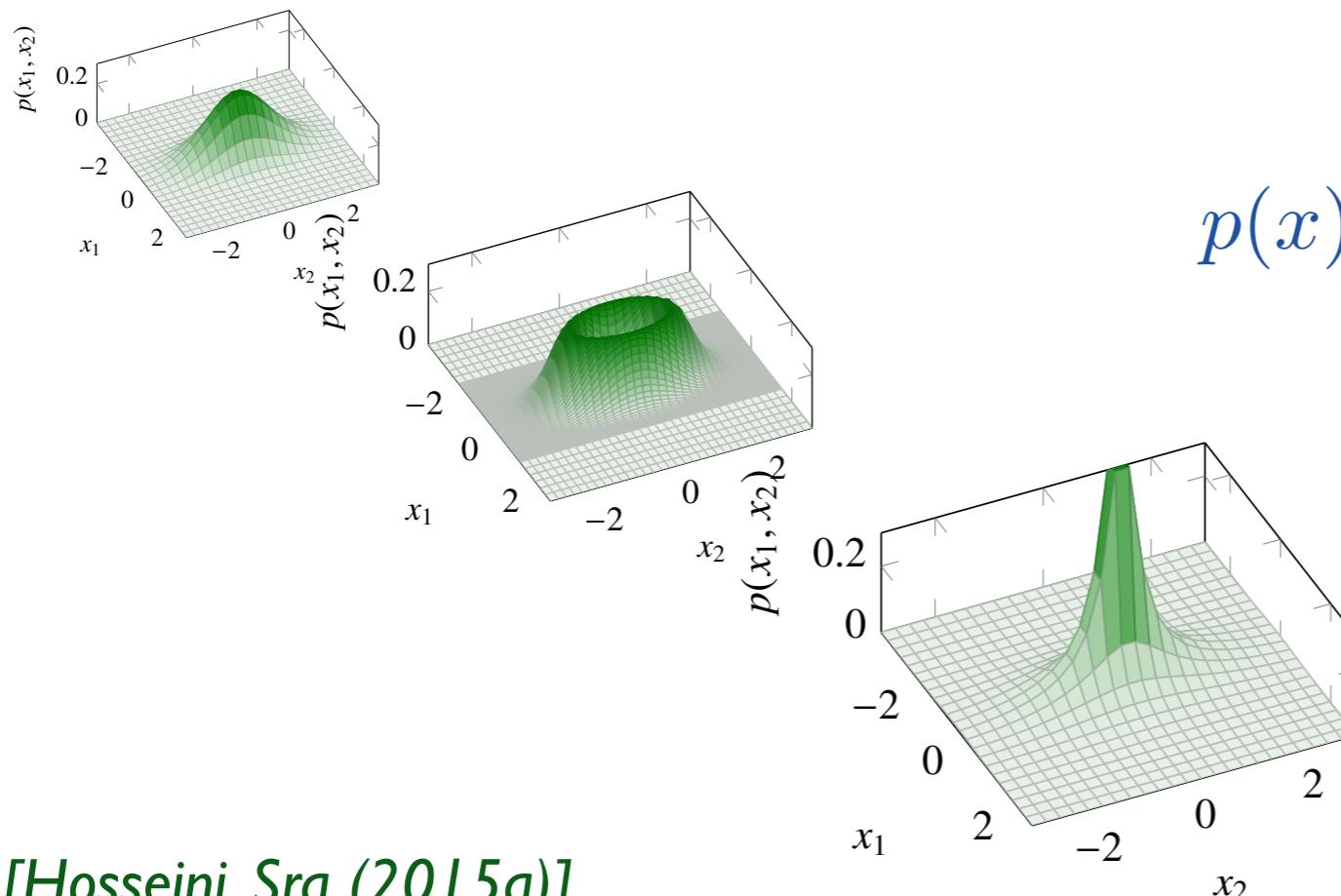
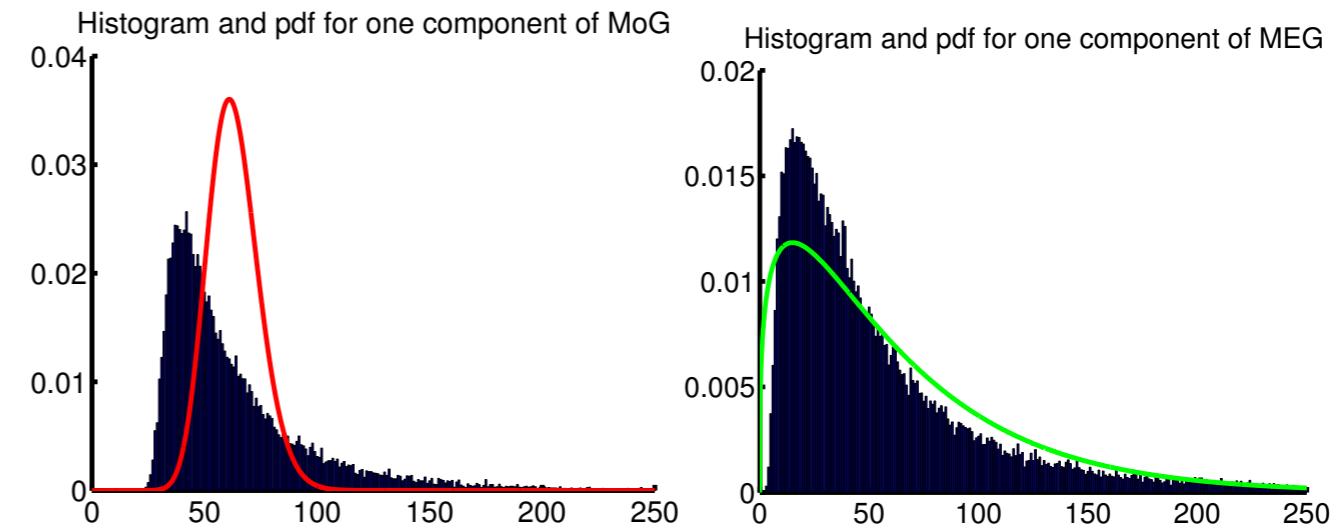
- ▶ Extract 200,000 training patches from 4167 images
- ▶ 10 sets of 100,000 test patches
- ▶ Log-transform intensities; add small amount of white noise



# M-estimators: Elliptical Distributions

## Natural Image Statistics

- ▶ Extract 200,000 training patches from 4167 images
- ▶ 10 sets of 100,000 test patches
- ▶ Log-transform intensities; add small amount of white noise



$$p(x) \propto \frac{(x^T \Sigma^{-1} x)^{a - \frac{d}{2}} e^{-\frac{a}{d} x^T \Sigma^{-1} x}}{\det(\Sigma)^{1/2}}$$

**Elliptically Contoured  
Distributions (ECD)**

[Hosseini, Sra (2015a)]

# M-estimators for ECDs

---

Given observations  $x_1, x_2, \dots, x_n$  find m.l.e. by solving

$$\begin{aligned} & \frac{n}{2} \log \det(\Sigma) - \left(a - \frac{d}{2}\right) \sum_{i=1}^n \log(x_i^T \Sigma^{-1} x_i) \\ & + \frac{a}{d} \text{trace}(\Sigma^{-1} \sum_i x_i x_i^T) \end{aligned}$$

---

[Hosseini, Sra (2015)]

# M-estimators for ECDs

Given observations  $x_1, x_2, \dots, x_n$  find m.l.e. by solving

$$\begin{aligned} & \frac{n}{2} \log \det(\Sigma) - \left(a - \frac{d}{2}\right) \sum_{i=1}^n \log(x_i^T \Sigma^{-1} x_i) \\ & + \frac{a}{d} \text{trace}(\Sigma^{-1} \sum_i x_i x_i^T) \end{aligned}$$

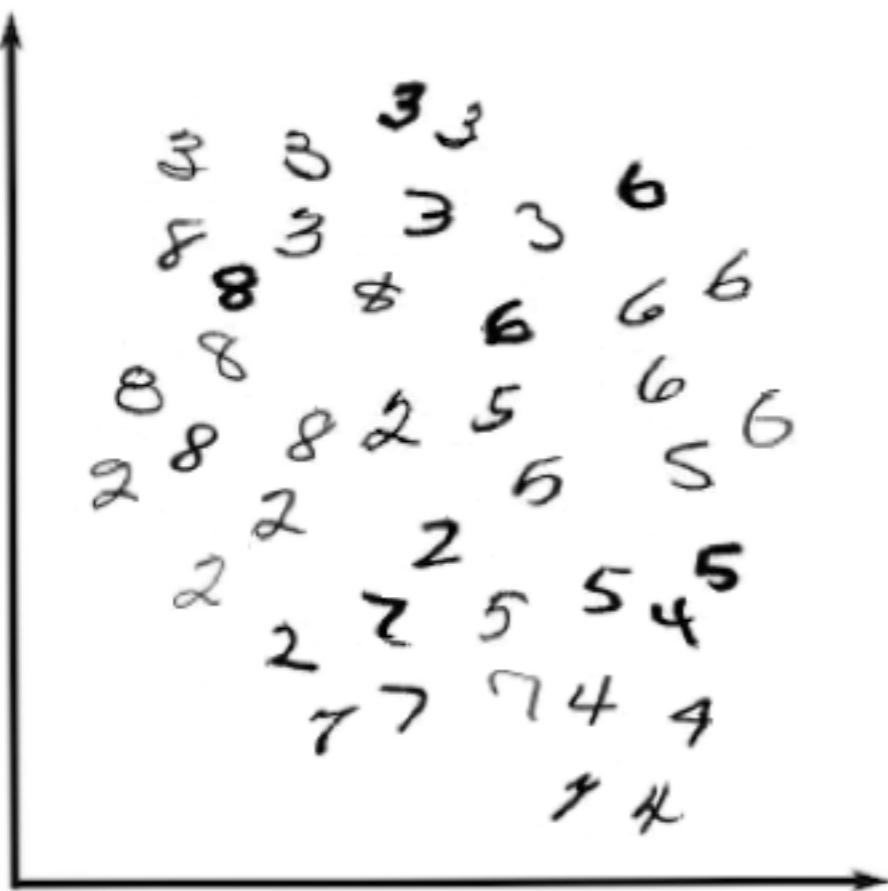
convex or nonconvex, but  
globally solvable due to g-convexity!

[Hosseini, Sra (2015)]

# Metric learning

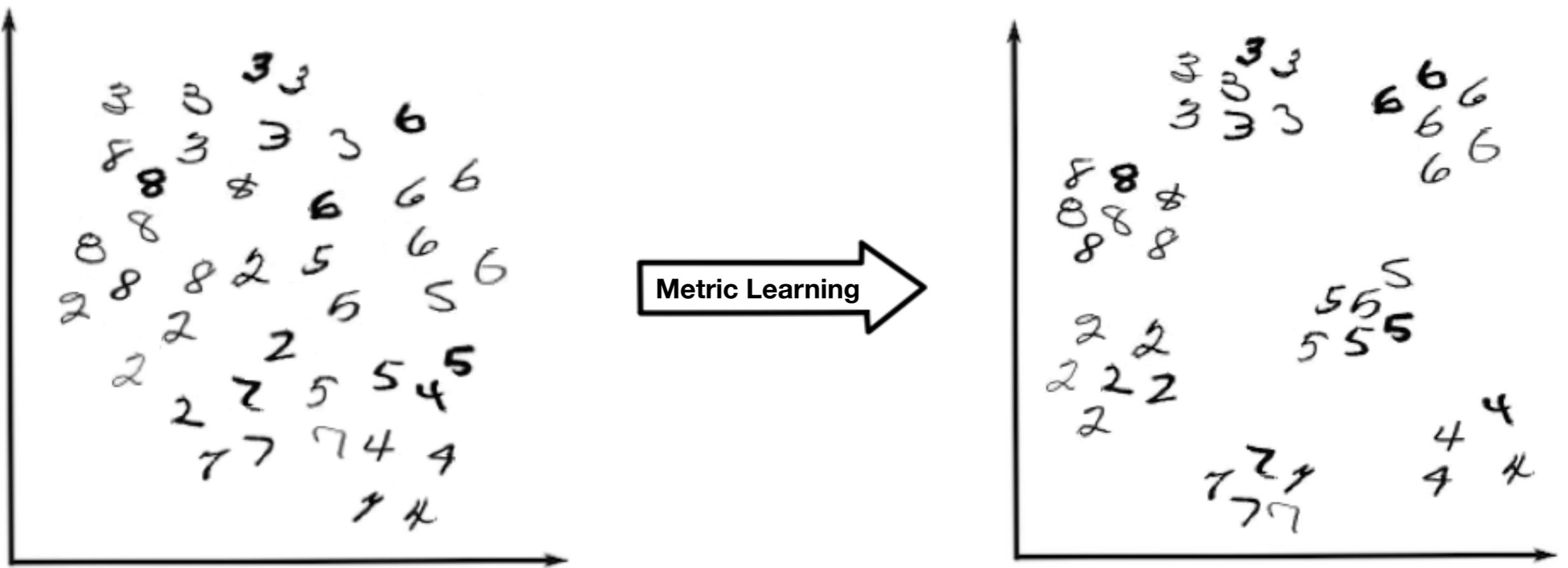
---

What does a metric learning method do?



# Metric learning

What does a metric learning method do?



# Euclidean metric learning

---

## *Pairwise constraints*

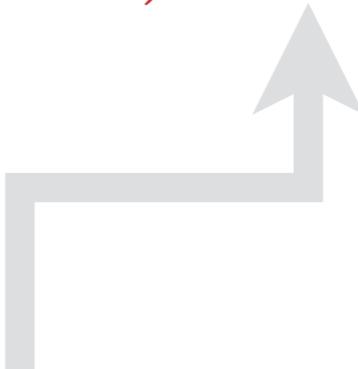
$$\mathcal{S} := \{(x_i, x_j) \mid x_i \text{ and } x_j \text{ are in the same class}\}$$

$$\mathcal{D} := \{(x_i, x_j) \mid x_i \text{ and } x_j \text{ are in different classes}\}$$

## *Goal*

*given pairwise constraints learn Mahalanobis distance*

$$d_A(x, y) := (x - y)^T A (x - y)$$



*Positive definite matrix **A***

# Metric learning methods

---

**MMC**

[Xing, Jordan, Russell, Ng 2002]

**LMNN**

[Weinberger, Saul 2005]

**ITML**

[Davis, Kulis, Jain, Sra, Dhillon 2007]

# Metric learning methods

---

**MMC**

[Xing, Jordan, Russell, Ng 2002]

**LMNN**

[Weinberger, Saul 2005]

**ITML**

[Davis, Kulis, Jain, Sra, Dhillon 2007]

$$\min_{\mathbf{A} \succeq 0} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j)$$

such that  $\sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} \sqrt{d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j)} \geq 1$

$$\min_{\mathbf{A} \succeq 0} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} \left[ (1 - \mu) d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) + \mu \sum_l (1 - y_{il}) \xi_{ijl} \right]$$

$$d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_l) - d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) \geq 1 - \xi_{ijl}$$
$$\xi_{ijl} \geq 0$$

$$\min_{\mathbf{A} \succeq 0} D_{\text{ld}}(\mathbf{A}, \mathbf{A}_0)$$

such that  $d_{\mathbf{A}}(\mathbf{x}, \mathbf{y}) \leq u, \quad (\mathbf{x}, \mathbf{y}) \in \mathcal{S},$   
 $d_{\mathbf{A}}(\mathbf{x}, \mathbf{y}) \geq l, \quad (\mathbf{x}, \mathbf{y}) \in \mathcal{D}$

$$D_{\text{ld}}(\mathbf{A}, \mathbf{A}_0) := \text{tr}(\mathbf{A}\mathbf{A}_0^{-1}) - \log \det(\mathbf{A}\mathbf{A}_0^{-1}) - d$$

# Metric learning methods

**MMC**

[Xing, Jordan, Russell, Ng 2002]

**LMNN**

[Weinberger, Saul 2005]

**ITML**

[Davis, Kulis, Jain, Sra, Dhillon 2007]

*tons of other methods!*

$$\min_{\mathbf{A} \succeq 0} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j)$$

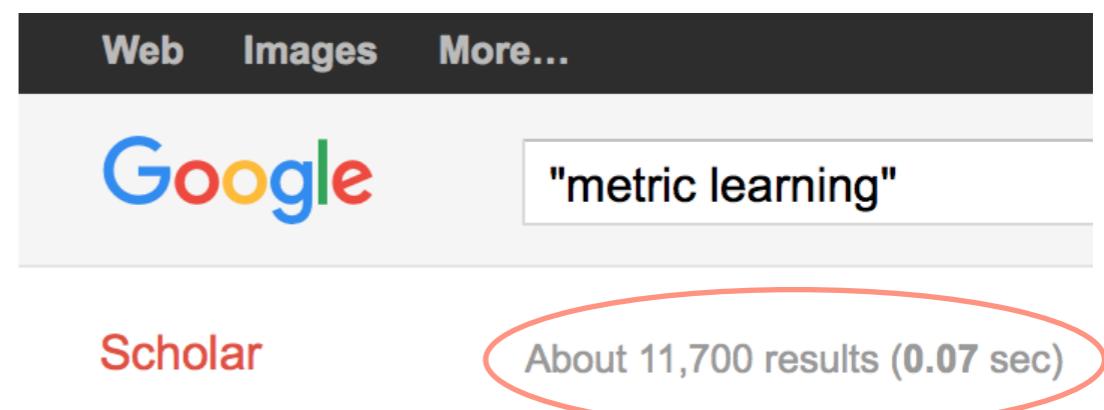
such that  $\sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} \sqrt{d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j)} \geq 1$

$$\min_{\mathbf{A} \succeq 0} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} \left[ (1 - \mu) d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) + \mu \sum_l (1 - y_{il}) \xi_{ijl} \right]$$
$$d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_l) - d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) \geq 1 - \xi_{ijl}$$
$$\xi_{ijl} \geq 0$$

$$\min_{\mathbf{A} \succeq 0} D_{\text{ld}}(\mathbf{A}, \mathbf{A}_0)$$

such that  $d_{\mathbf{A}}(\mathbf{x}, \mathbf{y}) \leq u, \quad (\mathbf{x}, \mathbf{y}) \in \mathcal{S},$   
 $d_{\mathbf{A}}(\mathbf{x}, \mathbf{y}) \geq l, \quad (\mathbf{x}, \mathbf{y}) \in \mathcal{D}$

$$D_{\text{ld}}(\mathbf{A}, \mathbf{A}_0) := \text{tr}(\mathbf{A}\mathbf{A}_0^{-1}) - \log \det(\mathbf{A}\mathbf{A}_0^{-1}) - d$$



# Related work

---

- FlatGeo *[Meyer, Bonnabel, Sepulchre '11]*

$$\min_{\mathbf{A} \succeq 0} \quad \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} \max(0, l - d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j))^2 + \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} \max(0, d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) - u)^2$$

- LMNN *[Weinberger, Saul '05]*; and 100s of other methods out there!

# Related work

---

- FlatGeo [Meyer, Bonnabel, Sepulchre '11]

$$\min_{\mathbf{A} \succeq 0} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} \max(0, l - d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j))^2 + \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} \max(0, d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) - u)^2$$

- LMNN [Weinberger, Saul '05]; and 100s of other methods out there!

- **Main concerns**

They does not scale well to the large problems, with respect to:

1. The number of constraints
2. The dimensionality

# A simple new way for metric learning

---

*Euclidean idea*

$$\min_{\mathbf{A} \succeq 0} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) - \lambda \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j)$$

# A simple new way for metric learning

---

*Euclidean idea*

$$\min_{\mathbf{A} \succeq 0} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) - \lambda \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j)$$

# A simple new way for metric learning

---

*Euclidean idea*

$$\min_{\mathbf{A} \succeq 0} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) - \lambda \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j)$$

*New idea*

$$\min_{\mathbf{A} \succeq 0} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) + \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} d_{\mathbf{A}^{-1}}(\mathbf{x}_i, \mathbf{x}_j)$$

# A simple new way for metric learning

---

*Euclidean idea*

$$\min_{\mathbf{A} \succeq 0} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) - \lambda \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j)$$

*New idea*

$$\min_{\mathbf{A} \succeq 0} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) + \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} d_{\mathbf{A}^{-1}}(\mathbf{x}_i, \mathbf{x}_j)$$

*Equivalently solve*

# A simple new way for metric learning

*Euclidean idea*

$$\min_{\mathbf{A} \succeq 0} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) - \lambda \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j)$$

*New idea*

$$\min_{\mathbf{A} \succeq 0} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) + \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} d_{\mathbf{A}^{-1}}(\mathbf{x}_i, \mathbf{x}_j)$$

*Equivalently solve*

$$\min_{\mathbf{A} \succ 0} h(\mathbf{A}) := \text{tr}(\mathbf{A}\mathbf{S}) + \text{tr}(\mathbf{A}^{-1}\mathbf{D})$$

$$\mathbf{S} := \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T,$$

$$\mathbf{D} := \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T$$

[Habibzadeh, Hosseini, Sra, ICML 2016]

# A simple new way for metric learning

---

**Closed form solution!**

$$\nabla h(A) = 0 \iff S - A^{-1}DA^{-1} = 0$$

# A simple new way for metric learning

Closed form solution!

$$X \#_t Y := X^{\frac{1}{2}} (X^{-\frac{1}{2}} Y X^{-\frac{1}{2}})^t X^{\frac{1}{2}}$$

$$\nabla h(A) = 0 \iff S - A^{-1} D A^{-1} = 0$$

$$A = S^{-1} \#_{\frac{1}{2}} D$$

# A simple new way for metric learning

Closed form solution!

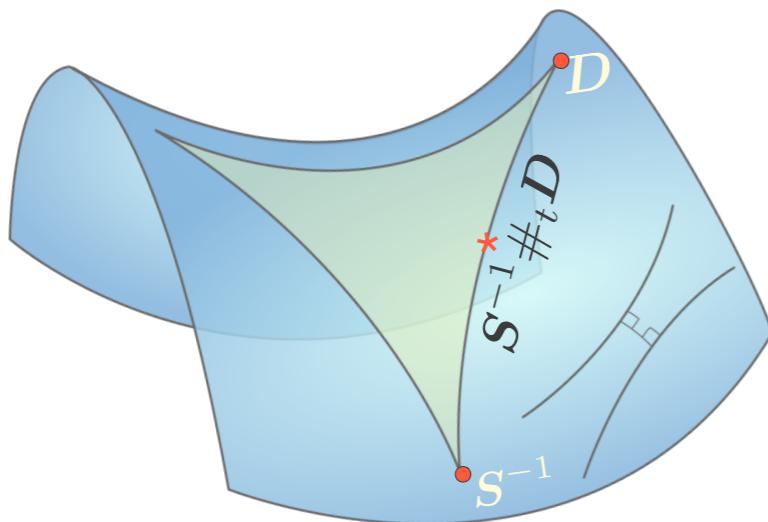
$$X \#_t Y := X^{\frac{1}{2}} (X^{-\frac{1}{2}} Y X^{-\frac{1}{2}})^t X^{\frac{1}{2}}$$

$$\nabla h(A) = 0 \iff S - A^{-1} D A^{-1} = 0$$

$$A = S^{-1} \#_{\frac{1}{2}} D$$

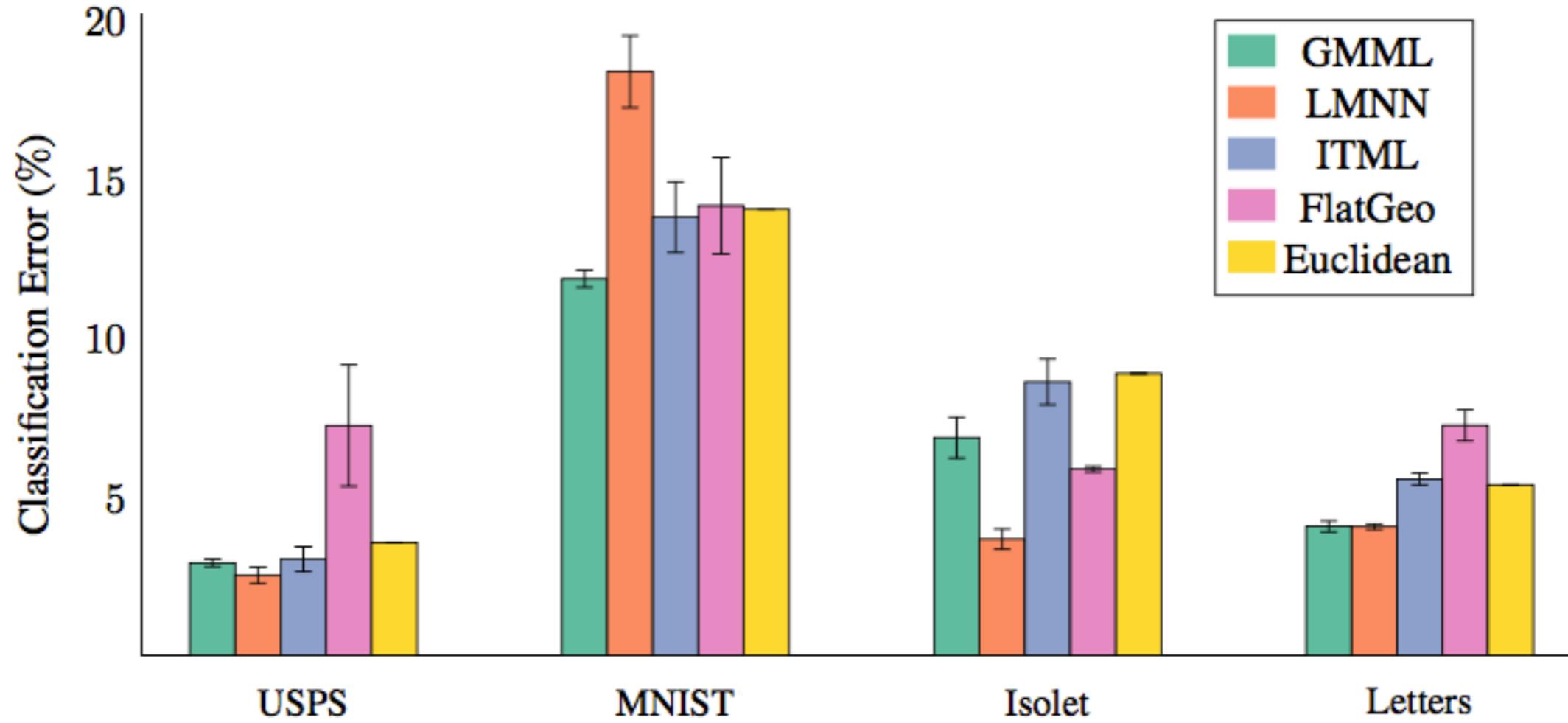
More generally

$$\min_{A \succ 0} (1-t)\delta_R^2(S^{-1}, A) + t\delta_R^2(D, A)$$



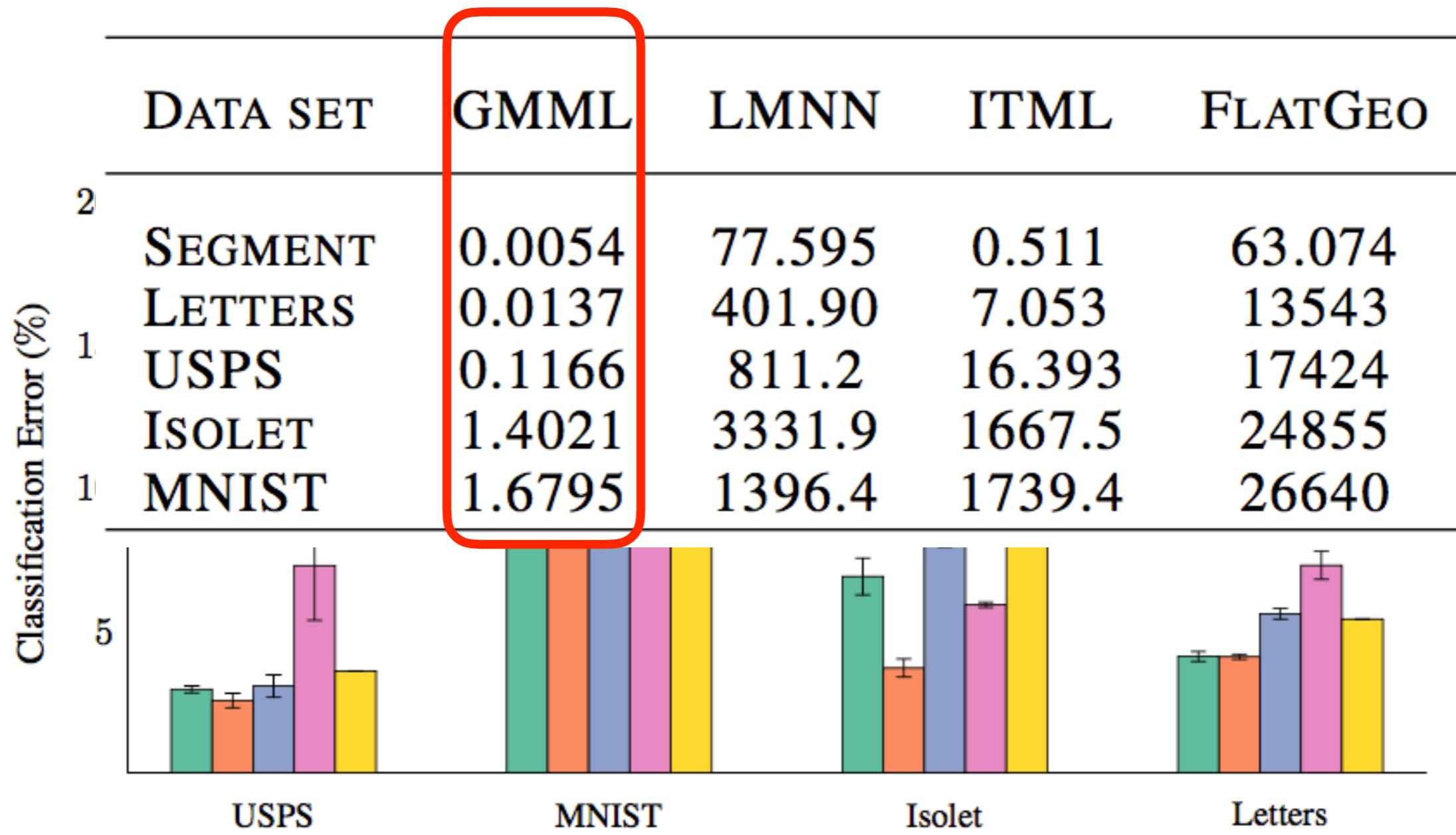
$$S^{-1} \#_t D$$

# Experiments



[Habibzadeh, Hosseini, Sra ICML 2016]

# Experiments



[Habibzadeh, Hosseini, Sra ICML 2016]

# Brascamp-Lieb Constant

---

# Brascamp-Lieb Constant

---

$$\int_{\mathbb{R}^n} \prod_{i=1}^m f_i(B_i x)^{p_i} dx \leq D^{-1/2} \prod_{i=1}^m \left( \int_{\mathbb{R}^{n_i}} f_i(y) dy \right)^{p_i}$$

---

$$p_i > 0, f_i \geq 0 \quad \sum_{i=1}^m p_i n_i = n$$

powerful inequality; includes Hölder, Loomis-Whitney, Young's, many others!

---

# Brascamp-Lieb Constant

---

$$\int_{\mathbb{R}^n} \prod_{i=1}^m f_i(B_i x)^{p_i} dx \leq D^{-1/2} \prod_{i=1}^m \left( \int_{\mathbb{R}^{n_i}} f_i(y) dy \right)^{p_i}$$

$$D := \inf \left\{ \frac{\det(\sum_i p_i B_i^* X_i B_i)}{\prod_i (\det X_i)^{p_i}} \mid X_i \succ 0, n_i \times n_i, \right\}$$

---

$$p_i > 0, f_i \geq 0 \quad \sum_{i=1}^m p_i n_i = n$$

powerful inequality; includes Hölder, Loomis-Whitney, Young's, many others!

---

# Brascamp-Lieb constant

---

$$\min_{X_1, \dots, X_m \succ 0} \log \det \left( \sum_i p_i B_i^* X_i B_i \right) - \sum_i p_i \log \det X_i$$

- Arises in geometric complexity theory:  
*[Garg, Gurvits, Oliveira, Wigderson; Jul 2016]*

# Brascamp-Lieb constant

---

$$\min_{X_1, \dots, X_m \succ 0} \log \det \left( \sum_i p_i B_i^* X_i B_i \right) - \sum_i p_i \log \det X_i$$

- Arises in geometric complexity theory:  
*[Garg, Gurvits, Oliveira, Wigderson; Jul 2016]*

*Exercise*

**Prove this is a g-convex opt problem**

# Brascamp-Lieb constant

---

$$\min_{X_1, \dots, X_m \succ 0} \log \det \left( \sum_i p_i B_i^* X_i B_i \right) - \sum_i p_i \log \det X_i$$

- Arises in geometric complexity theory:  
*[Garg, Gurvits, Oliveira, Wigderson; Jul 2016]*

## Exercise

**Prove this is a g-convex opt problem**

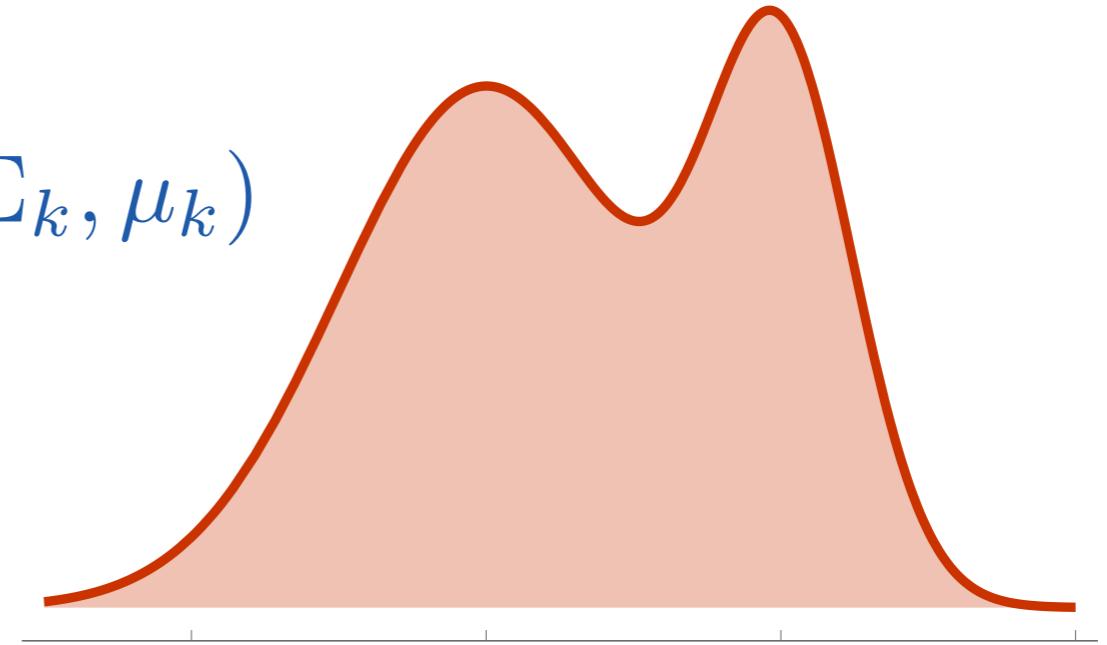
- G-convexity yields transparent algorithms & complexity analysis for global optimum.

# Gaussian mixture models

---

$$p_{\text{mix}}(x) := \sum_{k=1}^K \pi_k p_{\mathcal{N}}(x; \Sigma_k, \mu_k)$$

$$\max \prod_i p_{\text{mix}}(x_i)$$



*Expectation maximization (EM): default choice*

$$p_{\mathcal{N}}(x; \Sigma, \mu) \propto \frac{1}{\sqrt{\det(\Sigma)}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

# Gaussian mixture models

---

- **Nonconvex** – difficult, possibly several local optima
- **GMMs** – Recent surge of theoretical results
- **In Practice** – EM still default choice

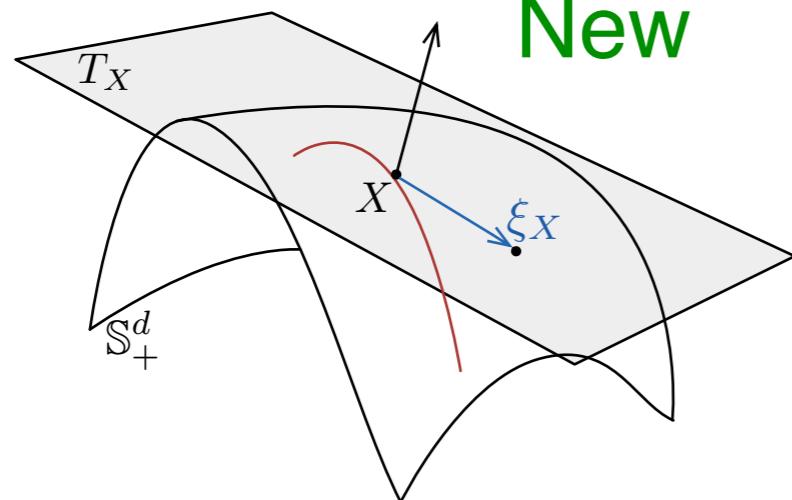
**Difficulty:** Positive definiteness constraint on  $\Sigma_k$

# Gaussian mixture models

- **Nonconvex** – difficult, possibly several local optima
- **GMMs** – Recent surge of theoretical results
- **In Practice** – EM still default choice

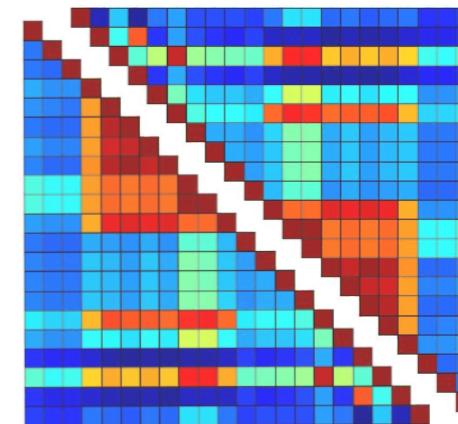
**Difficulty:** Positive definiteness constraint on  $\Sigma_k$

Geometric opt



Unconstrained, Cholesky  
Folklore

$$LL^T$$



[Hosseini, Sra NIPS 2015]

# Failure of “obvious” LL<sup>T</sup>

sep.	EM	CG-LL <sup>T</sup>
0.2	52s // 12.7	614s // 12.7
1	160s // 13.4	435s // 13.5
5	72s // 12.8	426s // 12.8

$$\|\mu_i - \mu_j\| \geq \text{sep} \max_{ij}\{\text{tr}\Sigma_i, \text{tr}\Sigma_j\}$$

*d=20  
simulation*

# Failure of Riemannian optimization

K	EM	Riem-CG
2	17s // 29.28	947s // 29.28
5	202s // 32.07	5262s // 32.07
10	2159s // 33.05	17712s // 33.03



[manopt.org](http://manopt.org)

*Riemannian opt. toolbox*

*d=35  
images  
dataset*

# What's wrong?

---



# What's wrong?

---



**log-likelihood for one component**

$$-\frac{n}{2} \log \det \Sigma - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

# What's wrong?

---



**log-likelihood for one component**

$$-\frac{n}{2} \log \det \Sigma - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

Euclidean convex problem  
**Not geodesically convex**

# What's wrong?



**log-likelihood for one component**

$$-\frac{n}{2} \log \det \Sigma - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

Euclidean convex problem  
**Not geodesically convex**



**Reformulate as g-convex**

$$y_i = \begin{bmatrix} x_i \\ 1 \end{bmatrix} \quad S = \begin{bmatrix} \Sigma + \mu\mu^T & \mu \\ \mu^T & 1 \end{bmatrix}$$
$$\max_{S\succ 0} \widehat{\mathcal{L}}(S) := \sum_{i=1}^n \log q_N(y_i; S),$$

**Thm.** The modified log-likelihood is g-convex. Local max of modified mixture LL is local max of original.

# Success of geometric optimization

K	EM	Riem-CG	L-RBFGS
2	17s // 29.28	<b>18s</b> // 29.28	<b>14s</b> // 29.28
5	202s // 32.07	<b>140s</b> // 32.07	<b>117s</b> // 32.07
10	2159s // 33.05	<b>1048s</b> // 33.06	<b>658s</b> // 33.06

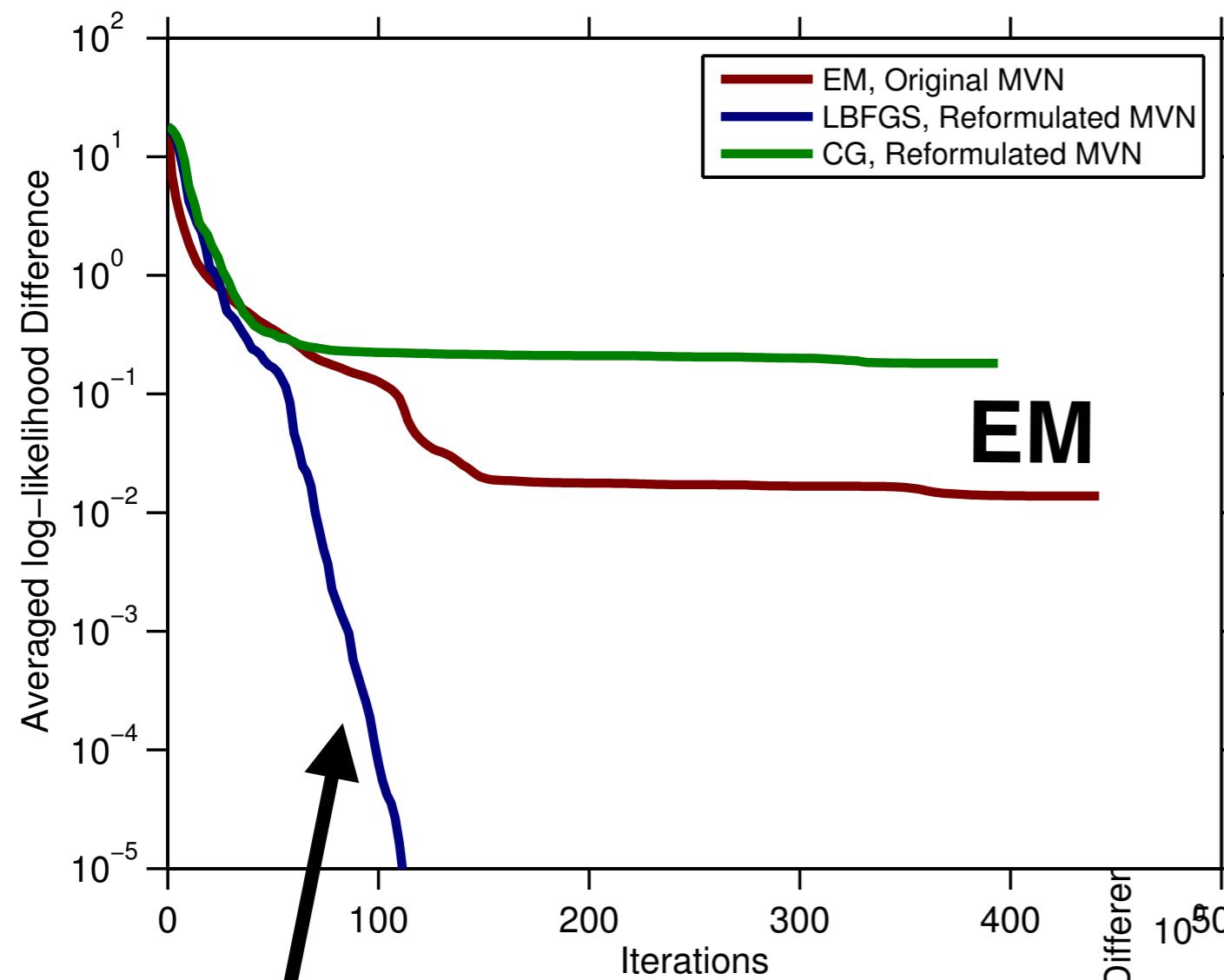
Riem-CG (*manopt*) savings:

947 → **18**; 5262 → **140**; 17712 → **1048**

*d=35  
images  
dataset*



[github.com/utvisionlab/mixest](https://github.com/utvisionlab/mixest)

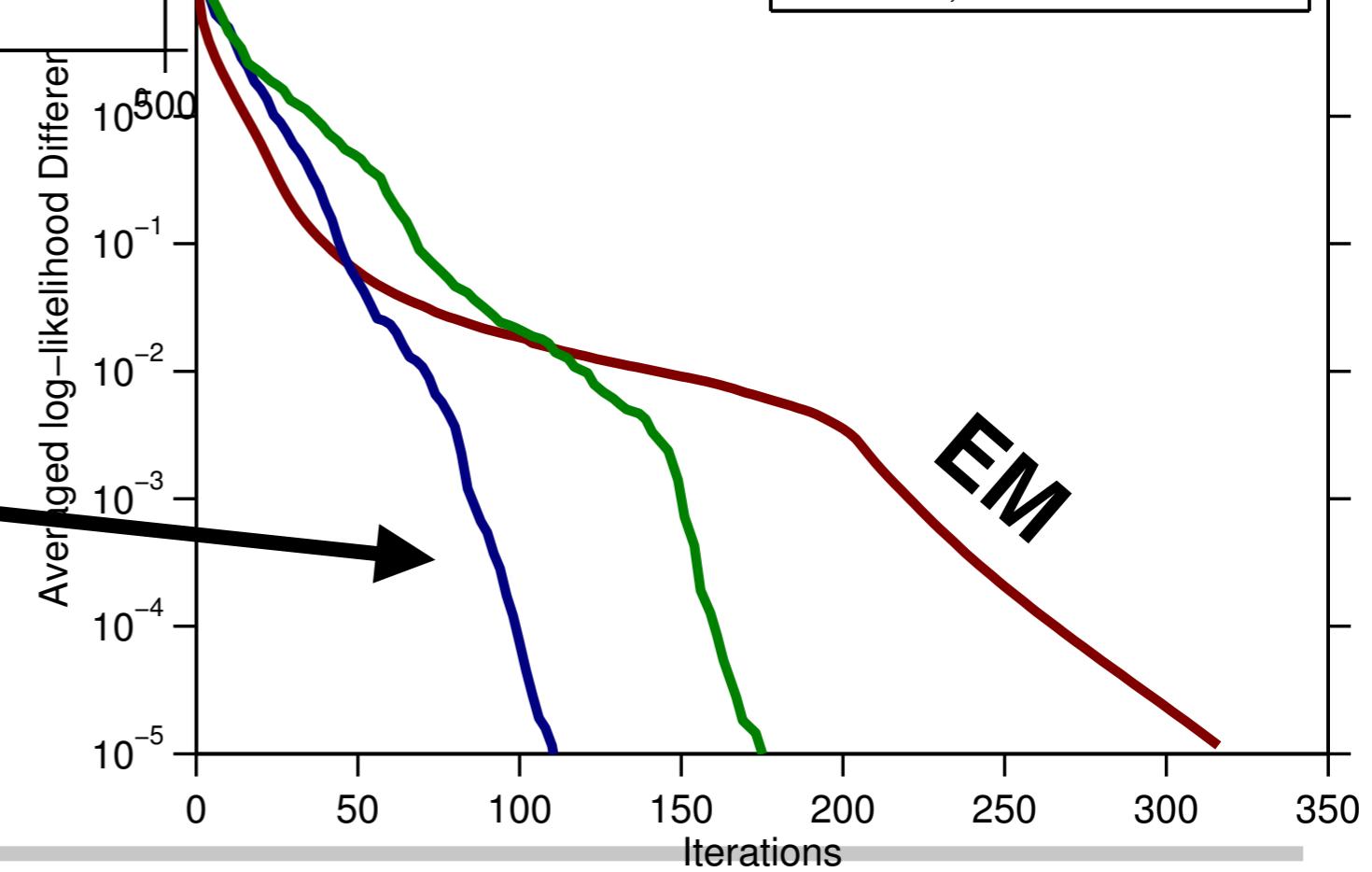


**corel, K=7, d=57**

**EM**

**Our method**

**year, K=9, d=90**



# Extensions

---

- Large-scale mixture models
  - Stochastic optimization
  - Parallel & distributed
- Easily use priors (Wishart and beyond)
- Allows additional constraints
- Non Gaussian mixtures



[github.com/utvisionlab/mixest](https://github.com/utvisionlab/mixest)



[suvrit.de/gopt.html](http://suvrit.de/gopt.html)

# Extensions

---

- Large-scale mixture models
  - Stochastic optimization
  - Parallel & distributed
- Easily use priors (Wishart and beyond)
- Allows additional constraints
- Non Gaussian mixtures



[github.com/utvisionlab/mixest](https://github.com/utvisionlab/mixest)



[suvrit.de/gopt.html](http://suvrit.de/gopt.html)

# Extensions

---

- Large-scale mixture models
    - Stochastic optimization
    - Parallel & distributed
  - Easily use priors (Wishart and beyond)
  - Allows additional constraints
  - Non Gaussian mixtures
- 
- 



[github.com/utvisionlab/mixest](https://github.com/utvisionlab/mixest)



[suvrit.de/gopt.html](http://suvrit.de/gopt.html)

# Extensions

---

- Large-scale mixture models
  - Stochastic optimization 
  - Parallel & distributed
- Easily use priors (Wishart and beyond) 
- Allows additional constraints
- Non Gaussian mixtures 



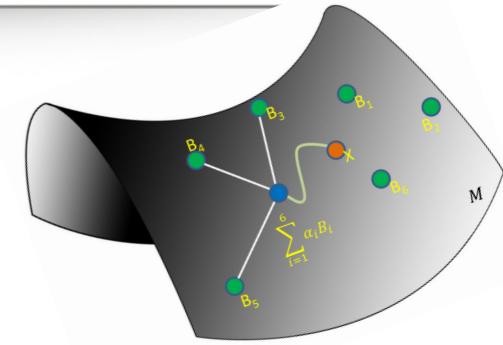
[github.com/utvisionlab/mixest](https://github.com/utvisionlab/mixest)



[suvrit.de/gopt.html](http://suvrit.de/gopt.html)

# Riemannian finite-sum problems

$$\min_{x \in \mathcal{M}} \quad f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x)$$



- $\mathcal{M}$  is a Riemannian manifold
- g-convex and g-nonconvex ‘f’ allowed
- First global complexity results for stochastic methods on Riemannian manifolds
- Riemannian SVRG

[Zhang, Reddi, Sra, NIPS 2016]

# Riemannian SGD with Arbitrary Retraction

- (i) The function satisfies the Lipschitz growth bound

$$f(\text{Ret}_x(\xi)) \leq f(x) + \langle \nabla f(x), \xi \rangle + \frac{L}{2} \|\xi\|^2.$$

- (ii) The stochastic gradients in all iterations are unbiased, i.e.,

$$\mathbb{E}[\nabla f_{i_t}(x_t) - \nabla f(x_t)] = 0.$$

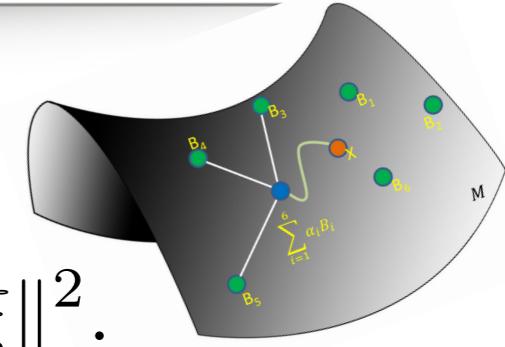
- (iii) The stochastic gradients have bounded variance, so that

$$\mathbb{E}[\|\nabla f_{i_t}(x_t) - \nabla f(x_t)\|^2] \leq \sigma^2, \quad 0 \leq \sigma < \infty.$$

- **We can proof the following convergence result**

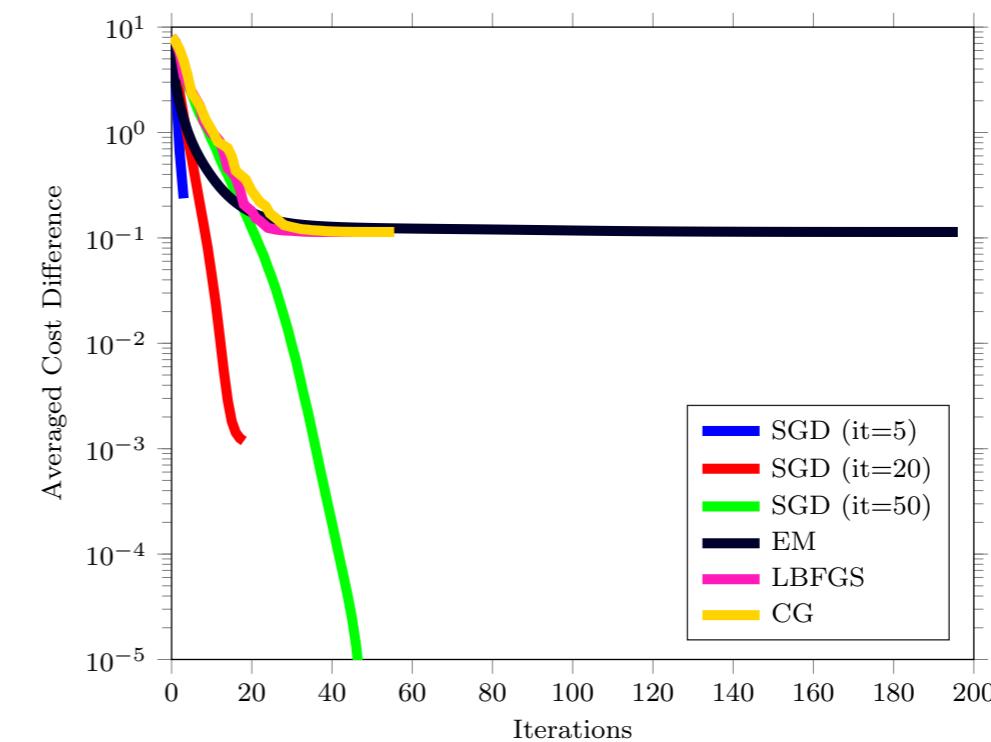
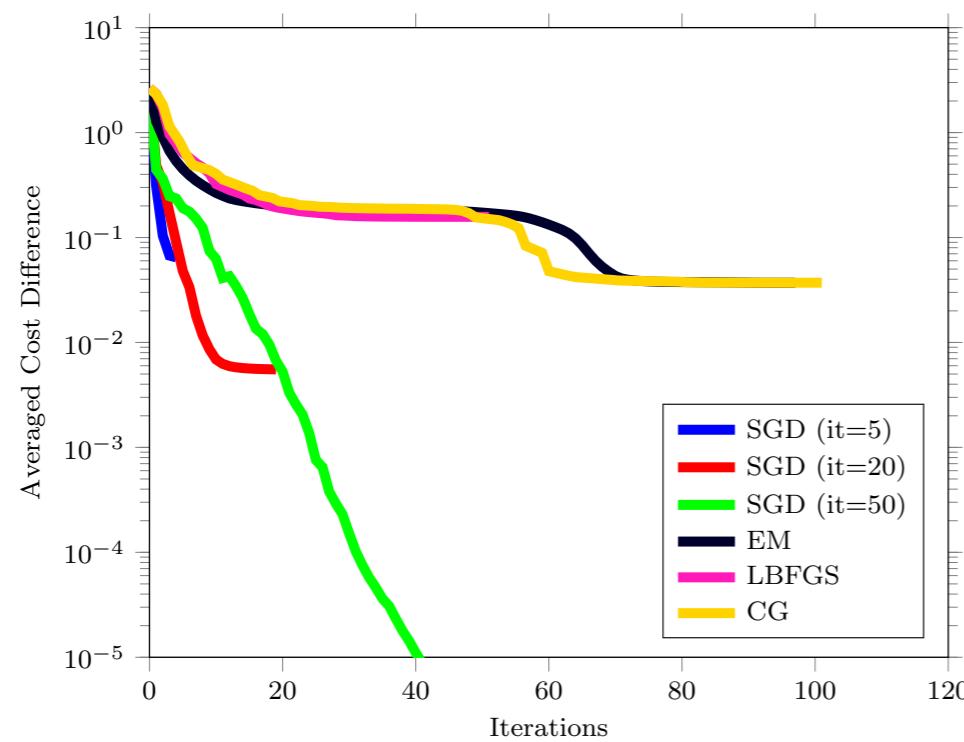
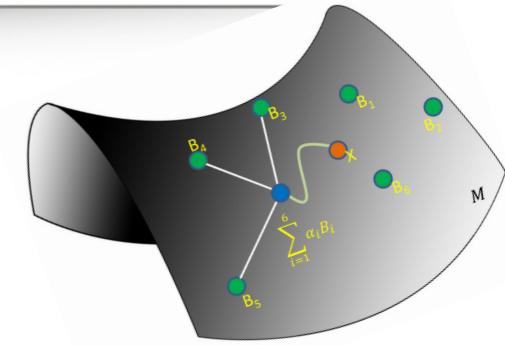
$$\mathbb{E}[\|\nabla f(x_a)\|^2] \leq \frac{2L\Delta_1}{T} + (c + c^{-1}\Delta_1) \frac{L\sigma}{\sqrt{T}} = \mathcal{O}\left(\frac{1}{T}\right) + \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$$

[Hosseini, Sra, MAPR 2019]



# Riemannian SGD with Arbitrary Retraction

- The approach successfully applied to the case of GMMs
- To show that gradients are bounded, we proved a key theorem that iterations stay within a compact set



[Hosseini, Sra, MAPR]

In particular, we study the **global complexity**  
of **first-order g-convex optimization**

## Global Complexity

**Gradient Descent**

**Stochastic Gradient Descent**

**Coordinate Descent**

**Accelerated Gradient Descent**

**Fast Incremental Gradient**

... ...

$$\mathbb{E}[f(x_a) - f(x^*)] \leq ?$$

Convex Optimization

G-Convex Optimization

Convergence rates depend on lower bounds  
on the sectional curvature

## (Sub)gradient

Lipschitz

Strongly convex / smooth

Strongly convex & smooth

## Stochastic (sub)gradient

convex

g-convex

$$O\left(\sqrt{\frac{1}{t}}\right)$$

$$O\left(\frac{1}{t}\right)$$

$$O\left((1 - \frac{\mu}{L_g})^t\right)$$

$$O\left(\sqrt{\frac{\zeta_{\max}}{t}}\right)$$

$$O\left(\frac{\zeta_{\max}}{t}\right)$$

$$O\left((1 - \min\left\{\frac{1}{\zeta_{\max}}, \frac{\mu}{L_g}\right\})^t\right)$$

... ...

See paper for other interesting results [Zhang, Sra, COLT 2016]

$$\zeta_{\max} \triangleq \frac{\sqrt{|\kappa_{\min}|} D}{\tanh\left(\sqrt{|\kappa_{\min}|} D\right)}$$

# Invitation for more!

---

- Other ML application
- Faster stochastic optimization
- Beyond Riemannian manifolds
- Faster gradient based algorithms
- Handling more complex constraints
- Lower bound theory
- ...

# Invitation for more!

---

- Other ML application
- Faster stochastic optimization
- Beyond Riemannian manifolds
- Faster gradient based algorithms
- Handling more complex constraints
- Lower bound theory
- ...

# Thanks!