

# **Nonlinear Signal Processing (2004-2005)**

Course Overview

Instituto Superior Técnico, Lisbon, Portugal

João Xavier

`{jxavier}@isr.ist.utl.pt`

## Outline

- **Motivation:** Signal Processing & Related Applications of Differential Geometry
  - ▷ Optimization
  - ▷ Kendall's theory of shapes
  - ▷ Random Matrix Theory
    - ◇ Coherent Capacity of Multi-Antenna Systems
  - ▷ Information Geometry
  - ▷ Geometrical Interpretation of Jeffreys' Prior
  - ▷ Performance Bounds for Constrained or Non-Identifiable Parametric Estimation
  
- **Course's Table of Contents**
  - ▷ Topological manifolds
  - ▷ Differentiable manifolds
  - ▷ Riemannian manifolds

## Outline

- Bibliography**

- ▷ Recommended textbooks
- ▷ Additional material (short notes on specialized topics)

- Grading**

- Discussion, questions, etc**

## Applications of DG: Optimization

- **Unconstrained minimization problem:**

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

- **Iterative line search:**

given initial point  $\mathbf{x}_0$

for  $k = 0, 1, \dots$

choose descent direction  $\mathbf{d}_k$

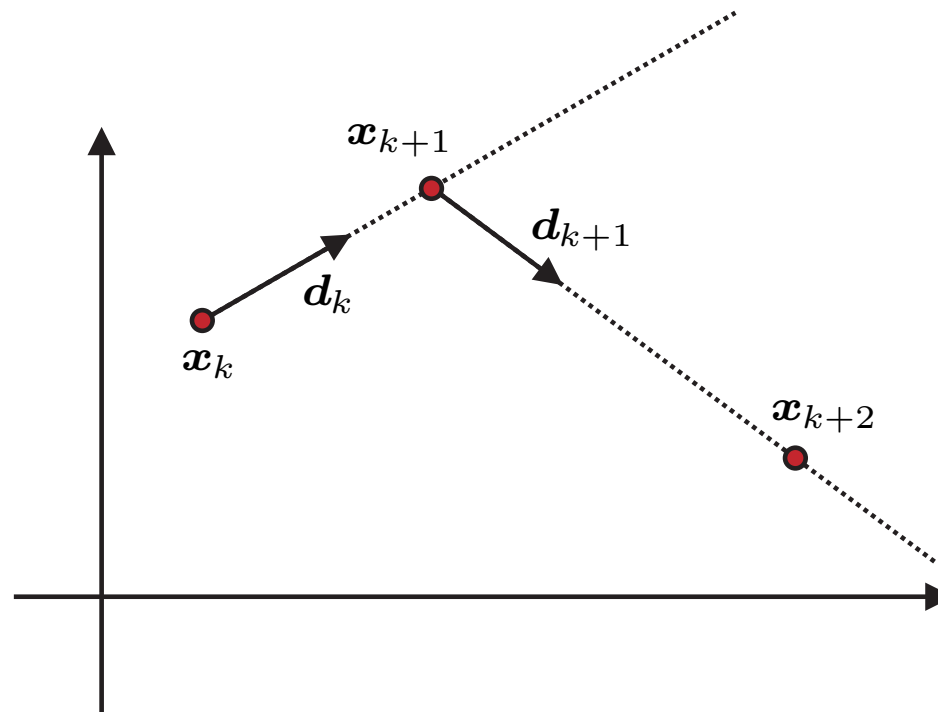
solve  $t^* = \arg \min_{t \geq 0} f(\mathbf{x}_k + t\mathbf{d}_k)$

$\mathbf{x}_{k+1} = \mathbf{x}_k + t^* \mathbf{d}_k$

end

## Applications of DG: Optimization

□ Sketch:



□ Descent direction:  $d_{\text{grad}} = -\nabla f(\mathbf{x}_k)$ ,  $d_{\text{newton}} = -[\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$

## Applications of DG: Optimization

- **Constrained minimization problem:**

$$\mathbf{x}^* = \arg \min_{\mathbf{h}(\mathbf{x})=0} f(\mathbf{x})$$

- **Iterative line search with projected gradient:**

given initial point  $\mathbf{x}_0$

for  $k = 0, 1, \dots$

compute  $\mathbf{d}_k = \mathbf{\Pi}(-\nabla f(\mathbf{x}_k))$

solve  $t^* = \arg \min_{t \geq 0} f(\mathbf{x}_k + t\mathbf{d}_k)$

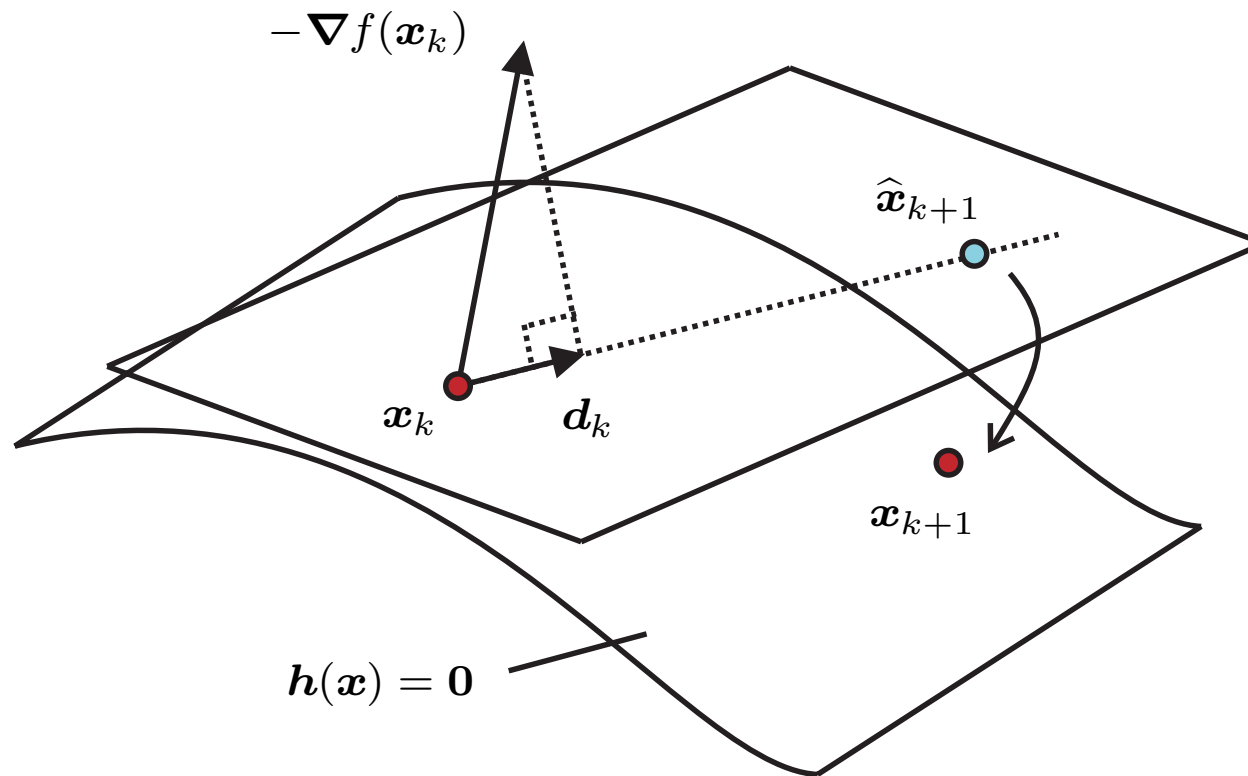
$\hat{\mathbf{x}}_{k+1} = \mathbf{x}_k + t^*\mathbf{d}_k$

return to the constraint surface  $\mathbf{x}_{k+1} = \arg \min_{\mathbf{h}(\mathbf{x})=0} \|\mathbf{x} - \hat{\mathbf{x}}_{k+1}\|^2$

end

## Applications of DG: Optimization

□ Sketch:



## Applications of DG: Optimization

□ Differential geometry enables a descent algorithm with feasible iterates

□ **Iterative geodesic search:**

given initial point  $\mathbf{x}_0$

for  $k = 0, 1, \dots$

choose descent direction  $\mathbf{d}_k$

solve  $t^* = \arg \min_{t \geq 0} f(\gamma_k(t))$

( $\gamma_k(t)$  = geodesic emanating from  $\mathbf{x}_k$  in the direction  $\mathbf{d}_k$ )

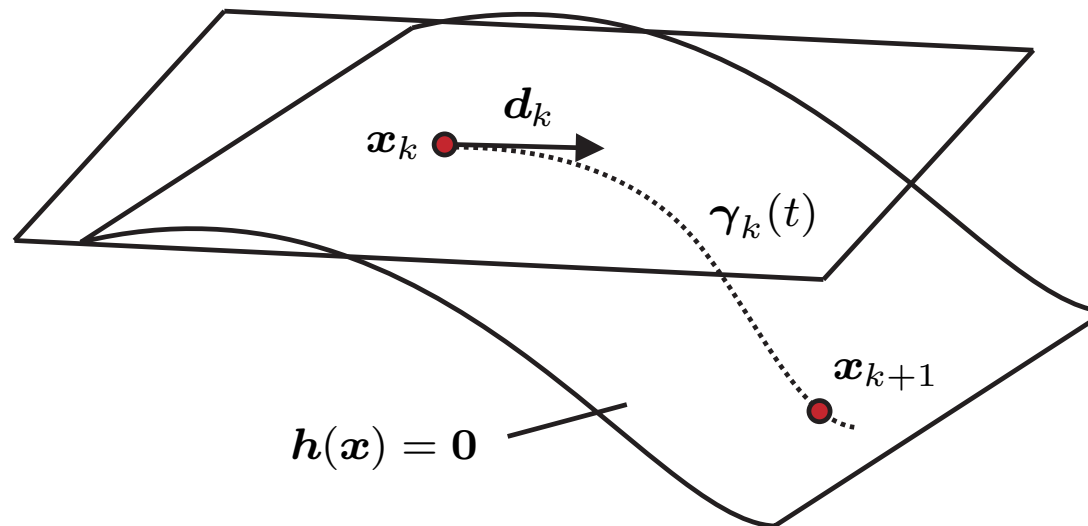
$\mathbf{x}_{k+1} = \gamma_k(t^*)$

end



## Applications of DG: Optimization

□ Sketch:



- **Descent direction:** generalizations of  $d_{\text{grad}}$  and  $d_{\text{newton}}$  are available
- Theory works for abstract spaces (e.g. projective spaces)

## Applications of DG: Optimization

□ **Example:** Signal model

$$\mathbf{y}[t] = \mathbf{Q}\mathbf{x}[t] + \mathbf{w}[t] \quad t = 1, 2, \dots, T$$

$\mathbf{Q}$ : orthogonal matrix ( $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}_N$ ),  $\mathbf{x}[t]$ : known and  $\mathbf{w}[t] \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{C})$

□ **Maximum-Likelihood Estimate:**

$$\mathbf{Q}^* = \arg \max_{\mathbf{Q} \in \mathbb{O}(N)} p(\mathbf{Y}; \mathbf{Q})$$

▷  $\mathbb{O}(N)$  = group of  $N \times N$  orthogonal matrices

▷  $\mathbf{Y} = [\mathbf{y}[1] \ \mathbf{y}[2] \ \cdots \ \mathbf{y}[T]]$  and  $\mathbf{X} = [\mathbf{x}[1] \ \mathbf{x}[2] \ \cdots \ \mathbf{x}[T]]$

## Applications of DG: Optimization

□ **Optimization problem:** Orthogonal Procrustes rotation

$$\begin{aligned} Q^* &= \arg \min_{Q \in \mathbb{O}(N)} \|Y - QX\|_{C^{-1}}^2 \\ &= \arg \min_{Q \in \mathbb{O}(N)} \text{tr} \left\{ Q^T C^{-1} Q \hat{R}_{xx} \right\} - \text{tr} \left\{ Q^T C^{-1} \hat{R}_{yx} \right\} \end{aligned}$$

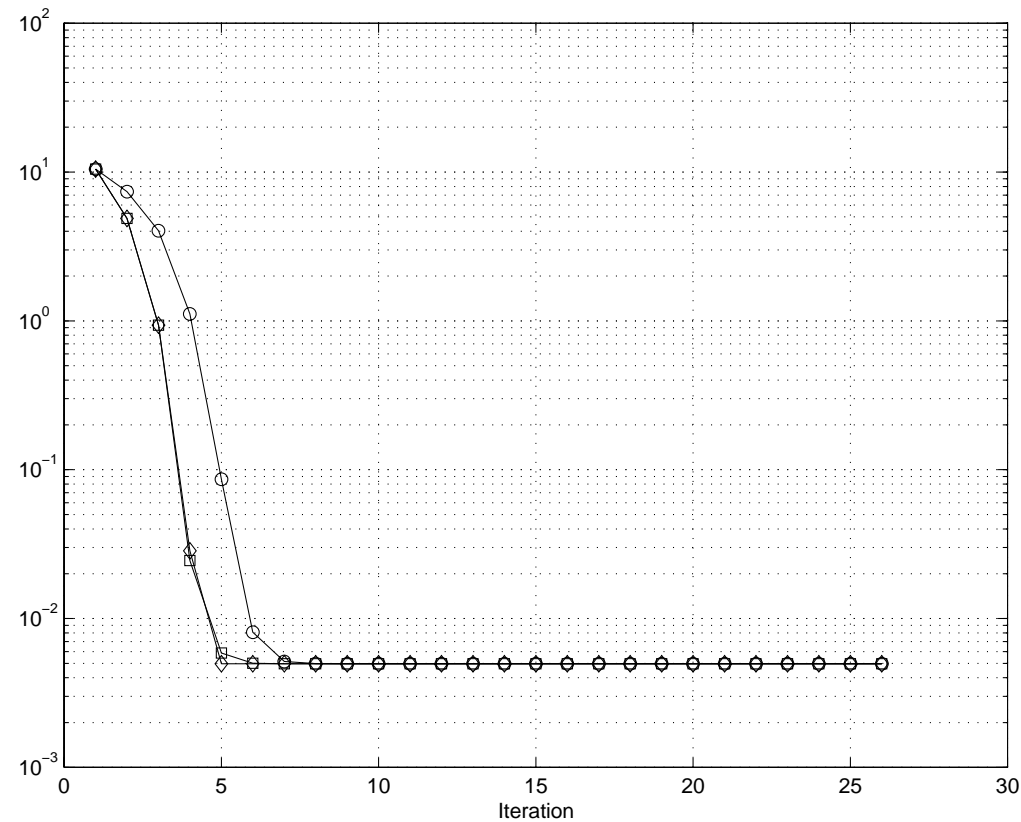
$$\triangleright \hat{R}_{yx} = \frac{1}{T} \sum_{t=1}^T \mathbf{y}[t] \mathbf{x}[t]^T \text{ and } \hat{R}_{xx} = \frac{1}{T} \sum_{t=1}^T \mathbf{x}[t] \mathbf{x}[t]^T$$

□ **Note:** the eigenstructure of  $C$  controls the Hessian of the objective

$$\kappa(C^{-1}) = \frac{\lambda_{\max}(C^{-1})}{\lambda_{\min}(C^{-1})} \text{ condition number of } C^{-1}$$

## Applications of DG: Optimization

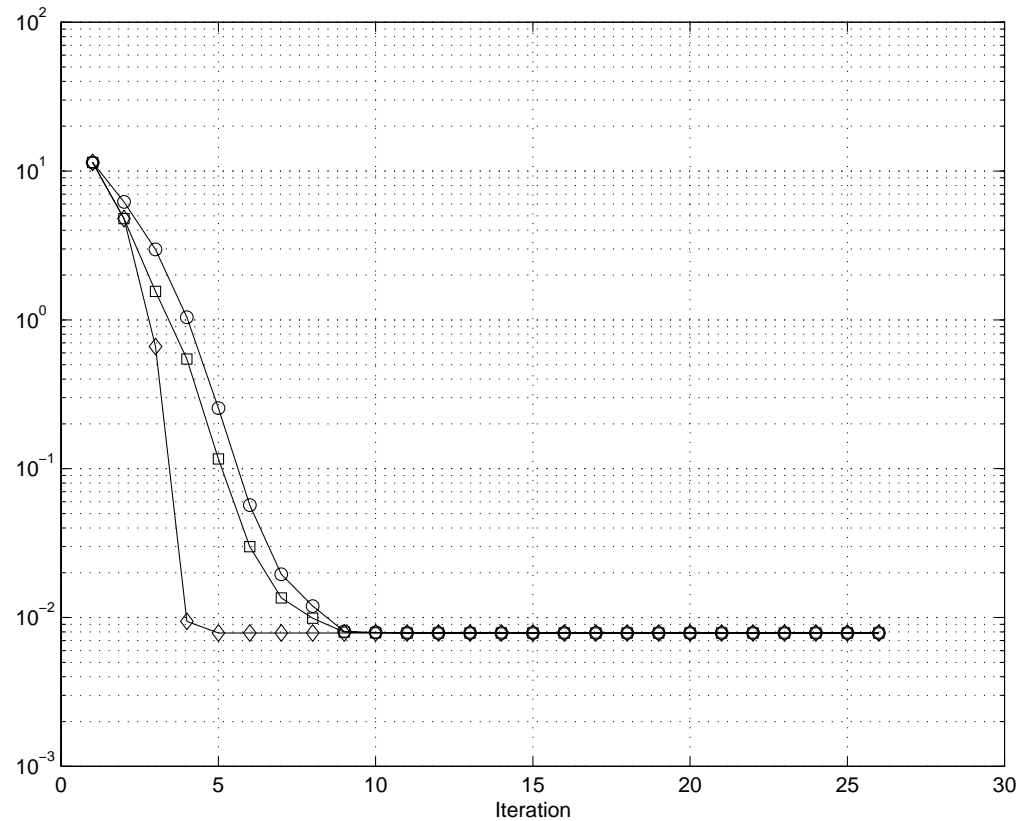
□ **Example:**  $N = 5$ ,  $T = 100$ ,  $\mathbf{C} = \text{diag}(1, 1, 1, 1, 1)$ ,  $\kappa(\mathbf{C}^{-1}) = 1$



○=projected gradient □=gradient geodesic descent ◇=Newton geodesic descent

## Applications of DG: Optimization

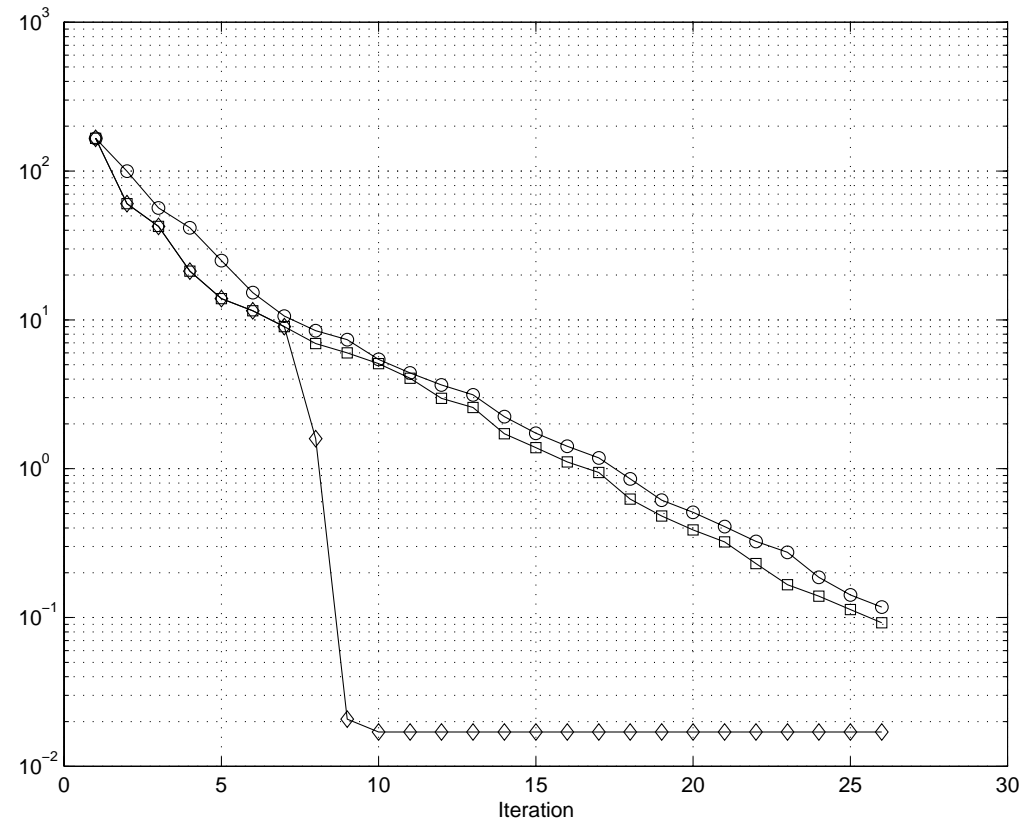
□ **Example:**  $N = 5$ ,  $T = 100$ ,  $\mathbf{C} = \text{diag}(0.2, 0.4, 0.6, 0.8, 1)$ ,  $\kappa(\mathbf{C}^{-1}) = 5$



○=projected gradient □=gradient geodesic descent ◇=Newton geodesic descent

## Applications of DG: Optimization

□ **Example:**  $N = 5$ ,  $T = 100$ ,  $\mathbf{C} = \text{diag}(0.02, 0.05, 0.14, 0.37, 1)$ ,  $\kappa(\mathbf{C}^{-1}) = 50$



○=projected gradient □=gradient geodesic descent ◇=Newton geodesic descent

## Applications of DG: Optimization

- **Important:** Following geodesics is not necessarily optimal. See:  
“Optimization algorithms exploiting unitary constraints”, J. Manton, IEEE Trans. on Signal Processing, vol. 50, no. 3, pp. 635–650, March 2002

## Applications of DG: Optimization

### □ Bibliography:

- ◇ “The geometry of weighted low-rank approximations”, J. Manton *et al.*, IEEE Trans. on Signal Processing, vol. 51, no. 2, pp. 500–514, February 2003
- ◇ “Efficient algorithms for inferences on Grassmann manifolds”, K. Gallivan *et al.*, Proc. 12th IEEE Workshop Statistical Signal Processing, 2003
- ◇ “Adaptive eigenvalue computations using Newton’s method on the Grassmann manifold”, E. Lundstrom *et al.*, SIAM J. Matrix Anal. Appl., vol. 23, no. 3, pp. 819–839, 2002
- ◇ “A Grassmann-Rayleigh quotient iteration for computing invariant subspaces”, P. Absil *et al.*, SIAM Review, vol. 44, no. 1, pp. 57–73, 2002
- ◇ “Algorithms on the Stiefel manifold for joint diagonalization”, M. Nikpour *et al.*, IEEE Int. Conf. on Acoust. Speech and Signal Proc. (ICASSP), vol. 2, pp. 1481–1484, 2002
- ◇ “Optimization algorithms exploiting unitary constraints”, J. Manton, IEEE Trans. on Signal Processing, vol. 50, no. 3, pp. 635–650, March 2002
- ◇ “Contravariant adaptation on structured matrix spaces”, T. Moon and J. Gunther, Signal Processing, 82, pp. 1389–1410, 2002



## Applications of DG: Optimization

### □ Bibliography (cont.):

- ◇ “The geometry of the Newton method on non-compact Lie groups”, R. Mahony and J. Manton, *Journal of Global Optimization*, vol. 23, pp. 309–327, 2002.
- ◇ “Prior knowledge and preferential structures in gradient descent learning algorithms”, R. Mahony and Williamson, *Journal of Machine Learning Research*, pp. 311–355, 2001.
- ◇ “Precoder assisted channel estimation in complex projective space”, J. Manton, *IEEE 3rd Workshop on Sig. Proc. Advanc. on Wir. Comm. (SPAWC)*, pp. 348–351, 2001
- ◇ “Optimization on Riemannian manifold”, *IEEE Proc. 38th conference on Decision and Control*, pp. 888–893, Dec. 1999.
- ◇ “Optimum phase-only adaptive nulling”, S. Smith, *IEEE Trans. on Signal Processing*, vol. 47, no. 7, pp. 1835–1843, July 1999
- ◇ “Motion estimation in computer vision: optimization on Stiefel manifolds”, Y. Ma *et al*, *IEEE Proc. 38th conference on Decision and Control*, vol. 4, pp. 3751–3756, Dec. 1998
- ◇ “The geometry of algorithms with orthogonality constraints”, A. Edelman *et al.*, *SIAM J. Matrix Anal. Appl.*, vol. 20, no. 2, pp. 303–353, 1998

## Applications of DG: Optimization

### □ Bibliography (cont.):

- ◇ “Optimal motion from image sequences: a Riemannian viewpoint”, Y. Ma *et al*, Electronic Research Lab Memorandum, UC Berkeley, 1998
- ◇ “Optimization techniques on Riemannian manifolds”, S. Smith, Fields Institute Communications, vol. 3, pp. 113–136, 1994
- ◇ “Optimization and Dynamical Systems”, U. Helmke and J. Moore, Springer-Verlag, 1994
- ◇ “Geometric optimization methods for adaptive filtering”, S. Smith, PhD Thesis, Harvard University, 1993
- ◇ “Constrained optimization along geodesics”, C. Botsaris, J. Math. Anal. Appl., vol. 79, pp. 295–306, 1981

# Applications of DG: Kendall's theory of shapes

Image 1

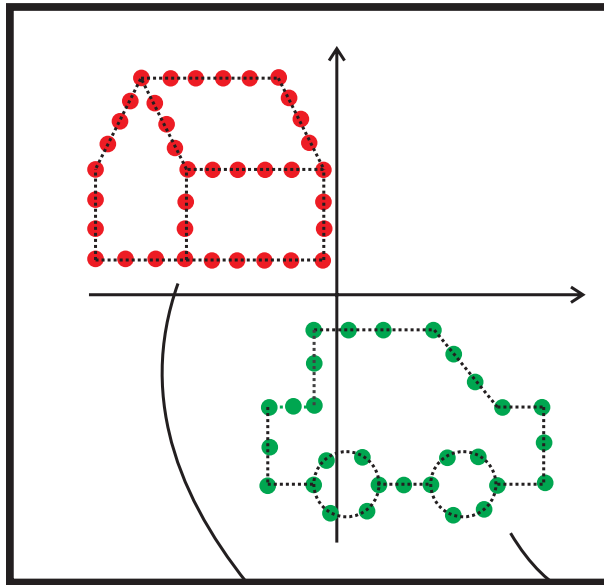
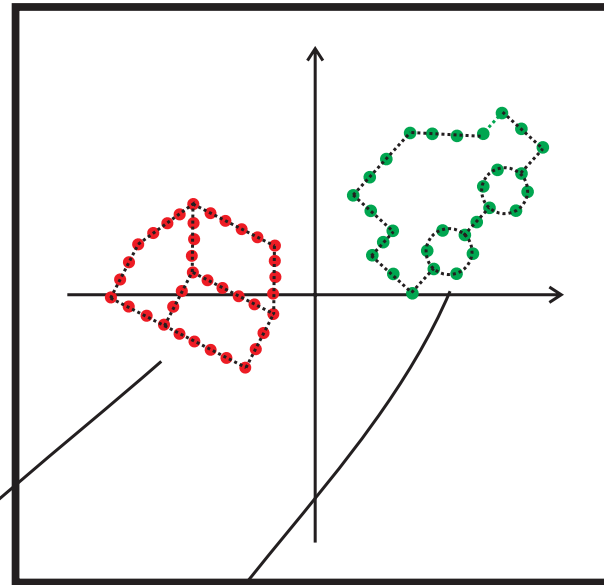


Image 2



Quotient space [manifold]  
Database of shapes

- ▷ (invariant) shape recognition
- ▷ morphing one shape into another
- ▷ statistics ("mean" shape, clustering)

## Applications of DG: Kendall's theory of shapes

### □ Bibliography:

- ◇ "Multivariate shape analysis", I. Dryden and K. Mardia, *Sankhya: The Indian Journal of Statistics*, 55, pp. 460–480, 1993
- ◇ "Procrustes methods in the statistical analysis of shape", C. Goodall, *J. R. Statist. Soc. B*, 53, no.2, pp. 285–339, 1991
- ◇ "A survey of the statistical theory of shapes", D. Kendall, *Statist. Sci.*, 4, pp. pp. 87–120, 1989
- ◇ "Shape manifolds, Procrustean metrics and complex projective spaces", D. Kendall, *Bull. London Math. Soc.*, 16, pp. 81–121, 1984
- ◇ "Directional Statistics", K. Mardia and P. Jupp, *Wiley Series in Probability and Statistics*

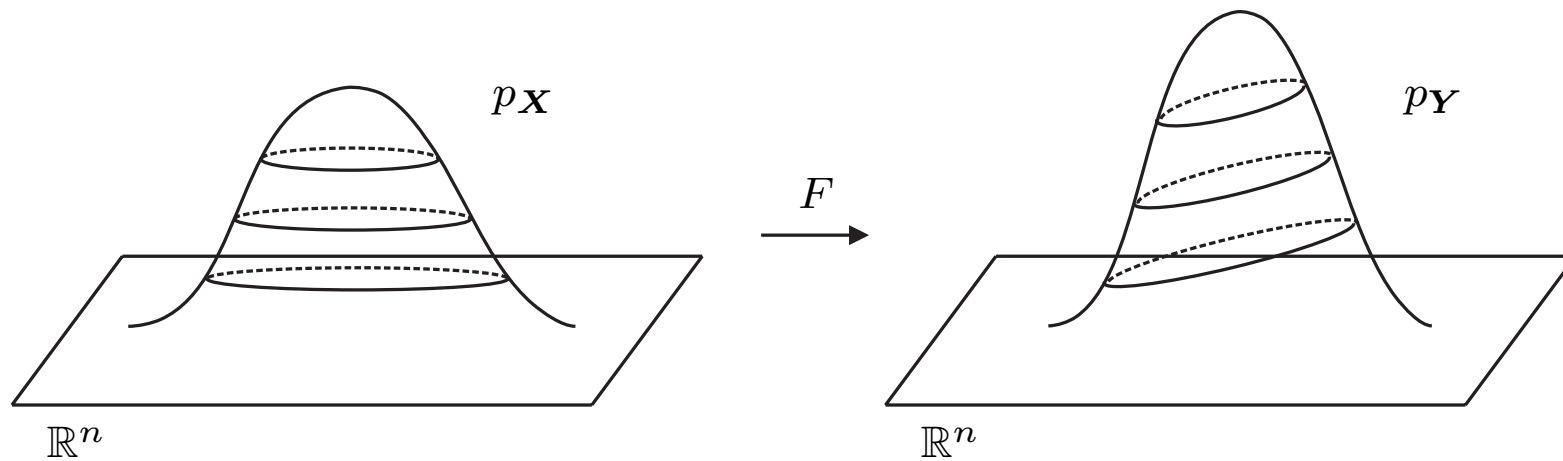
## Applications of DG: Random Matrix Theory

□ **Basic statistics:** transformation of random objects in Euclidean spaces

$$\left\{ \begin{array}{l} \mathbf{x} \text{ is a random vector in } \mathbb{R}^n \\ \mathbf{x} \sim p_{\mathbf{X}}(\mathbf{x}) \\ F : \mathbb{R}^n \rightarrow \mathbb{R}^n \text{ smooth, bijective} \\ \mathbf{y} = F(\mathbf{x}) \end{array} \right.$$

$\Rightarrow$

$$\mathbf{y} \sim p_{\mathbf{Y}}(\mathbf{y}) = p_{\mathbf{X}}(F^{-1}(\mathbf{y})) J(\mathbf{y})$$
$$J(\mathbf{y}) = \frac{1}{\det(DF(F^{-1}(\mathbf{y})))}$$

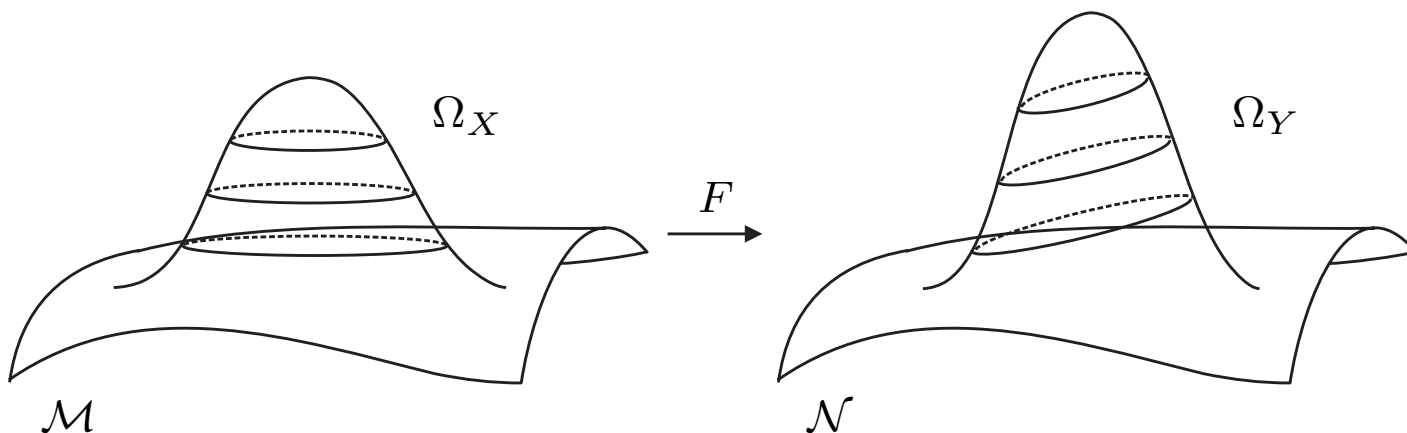


## Applications of DG: Random Matrix Theory

□ **Generalization:** transformation of random objects in manifolds  $\mathcal{M}, \mathcal{N}$

$$\left\{ \begin{array}{l} x \text{ is a random point in } \mathcal{M} \\ x \sim \Omega_X \text{ (exterior form)} \\ F : \mathcal{M} \rightarrow \mathcal{N} \text{ smooth, bijective} \\ y = F(x) \end{array} \right. \Rightarrow y \sim \Omega_Y = \dots$$

The answer is provided by the calculus of exterior differential forms



## Applications of DG: Random Matrix Theory

□ **Example A:** decoupling a random vector in amplitude and direction

$$\left\{ \begin{array}{l} \mathcal{M} = \mathbb{R}^n - \{\mathbf{0}\} = \{\mathbf{x} : \mathbf{x} \neq \mathbf{0}\} \\ \mathcal{N} = \mathbb{R}^+ \times \mathbb{S}^{n-1} = \{(R, \mathbf{u}) : R > 0, \|\mathbf{u}\| = 1\} \\ (R, \mathbf{u}) = F(\mathbf{x}) = \left( \|\mathbf{x}\|, \frac{\mathbf{x}}{\|\mathbf{x}\|} \right) \\ \mathbf{x} \sim p_{\mathbf{X}}(\mathbf{x}) \end{array} \right. \Rightarrow p(R, \mathbf{u}) = p_{\mathbf{X}}(R\mathbf{u}) R^{n-1}$$

□ **Example B:** decoupling a random matrix through the polar decomposition

$$\left\{ \begin{array}{l} \mathcal{M} = \text{GL}(n) = \{\mathbf{X} \in \mathbb{R}^{n \times n} : |\mathbf{X}| \neq 0\} \\ \mathcal{N} = \mathbb{P}(n) \times \mathbb{O}(n) = \{(\mathbf{P}, \mathbf{Q}) : \mathbf{P} \succ \mathbf{0}, \mathbf{Q}^T \mathbf{Q} = \mathbf{I}_n\} \\ (\mathbf{P}, \mathbf{Q}) = F(\mathbf{X}) \Leftrightarrow \mathbf{X} = \mathbf{P}\mathbf{Q} \\ \mathbf{X} \sim p_{\mathbf{X}}(\mathbf{X}) \end{array} \right. \Rightarrow p(\mathbf{P}, \mathbf{Q}) = \dots (\text{known})$$

## Applications of DG: Random Matrix Theory

□ **Example C:** decoupling a random symmetric matrix by eigendecomposition

$$\left\{ \begin{array}{l} \mathcal{M} = \mathbb{S}(n) = \{ \mathbf{X} \in \mathbb{R}^{n \times n} : \mathbf{X} = \mathbf{X}^T \} \\ \mathcal{N} = \mathbb{O}(n) \times \mathbb{D}(n) = \{ (\mathbf{Q}, \mathbf{\Lambda}) : \mathbf{Q}^T \mathbf{Q} = \mathbf{I}_n, \mathbf{\Lambda} : \text{diag} \} \\ (\mathbf{Q}, \mathbf{\Lambda}) = F(\mathbf{X}) \Leftrightarrow \mathbf{X} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T \\ \mathbf{X} \sim p_{\mathbf{X}}(\mathbf{X}) \end{array} \right. \Rightarrow p(\mathbf{Q}, \mathbf{\Lambda}) = \dots (\text{known})$$

□ **Many other examples...** (e.g. Cholesky, QR, LU, SVD)



## Applications of DG: Random Matrix Theory

### □ Bibliography:

- ◇ “Matrix Variate Distributions”, A. Gupta, Chapman & Hall, 1999
- ◇ “Jacobians of Matrix Transformations and Functions of Matrix Argument”, A. Mathai, World Scientific, 1997
- ◇ “Random Matrices”, M. Mehta, Academic Press, 1991
- ◇ “Eigenvalues and Condition Numbers of Random Matrices”, A. Edelman, PhD Thesis, Massachusetts Institute of Technology, 1989
- ◇ “Multivariate Calculation”, R. Farrell, Springer-Verlag, 1985
- ◇ “Aspects of Multivariate Statistical Theory”, R. Muirhead, John Wiley & Sons, 1982
- ◇ “Distributions of matrix variates and latent roots derived from normal samples”, A. James, Annals of Math. Statistics, vol. 35, pp. 475–501, 1964

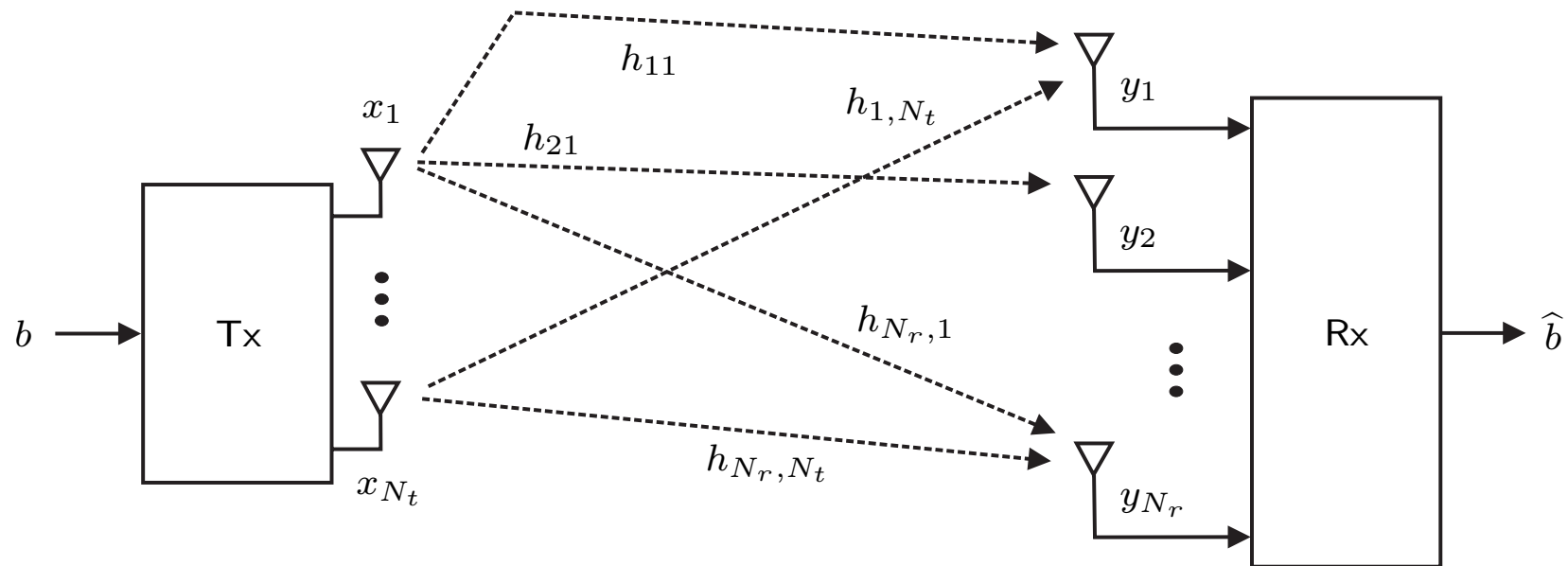
## RMT and DG concepts in signal processing

### □ Bibliography (only a small sample):

- ◇ “Random Matrix Theory and Wireless Communications”, A. Tulino and S. Verdú, Now Publishers Inc., 2004.
- ◇ “Grassmann-based signal design for non-coherent reception”, I. Kammoun and J. C. Belfiore, *Signal Processing Advances in Wireless Communications, 2003*, SPAWC 2003, 4th IEEE Workshop 2003 (pp.507–511)
- ◇ “Communication on the Grassmann manifold: a geometric approach to the noncoherent multiple-antenna channel”, L. Zheng and D. Tse, *IEEE Transactions on Information Theory*, vol. 48, no. 2, pp. 359–383, February 2002.
- ◇ A. Srivastava, “A Bayesian approach to geometric subspace estimation,” *IEEE Transactions on Signal Processing*, vol. 48, no. 5, pp. 1390–1400, May 2000.

## Applications of RMT: Coherent Capacity of Multi-Antenna Systems

- **Scenario:** point-to-point single-user communication with multiple Tx antennas



## Applications of RMT: Coherent Capacity of Multi-Antenna Systems

□ **Data model:**  $\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{n}$  with  $\mathbf{y}, \mathbf{n} \in \mathbb{C}^{N_r}$ ,  $\mathbf{H} \in \mathbb{C}^{N_r \times N_t}$ ,  $\mathbf{x} \in \mathbb{C}^{N_t}$

◇  $N_t$  = number of Tx antennas

◇  $N_r$  = number of Rx antennas

Assumption:  $n_i \stackrel{\text{iid}}{\sim} \mathcal{CN}(0, 1)$

□ **Decoupled data model:**

◇ SVD:  $\mathbf{H} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^H$  with  $\mathbf{U} \in \mathbb{U}(N_r)$ ,  $\mathbf{V} \in \mathbb{U}(N_t)$ ,  $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_f, \mathbf{0})$ ,

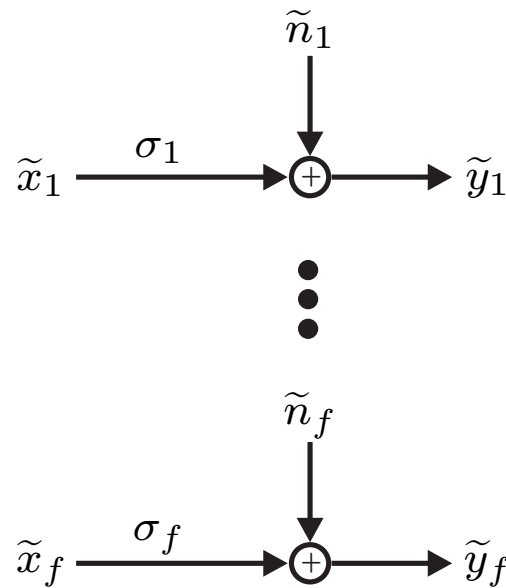
$(\sigma_1, \dots, \sigma_f)$  = nonzero singular values of  $\mathbf{H}$ ,  $f = \min\{N_r, N_t\}$

◇ Transform the data:  $\tilde{\mathbf{y}} = \mathbf{U}^H \mathbf{y}$ ,  $\tilde{\mathbf{x}} = \mathbf{V}^H \mathbf{x}$  and  $\tilde{\mathbf{n}} = \mathbf{U}^H \mathbf{n}$

◇ Equivalent diagonal model:  $\tilde{\mathbf{y}} = \mathbf{\Sigma}\tilde{\mathbf{x}} + \tilde{\mathbf{n}}$

## Applications of RMT: Coherent Capacity of Multi-Antenna Systems

- **Interpretation:** The matrix channel  $\mathbf{H}$  is equivalent to  $f$  parallel scalar channels



## Applications of RMT: Coherent Capacity of Multi-Antenna Systems

□ **Assumption:**  $\mathbf{H}$  is random and known only at the Rx

□ **Channel capacity:**

$$C = \max_{p(\mathbf{x}), \mathbf{E}\{\|\mathbf{x}\|^2\} \leq P} I(\mathbf{x}; (\mathbf{y}, \mathbf{H}))$$

$I$  = mutual information

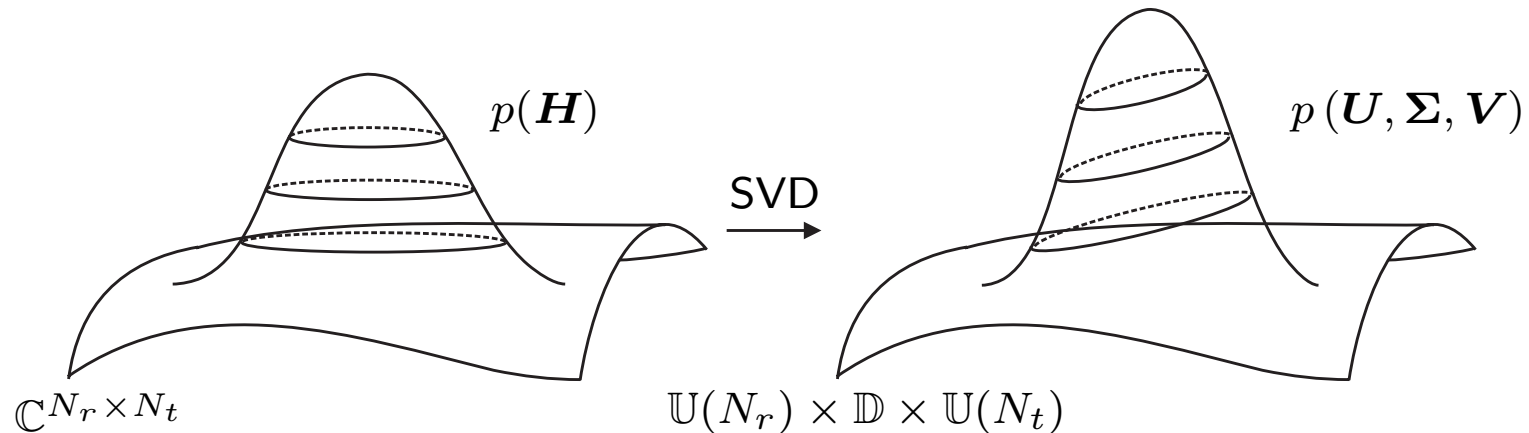
□ **Solution:**

$$C = \mathbf{E}_{\mathbf{H}} \left\{ \sum_{i=1}^f \log \left( 1 + (P/N_t) \sigma_i^2 \right) \right\}$$

Recall:  $(\sigma_1, \dots, \sigma_f)$  = nonzero singular values of  $\mathbf{H}$ ,  $f = \min \{N_r, N_t\}$

## Applications of RMT: Coherent Capacity of Multi-Antenna Systems

- $\mathbf{H}$  is random and  $\mathbf{H} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^H$  (SVD)



- **Capacity:** when  $[\mathbf{H}_{ij}] \stackrel{\text{iid}}{\sim} \mathcal{CN}(0, 1)$

$$C = \int_0^\infty \log(1 + (P/N_t)\lambda) \sum_{k=0}^{f-1} \frac{k!}{(k + g - f)!} (L_k^{g-f}(\lambda))^2 \lambda^{g-f} e^{-\lambda} d\lambda$$

$g = \max\{N_r, N_t\}$  and  $L_j^i = \text{Laguerre polynomials}$

## Applications of RMT: Coherent Capacity of Multi-Antenna Systems

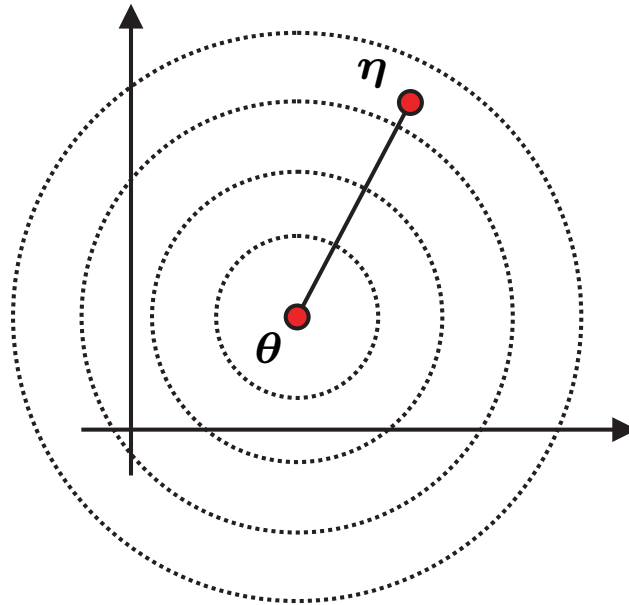
### □ Bibliography:

- ◇ “Keyholes, correlations and capacities of multielement transmit and receive antennas”, D. Chizhik, *IEEE Trans. Wireless Comm.*, vol. 1, pp. 361–368, April 2002
- ◇ “Capacity scaling in MIMO wireless systems under correlated fading”, C. Chuah, *IEEE Trans. Information Th.*, vol. 48, pp. 637–650, March 2002
- ◇ “Capacity of mobile multiple-antenna communication link in Rayleigh flat-fading”, T. Marzetta *et al.*, *IEEE Trans. Information Th.*, vol. 45, no. 1, pp. 139–157, January 1999
- ◇ “On limits of wireless communications in fading environment when using multiple antennas”, G. Foschini and M. Gans, *Wireless Personal Communications*, vol. 6, no.3, pp. 311–355, 1998
- ◇ “Layered space-time architecture for wireless communication in a fading environment when using multi-element antennas”, G. Foschini, *Bell Labs Technical Journal*, vol. 1, no. 2, pp. 41–59, 1996
- ◇ “Capacity of multi-antenna Gaussian channels”, I. Telatar, AT&T Bell Labs, Internal Technical Memorandum, 1995



## Applications of DG: Information Geometry

- **Problem:** Given a parametric statistical family  $\mathcal{F} = \{p(\mathbf{x}; \boldsymbol{\theta}) : \boldsymbol{\theta} \in \Theta\}$  assign a distance function  $d : \Theta \times \Theta \rightarrow \mathbb{R}$
- **Example:**  $\mathcal{F} = \{p(\mathbf{x}; \boldsymbol{\theta}) \sim \mathcal{N}(\boldsymbol{\theta}, \boldsymbol{\Sigma}) : \boldsymbol{\theta} \in \Theta = \mathbb{R}^n\}$  (note:  $\boldsymbol{\Sigma}$  is fixed)  
Naive choice (Euclidean distance):  $d(\boldsymbol{\theta}, \boldsymbol{\eta}) = \|\boldsymbol{\theta} - \boldsymbol{\eta}\|$

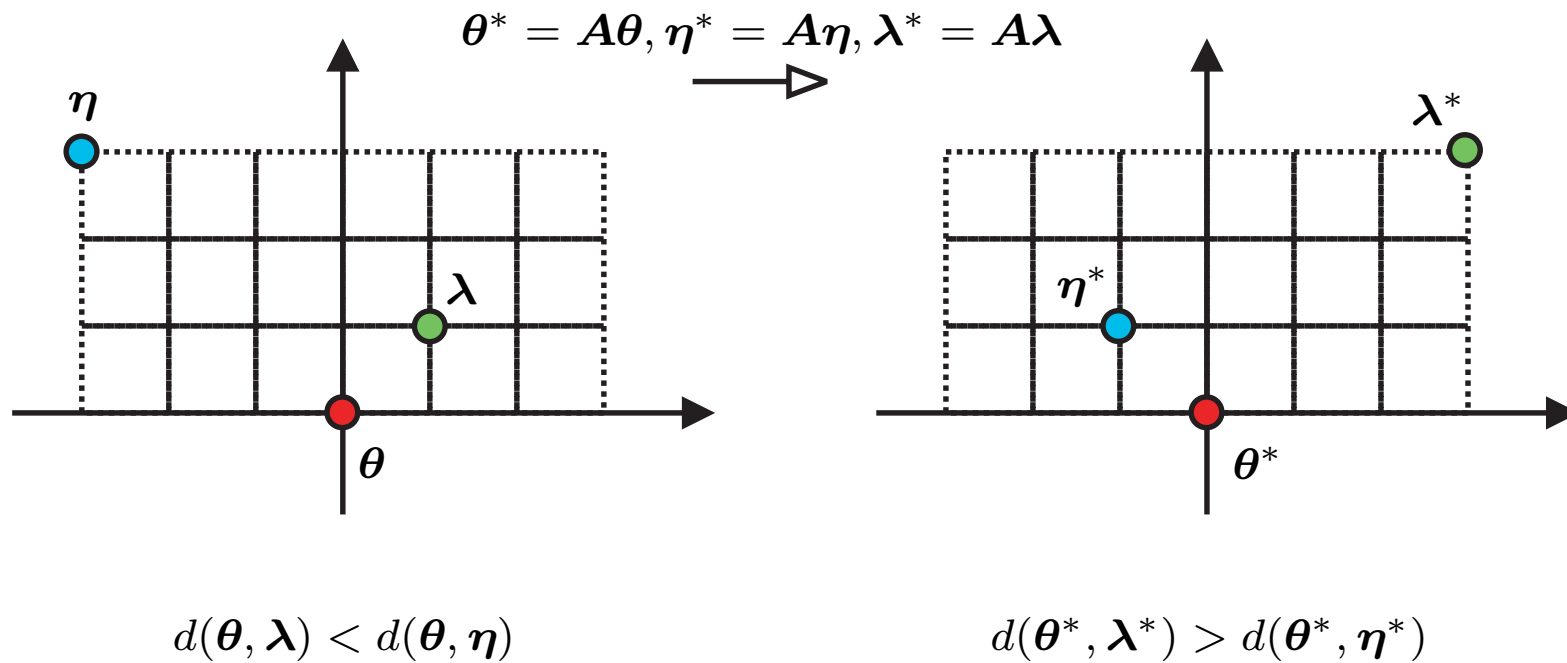


- This method does not produce “intrinsic” distances (parameter invariant)

## Applications of DG: Information Geometry

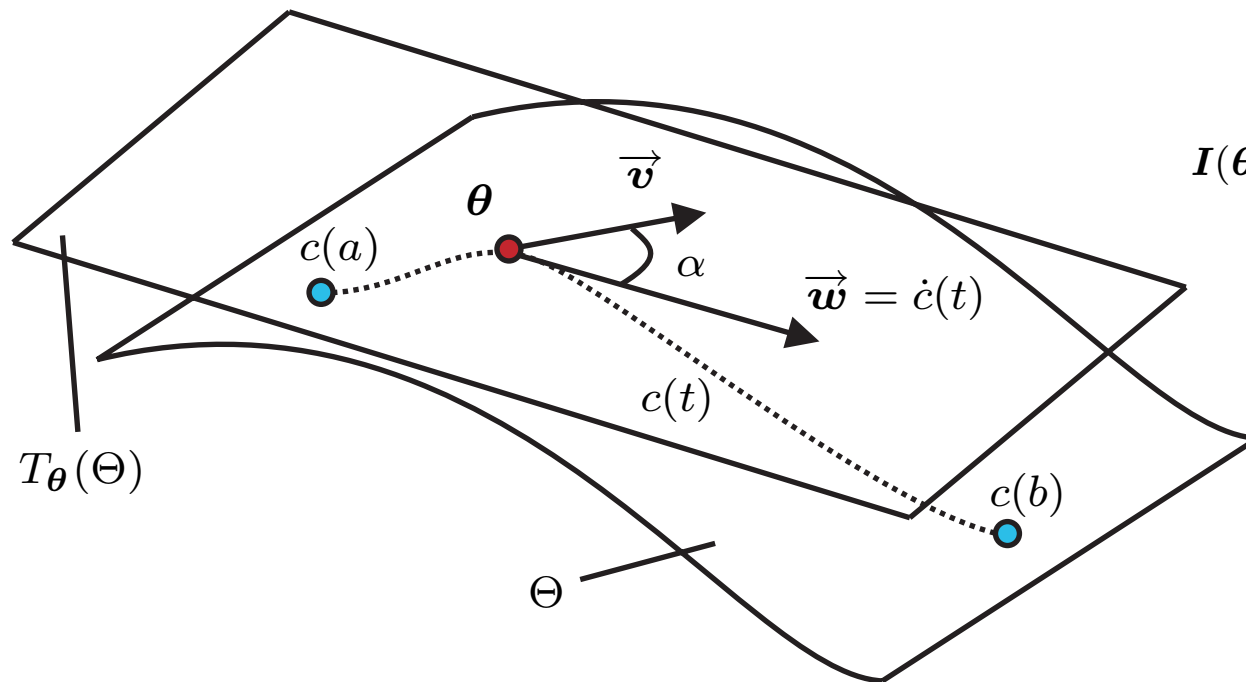
□ With  $\theta^* = A\theta$ :  $\mathcal{F} = \{p(\mathbf{x}; \theta^*) \sim \mathcal{N}(A^{-1}\theta^*, \Sigma) : \theta^* \in \Theta^* = \mathbb{R}^n\}$

□ **Example:**  $\theta = (0, 0)$ ,  $\eta = (-3, 3)$ ,  $\lambda = (1, 1)$ ,  $A = \begin{bmatrix} 5/3 & 4/3 \\ 4/3 & 5/3 \end{bmatrix}$



## Applications of DG: Information Geometry

- Rao suggested the information metric to obtain distances between pdf's
- **Differential geometric interpretation:** The Fisher Information Matrix is adopted as the Riemannian tensor on  $\Theta$



$$\langle \vec{v}, \vec{w} \rangle = \vec{v}^T I(\theta) \vec{w}$$

$$I(\theta) = -E_\theta \{ \nabla_\theta^2 \log p(\mathbf{x}; \theta) \}$$

$$|\vec{v}| = \sqrt{\langle \vec{v}, \vec{v} \rangle}$$

$$\alpha = \frac{\langle \vec{v}, \vec{w} \rangle}{|\vec{v}| |\vec{w}|}$$

$$\text{length}(c) = \int_a^b |\dot{c}(t)| dt$$

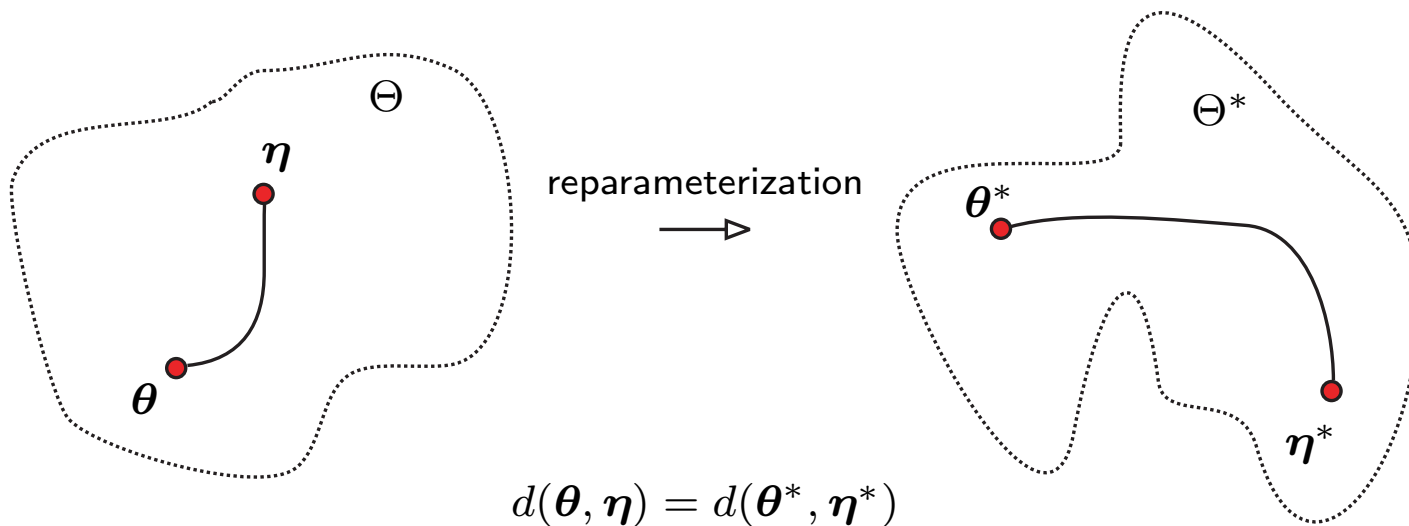
- **Insight:** A parametric statistical family is an autonomous geometrical object

## Applications of DG: Information Geometry

□ **Information distance:**

$$d(\boldsymbol{\theta}, \boldsymbol{\eta}) = \inf \{ \text{length}(c) : c \text{ is a curve on } \Theta \text{ connecting } \boldsymbol{\theta} \text{ to } \boldsymbol{\eta} \}$$

□ The information distance is invariant to reparameterizations



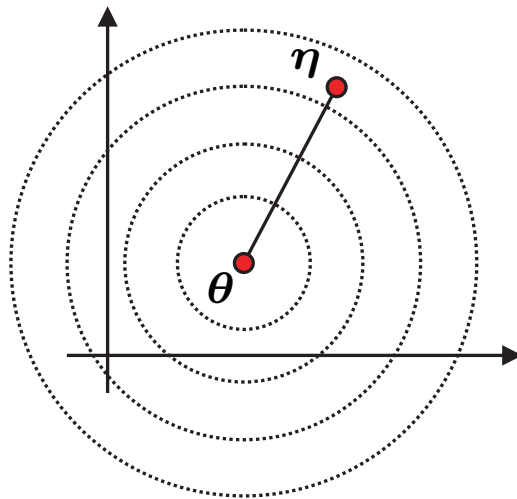
□ **Link with Kullback-Leibler distance:**  $d_{\text{KL}}(\boldsymbol{\theta}, \boldsymbol{\eta}) = \frac{1}{2} d(\boldsymbol{\theta}, \boldsymbol{\eta})^2 + O(d(\boldsymbol{\theta}, \boldsymbol{\eta})^3)$

## Applications of DG: Information Geometry

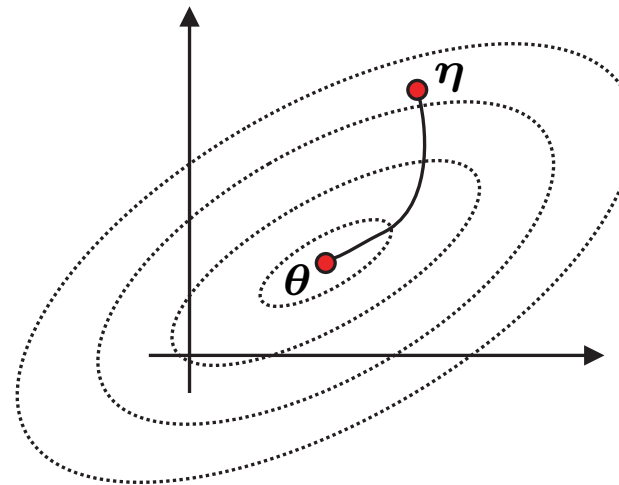
□ Some examples:

$$\diamond \mathcal{F} = \{p(\mathbf{x}; \boldsymbol{\theta}) \sim \mathcal{N}(\boldsymbol{\theta}, \boldsymbol{\Sigma}) : \boldsymbol{\theta} \in \Theta = \mathbb{R}^n\} \quad (\boldsymbol{\Sigma} \text{ is fixed})$$

$$d(\boldsymbol{\theta}, \boldsymbol{\eta}) = \sqrt{(\boldsymbol{\theta} - \boldsymbol{\eta})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\theta} - \boldsymbol{\eta})} \quad [\text{Mahalanobis distance}]$$



Euclidean distance (geodesic)

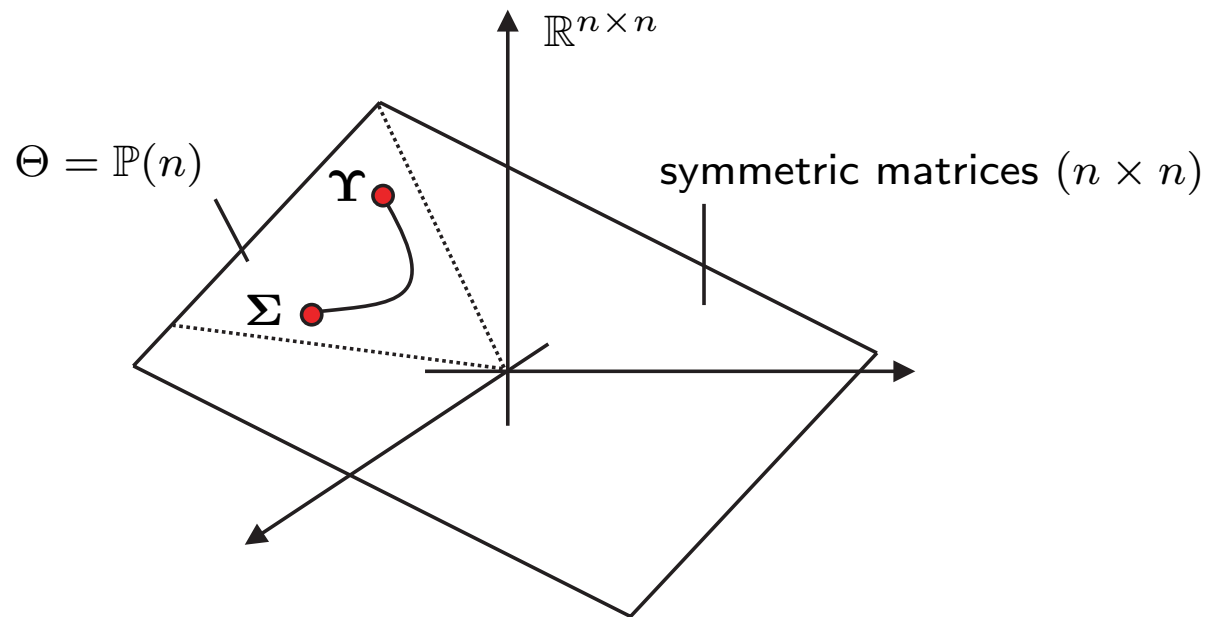


Information distance (geodesic)

## Applications of DG: Information Geometry

◇  $\mathcal{F} = \{p(\mathbf{x}; \Sigma) \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma) : \Sigma \in \Theta = \mathbb{P}(n)\} \quad (\boldsymbol{\mu} \text{ is fixed})$

$$d(\Sigma, \Upsilon) = \sqrt{\frac{1}{2} \sum_{i=1}^n (\log \lambda_i)^2}, \quad (\lambda_1, \dots, \lambda_n) = \text{generalized eigenvalues of } (\Sigma, \Upsilon)$$



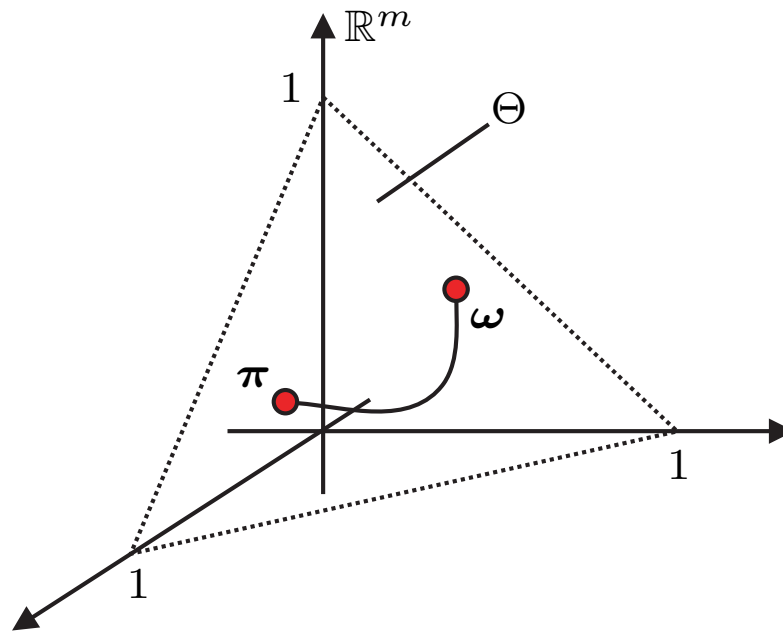
□ Recall:  $\mathbb{P}(n) = \text{set of } n \times n \text{ positive definite matrices}$

## Applications of DG: Information Geometry

◇  $\mathcal{F} = \{p(\mathbf{x}; \boldsymbol{\pi}) \sim \text{multinomial}(n, \boldsymbol{\pi}) : \boldsymbol{\pi} \in \Theta = \text{simplex}(\mathbb{R}^m)\}$

$\mathbf{x} = (x_1, \dots, x_m) \in \mathbb{N}^m, \sum_{i=1}^m x_i = n, \boldsymbol{\pi} = (\pi_1, \dots, \pi_m), \sum_{i=1}^m \pi_i = 1$

$$p(\mathbf{x}; \boldsymbol{\pi}) = \frac{n!}{x_1! \cdots x_m!} \pi_1^{x_1} \cdots \pi_m^{x_m} \quad d(\boldsymbol{\pi}, \boldsymbol{\omega}) = 2\sqrt{n} \arccos \left( \sum_{i=1}^m \pi_i \omega_i \right)$$



## Applications of DG: Information Geometry

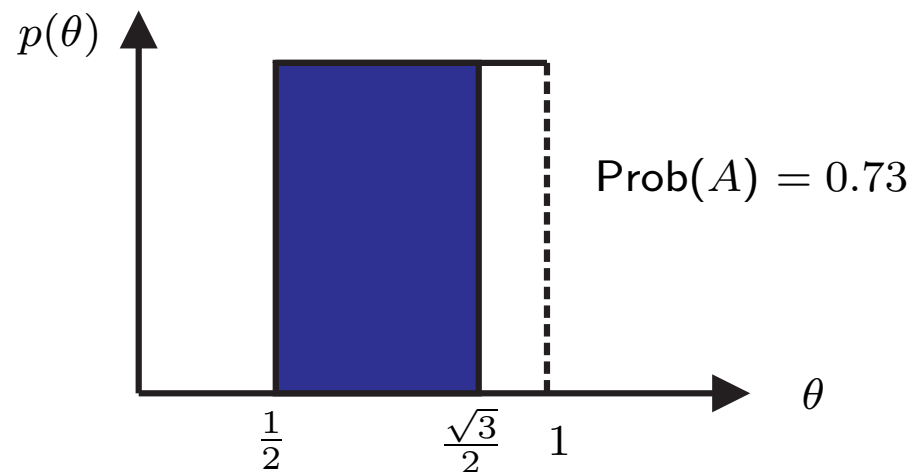
### □ Bibliography:

- ◇ “Differential Geometry and Statistics”, M. Murray *et al.*, Chapman & Hall, 1993
- ◇ “The geometry of asymptotic inference”, R. Kass, *Statistical Science*, vol. 4, no. 3, pp. 188–234, 1989
- ◇ “Differential Geometry in Statistical Inference”, S. Amari *et al.*, Institute of Mathematical Statistics, Lecture Notes, 1987
- ◇ “The role of differential geometry in statistical theory”, O.E. Barndorff-Nielsen *et al.*, *International Statistical Review*, 54, pp. 83–96, 1986
- ◇ “Information and accuracy attainable in the estimation of statistical parameters”, C. Rao, *Bull. Calcutta Math. Soc.*, 37, pp. 81–91, 1945



## Applications of DG: Geometrical Interpretation of Jeffreys' Prior

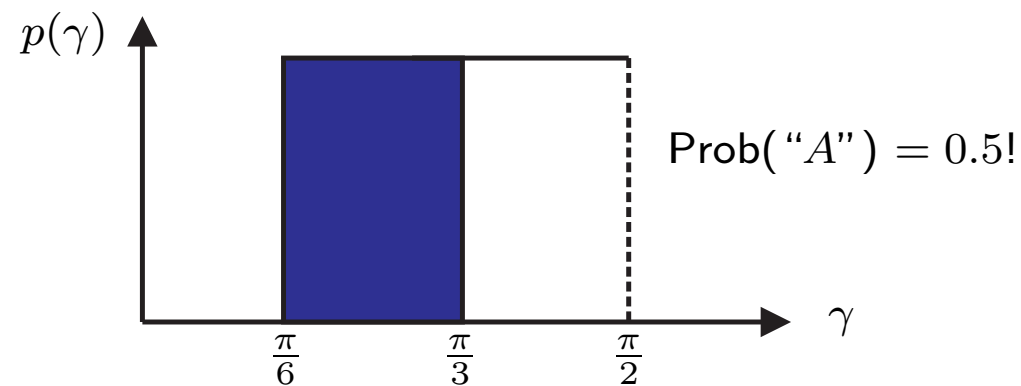
- **Problem:** Given a parametric statistical family  $\mathcal{F} = \{p(x; \theta) : \theta \in \Theta\}$  assign a non-informative prior  $p(\theta)$  for the parameter  $\theta$
- **Example:**  $\mathcal{F} = \{p(x; \theta) \sim \mathcal{N}(0, \theta^2) : \theta \in \Theta = (1/2, 1)\}$   
Naive choice (uniform distribution):



- This method does not produce “intrinsic” priors (parameter invariant)

## Applications of DG: Geometrical Interpretation of Jeffreys' Prior

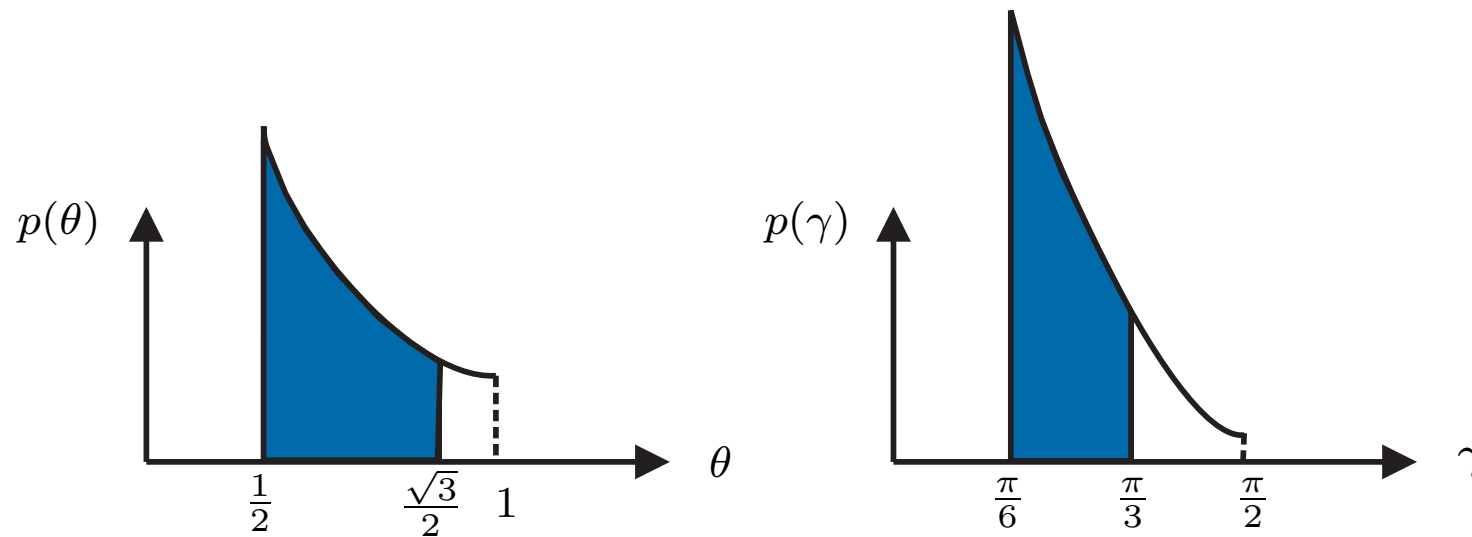
□ With  $\theta = \sin(\gamma)$ :  $\mathcal{F} = \{p(x; \gamma) \sim \mathcal{N}(0, \sin^2(\gamma)) : \gamma \in \Gamma = (\pi/6, \pi/2)\}$



□ Jeffreys' prior:  $p(\boldsymbol{\theta}) \propto \sqrt{\det(\mathbf{I}(\boldsymbol{\theta}))}$  where  $\mathbf{I}(\boldsymbol{\theta})$  is the Fisher information matrix

## Applications of DG: Geometrical Interpretation of Jeffreys' Prior

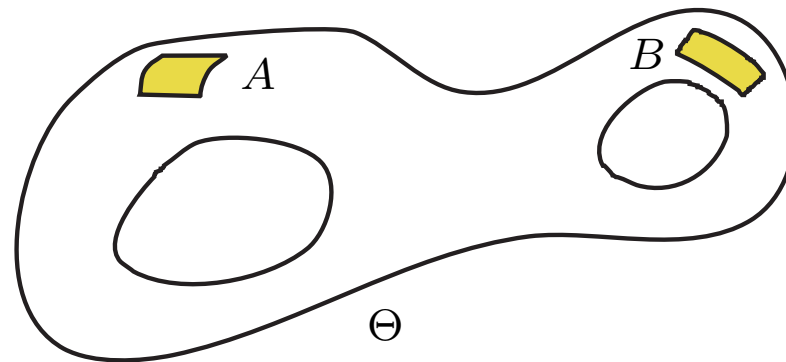
□ For the current example:  $p(\theta) \propto \frac{1}{\theta}$  and  $p(\gamma) \propto \cotg(\gamma)$



$$\text{Prob}(A) = \text{Prob}("A") = 0.79$$

## Applications of DG: Geometrical Interpretation of Jeffreys' Prior

- **Differential geometric interpretation:** Jeffreys' prior is simply the Riemannian volume element induced by the Fisher metric!
- **Insight:** A parametric statistical family is an autonomous geometrical object carrying its own “uniform” prior (applies equal mass to sets of equal area)



$$\text{Area}(A) = \text{Area}(B) \Rightarrow \text{Prob}(\boldsymbol{\theta} \in A) = \text{Prob}(\boldsymbol{\theta} \in B)$$

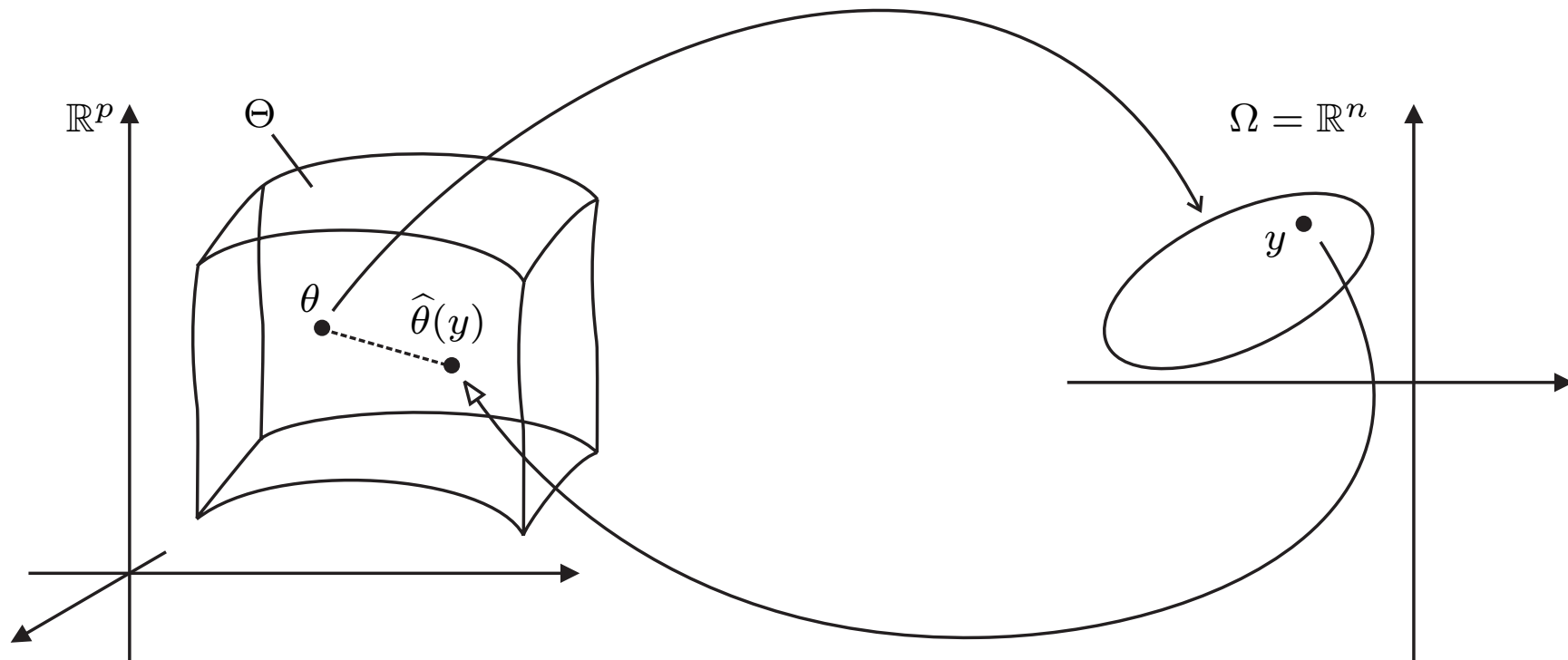
## Applications of DG: Geometrical Interpretation of Jeffreys' Prior

### □ Bibliography:

- ◇ “The geometry of asymptotic inference”, R. Kass, *Statistical Science*, vol. 4, no. 3, pp. 188–234, 1989
- ◇ “Differential Geometry in Statistical Inference”, S. Amari *et al.*, Institute of Mathematical Statistics, Lecture Notes, 1987
- ◇ “The role of differential geometry in statistical theory”, O.E. Barndorff-Nielsen *et al.*, *International Statistical Review*, 54, pp. 83–96, 1986
- ◇ “Theory of Probability”, 3rd ed., H. Jeffreys, Oxford University, 1961
- ◇ “An invariant form for the prior probability in estimation problems”, H. Jeffreys, *Proc. Royal Soc. London Ser. A*, 196, pp. 453–461, 1946

## Application of DG: bounds

□ Classical Euclidean setup:

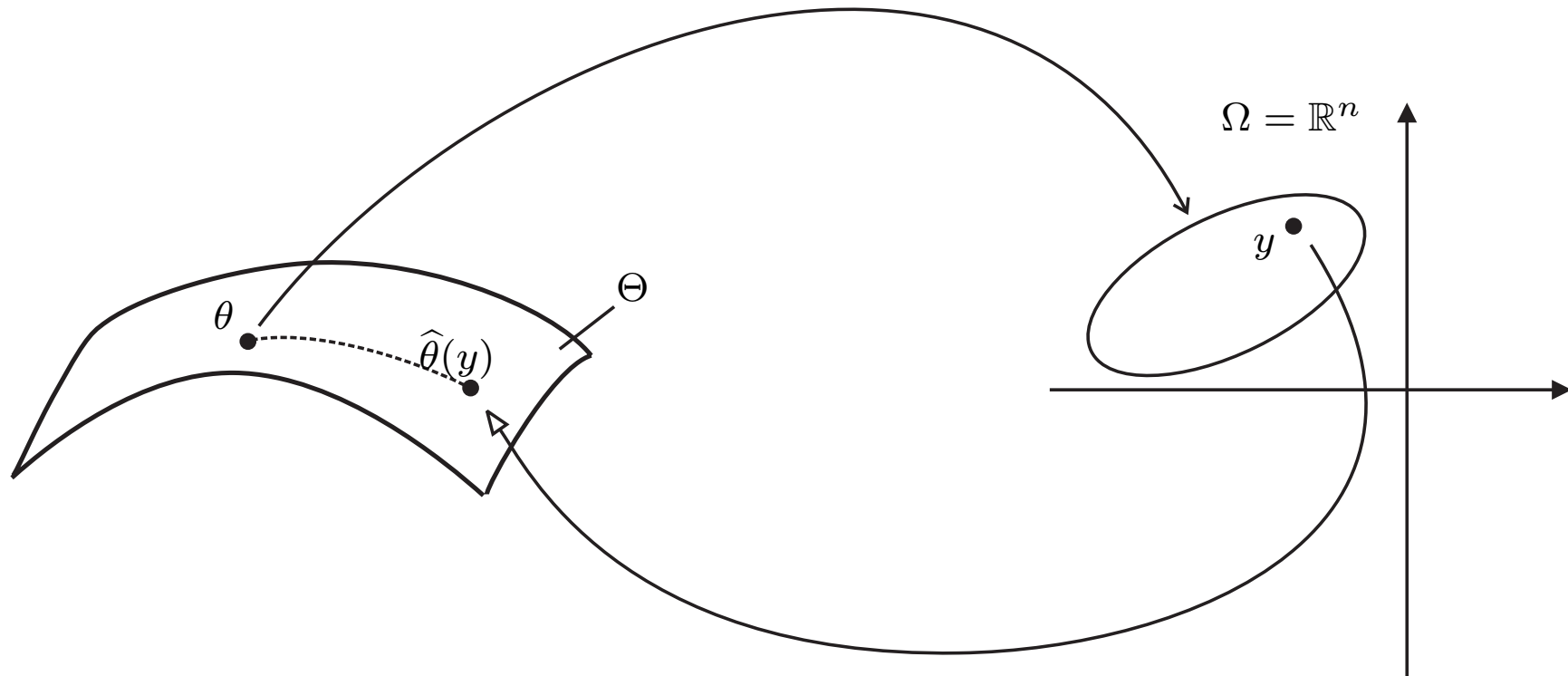


□ Cramér-Rao Bound (CRB):

$$\text{var}_{\theta}(\hat{\theta}) = \mathbb{E}_{\theta} \left\{ d(\theta, \hat{\theta}(Y))^2 \right\} \geq \text{tr}(I_{\theta}^{-1}) \quad (I_{\theta} = \text{Fisher matrix})$$

## Application of DG: bounds

□ Riemannian setup:



□ Intrinsic Variance Lower Bound (IVLB):

$$\text{var}_{\theta}(\hat{\theta}) = \mathbb{E}_{\theta} \left\{ d(\theta, \hat{\theta}(Y))^2 \right\} \geq ?$$

## Applications of DG: bounds

□ **Theorem (IVLB).** Suppose:

- ▷ The sectional curvature of  $\Theta$  is upper bounded by  $C \geq 0$
- ▷ + some technical conditions

Then,

$$\text{var}_\theta(\hat{\theta}) \geq \begin{cases} \lambda_\theta & , \text{ if } C = 0 \\ \frac{\lambda_\theta C + 1 - \sqrt{2\lambda_\theta C + 1}}{C^2 \lambda_\theta / 2} & , \text{ if } C > 0 \end{cases}$$

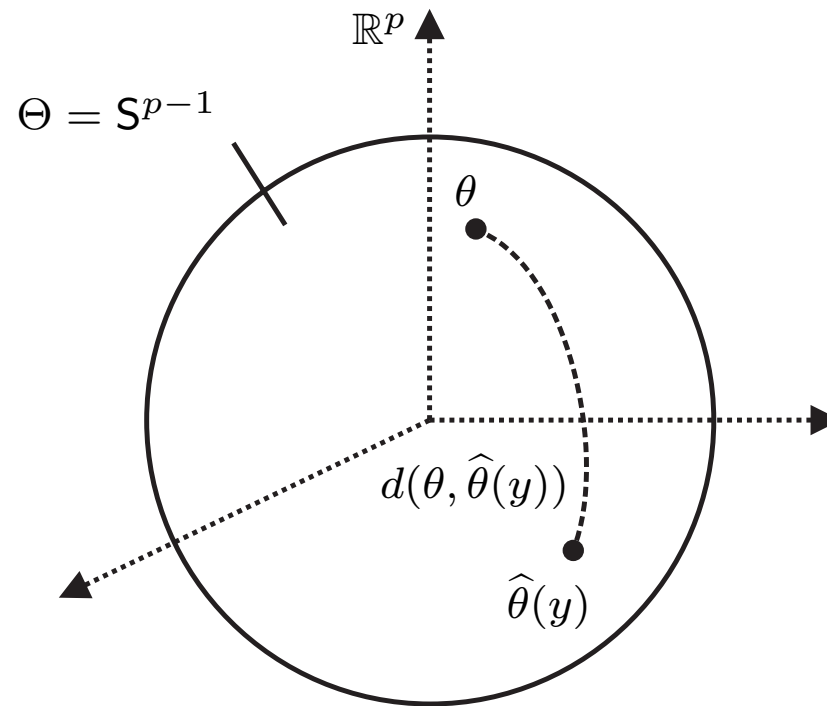
where:

- ▷  $\lambda_\theta = \text{tr}(I_\theta^{-1})$  ( $I_\theta =$  Fisher tensor )



## Example: inference on $S^{p-1}$

□  $S^{p-1} = \{x \in \mathbb{R}^p : \|x\| = 1\}$  is the unit-sphere in  $\mathbb{R}^p$



□ **Geometry of  $\Theta$ :**  $d(\theta, \hat{\theta}(y)) = \text{acos}(\theta^T \hat{\theta}(y))$  and  $C = 1$

## Example: inference on $S^{p-1}$

□ **Observation:**  $y = \theta + w \in \mathbb{R}^p$  ( $p = 10$ )

▷  $\theta \in \Theta = S^{p-1}$

▷  $w \sim \mathcal{N}(0, \sigma^2 I_p)$

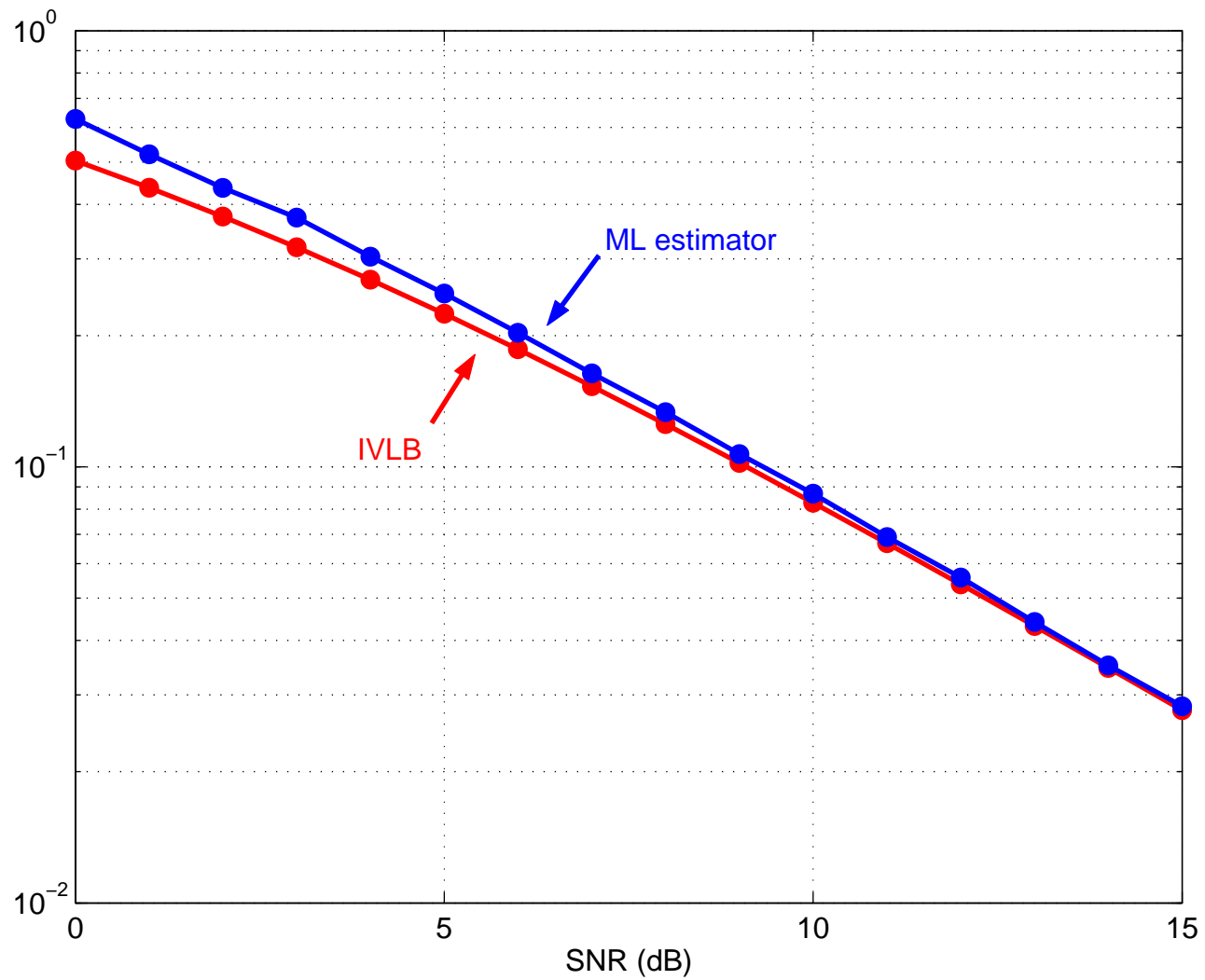
□ Maximum-likelihood estimator:

$$\hat{\theta}(y) = \frac{y}{\|y\|}$$

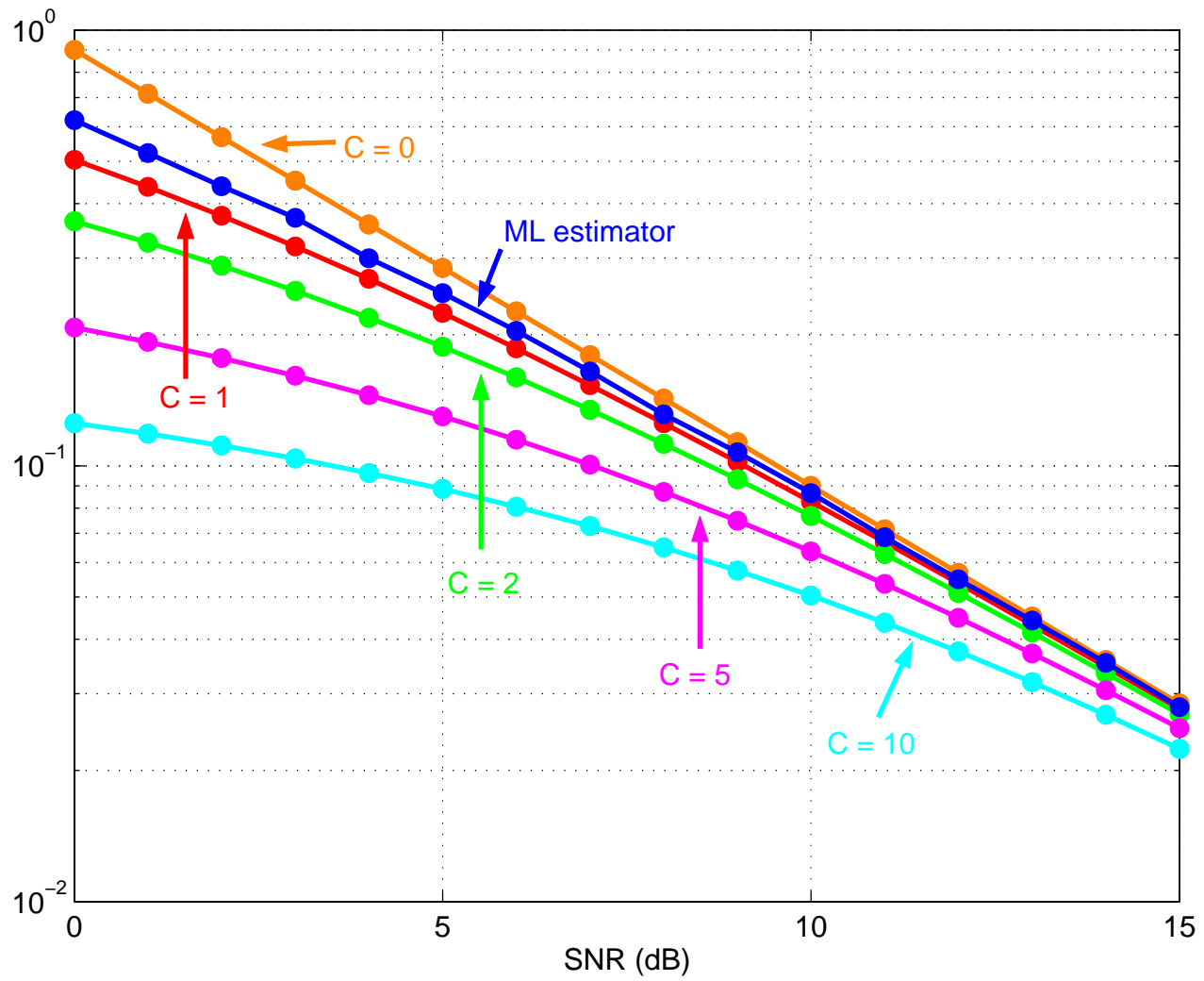
□ Signal-to-noise ratio:

$$\text{SNR} = \frac{\mathbb{E} \left\{ \|\theta\|^2 \right\}}{\mathbb{E} \left\{ \|w\|^2 \right\}} = \frac{1}{p \sigma^2}$$

Example: inference on  $S^{p-1}$



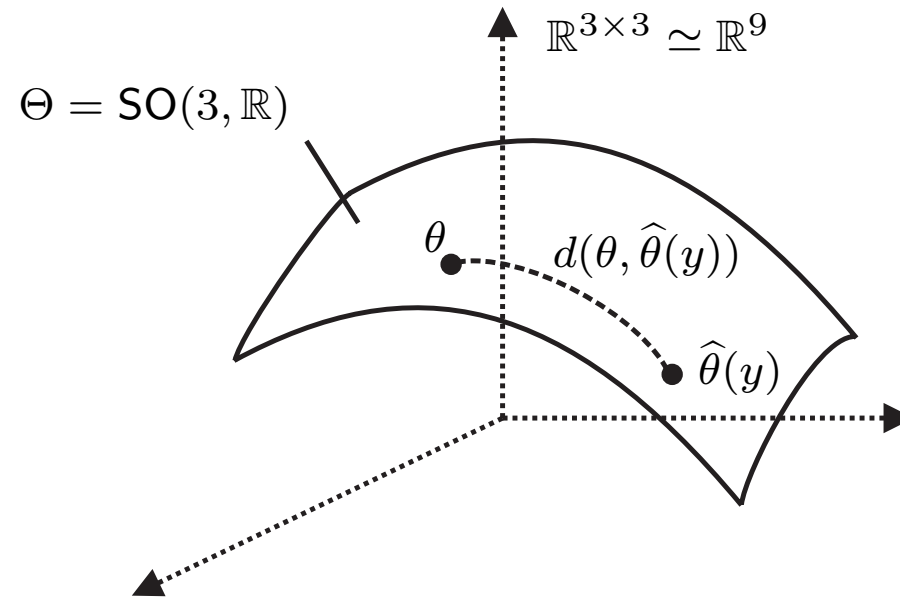
# Example: inference on $S^{p-1}$



## Example: inference on $SO(3, \mathbb{R})$

□  $SO(3, \mathbb{R})$  is the special orthogonal group:

$$SO(3, \mathbb{R}) = \left\{ Q \in \mathbb{R}^{3 \times 3} : Q^T Q = I_3, \det(Q) = 1 \right\}$$



□ **Geometry of  $\Theta$ :**  $d(\theta, \hat{\theta}(y)) = \sqrt{2} \operatorname{acos}(0.5[\operatorname{tr}(\theta^T \hat{\theta}(y)) - 1])$  and  $C = 1/8$

## Example: inference on $SO(3, \mathbb{R})$

□ **Observation:**  $Y = \theta X + W \in \mathbb{R}^{3 \times k}$  ( $k = 10$ )

▷  $\theta \in \Theta = SO(3, \mathbb{R})$ : unknown rotation matrix [Procrustean analysis]

▷  $X = [x_1 \ x_2 \ \cdots \ x_k]$ : constellation of known  $k$  landmarks in  $\mathbb{R}^3$  ( $XX^T = I_3$ )

▷  $W = [w_1 \ w_2 \ \cdots \ w_k]$ ,  $w_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2 I_3)$ : additive observation noise

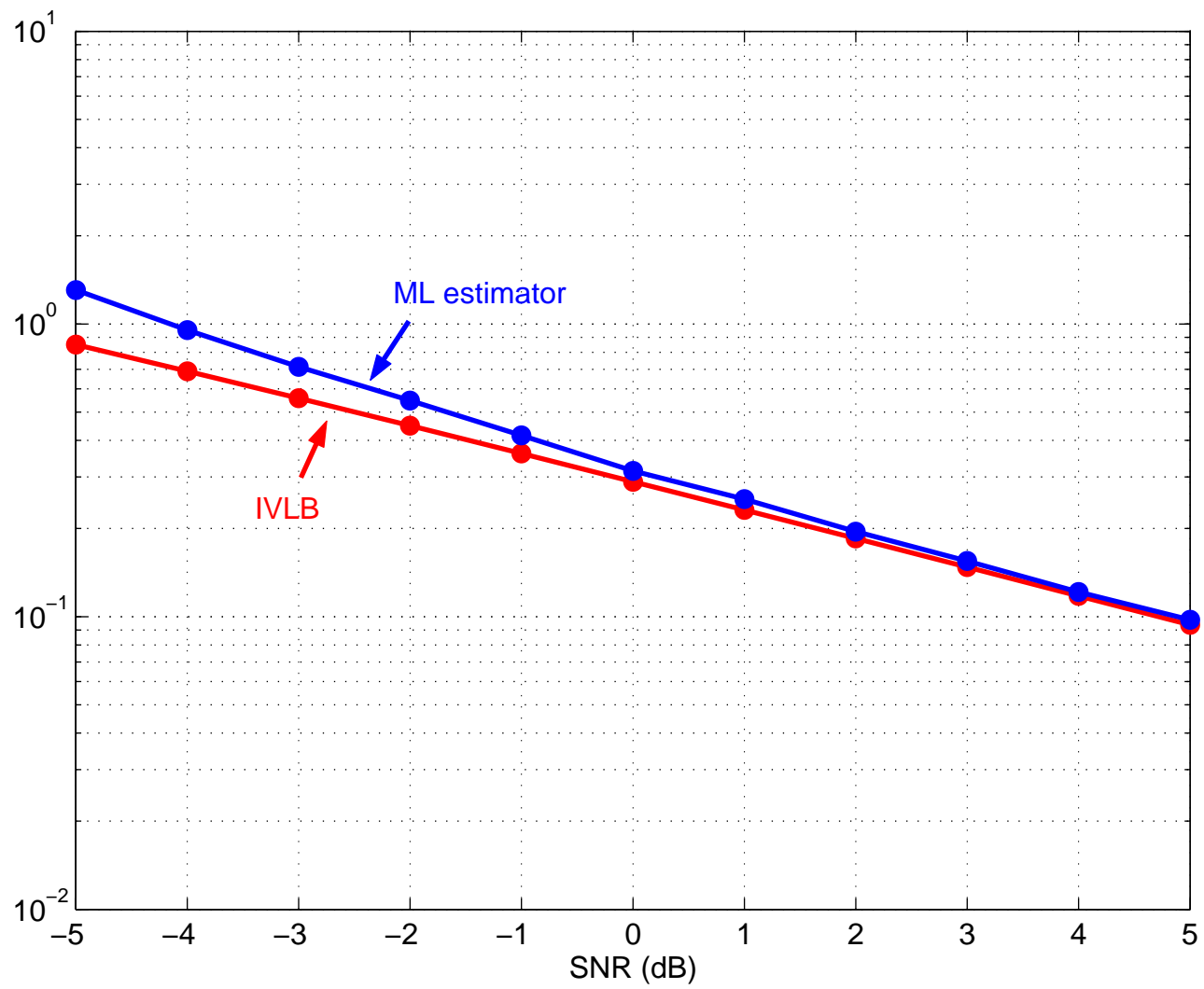
□ Maximum-likelihood estimator:

$$\hat{\theta}(Y) = \cdots \text{ (closed - form)}$$

□ Signal-to-noise ratio:

$$\text{SNR} = \frac{\mathbb{E} \left\{ \|\theta X\|^2 \right\}}{\mathbb{E} \left\{ \|W\|^2 \right\}} = \frac{1}{k \sigma^2}$$

# Example: inference on $SO(3, \mathbb{R})$



## Applications of DG: Bounds

### □ Bibliography:

- ◇ “Covariance, subspace, and intrinsic Cramér-Rao bounds,” S. Smith, IEEE Trans. on Signal Proc., vol. 53, no.5, May 2005
- ◇ “Intrinsic variance lower bound (IVLB): an extension of the Cramér-Rao bound to Riemannian manifolds”, J. Xavier and V. Barroso, IEEE Int. Conf. on Acoust., Sp. and Sig. Proc. (ICASSP), March 2005
- ◇ “The Riemannian geometry of certain parameter estimation problems with singular Fisher matrices”, J. Xavier and V. Barroso, IEEE Int. Conf. on Acoust., Sp. and Sig. Proc. (ICASSP), May 2004
- ◇ “Hilbert-Schmidt lower bounds for estimators on matrix Lie groups for ATR”, U. Grenander *et al.*, IEEE Trans. on Patt. Anal. and Mach. Intell., vol. 20, no. 8, pp. 790–801, August 1998
- ◇ “On the Cramér-Rao bound under parametric constraints”, P. Stoica *et al.*, IEEE Sig. Proc. Lett., vol. 5, no. 7, pp. 177–179, July 1998
- ◇ “Intrinsic analysis of statistical estimation”, J. Oller *et al.*, The Annals of Stat., vol. 23, no. 5, pp. 1562–1581, 1995
- ◇ “A Cramér-Rao type lower bound for estimators with values in a manifold”, H. Hendricks, Journal of Multivar. Anal., no. 38, pp. 245–261, 1991




## Course's Table of Contents

**Three main topics:**

- ▷ Topological manifolds
- ▷ Differentiable manifolds
- ▷ Riemannian manifolds

**Three layers of structure:**



Riemannian structure	Length of curves ; Geodesics ; Distance ; Connections ; etc
Differentiable structure	Tangent vectors; Smooth maps; Tensors; Integration ; etc
Topological structure	Boundary of sets; Convergent sequences; Continuous maps ; etc
Plain set	

## Course's Table of Contents

### **Topological manifolds:** "Introduction to Topological Manifolds", J. Lee, Springer-Verlag

- ◇ Ch.2: Topological spaces
- ◇ Ch.3: New spaces from old
- ◇ Ch.4: Connectedness and compactness

### **Smooth manifolds:** "Introduction to Smooth Manifolds", J. Lee, Springer-Verlag

- ◇ Ch.2: Smooth maps
- ◇ Ch.3: The tangent bundle
- ◇ Ch.5: Submanifolds
- ◇ Ch.7: Lie group actions
- ◇ Ch.8: Tensors
- ◇ Ch.9: Differential forms
- ◇ Ch.10: Integration on manifolds

## Course's Table of Contents

- **Riemannian manifolds:** “Riemannian Manifolds”, J. Lee, Springer-Verlag
  - ◇ Ch.3: Definitions and examples of Riemannian metrics
  - ◇ Ch.4: Connections
  - ◇ Ch.5: Riemannian geodesics

## Bibliography for the Course

### □ Topological manifolds

- ◇ “Introduction to Topological Manifolds”, J. Lee, Springer-Verlag, 2000
- ◇ “Introduction to Topology and Modern Analysis”, G. Simmons, 1963

### □ Smooth manifolds

- ◇ “Introduction to Smooth Manifolds”, J. Lee, Springer-Verlag, 2002
- ◇ “An Introduction to Differentiable Manifolds and Riemannian Geometry”, 2nd ed., W. Boothby, Academic Press, 1986
- ◇ “Manifolds, Tensor Analysis and Applications”, R. Abraham *et al.*, Springer-Verlag, 1988
- ◇ “A Comprehensive Introduction to Differential Geometry”, vol. I, M. Spivak, Publish or Perish, 1979
- ◇ “Lectures on Differential Geometry”, S. Chern, W. Chern and K. Lam, World Scientific, 1999

### □ Riemannian manifolds

- ◇ “Riemannian Manifolds”, J. Lee, Springer-Verlag
- ◇ “Riemannian Geometry”, M. Carmo, Birkhauser, 1992

## Bibliography

### Other references (introductory):

- ◇ “Differential Forms with Applications to the Physical Sciences”, H. Flanders, Dover, 1963
- ◇ “Differential Forms with Applications”, M. Carmo, Springer-Verlag, 1994

### Other references (advanced):

- ◇ “Riemannian Geometry”, S. Gallot, D. Hulin and J. Lafontaine, Springer-Verlag, 1987
- ◇ “A Comprehensive Introduction to DG”, vol.II-V, M. Spivak, Publish or Perish, 1979
- ◇ “Riemannian Geometry: A Modern Introduction”, I. Chavel, Cambridge Press, 1993
- ◇ “Riemannian Geometry and Geometric Analysis”, J. Jost, Springer-Verlag, 1998
- ◇ “Foundations of Differential Geometry”, vol. I-II, S. Kobayashi and K. Nomizu, Wiley 1969
- ◇ “DG, Lie Groups and Symmetric Spaces”, S. Helgason, Academic Press, 1978

### Many others. . .

## Grading

**Grade** = Homework (60%) + Project (40%)

**Homeworks:**

#	Received	Due
1	March, 29	April, 19
2	April, 19	May, 10
3	May, 10	May, 31
4	May, 31	June, 21

**Project (individual):** A paper will be assigned for each student to study

Output: public presentation of the paper

Start: May, 10 End: July, 31

**Discussion, questions, etc**