



Statistical Theories and Machine Learning Using Geometric Methods

メタデータ	言語: en 出版者: Osaka Central Advanced Mathematical Institute (OCAMI) Osaka Metropolitan University 公開日: 2024-05-27 キーワード (Ja): キーワード (En): Lie group, information geometry, exponential families, machine learning, geometric analysis, kernel methods 作成者: Tojo, Koichi, Nakashima, Hideto, Konno, Yoshihiko, Ishi, Hideyuki, Fukumizu, Kenji メールアドレス: 所属:
URL	https://doi.org/10.24729/0002000871

Osaka Central Advanced Mathematical Institute (OCAMI)
Osaka Metropolitan University
MEXT Joint Usage/Research Center on Mathematics and Theoretical Physics

OCAMI Reports Vol. 2 (2024)

Statistical Theories and Machine Learning Using Geometric Methods

Organized by
Koichi Tojo
Hideto Nakashima
Yoshihiko Konno
Hideyuki Ishi
Kenji Fukumizu

December 14–15, 2023

Abstract

This workshop was held on December 14–15, 2023 in order to connect researchers in several fields, in particular Statistics, Machine Learning and Mathematics, and to share problems and researches in these fields interdisciplinary.

2020 Mathematics Subject Classification.
20G05, 22F30, 53B12, 60E05, 62E10, 62H12, 62J07

Key words and Phrases.
Lie group, information geometry, exponential families,
machine learning, geometric analysis, kernel methods

Preface

This is a proceedings of the international workshop “Statistical Theories and Machine Learning Using Geometric Methods” held from December 14th to December 15th in 2023. This workshop aimed to connect researchers in several fields, in particular Statistics, Machine Learning and Mathematics, and to share problems and researches in these fields in an interdisciplinary manner.

This workshop was supported by Osaka Central Advanced Mathematical Institute (MEXT Promotion of Distinctive Joint Research Center Program JPMXP0723833165), Osaka Metropolitan University.

This workshop was held in a hybrid format. Domestic speakers were gathered in Academic Extension Center (Osaka Metropolitan University), and overseas speakers participated by Zoom. We had 11 talks, 9 of which were from Japan and the others were from abroad, and 51 people participated in this workshop.

Organizers

Koichi Tojo

RIKEN Center for Advanced Intelligence Project, Nihonbashi 1-chome Mitsui Building,
15th floor, 1-4-1 Nihonbashi, Chuo-ku, Tokyo 103-0027, Japan

Email address: `koichi.tojo@riken.jp`

Hideto Nakashima

Research Center for Statistical Machine Learning, The Institute of Statistical Mathematics,
10-3 Midori-cho, Tachikawa, Tokyo 190-8562, Japan

Email address: `hideto@ism.ac.jp`

Yoshihiko Konno

Department of Mathematics, Osaka Metropolitan University, 1-1, Gakunen-cho, Naka-ku,
Sakai-shi, 599-8531, Japan

Email address: `konno@omu.ac.jp`

Hideyuki Ishi

Department of Mathematics, Osaka Metropolitan University, 3-3-138, Sugimoto, Sumiyoshi-
ku, Osaka, 558-8585, Japan

Email address: `hideyuki-ishi@omu.ac.jp`

Kenji Fukumizu

Research Center for Statistical Machine Learning, The Institute of Statistical Mathematics,
10-3 Midori-cho, Tachikawa, Tokyo 190-8562, Japan

Email address: `fukumizu@ism.ac.jp`

Contents

Hiroto Inoue	
<i>Mean-variance joint statistic valued in a real Siegel domain</i>	1
Eren Mehmet Kiral	
<i>Bayesian Learning with Lie Groups</i>	17
Hajime Fujita	
<i>The generalized Pythagorean theorem on the compactifications of certain dually flat spaces via toric geometry</i>	31
Atsumi Ohara	
<i>Doubly autoparallel structure and curvature integrals: An application to iteration complexity analysis of convex optimization</i>	53
Adam Chojecki, Hideyuki Ishi	
<i>Uncovering Data Symmetries: Estimating Covariance Matrix in High-Dimensional Setting With ‘gips’ R Package</i>	84
Tomasz Skalski	
<i>Maximum likelihood estimation for discrete exponential families, its geometry and combinatorics</i>	99
Tomonari Sei	
<i>Some open problems on minimum information dependence models</i>	113
Tomonari Sei, Ushio Tanaka	
<i>Stein identity, Poincare inequality and exponential integrability on a metric measure space¹</i>	130
Hikaru Watanabe	
<i>Infinite dimensional parameterized measure models</i>	131
Eiki Shimizu	
<i>Neural-Kernel Conditional Mean Embeddings</i>	148
Satoshi Kuriki	
<i>Bonferroni method and tube method for heavy-tailed distributions¹</i>	169
Program	170

¹Slides are not included in this report

Mean-variance joint statistic valued in a real Siegel domain

Hiroto Inoue

Nishinippon Institute of Technology

Abstract

Let \bar{x} and s_x^2 be the mean and the variance of given data $x = (x_1, \dots, x_n) \in \mathbb{R}^n$. The main idea of this talk is to characterize the pair of statistics $x \mapsto (\bar{x}, s_x^2)$ as a quotient map $\mathbb{R}^n \rightarrow \mathbb{R}^n/H$ with the $n-1$ dimensional orthogonal group $H = O(n-1)$. This is realized by a combination of a quadratic form and the Cholesky decomposition, or the QR decomposition for matrices. Since the quotient space \mathbb{R}^n/H can be identified with the affine group $G = \mathbb{R}_+ \ltimes \mathbb{R}$ or a real Siegel domain Ω , the joint statistics (\bar{x}, s_x^2) is considered to be valued in G or Ω .

From this characterization, we find some fundamental properties and advantages of the joint statistic (\bar{x}, s_x^2) as follows.

- (\bar{x}, s_x^2) is invariant under the actions of H , especially under the permutations between x_1, \dots, x_n . More precisely, it is the maximal invariant statistic in the context of statistical inference.
- The quotient map is equivariant under the affine transformation, which is consequently followed by the transformation rule $ax + b \mapsto (a\bar{x} + b, a^2 s_x^2)$. Moreover, under the assumption of the normal distributions for the data $x \in \mathbb{R}^n$, the parameter estimators and test statistics are written in the product and inverse in the affine group G .
- The probability density functions (PDFs) of each statistic which is a function of the joint statistic (\bar{x}, s_x^2) can be calculated by using the integral formula on the space \mathbb{R}^n/H in terms of the Haar measure on G . The familiar distributions such as the χ^2 , t and F distributions are derived as the marginal ones. It enables us to describe the simultaneous hypothesis testing for the pair of parameters (μ, σ^2) of the normal distribution.

We notice that the PDF of the joint statistic (\bar{x}, s_x^2) is obtained from the Bartlett decomposition of the classical Wishart distribution restricted onto Ω . Then it might be meaningful to consider the same kind of statistic valued in other homogeneous cones or Siegel domains as our joint statistic.

Mean-variance joint statistic valued in a real Siegel domain

Hiroto Inoue

Nishinippon Institute of Technology

OCAMI Workshop "Statistical Theories and Machine Learning
Using Geometric Methods" Dec. 14, 2023

Mean-variance joint statistic valued in a real Siegel domain

Basic idea

For a set of real data $x_1, x_2, \dots, x_n \in \mathbb{R}$, we consider a $2 \times n$ matrix

$$x = \begin{pmatrix} x_1 & \cdots & x_n \\ 1 & \cdots & 1 \end{pmatrix} \in \mathbb{R}^{2 \times n}.$$

Consider the QR decomposition of x in the following form

$$x = RU; \quad R = \begin{pmatrix} s & t \\ 0 & 1 \end{pmatrix}, \quad s > 0, \quad UU^T = nI_2.$$

Then, it can be checked that

$$s = s_x := \left\{ \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})^2 \right\}^{\frac{1}{2}}, \quad t = \bar{x} := \frac{1}{n} \sum_{k=1}^n x_k.$$

Mean-variance joint statistic valued in a real Siegel domain

By the definition, the factor $R = R(x)$ has the following properties.

■ **Equivariance**

$$R(gx) = gR(x), \quad \forall g = \begin{pmatrix} a & b \\ 0 & 1 \end{pmatrix}, \quad a > 0.$$

This is equivalent to $(s_{ax+b}, \overline{ax+b}) = (as_x, a\bar{x} + b)$.

■ **Invariance**

$$R(xh) = R(x), \quad \forall h \in O(n) \text{ s.t. } (1, \dots, 1)h = (1, \dots, 1).$$

Especially, $R(x)$ is symmetric in x_1, \dots, x_n .

▷ We will obtain the **test statistics** and their **PDFs** in a group theoretical way.

3 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Mean-Variance Joint Statistic

1 Mean-Variance Joint Statistic

- Definition
- Decomposition of volume element

2 Probability distributions on G_0

3 Hypothesis test of normal distributions

- One-sample testing
- Two-sample testing

4 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Mean-Variance Joint Statistic

└ Definition

Setting

- **Data space** : a set $E_n = E_{n,r}$ ($n \geq 2, r \geq 1$) of matrices s.t.

$$E_n = \left\{ x = \begin{pmatrix} x_1 & \cdots & x_n \\ 1 & \cdots & 1 \end{pmatrix} \in \mathbb{R}^{(r+1) \times n} \mid \text{rank } x = r+1 \right\}^1$$

- **Range of statistic** : a region $\Omega_n \subset \text{Sym}_{r+1}^+(\mathbb{R})$,

$$\Omega_n = \left\{ \mathbf{w} = (w_{ij}) \in \text{Sym}_{r+1}^+(\mathbb{R}) \mid w_{r+1,r+1} = n \right\}$$

- **Quadratic statistic** : a map $\mathcal{Q} : E_n \rightarrow \Omega_n$,

$$\mathcal{Q}(x) := xx^T \quad (x \in E_n)$$

¹The condition of rank can be omitted if we ignore the zero-measure set in the following argument.

5 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Mean-Variance Joint Statistic

└ Definition

- **Transformation group** : a group $G_0 = B \ltimes \mathbb{R}^r$ by

$$(s, t) \cdot (a, b) := (sa, sb + t) \quad ((s, t), (a, b) \in B \times \mathbb{R}^r),$$

where

$$B = \left\{ s \in \text{GL}_r(\mathbb{R}) \mid \begin{array}{ll} s_{ij} = 0 & (i < j), \\ s_{ii} > 0 & (i = 1, \dots, r) \end{array} \right\},$$

the group of the upper triangular matrices with positive diagonals.

6 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Mean-Variance Joint Statistic

└ Definition

■ Actions of G_0 on E_n, Ω_n by

$$G_0 \curvearrowright E_n : (s, t)x = \begin{pmatrix} s & t \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 & \cdots & x_n \\ 1 & \cdots & 1 \end{pmatrix} \quad (x \in E_n)$$

$$G_0 \curvearrowright \Omega_n : (s, t)\mathbf{w} = \begin{pmatrix} s & t \\ 0 & 1 \end{pmatrix} \mathbf{w} \begin{pmatrix} s & t \\ 0 & 1 \end{pmatrix}^T \quad (\mathbf{w} \in \Omega_n)$$

Notice that $G_0 \curvearrowright \Omega_n$ is **simply transitive** due to the Cholesky decomposition.

7 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Mean-Variance Joint Statistic

└ Definition

Proposition 1.1 (QR decomposition)

Any $x \in E_n$ is uniquely decomposed as follows.

$$x = gu, \quad g \in G_0, \quad u \in \mathcal{Q}^{-1}(ne). \quad (1)$$

Here $e \in \Omega_n$ is the identity matrix.

So the element g is uniquely determined from x and we denote it by $g = R(x)$.

Proposition 1.2

For $x \in E_n$, it holds that $R(x) = (s_x, \bar{x})$, where

$$s_x s_x^T = \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})(x_k - \bar{x})^T, \quad \bar{x} = \frac{1}{n} \sum_{k=1}^n x_k \quad (2)$$

8 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Mean-Variance Joint Statistic

└ Definition

By the definition of R , we see again

■ G_0 -Equivariance

$$R(gx) = gR(x), \quad \forall g \in G_0.$$

■ $O(n-1)$ -Invariance

$$R(xh) = R(x), \quad \forall h \in O(E_n),$$

where

$$O(E_n) = \{h \in O(n) \mid E_n h = E_n\} \cong O(n-1).$$

$$\begin{array}{ccc} E_n & \xrightarrow{\cong} & G_0 \times \mathcal{Q}^{-1}(ne) \\ \downarrow \mathcal{Q} & \searrow R & \downarrow \text{pr} \\ \Omega_n & \xrightarrow{\cong} & G_0 \end{array}$$

9 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Mean-Variance Joint Statistic

└ Decomposition of volume element

Notations for volume elements For $\mathbf{s} = (s_1, \dots, s_r) \in \mathbb{R}^r$ and $\mathbf{w} \in \Omega_n$, we put²

$$\Delta_{\mathbf{s}}(\mathbf{w}) := \Delta_1(\mathbf{w})^{s_1-s_2} \Delta_2(\mathbf{w})^{s_2-s_3} \dots \Delta_r(\mathbf{w})^{s_r},$$

$$\Delta_k(\mathbf{w}) = \det(\mathbf{w}_{ij})_{r-k+1 \leq i, j \leq r+1}.$$

We notice that

$$\Delta_{\mathbf{s}}(ge) = \prod_{i=1}^r (s_{ii}^2)^{s_{r-i+1}} \quad (g = (s, t) \in G_0).$$

For convenience of the index, we let

$$[r] = (1, \dots, r), \quad \mathbf{s} + \lambda = \mathbf{s} + (\lambda, \dots, \lambda) \quad (\mathbf{s} \in \mathbb{R}^r, \lambda \in \mathbb{R}).$$

²Notations by Faraut-Korányi [1]

Mean-variance joint statistic valued in a real Siegel domain

└ Mean-Variance Joint Statistic

└ Decomposition of volume element

We take a left Haar measure $d\mu(g)$ on G_0 as

$$d\mu(g) = \Delta_{-\frac{1}{2}[r]}(ge) \prod_{i=1}^r s_{ii}^{-1} ds_{ii} \prod_{i < j}^r ds_{ij} \prod_{i=1}^r dt_i, \quad g = (s, t).$$

It can be checked the following transformation rule.

$$\begin{cases} d\mu(h_1 g h_2) = \Delta_G(h_2) d\mu(g) & (h_1, h_2 \in G_0), \\ d\mu(g^{-1}) = \Delta_G(g)^{-1} d\mu(g), \end{cases}$$

where $\Delta_G(g) = \Delta_{\frac{1}{2}r-[r]}(ge)$.

11 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Mean-Variance Joint Statistic

└ Decomposition of volume element

Next we denote by dx the Lebesgue measure on $E_n \subset \mathbb{R}^{rn}$;

$$dx = \prod_{k=1}^n dx_k, \quad x = \begin{pmatrix} x_1 & \cdots & x_n \\ 1 & \cdots & 1 \end{pmatrix}$$

We notice the transformation rule

$$d(hx) = |h|^n dx \quad (h \in G_0), \quad \text{where } |h| = \Delta_{+\frac{1}{2}}(he).$$

12 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Mean-Variance Joint Statistic

└ Decomposition of volume element

Proposition 1.3

Under the relation $x = gu$ in (1), there exists a volume element du on $Q^{-1}(ne)$ such that

$$dx = |g|^n d\mu(g) du. \quad (3)$$

The integral of du is evaluated as follows.

$$\int_{Q^{-1}(ne)} du = \frac{\sqrt{2\pi}^{rn}}{c_n}, \quad (4)$$

$$c_n = \frac{1}{2^r} \left(\frac{n}{2}\right)^{-\frac{1}{2}rn} \sqrt{\pi}^{\frac{1}{2}r(r+1)} \prod_{i=1}^r \Gamma\left(\frac{1}{2}(n-i)\right).$$

Here $\Gamma(x)$ be the Gamma function.

13 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Probability distributions on G_0

1 Mean-Variance Joint Statistic

- Definition
- Decomposition of volume element

2 Probability distributions on G_0

3 Hypothesis test of normal distributions

- One-sample testing
- Two-sample testing

14 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Probability distributions on G_0

Let $\Phi_G(g)$ denote the PDF of random variables G on G_0 ;

$$P(G \in \mathcal{V}) = \int_{\mathcal{V} \subset G_0} \Phi_G(g) d\mu(g).$$

We define an expectation $E[G] \in G_0$ by

$$E[G]e = \int_{G_0} (he) \Phi_G(h) d\mu(h).$$

Proposition 2.1

For a random variable G on G_0 and $h \in G_0$, it holds that

$$E[hG] = hE[G].$$

15 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Probability distributions on G_0

Distributions on G_0 derived from normal distributions

Suppose \mathbf{X} is a random variable on E_n and

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 & \cdots & \mathbf{X}_n \\ 1 & \cdots & 1 \end{pmatrix}, \quad \mathbf{X}_i \sim N_r(t, ss^T) \quad \text{i.i.d..}$$

We denote it by $\mathbf{X} \sim N_r(h)^n$, where $h = (s, t) \in G_0$.

■ Parameter change

$$g \in G_0, \mathbf{X} \sim N_r(h)^n \Rightarrow g\mathbf{X} \sim N_r(gh)^n.$$

■ Unbiased estimator

$$\mathbf{X} \sim N_r(h)^n \Rightarrow E[R(\mathbf{X})] = h.$$

16 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Probability distributions on G_0

The PDF of $\mathbf{X} \sim N_r(g)^n$ is written as

$$f_{\mathbf{X}}(x) = \frac{1}{\sqrt{2\pi}^{rn}} e^{-\frac{1}{2}\text{tr}\mathcal{Q}(g^{-1}x) + \frac{1}{2}n}|g|^{-n}.$$

The PDF of $R(\mathbf{X})$ is obtained as a marginal distribution.

Proposition 2.2

If $\mathbf{X} \sim N_r(e)^n$, then $\Phi_{R(\mathbf{X})}(h)$ is given by

$$\Phi_n(h) := \frac{1}{c_n} e^{-\frac{1}{2}n\text{tr}(he) + \frac{1}{2}n}|h|^n. \quad (5)$$

17 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Probability distributions on G_0

Proposition 2.3

If $\mathbf{X} \sim N_r(e)^n$, $\mathbf{Y} \sim N_r(e)^m$ independently, then $\Phi_{R(\mathbf{X})^{-1}R(\mathbf{Y})}(h)$ is given as

$$\Phi_{n,m}(h) := \frac{c_{n+m}}{c_n c_m} |h|^m \Delta_{[r] - \frac{1}{2}(n+m+r)} \left(\frac{n}{n+m} e + \frac{m}{n+m} h e \right). \quad (6)$$

Sketch of proof. It is obtained by calculating the convolution of the PDFs

$$\Phi_{R(\mathbf{X})^{-1}R(\mathbf{Y})}(h) = \int_{G_0} \Phi_n(g) \Phi_m(gh) d\mu(g).$$

□

18 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Probability distributions on G_0

Marginal distributions on \mathbb{R}_+, \mathbb{R} when $r = 1$

When $r = 1$, for $h = (s, t) \in G_0$ we have

$$\Phi_{R(\mathbf{X})}(h) = \frac{1}{c_n} e^{-\frac{1}{2}n(s^2+t^2)} s^n, \quad c_n = \frac{\sqrt{\pi}}{2} \left(\frac{n}{2}\right)^{-\frac{n}{2}} \Gamma\left(\frac{n-1}{2}\right).$$

Corollary 2.1

For $R(\mathbf{X}) = (S_X, \overline{\mathbf{X}})$, it holds that

$$\begin{cases} \Phi_{S_X}(s) = \frac{\sqrt{2\pi/n}}{c_n} e^{-\frac{1}{2}ns^2} s^{n-2} & : \chi^2 \text{ dist.} \\ \Phi_{\overline{\mathbf{X}}}(t) = \frac{1}{\sqrt{2\pi/n}} e^{-\frac{1}{2}nt^2} & : \text{normal dist.} \end{cases}$$

19 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Probability distributions on G_0

We have

$$\Phi_{R(\mathbf{X})^{-1}}(h) = \Phi_n(h^{-1}) \Delta_G(h)^{-1} = \frac{1}{c_n} e^{-\frac{1}{2s^2}n(1+t^2)} s^{-n+1}$$

Corollary 2.2

For $R(\mathbf{X})^{-1} = (S_X^{-1}, -S_X^{-1}\overline{\mathbf{X}})$, it holds that

$$\begin{cases} \Phi_{S_X^{-1}}(s) = \frac{\sqrt{\pi/n}}{c_n} e^{-\frac{n}{2s^2}} s^{-n-1} & : \text{inv. } \chi^2 \text{ dist.} \\ \Phi_{-S_X^{-1}\overline{\mathbf{X}}}(t) = \frac{1}{B\left(\frac{n-1}{2}, \frac{1}{2}\right)} (1+t^2)^{-\frac{n}{2}} & : t \text{ dist.} \end{cases}$$

20 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Probability distributions on G_0

We have

$$\Phi_{R(\mathbf{X})^{-1}R(\mathbf{Y})}(h) = \frac{c_{n+m}}{c_n c_m} s^m \left(\frac{n}{n+m} + \frac{m}{n+m} s^2 + \frac{nm}{(n+m)^2} t^2 \right)^{-\frac{n+m-1}{2}}.$$

Let $B(x, y)$ be the beta function.

Corollary 2.3

For $R(\mathbf{X})^{-1}R(\mathbf{Y}) = (S_{\mathbf{X}}^{-1}S_{\mathbf{Y}}, S_{\mathbf{X}}^{-1}(\bar{\mathbf{Y}} - \bar{\mathbf{X}}))$, it holds that

$$\begin{cases} \Phi_{S_{\mathbf{X}}^{-1}S_{\mathbf{Y}}}(s) = \frac{2n^{\frac{n-1}{2}} m^{\frac{m-1}{2}}}{B\left(\frac{n-1}{2}, \frac{m-1}{2}\right)} s^{m-2} \left(\frac{1}{n+ms^2} \right)^{\frac{n+m}{2}-1} & : F \text{ dist.} \\ \Phi_{S_{\mathbf{X}}^{-1}(\bar{\mathbf{Y}} - \bar{\mathbf{X}})}(t) = \frac{1}{B\left(\frac{n-1}{2}, \frac{1}{2}\right)} \left(\frac{m}{n+m} \right)^{\frac{1}{2}} \left(1 + \frac{m}{n+m} t^2 \right)^{-\frac{n}{2}} & : t \text{ dist.} \end{cases}$$

21 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Probability distributions on G_0

On the orthogonal factor $U(x)$, when $r = 1$

For $x \in E_n$, we write its decomposition as

$$x = R(x)U(x), \quad U(x) \in \mathcal{Q}^{-1}(ne).$$

For $x, y \in E_n$, the determinant of transposed product

$$r_{xy} := \det(U(x)U(y)^T)$$

is the [correlation coefficient](#) of x and y .

22 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Probability distributions on G_0

Assume that $\mathbf{X} \sim N_r(e)^n$. Then,

- $R(\mathbf{X})$ and $U(\mathbf{X})$ are independent of each other.
 - $U(\mathbf{X}) = ((x_i - \bar{x})/\sigma_x)$ follows the uniform distribution on a sphere S^{n-1} .³
- ⇒ For $y \in E_n$, $\mathbf{R} = \det(U(\mathbf{X})U(y)^T)$ has a following PDF

$$f(r) = c(1 - r^2)^{\frac{n}{2}-1} \quad (-1 \leq r \leq 1), \quad c > 0.$$

⇒ $\mathbf{T} := \frac{\sqrt{n-2}\mathbf{R}}{\sqrt{1-\mathbf{R}^2}}$ follows the t distribution (used in correlation analysis).

³Each $(x_i - \bar{x})/\sigma_x$ is called an ancillary statistic.

Mean-variance joint statistic valued in a real Siegel domain

└ Hypothesis test of normal distributions

1 Mean-Variance Joint Statistic

- Definition
- Decomposition of volume element

2 Probability distributions on G_0

3 Hypothesis test of normal distributions

- One-sample testing
- Two-sample testing

Mean-variance joint statistic valued in a real Siegel domain

└ Hypothesis test of normal distributions

└ One-sample testing

One-sample testing

We consider one-sample test for a normal distribution $N_1(h)$,

$$H_0 : h = h_0$$

with a set of sample data $\mathbf{X} \sim N_1(h)^n$.

- We define a G_0 -valued **test statistic** U as

$$U = h_0^{-1} R(\mathbf{X}).$$

Now H_0 implies $\Phi_U(g) = \Phi_n(g)$.

- We take an **acceptance region** \mathcal{V} : a neighbor of the identity element $e \in G_0$ such that

$$\int_{\mathcal{V}} \Phi_n(g) d\mu(g) = 1 - \alpha \quad (\alpha \text{ is a significance level})$$

25 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Hypothesis test of normal distributions

└ One-sample testing

It realizes the following testings depending on the choice of \mathcal{V} .

- $\mathcal{V} = \mathcal{I} \times \mathbb{R}$:

$$U \in \mathcal{V} \Leftrightarrow \frac{S_{\mathbf{X}}}{\sigma_0} \in \mathcal{I} \quad (\chi^2\text{-test for } \sigma^2)$$

- $\mathcal{V} = \mathbb{R}_+ \times \mathcal{J}$:

$$U \in \mathcal{V} \Leftrightarrow \frac{\overline{\mathbf{X}} - \mu_0}{\sigma_0} \in \mathcal{J} \quad (Z\text{-test for } \mu \text{ with known } \sigma^2)$$

- $\mathcal{V} = (\mathbb{R}_+ \times \mathcal{J})^{-1} := \{h \in G_0; h^{-1} \in \mathbb{R}_+ \times \mathcal{J}\}$

$$U \in \mathcal{V} \Leftrightarrow -\frac{\overline{\mathbf{X}} - \mu_0}{S_{\mathbf{X}}} \in \mathcal{J} \quad (t\text{-test for } \mu \text{ with unknown } \sigma^2)$$

Here $R(\mathbf{X}) = (S_{\mathbf{X}}, \overline{\mathbf{X}})$, $h_0 = (\sigma_0, \mu_0)$.

26 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Hypothesis test of normal distributions

└ Two-sample testing

Two-sample testing

We consider a two-sample test for normal distributions $N_1(h_x)$, $N_1(h_y)$,

$$H_0 : h_x = h_y$$

with two sets of data $\mathbf{X} \sim N_1(h_x)^n$, $\mathbf{Y} \sim N_1(h_y)^m$.

- We define a G_0 -valued **test statistic** U as

$$U = R(\mathbf{X})^{-1}R(\mathbf{Y}).$$

Now H_0 implies $\Phi_U(g) = \Phi_{n,m}(g)$.

- We take an **acceptance region** \mathcal{V} : a neighbor of the identity element $e \in G_0$ such that

$$\int_{\mathcal{V}} \Phi_{n,m}(g) d\mu(g) = 1 - \alpha \quad (\alpha \text{ is a significance level})$$

27 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Hypothesis test of normal distributions

└ Two-sample testing

Depending on the choice of \mathcal{V} , it realizes the following testings.

- $\mathcal{V} = \mathcal{I} \times \mathbb{R}$:

$$U \in \mathcal{V} \Leftrightarrow \frac{S_Y}{S_X} \in \mathcal{I} \quad (\text{F-test for the ratio } \sigma_x^2/\sigma_y^2)$$

- $\mathcal{V} = \mathbb{R}_+ \times \mathcal{J}$:

$$U \in \mathcal{V} \Leftrightarrow \frac{\overline{Y} - \overline{X}}{S_X} \in \mathcal{J} \quad (\text{t-test for the difference } \mu_x - \mu_y)^4$$

Here $R(\mathbf{X}) = (S_X, \overline{X})$, $R(\mathbf{Y}) = (S_Y, \overline{Y})$.

⁴This statistic is however not common for testing the difference $\mu_x - \mu_y$.

Mean-variance joint statistic valued in a real Siegel domain

└ Hypothesis test of normal distributions

└ Two-sample testing

Remarks

- We notice the similarity of $Q(\mathbf{X}), R(\mathbf{X})$ to the [Wishart matrix](#) and its [Bartlett decomposition](#).
- Ω_n , the range of $Q(\mathbf{X})$, is isomorphic to the following convex domain D ,

$$\Omega_n \cong D = \left\{ (\Sigma, \mu) \in \text{Sym}_n^+(\mathbb{R}) \times \mathbb{R}^n \mid \Sigma - \mu\mu^T \in \text{Sym}_n^+(\mathbb{R}) \right\}$$

It is one of so-called [real Siegel domains](#).

▷ [Future work](#) Take other data space E to define a similar statistic $\tilde{Q} : E^{\oplus n} \rightarrow \Omega$ valued in a convex domain Ω .

29 / 30

Mean-variance joint statistic valued in a real Siegel domain

└ Hypothesis test of normal distributions

└ Two-sample testing

References



J. Faraut and A. Korányi.

Analysis on symmetric cones.

Oxford mathematical monographs. Clarendon Press, Oxford University Press, 1994.

30 / 30

Bayesian Learning with Lie Groups

Eren Mehmet Kiral

Abstract

There is uncertainty in the data, and there can be more than one explanation for the same data. Two models can distinguish between a cat picture and a dog picture but one may focus on the fur texture, while the other may make its decisions mostly based on the pointedness of the ears. These can be equivalently good explanations of the data. In this talk I will talk about how we can learn multiple models “at the same time”, by introducing some uncertainty to our model parameters. I will present a way of updating probability distributions on the model parameters using Lie Groups, and how desired structures can be preserved by the action of a Lie group.

This is joint work with Thomas Möllenhoff and M. Emtiyaz Khan. It was presented in the AISTATS 2023 conference.

Lie Group Bayesian Learning Rule

E. Mehmet Kiral, joint w. Thomas Möllenhoff and M. Emtiyaz Khan.

Dec 14, 2023 OCAMI Workshop
Statistical Theories and Machine Learning Using Geometric Methods



The work is supported mainly by the Bayes-duality project, JST CREST Grant Number JPMJCR2112.

Presented at **AISTATS 2023**, arxiv # 2303.04397

Navigation icons: back, forward, search, etc.

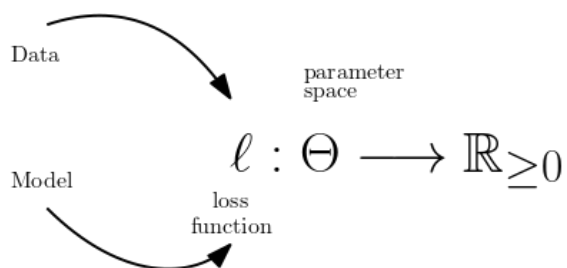
Kiral, Möllenhoff, Khan

Lie Group Bayesian Learning Rule

Dec 2023

1 / 25

The classical and Bayesian learning setups



Navigation icons: back, forward, search, etc.

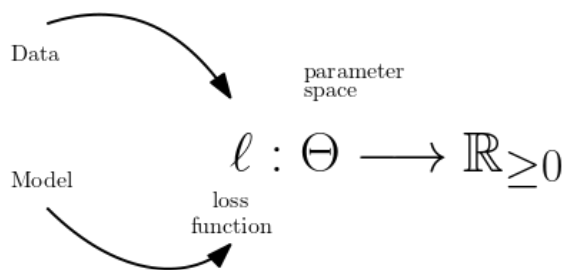
Kiral, Möllenhoff, Khan

Lie Group Bayesian Learning Rule

Dec 2023

2 / 25

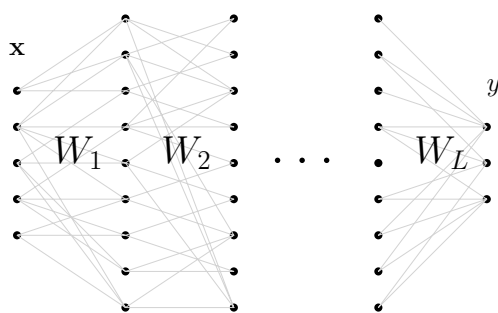
The classical and Bayesian learning setups



Classically: find $\theta^* \in \Theta$ minimizing ℓ .

Bayesian : find a distribution $q \in \mathcal{P}(\Theta)$

Example of loss function



$$f(\mathbf{x}, \theta) = W_L(\sigma(W_{L-1}(\cdots \sigma(W_1 \mathbf{x}))))$$

e.g. $\sigma : x_i \mapsto \max\{0, x_i\}$ componentwise. The tunable parameters are

$$\theta = (W_1, W_2, \dots, W_L) \in \mathbb{R}^P = \Theta.$$

For data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1, \dots, N}$ we search for a parameter $\theta \in \Theta$ s.t. $f(\mathbf{x}_i, \theta) \approx y_i$ for all i .

Let $c(\bar{y}, y) = \|y - \bar{y}\|^2$, say. Then using gradient, minimize,

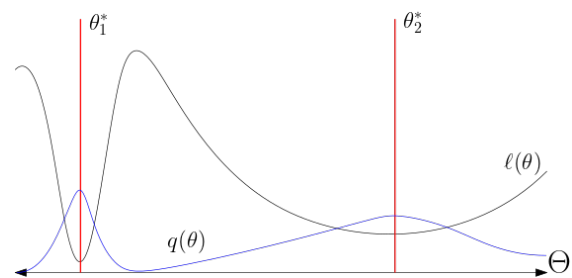
$$\ell(\theta) = \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}} \underbrace{c(f(\mathbf{x}_i, \theta), y_i)}_{\ell_i(\theta)} + \underbrace{\frac{1}{2} \lambda \|\theta\|^2}_{R(\theta)}.$$

Classical vs. Bayesian learning

The loss function is highly nonconvex. Usually

$$\ell(\theta) = \sum_{i=1}^N \ell_i(\theta) + R(\theta)$$

where $\ell_i(\theta)$ is the loss contribution from the i^{th} data point and $R(\theta)$ regularizer.



θ_1^* and θ_2^* are both equally valid explanations of the same data.

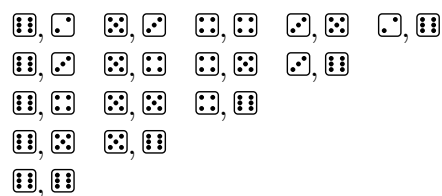
A distribution over the data considers both explanations “*at the same time*”.

Betting it all on one outcome

Say two dice are thrown and I tell you that the sum is greater than 7. satisfies this.

We could say the result was definitely .

But there are a total of 15 possibilities



It is much more sensible to say it is one of these 15 outcomes, with equal probability.

(principle of indifference, principle of maximum entropy)

The Bayesian Learning Problem

$\ell(\theta)$, a loss function on model parameters $\theta \in \Theta$. Pick a base measure ν on Θ . We solve

$$q_* \in \arg \min_{q \in \mathcal{Q}} \mathbb{E}_q[\ell] - \tau \mathcal{H}(q)$$

for some family of distributions $\mathcal{Q} \subseteq \mathcal{P}_\nu(\Theta) = \{q(\theta)d\nu(\theta)\}$ on the parameters.

- The expectation $\mathbb{E}_q[\ell] = \int_\Theta \ell(\theta)q(\theta)d\nu(\theta)$ prefers regions with low loss.
- The entropy $\mathcal{H}_\nu(q) = - \int_\Theta q(\theta) \log q(\theta)d\nu(\theta)$ prefers a higher spread of q .
- The temperature $\tau > 0$ is a balancing term.

Constrained maximization: Statistical mechanics interpretation

Assume θ to be a kind of “microstate” with energy level $\ell(\theta)$. So Θ is some “state space”.

Statistical mechanics: Assume a distribution of the microstates (across “particles”) maximizing entropy, constrained to have expected energy $\leq E_0$.

Lagrange multiplier $\beta \geq 0$:

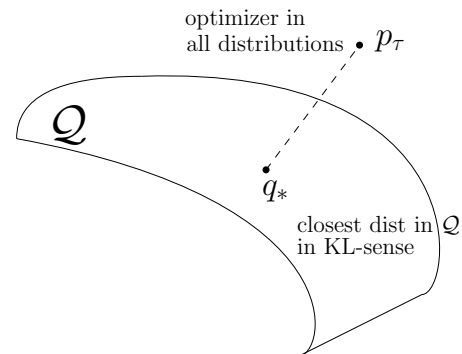
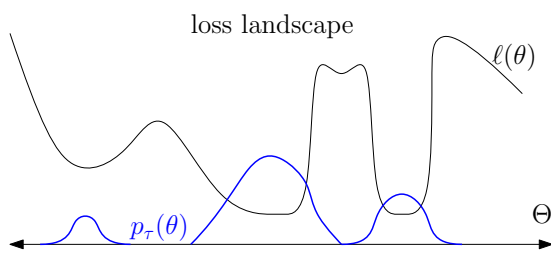
$$\arg \min_{q \in \mathcal{P}_\nu(\Theta)} -\mathcal{H}_\nu(q) + \beta(\mathbb{E}_{q d\nu}[\ell] - E_0) = \arg \min_{q \in \mathcal{P}_\nu(\Theta)} \mathbb{E}_{q d\nu}[\ell] - \frac{1}{\beta} \mathcal{H}_\nu(q)$$

$\tau = \frac{1}{\beta}$ corresponds to the thermodynamical notion of temperature.

The exact posterior.

If $\mathcal{Q} = \mathcal{P}_\nu(\Theta)$ then there is a unique minimizer $p_\tau(\theta) \propto e^{-\frac{1}{\tau}\ell(\theta)}$:

$$\arg \min_{q \in \mathcal{Q}} \mathbb{E}_{q \text{d}\nu}[\ell] - \tau \mathcal{H}(q) = \arg \min_{q \in \mathcal{Q}} \mathbb{D}_\nu(q \| p_\tau).$$



Minimize the objective $\mathcal{E}(q) := \mathbb{D}(q \| p_\tau)$ for $q \in \mathcal{Q}$...

...an approximate Bayesian solution.

Previously... BLR: The Bayesian Learning Rule, [KR21]¹

[KR21] take \mathcal{Q} as exponential families are $q_\lambda(\theta) \propto e^{-\lambda^\top T(\theta)}$.
 $T : \Theta \rightarrow \mathbb{R}^d$ is a *sufficient statistic*, λ are *natural parameters*.

$$\lambda \leftarrow \lambda - \alpha F(\lambda)^{-1} \nabla_\lambda \mathcal{E}(q_\lambda)$$

BLR is Natural Gradient Descent on λ parameters.

Gaussians, Exponential distributions, Gamma, inverse Gamma, Wishart, von-Mises, etc.

$\alpha > 0$ step size
 $F(\lambda)$ the Fisher matrix

- Issue 1** The candidates \mathcal{Q} is required to be an exponential family,
- Issue 2** Not every λ is allowed as a natural parameter, and the linear update rule could overshoot the constraints.
- Issue 3** Computing $\nabla_\lambda \mathcal{E}(q_\lambda)$ is not efficient in general but for special exponential families.

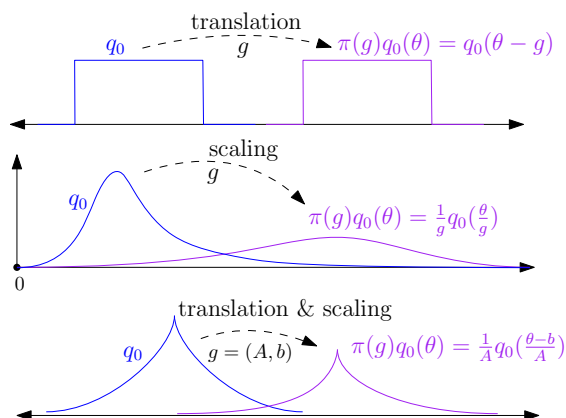
¹[KR21]: Khan, M. E. and Rue H., *The Bayesian Learning Rule*

Parametrizing \mathcal{Q} by groups.

Group G acting on the parameter set Θ , also acts on distributions on Θ .
 \mathcal{Q} is formed as the orbit of such an action for any base distribution q_0 :

$$\mathcal{Q} = \{\pi(g)q_0 : g \in G\}.$$

Call $q_g = \pi(g)q_0$.



- $G = (\mathbb{R}, +)$, $\Theta = \mathbb{R}$,

- $G = (\mathbb{R}_{>0}, \times)$, $\Theta = \mathbb{R}_{>0}$,

- $G = \text{Aff}(\mathbb{R}) = \mathbb{R}_{>0} \ltimes \mathbb{R}$, $\Theta = \mathbb{R}$

Navigation icons: back, forward, search, etc.

Optimization on the group

We now solve

$$\arg \min_{g \in G} \mathcal{E}(q_g) = \arg \min_{g \in G} \int_{\Theta} q_g \log \left(\frac{q_g}{e^{-\frac{1}{\tau} \ell}} \right)$$

Given $X \in \mathfrak{g} = T_e G$ the differential in the direction of X is

$$\left. \frac{d}{dt} \mathcal{E}(q_{ge^{tX}}) \right|_{t=0} = \underbrace{\left. \frac{d}{dt} \int_{\Theta} q_{ge^{tX}}(\theta) \frac{1}{\tau} \ell(\theta) d\nu(\theta) \right|_{t=0}}_{\text{data contribution}} + \underbrace{\left. \int_{\Theta} q_{ge^{tX}}(\theta) \log q_g(\theta) d\nu(\theta) \right|_{t=0}}_{\text{entropy contribution}}$$

The data contribution can be rewritten as

$$\int_{\Theta} q_g(\theta) (\nabla_{\theta} \ell(\theta))^{\top} (\text{Ad}_g(X) \cdot \theta) d\nu(\theta) \approx \frac{1}{K} \sum_{\substack{i=1 \\ \theta_i \sim q_g}}^K \nabla \ell(\theta_i)^{\top} (\text{Ad}_g(X) \cdot \theta_i)$$

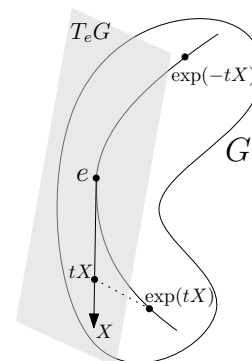
Navigation icons: back, forward, search, etc.

Important aspects of a Lie-Group.

A Lie group is also a **smooth manifold** and **directions** makes sense.
Can take derivatives of paths on the group, giving tangent vectors.

There is a special map $\exp : T_e G \rightarrow G$, which works as a global **retraction**. Multiplying the whole path by a $g \in G$ carries it to a neighborhood of g . So, directions at g can also be parametrized by the tangents at identity.

But also other retractions can be created using the exponential map.

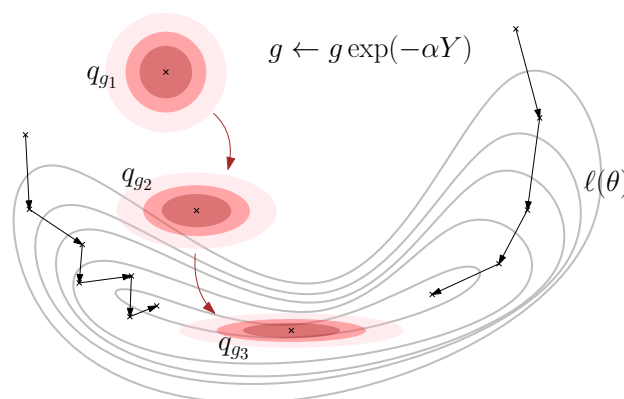


A rich mine of examples of Lie Groups are **matrix groups** (e.g. invertible matrices $GL(n)$, orthogonal matrices $O(n)$, matrices with determinant one $SL(n)$) with matrix multiplication.

Navigation icons: back, forward, search, etc.

Classical Learning vs. Learning via Group

The *point based* gradient descent updates parameters: $\theta \leftarrow \theta - \alpha \nabla \ell(\theta)$
Bayesian Learning Rule(s) update the distribution over the parameters θ .



$Y \in T_e G$ is the direction of fastest ascent of $\mathcal{E}(q_g)$ w.r.t. the Fisher metric.

Navigation icons: back, forward, search, etc.

Solved issues

Issue 1 \mathcal{Q} is required to be an exponential family.

Solution Can choose q_0 freely and push it around with a group.

Issue 2 The updates could overshoot and leave the manifold.

Solution Closure of the group under operation keeps updates on \mathcal{Q} .

Issue 3 The gradient $\nabla_{\lambda} \mathcal{E}(q_{\lambda})$ can only be computed in special cases.

Solution Group action is the correct generality for reparametrization

$$\frac{d}{dg} \mathbb{E}_{q_g}[\ell] = \int_{\Theta} q_0(\theta) (\nabla_{\theta} \ell(g \cdot \theta))^{\top} \frac{dg \cdot \theta}{dg} d\theta$$

Also called *pathwise gradient estimators*²

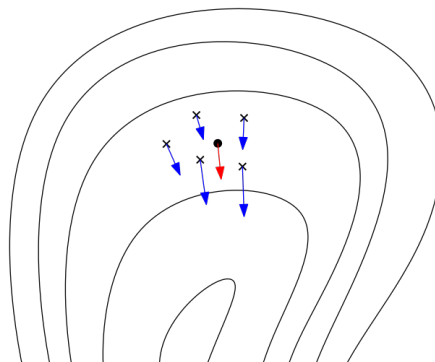
Bonus 1 The Fisher metric is invariant under translations by G .

Bonus 2 The tangent directions Y at each step lie in the same vector space $T_e G$, so they can be accumulated from previous steps.

²Mohamed et. al. *Monte carlo Gradient Estimation in Machine Learning* JMLR 2020

Specific Update Formulas: The Additive Group

$$g \in \mathbb{R}^P \text{ additive} \quad \Rightarrow \quad g \leftarrow g - \alpha \mathbb{E}_{q_g} [\nabla_{\theta} \ell]$$



Instead of going in the direction of the derivative **at** g , the direction is chosen by **consensus** with at points sampled from q_g .

Multiplicative and Affine Update Formulas

$$g \in \mathbb{R}_{>0} \text{ multiplicative} \implies$$

$$g \longleftarrow g \exp \left(-\alpha (\mathbb{E}_{q_g} [\theta \partial_{\theta} \ell] - \tau) \right)$$

 $(A, b) \in \text{Aff}(\mathbb{R})$ affine group \implies

$$b \longleftarrow b + \frac{c_X}{c_y} A \frac{\exp(-\alpha U) - 1}{U} V$$

$$A \longleftarrow A \exp(-\alpha U)$$

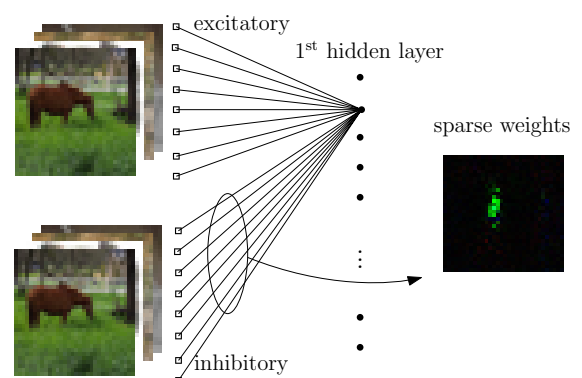
where $U = \mathbb{E}_{q_\theta}[(\theta - b)\partial_\theta \ell] - \tau$

$$V = A\mathbb{E}_{q_q}[\partial_\theta \ell]$$

Filters of the multiplicative group

Label nodes in a neural network “excitatory” or “inhibitory” like biology.

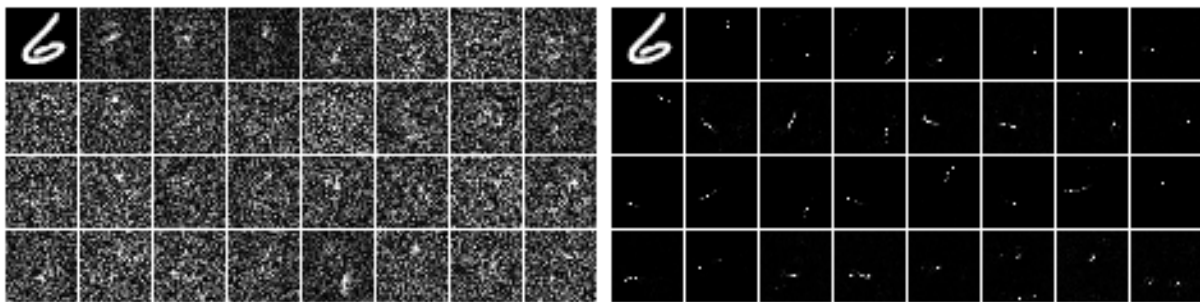
Magnitudes of the weights (in $\mathbb{R}_{>0}$) are the parameters (signs are fixed).



Multiplicative vs Additive filters

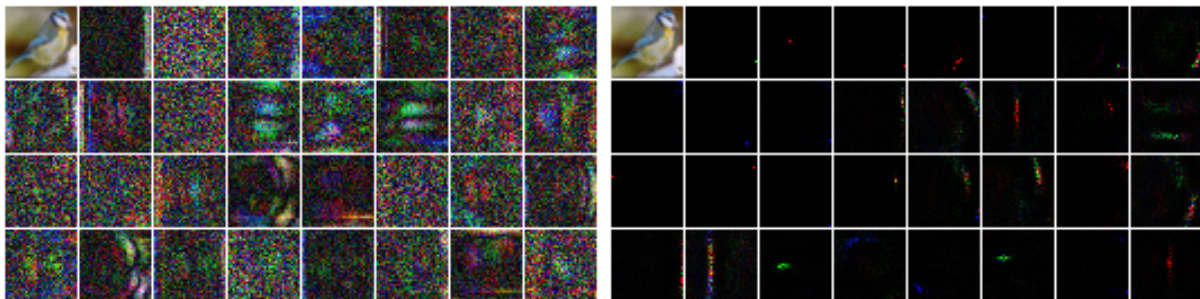
Model & Dataset	Method	Accuracy \uparrow (higher is better)	NLL \downarrow (lower is better)	ECE \downarrow (lower is better)
MNIST MLP	add.	98.38 ± 0.02	0.083 ± 0.001	0.012 ± 0.000
	mult.	98.59 ± 0.02	0.058 ± 0.001	0.006 ± 0.000
CIFAR-10 MLP	add.	58.85 ± 0.08	1.236 ± 0.002	0.085 ± 0.001
	mult.	59.19 ± 0.07	1.160 ± 0.001	0.026 ± 0.001

Additive rule is similar to SGD with momentum, multiplicative is different. They both learn.



Navigation icons: back, forward, search, etc.

The additive vs multiplicative filters for RGB images



Navigation icons: back, forward, search, etc.

Overview and Future work

- The Bayesian Learning Problem is a promising and fertile setting.
- The BLR of Khan & Rue both specialized to many existing algorithms and gave rise to many new and successful algorithms.
- Yet there are some issues with BLR such as: closure in the statistical manifold under updates, and calculation of derivatives w.r.t. distribution parameters.
- The group framework solves the closure problem by design and is able to very generally employ the **reparametrization trick**.
- Each new group would deserve an empirical study to investigate their learning behaviours (like multiplicative vs. additive)
- There may be implementation problems with arbitrary Lie groups, e.g. the exponential map may not always be feasible to compute, so approximations may be necessary.

Kiral, Möllenhoff, Khan

Lie Group Bayesian Learning Rule

Dec 2023

20 / 25

Teşekkürler

Kiral, Möllenhoff, Khan

Lie Group Bayesian Learning Rule

Dec 2023

21 / 25

Stiefel Manifold Update

Assume parameters are given as a matrix and want to preserve orthogonality of columns.

$$\Theta = \text{St}(n, m) = \{\theta \in \text{Mat}(n, m) : \theta^\top \theta = I_{m \times m}\}$$

The group $S = \text{SO}(n)$ preserves this manifold. And given a loss function $\ell : \Theta \rightarrow \mathbb{R}_{\geq 0}$

$$Y \in \mathfrak{so}(n) \text{ the update direction} \quad Y = \text{Skew} Y_0 = \frac{Y_0 - Y_0^\top}{2} \quad Y_0 = \mathbb{E}_{q_\Lambda}[\nabla \ell \theta^\top]$$

Here the distributions are parametrized by $\Lambda \in \text{Mat}(n, m)$

$$q_\Lambda(\theta) \propto e^{-\text{Tr}(\Lambda^\top \theta)}$$

and the update is given by

$$\Lambda \leftarrow e^{-\alpha Y} \Lambda \quad (\text{actually an efficient variation is used})$$

Koichi Tojo, Taro Yoshino's: "Harmonic Exponential Families".

G a Lie group $H \leq G$. Let ν be a relatively invariant measure on G $\pi : G \rightarrow \text{GL}(V)$ a representation of G . Let α be a 1-cocycle of π such that $\alpha|_H \equiv 0$. So $\alpha : G \rightarrow V$ satisfies

$$\alpha(gh) = \pi(g)\alpha(h) + \alpha(g) = \alpha(g). \quad \text{So } \alpha : \underbrace{G/H}_{:=\Theta} \rightarrow V$$

Let ν be a relatively invariant measure on Θ , meaning $\nu(gE) = \chi(g)\nu(E)$ for some homomorphism χ . Let $\lambda \in V^\vee$ s.t. $A(\lambda) = \log \int_\Theta e^{-\langle \lambda, \alpha(\theta) \rangle} d\nu(\theta) < \infty$. For such λ

$$q_\lambda(\theta) d\nu(\theta) := e^{-\langle \lambda, \alpha(\theta) \rangle - A(\lambda)} d\nu(\theta)$$

forms an exponential family satisfying

$$\frac{1}{\chi(g)} q_\lambda(g^{-1}\theta) := q_{\pi^\vee(g)\lambda}(\theta) \quad \text{where } \langle \pi^\vee(g)\lambda, v \rangle = \langle \lambda, \pi(g)v \rangle.$$

Why call it Bayesian?

Let $\ell(\theta) = \sum_{i=1}^N \ell_i(\theta) + R(\theta)$. Observe new data $(\mathbf{x}_{\text{new}}, y_{\text{new}})$ with loss contribution ℓ_{new} .

How to update p_τ ? Take $\tau = 1$

Bayes' rule is about conditional probabilities, and updating priors:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Interpret $e^{-\ell_i(\theta)}$ as the likelihood of **observing label y_i** given the model parameter θ and \mathbf{x}_i .

Interpret $\pi(\theta) \propto e^{-R(\theta)}$ as the prior on the parameters.

After one round of learning the posterior $p \propto e^{-\sum_i \ell_i} \pi$ is our prior belief about θ distribution. According to Bayes rule updated belief should be *after a new data* point.

$$p_{\text{updated}}(\theta) \propto e^{-\ell_{\text{new}}(\theta)} p(\theta).$$

This is also the optimizer if we had initially considered the loss function $\ell_{\text{updated}} = \ell + \ell_{\text{new}}$.

Exponential Families

Let $T : \Theta \rightarrow V$, called the sufficient statistic. Call

$$\Omega = \left\{ \lambda \in V^\vee : Z(\lambda) := \int_{\Theta} e^{-\langle \lambda, T(\theta) \rangle} d\nu(\theta) < \infty \right\}.$$

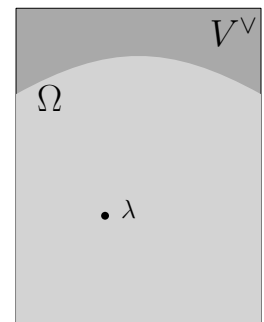
Then $q_\lambda(\theta) = \frac{1}{Z(\lambda)} e^{-\langle \lambda, T(\theta) \rangle}$ form an exponential family of distributions.

$$-\frac{\partial \ln Z}{\partial \lambda_i} = \int_{\Theta} T_i(\theta) \frac{1}{Z(\lambda)} e^{-\langle \lambda, T(\theta) \rangle} d\nu(\theta) = \mathbb{E}_{q_\lambda d\nu}[T_i] =: \mu_i$$

$$\frac{\partial^2 \ln Z}{\partial \lambda_i \partial \lambda_j} = \int_{\Theta} (T_i(\theta) - \mu_i)(T_j(\theta) - \mu_j) q_\lambda(\theta) d\nu(\theta)$$

$$= \mathbb{E}_{q_\lambda} \left[\left(\frac{\partial}{\partial \lambda_i} \log q_\lambda \right) \left(\frac{\partial}{\partial \lambda_j} \log q_\lambda \right) \right] =: F_{i,j}(\lambda) \quad \text{Fisher Matrix}$$

Example: If $T(\theta) = \begin{bmatrix} \theta \\ \theta^2 \end{bmatrix}$ then we get 1-D Gaussians $q_\lambda(\theta) \propto e^{-\lambda_1 \theta - \lambda_2 \theta^2}$ for $\lambda_2 > 0$.



The generalized Pythagorean theorem on the compactifications of certain dually flat spaces via toric geometry

Hajime Fujita (Japan Women's University)

The *dually flat space* is a fundamental object in information geometry, which appear as a geometric structure of a family of probability measures with “nice properties” such as exponential family. The dually flat space is a data consisting of a Riemannian manifold (X, g) and two flat connections (∇, ∇^*) which are dual to each other with respect to the metric g . The Riemannian metric g of the dually flat space is a Hesse metric, and hence, it admits a (local) potential function φ for a ∇ -affine coordinate and a dual potential ψ for a ∇^* -affine coordinate. The potential functions determine the *Bregman divergence* $D(\cdot|\cdot) : X \times X \rightarrow \mathbb{R}$. One important feature of the Bregman divergence is the *generalized Pythagorean theorem*, which is a fundamental tool in statistical inference.

Toric geometry is an area which have a position in the intersection of algebraic geometry, complex geometry and symplectic geometry. We focus on its complex and symplectic aspects, i.e., the Kähler structure. Any toric manifold has a structure of a singular Lagrangian torus fibration by its *moment map* and it associates a convex polytope as its image which is of spacial class called *Delzant polytopes*. There is a dictionary between toric manifolds and Delzant polytopes, and several geometric quantities of given toric manifold can be converted into them on the Delzant polytope. The description of the Riemannian metric on the toric manifold as in [1] tells us that the interior of the corresponding Delzant polytope has a natural structure of a dually flat space. This observation can be understood as a variant of Dombrowski's construction ([2]) which asserts a correspondence between Kähler manifold and dually flat space, and it leads to the notion of the *torification* introduced in [4].

In this talk we explain an extension of the dually flat structure of the Delzant polytope to its boundary following [3]. In particular we give a continuous extension of the Bregman divergence to the boundary so that the generalized Pythagorean theorem holds including the boundary points. A typical example is the complex projective space $\mathbb{C}P^n$ and the probability simplex Δ^n as its Delzant polytope. In the language of information geometry this setting is the family of probability measures on the finite set $[n+1] = \{1, \dots, n+1\}$ together with the expectation parameter. Our result offers an framework to handle the generalized Pythagorean theorem including the zero probabilities.

References

- [1] M. Abreu, *Kähler geometry of toric manifolds in symplectic coordinates*, Symplectic and contact topology: interactions and perspectives (Toronto, ON/Montreal, QC, 2001), 2003, pp. 1-24.
- [2] P. Dombrowski, *On the geometry of the tangent bundle*, J. Reine Angew. Math. 210 (1962), 73-88.
- [3] H. Fujita, *The generalized Pythagorean theorem on the compactifications of certain dually flat spaces via toric geometry*, arXiv:2305.08422v3, to appear in Information Geometry (DOI : 10.1007/s41884-023-00123-y).
- [4] M. Molitor, *Kähler toric manifolds from dually flat spaces*. arXiv:2109.04839v1.

The generalized Pythagorean theorem on the compactifications of certain dually flat spaces via toric geometry

Hajime Fujita
(Japan Women's university)

based on : arXiv:2305.08422 [math.SG], Information Geometry (published online).

Self-introduction & personal background of today's talk

My research area : symplectic geometry

(Geometric Quantization of Hamiltonian torus manifold using Dirac operators,
metric geometry of toric symplectic manifold)

[Submitted on 4 Mar 2020]

Distance functions on convex bodies and symplectic toric manifolds

Hajime Fujita, Yu Kitabeppu, Ayato Mitsuishi

In this paper we discuss three distance functions on the set of convex bodies in symplectic toric geometry. By using these observations, we derive a distance.

Comments: 19 pages

Subjects: **Metric Geometry (math.MG)**

[Submitted on 7 Jan 2020 (v1), last revised 1 Mar 2021 (this version, v3)]

Deformation of Dirac operators along orbits and quantization of non-compact Hamiltonian torus manifolds

Hajime Fujita

We give a formulation of a deformation of Dirac operator along orbits of a group action on a possibly non-compact manifold to get an equivariant index and a K-homology cycle representing the index. We apply this framework to non-compact Hamiltonian torus manifolds to define geometric quantization from the view point of index theory. We give two applications. The first one is a proof of a $[Q,R]=0$ type theorem, which can be regarded as a proof of the Vergne conjecture for Abelian case. The other is a Danilov-type formula for toric case in the non-compact setting, which shows that this geometric quantization is independent of the choice of polarization. The proofs are based on the localization of index to lattice points.

Comments: 27pages. Due to referee's comments several expositions are rewritten, and typos are corrected. Especially descriptions for non-abelian case are withdrawn. References uploaded. To appear in Canadian Journal of Mathematics

Subjects: **Differential Geometry (math.DG)**; K-Theory and Homology (math.KT); Symplectic Geometry (math.SG)

I am an armature of statistical theories and machine learning.

Two years ago I started studying information geometry just from curiosity.

I learned a connection between information geometry and Kähler geometry.

Last year I met Tojo-san (he was a part-time lecturer at my university) and

I attended his lecture on information geometry. At that time we had a conversation, and Tojo-san presented me a paper :

Self-introduction & personal background of today's talk

[Submitted on 10 Sep 2021]

Kähler toric manifolds from dually flat spaces

Mathieu Molitor

We present a correspondence between real analytic Kähler toric manifolds and dually flat spaces, similar to Delzant correspondence in symplectic geometry. This correspondence gives rise to a lifting procedure: if $f : M \rightarrow M'$ is an affine isometric map between dually flat spaces and if N and N' are Kähler toric manifolds associated to M and M' , respectively, then there is an equivariant Kähler immersion $N \rightarrow N'$. For example, we show that the Veronese and Segre embeddings are lifts of inclusion maps between appropriate statistical manifolds. We also discuss applications to Quantum Mechanics.

Subjects: **Differential Geometry (math.DG)**; Information Theory (cs.IT); Mathematical Physics (math-ph); Symplectic Geometry (math.SG)

In this paper a correspondence between a specific dually flat space and a toric Kähler manifold is discussed.

Molitor treated the “open dense part” of toric manifold. I tried to capture the behavior on the “compactification of it”.

As a result I got an application for the divergence on the boundary, the extension of “the extended Pythagorean theorem”.

Unfortunetely I do not have any statistical applications so far.

Please let me inform and have a diusscusion if you have any idea or insight!!

Contents

§1. Dually flat space and its example

§2. Dually flat structure on a convex polytope

§3. Delzant polytope and toric symplectic manifold

§4. The “generalized” generalized Pythagorean theorem

§5. Further discussions

§1. Dually flat space and its examples

↑
geometric structure which appears
in statistical model including
exponential family.

§1. Dually flat space and its example

Def A dually flat space is a data
 (X, h, ∇, ∇^*) consisting of

- (X, h) : Riemannian manifold.
- ∇ and ∇^* are dual to each other w.r.t. h .

$$V_1(h(V_2, V_3)) = h(\nabla_{V_1} V_2, V_3) + h(V_2, \nabla_{V_1}^* V_3)$$

$$\forall V_1, V_2, V_3 \in P(TX)$$

- ∇ and ∇^* are flat connections of TX .

$$\uparrow \text{torsion} \equiv 0, \text{curvature} \equiv 0$$

(h, ∇) is called a dually flat structure on X .

§1. Dually flat space and its example

Prop. (X, h, ∇, ∇^*) : dually flat

$\Rightarrow \begin{cases} \exists x = (x_1, \dots, x_n) : \text{(local) coordinate of } X \\ \exists \varphi = \varphi(x) : C^\infty\text{-function on } X. \end{cases}$

s.t. $\begin{cases} (1) \ x \text{ is } \nabla\text{-affine} ; \nabla_{\partial x_i} \partial x_j = 0 \ (\forall i, j) \\ (2) \ h = \text{Hess}_x \varphi ; h(\partial x_i, \partial x_j) = \frac{\partial^2 \varphi}{\partial x_i \partial x_j} \end{cases}$

Def. φ is called a potential of (h, ∇) .

§1. Dually flat space and its example

Prop. (X, h, ∇, ∇^*) : dually flat

$x = (x_1, \dots, x_n), \varphi = \varphi(x)$ as in (1) & (2).

$\Rightarrow \begin{cases} \exists y = (y_1, \dots, y_n) : \text{(local) coordinate of } X \\ \exists \psi = \psi(y) : C^\infty\text{-function on } X. \end{cases}$

s.t. $\begin{cases} (1) \ y \text{ is } \nabla^*\text{-affine} ; \nabla_{\partial y_i}^* \partial y_j = 0 \ (\forall i, j) \\ (2) \ y = \text{grad}_x \varphi \\ (3) \ \psi \text{ is the Legendre dual of } \varphi ; \\ \varphi(x) + \psi(y) - x \cdot y = 0 \end{cases}$

Def. y is called the dual affine coordinate.
 ψ is called the dual potential of φ .

§1. Dually flat space and its example

Def. For a dually flat space (X, h, ∇, ∇^*)

the two-variable function

$$D(\cdot, \cdot) : X \times X \longrightarrow \mathbb{R}$$

$$(\xi, \xi') \longmapsto D(\xi, \xi')$$

$$:= \varphi(x(\xi)) + \psi(y(\xi')) - \chi(\xi) \cdot y(\xi')$$

is called the (Bregman) divergence.

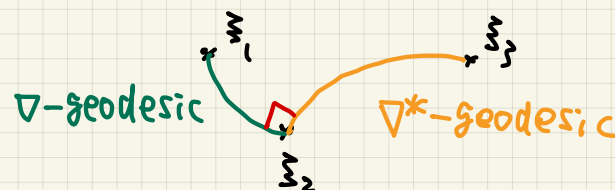
The divergence can be thought as a variant of the "square of the distance function" on (X, h) .

§1. Dually flat space and its example

Prop. $D(\xi, \xi') \geq 0$ for $\forall (\xi, \xi') \in X \times X$.
and $D(\xi, \xi') = 0 \iff \xi = \xi'$.

Thm. (The generalized Pythagorean theorem)

For $\xi_1, \xi_2, \xi_3 \in X$, if:



then $D(\xi_1, \xi_3) = D(\xi_1, \xi_2) + D(\xi_2, \xi_3)$.

§1. Dually flat space and its example

The following gives the example of our dually flat space.

Prop U : open subset in an affine space
 with the trivial flat connection ∇_0
 $\varphi: U \rightarrow \mathbb{R}$: smooth convex function
 $\text{Hess}\varphi = \nabla_0(d\varphi) > 0$
 $\Rightarrow (\text{Hess}\varphi, \nabla_0)$ is a dually flat str. on U
 with dual affine coord. $y = \text{grad}\varphi$.

Ex $U \subset \mathbb{R}^n$ with std. coordinate (x_1, \dots, x_n)
 $\Rightarrow \nabla_0 = \sum_i \frac{\partial}{\partial x_i} dx_i$, $\text{Hess}\varphi = \left(\frac{\partial^2 \varphi}{\partial x_i \partial x_j} \right)_{i,j}$

§1. Dually flat space and its example

Summary

Any dually flat space (X, h, ∇)
 \parallel
 $\text{Hess}\varphi$
 determines the divergence

$$\underline{D: X \times X \rightarrow \mathbb{R}.}$$

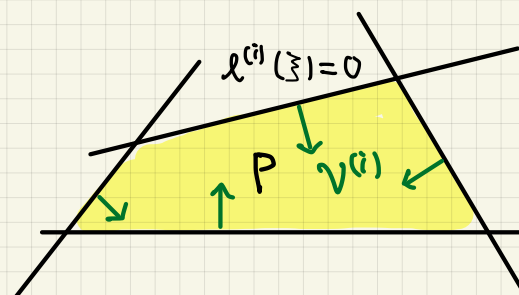
which satisfies the generalized Pythagorean thm.

§2. Dually flat structure on a convex polytope

§2. Dually flat structure on a convex polytope

P : convex polytope in \mathbb{R}^n defined by
 N - inequalities;

$$\ell^{(i)}(\xi) = \underbrace{v^{(i)}}_{\substack{\text{inward} \\ \text{normal vector}}} \cdot \xi + \overbrace{\lambda^{(i)}}^{\in \mathbb{R}} \geq 0 \quad (i=1, 2, \dots, N) \\ (\xi \in \mathbb{R}^n)$$



§2. Dually flat structure on a convex polytope

Def. Define $\varphi_P: P \rightarrow \mathbb{R}$ by

$$\varphi_P(\xi) := \sum_{i=1}^N \ell^{(i)}(\xi) \log \ell^{(i)}(\xi)$$

called the Guillemin potential of P

Prop φ_P is continuous on P and smooth convex on $\overset{\circ}{P} = \text{int}(P)$.

dually flat sp.
($\overset{\circ}{P}$, Hess φ_P)

$\ell^{(i)}(\xi) > 0$. ($\forall i$)

§2. Dually flat structure on a convex polytope

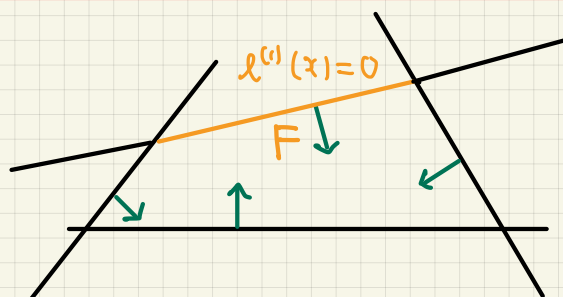
Take a face F of P determined by

$$\ell^{(1)}(\xi) = \dots = \ell^{(n-k)}(\xi) = 0 \quad (k = \dim F).$$

$$\ell^{(n-k+1)}(\xi), \dots, \ell^{(N)}(\xi) \geq 0$$

Define $\varphi_F: F \rightarrow \mathbb{R}$ by

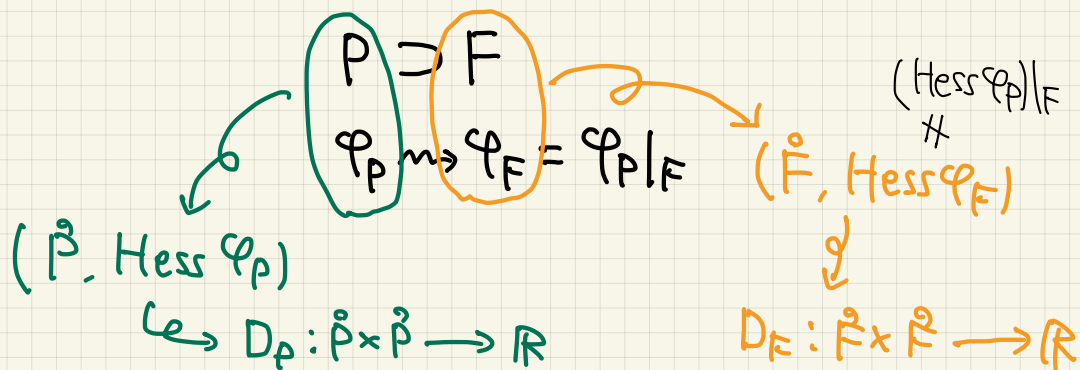
$$\varphi_F(\eta) = \sum_{i=n-k+1}^N \ell^{(i)}(\eta) \log \ell^{(i)}(\eta)$$



§2. Dually flat structure on a convex polytope

Prop $\varphi_F : F \rightarrow \mathbb{R}$ is continuous on F
and smooth convex on $\overset{\circ}{F} = \text{int}(F)$.

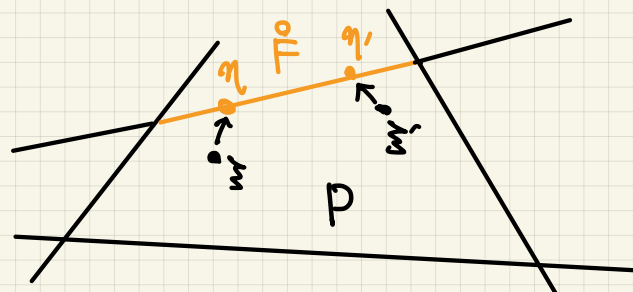
\hookrightarrow dually flat sp. $(\overset{\circ}{F}, \text{Hess}\varphi_F)$



§2. Dually flat structure on a convex polytope

Prop We have the following "continuity".

$$\lim_{\xi \rightarrow \eta} \lim_{\xi' \rightarrow \eta'} D_P(\xi, \xi') = D_F(\eta, \eta')$$



Rem. $\lim_{\xi \rightarrow \eta} \lim_{\xi' \rightarrow \eta'}$ is not valid.

§2. Dually flat structure on a convex polytope

Example $P: \xi_1 \geq 0, \xi_2 \geq 0, 1 - (\xi_1 + \xi_2) \geq 0$

Guillemin potential

$$\bullet \varphi_P(\xi) = \xi_1 \log \xi_1 + \xi_2 \log \xi_2 + (1 - \xi_1 - \xi_2) \log (1 - \xi_1 - \xi_2)$$

metric

$$\bullet \text{Hess } \varphi_P = \begin{pmatrix} \frac{1}{\xi_1} + \frac{1}{1 - \xi_1 - \xi_2} & \frac{1}{1 - \xi_1 - \xi_2} \\ \frac{1}{1 - \xi_1 - \xi_2} & \frac{1}{\xi_2} + \frac{1}{1 - \xi_1 - \xi_2} \end{pmatrix}$$

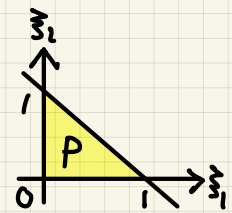
dual coord

$$\bullet y(\xi) = \text{grad } \varphi_P = \left(\log \frac{\xi_1}{1 - \xi_1 - \xi_2}, \log \frac{\xi_2}{1 - \xi_1 - \xi_2} \right) = (y_1, y_2)$$

divergence

$$\bullet D_P(\xi, \xi') = \xi_1 \log \frac{\xi_1}{\xi'_1} + \xi_2 \log \frac{\xi_2}{\xi'_2} + (1 - \xi_1 - \xi_2) \log \frac{1 - \xi_1 - \xi_2}{1 - \xi'_1 - \xi'_2}$$

$\xi = (\xi_1, \xi_2), \xi' = (\xi'_1, \xi'_2) \in P$



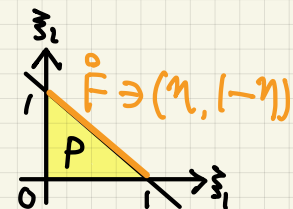
§2. Dually flat structure on a convex polytope

Example Take a face F as $\xi_1 + \xi_2 = 0$.
(continued)

$$\bullet \varphi_F(\eta) = \eta \log \eta + (1 - \eta) \log (1 - \eta)$$

$$\bullet \text{Hess } \varphi_F = \frac{1}{\eta} + \frac{1}{1 - \eta}$$

$$\bullet D_F(\eta, \eta') = \eta \log \frac{\eta}{\eta'} + (1 - \eta) \log \frac{1 - \eta}{1 - \eta'}$$



You can see:

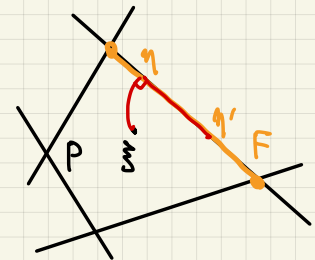
$$\lim_{\substack{\xi'_1 \rightarrow \eta \\ \xi'_2 \rightarrow 1 - \eta}} \lim_{\substack{\xi_1 \rightarrow \eta \\ \xi_2 \rightarrow 1 - \eta}} D_P(\xi, \xi') = D_F(\eta, \eta') \quad \text{for } \eta, \eta' > 0.$$

Question

Can we generalize
the generalized Pythagorean thm.
 among interior and boundary points?

My answer

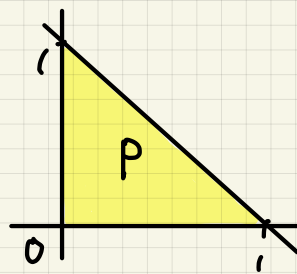
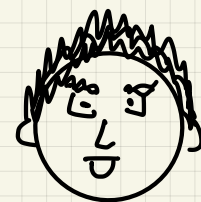
I did it when P : Delzant
 by using "toric geometry".



By the way :

Statistician

This P is the
 set of all probabilities
 on the three points
 $\{3\} = \{1, 2, 3\}$.
 D_P is the KL-divergence.

Symplectic geometer

This P is a
Delzant polytope
 = moment map image
 of $\mathbb{C}P^2$.

§3. Delzant polytope and toric Kähler manifold

§3. Delzant polytope and toric Kahler manifold

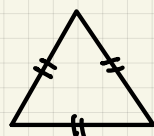
Def. $P \subset \mathbb{R}^n$: convex polytope is Delzant if:

- P is simple, i.e., each vertex has exactly n -edges.
- P is rational, i.e., each edge has directional vectors $\in \mathbb{Z}^n$.
- P is smooth, i.e., directional vectors at each vertex form a \mathbb{Z} -basis of \mathbb{Z}^n .

non-simple



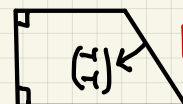
non-Delzant



Delzant



Delzant



§3. Delzant polytope and toric Kähler manifold

FACT (Delzant construction)

P : Delzant polytope in \mathbb{R}^n
 $\leadsto \exists M_P$: toric Kähler mfd.
 $\stackrel{M}{\cong} \dim_{\mathbb{R}} M_P = 2n$

What is a toric Kähler mfd?



Symplectic mfd with

- "nice" {
- torus action
 - Riemannian metric
 - complex structure
 - map $\mu: M_P \rightarrow P$

EXCUSE 😞

I do not give the precise definition of toric symplectic/Kähler manifold because of the constraint of time.

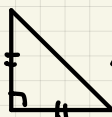
§3. Delzant polytope and toric Kähler manifold

FACT (Delzant construction)

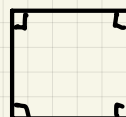
P : Delzant polytope
 $\leadsto \exists M_P$: toric Kähler mfd.
 $\stackrel{M}{\cong} \dim_{\mathbb{R}} M_P = 2n$

In fact, the converse direction holds. ($\mu(M) = P$)

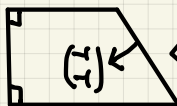
Ex.



$$\iff M_P = \mathbb{C}P^2$$



$$\iff M_P = \mathbb{C}P^1 \times \mathbb{C}P^1$$

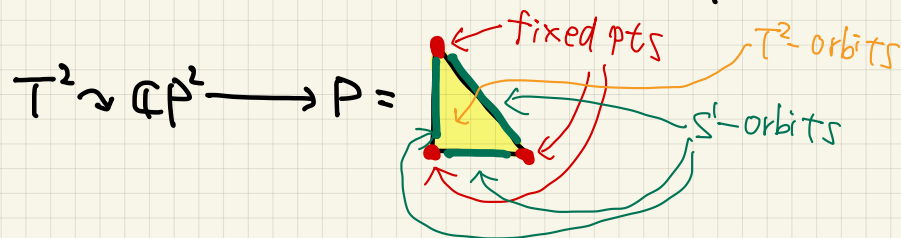


$$\iff M_P = \text{Hirzebruch surface}$$

§3. Delzant polytope and toric Kahler manifold

Fundamental relation between $M = M_P$ & P .

- (1) The map $\mu: M \rightarrow P$ is called the moment map which is in fact the quotient map.



- (2) $\mu: M \rightarrow P$ is a singular Lagrangian torus fibration.

- (3) \hat{P} parametrizes the principal orbits.

§3. Delzant polytope and toric Kahler manifold

Ex

$$\begin{array}{c}
 (U(1))^2 = T^2 \simeq M = \mathbb{CP}^2 = (\mathbb{C}^3 \setminus 0) / \mathbb{C}^\times \xrightarrow{\mu} P \subset \mathbb{R}^2 \\
 \downarrow (t_1, t_2) \quad \downarrow [z_1, z_2, z_3] \quad \downarrow \frac{(|z_1|^2, |z_2|^2)}{|z_1|^2 + |z_2|^2 + |z_3|^2} \\
 \downarrow [t_1 z_1, t_2 z_2, z_3] \\
 \end{array}$$

$\begin{array}{c} \uparrow \xi_1 \\ \uparrow \xi_2 \end{array}$

(1)

$F: \xi_1 + \xi_2 = 1$

$\leadsto M_F = \mu^{-1}(F)$

$= \{[z_1, z_2, z_3] \mid z_3 = 0\}$

$\cong \mathbb{CP}^1$ (4)

$F \leftarrow P_F = [0, 1] \subset \mathbb{R}$

$(\eta, 1-\eta) \leftarrow \eta$

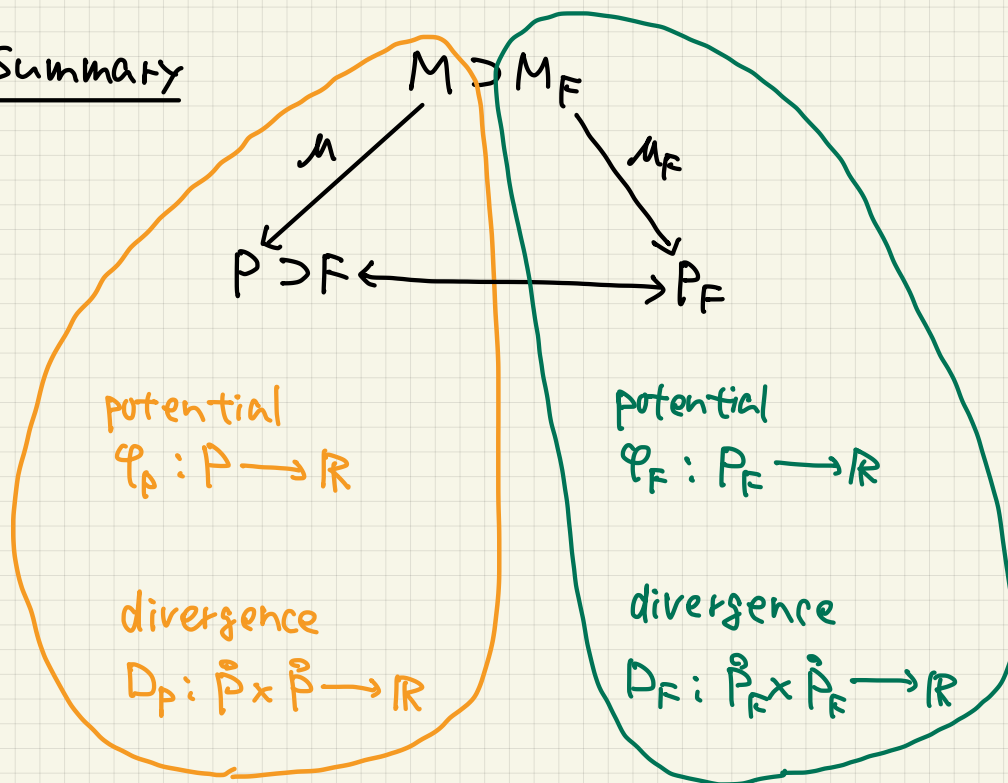
(6)

$M \xleftarrow{\mu} P$

$\downarrow \quad \downarrow$

$[\sqrt{\xi_1}, \sqrt{\xi_2}, \sqrt{1-\xi_1-\xi_2}] \leftarrow (\xi_1, \xi_2)$

§3. Delzant polytope and toric Kahler manifold

Summary

§4. The “generalized” generalized Pythagorean theorem

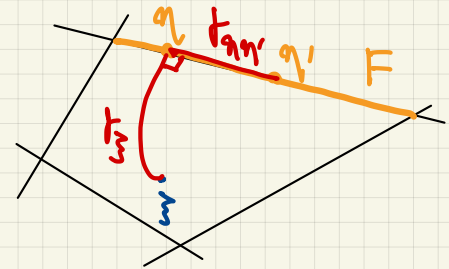
§4. The "generalized" generalized Pythagorean theorem

Thm. P : Delzant polytope. $\Leftrightarrow M$: toric Kähler

F : face. Take $\xi \in \overset{\circ}{P}$, $\eta, \eta' \in \overset{\circ}{F}$.

Let Γ_ξ and $\Gamma_{\eta\eta'}$ be ∇^* , ∇ -geodesics s.t.

- $\lim_{t \rightarrow \infty} S(\Gamma_\xi(t)) = S(\eta) \in M$
- $\lim_{t \rightarrow \infty} \frac{d}{dt} S(\Gamma_\xi(t)) \perp \frac{d}{dt} S(\Gamma_{\eta\eta'}(t))$
at $S(\eta) \in M_F \subset M$

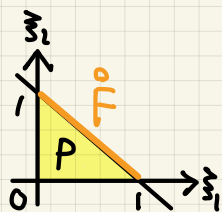


Then: $D'_F(\eta', \xi) = D_F(\eta', \eta) + D'_F(\eta, \xi)$
 $\left(\lim_{\xi \rightarrow \eta'} D(\xi, \xi) \right)$

§4. The "generalized" generalized Pythagorean theorem

\overline{F}_x

$\xi = (a, b) \in \overset{\circ}{P}$



∇^* -geodesic Γ_ξ

\parallel
affine map w.r.t. g

$$y = \text{grad } \varphi = \left(\log \frac{\xi_1}{1 - \xi_1 - \xi_2}, \log \frac{\xi_2}{1 - \xi_1 - \xi_2} \right)$$

$$x = x(y) = \left(\frac{e^{y_1}}{1 + e^{y_1} + e^{y_2}}, \frac{e^{y_2}}{1 + e^{y_1} + e^{y_2}} \right)$$

$$t \mapsto \left(\log \frac{a}{1-a-b} + t v_1, \log \frac{b}{1-a-b} + t v_2 \right) =: \overline{F}_\xi(t)$$

$$\Gamma_\xi(t) = x(\overline{F}_\xi(t)) \quad (\leftarrow \nabla^*\text{-geod. in } x\text{-coord.})$$

$$= \frac{1}{1-a-b + a e^{t v_1} + b e^{t v_2}} (e^{t v_1}, e^{t v_2})$$

$$\xrightarrow{t \rightarrow \infty} \left(\frac{a}{a+b}, \frac{b}{a+b} \right) \in \overset{\circ}{F}$$

if $v_1 = v_2 > 0$

§4. The "generalized" generalized Pythagorean theorem

Ex
(continued)

$$S: P \longrightarrow \mathbb{CP}^2$$

$$(\xi_1, \xi_2) \longmapsto [\sqrt{\xi_1}, \sqrt{\xi_2}, \sqrt{1-\xi_1-\xi_2}]$$

$$S(t_j(t)) = \left[\sqrt{\frac{a e^{t v_1}}{A}}, \sqrt{\frac{b e^{t v_1}}{A}}, \sqrt{\frac{1-a-b}{A}} \right]$$

$$= [\sqrt{a}, \sqrt{b}, e^{-\frac{t}{2} v_1} \sqrt{1-a-b}] \in \mathbb{CP}^2.$$

perpendicular to $\bigcup M_F \cong \mathbb{CP}^1$

$\xi = (a, b) \in \beta,$
 $\eta = \left(\frac{a}{a+b}, \frac{b}{a+b} \right) \in F,$
 $\eta' = (\eta', 1-\eta') \in F.$

satisfy the assumptions in the theorem

§4. The "generalized" generalized Pythagorean theorem

Ex
(continued) By the direct computations we have:

$$\left(\begin{array}{l} \cdot D'_F(\eta', \xi) = \dots = \eta' \log \frac{\eta'}{a} + (1-\eta') \log \frac{1-\eta'}{b} \\ \cdot D_F(\eta', \eta) = \dots = \eta' \log \frac{\eta'}{a} + (1-\eta') \log \frac{1-\eta'}{b} - \log(a+b) \\ \cdot D'_F(\eta, \xi) = \dots = \log(a+b) \end{array} \right.$$

$$\underline{D'_F(\eta', \xi) = D_F(\eta', \eta) + D'_F(\eta, \xi)} \quad \text{holds !!}$$

§5. Further discussions

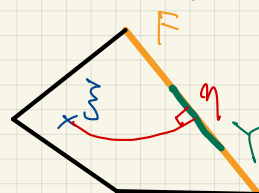
§5. Further discussions

Toward the "Projection thm".

For $\xi \in \mathring{P}$ and $Y \subset F$: submfd.

$\eta \in \text{Crit}(f_\xi) \stackrel{??}{\iff}$ orthogonality of the
 ∇^* -geod. and Y

$$\left(\begin{array}{ccc} f_\xi : Y & \longrightarrow & \mathbb{R} \\ \downarrow \eta & & \downarrow \\ \eta & \longmapsto & D'_F(\eta, \xi) \end{array} \right)$$

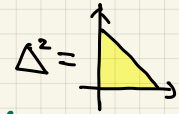


Rem \mathbb{CP}^2 (\mathbb{CP}^n) case is confirmed.

§5. Further discussions

Applications to statistical inference?

\mathbb{CP}^n Some P has an interpretation
as probability densities on a finite set.



$$\Delta^n = \{ \xi = (\xi_1, \dots, \xi_n) \in \mathbb{R}^n \mid \forall \xi_i \geq 0, \sum_{i=1}^n \xi_i \leq 1 \}$$

\rightarrow prob. densities on $[n+1] = \{1, \dots, n+1\}$
with the expectation parameter ξ .

- $D_{\Delta^n} = D_{KL} \leftarrow$ The Pythagorean thm including 0-probabilities.
- $\Delta^{n_1} \times \Delta^{n_2} \leftrightarrow$ a subfamily on $[n_1 n_2]$ (mixture)

§5. Further discussions

Relation with Nakajima-Ohmoto's theory?**The dually flat structure for singular models**

Naomichi Nakajima, Toru Ohmoto

The dually flat structure introduced by Amari-Nagaoka is highlighted in information geometry and related fields. In practical applications, however, the underlying pseudo-Riemannian metric may often be degenerate, and such an excellent geometric structure is rarely defined on the entire space. To fix this trouble, in the present paper, we propose a novel generalization of the dually flat structure for a certain class of singular models from the viewpoint of Lagrange and Legendre singularity theory - we introduce a quasi-Hessian manifold endowed with a possibly degenerate metric and a particular symmetric cubic tensor, which exceeds the concept of statistical manifolds and is adapted to the theory of (weak) contrast functions. In particular, we establish Amari-Nagaoka's extended Pythagorean theorem and projection theorem in this general setup, and consequently, most of applications of these theorems are suitably justified even for such singular cases. This work is motivated by various interests with different backgrounds from Frobenius structure in mathematical physics to Deep Learning in data science.

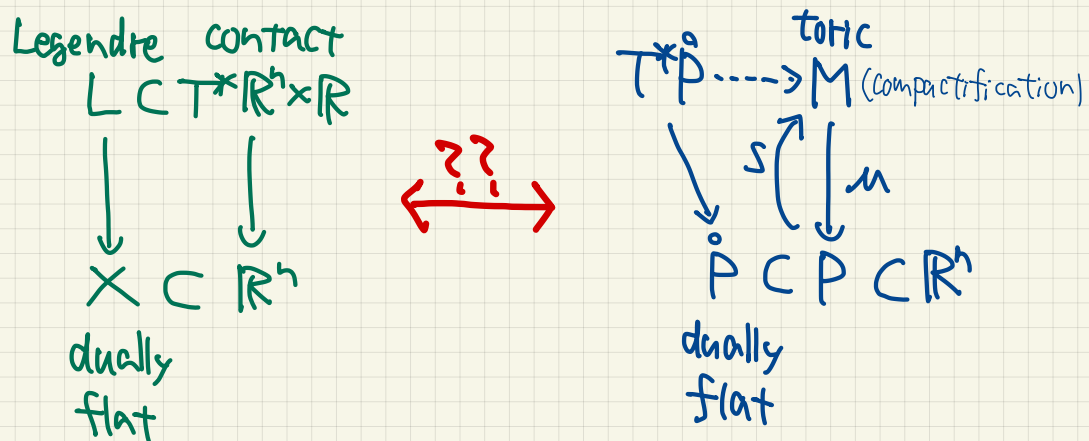
Comments: 29pages, 5figures
Subjects: **Differential Geometry (math.DG)**; Mathematical Physics (math-ph); Statistics Theory (math.ST)
MSC classes: 53B12 (Primary)

They gave a reformulation of
dually flat spaces so that they
can handle singular models.

§5. Further discussions

Relation with Nakajima-Ohmoto's theory?

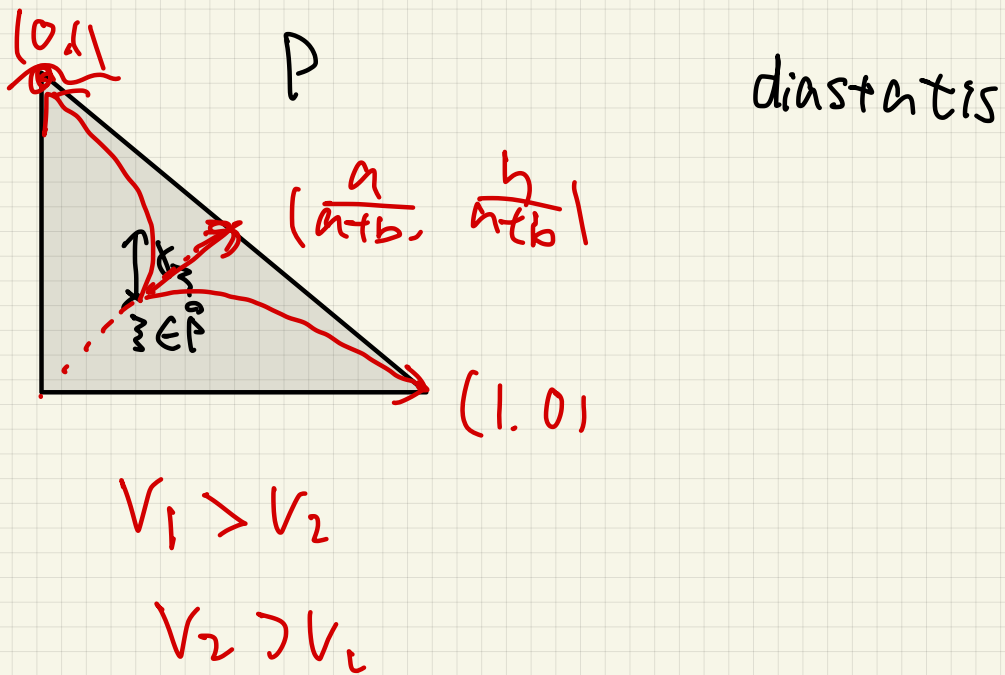
The basis of their theory are
contact manifold and Legendre submfd.



Thank you for
your attension!!

Summary

- Each convex polytope P has a natural dually flat str. by the Guillemin potential φ_P
- We can extend D_P to any face F .
- If P is Delzant, then we can generalize the Pythagorean thm. to F .



Doubly autoparallel structure and curvature integrals: An application to iteration complexity analysis of convex optimization

Atsumi Ohara

University of Fukui, Japan

ohara@fuee.u-fukui.ac.jp

Abstract

On a statistical manifold, we can define autoparallel submanifolds and path integrals of the second fundamental forms (curvature integrals) for its primal and dual affine connections, respectively. A submanifold is called doubly autoparallel if it is simultaneously autoparallel with respect to the both connections.

In this presentation we first discuss common properties of such submanifolds. In particular we next give an algebraic characterization of them in Jordan algebras and show their applications. Further, we exhibit that both curvature integrals induced from dually flat structure are interestingly related to an unexpected quantity, i.e., iteration-complexity of the interior-point algorithms for convex optimization defined on a submanifold that is *not* doubly autoparallel.

This is a joint work with Hideyuki Ishi and Takashi Tsuchiya.

References

- [1] A. Ohara, H. Ishi and T. Tsuchiya, Doubly autoparallel structure and curvature integrals: An application to iteration complexity for solving convex programs, *Information Geometry*, 2023. <https://doi.org/10.1007/s41884-023-00116-x> (Open Access)

Doubly autoparallel structure and curvature integrals

- An application to iteration-complexity analysis of convex optimization -

Atsumi Ohara University of Fukui

Joint work with
H. Ishi (Osaka Metropolitan Univ.),
T. Tsuchiya (GRIPS)

Statistical Theories and Machine Learning
Using Geometric Methods
December 14-15, 2023 @Osaka Metropolitan Univ.

1

Introduction

- Doubly autoparallel (DA) submanifold
 - natural notion in [information geometry](#) with symmetry
 - sometimes appears but much attention has not been paid
 - important applications in statistics and optimization
 - convex optimization on symmetric cones (e.g. SemiDef. Prog.)
 - MLE of structured covariance matrices
 - means on symmetric cones
 - and so on.

2

Introduction

- A goal of the presentation is to demonstrate
 - Characterization of DA submfd's in a symmetric cone
 - If the feasible region of SDP is DA
 - an explicit formula for the optimal solutions
 - If the feasible region of SDP is **not** DA
 - **curvature integral** evaluates an iteration-complexity to obtain the optimal solutions with IP algorithm

3

Outline

- Doubly autoparallel submanifolds
- Preliminaries: Dually flat structure on a symmetric cone
- Characterization of DA submfd's in a symmetric cone
 - Applications to
 - SemiDef. Prog. (SDP) - MLE for structured cov. mtx.
 - explicit formulas for an optimal solution for them
- **Non DA case**: curvature integrals
 - Iteration-complexity analysis for **conic linear program** (an extension of SDP) via an Interior-Point algorithm
 - Application to primal-dual path following methods
- Concluding remark

4

Information geometry [Amari & Nagaoka 00]

Def. Statistical manifold: (S, g, ∇, ∇^*)

$$Xg(Y, Z) = g(\nabla_X Y, Z) + g(Y, \nabla_X^* Z)$$

X, Y and Z : arbitrary vector fields on S

- ★ g : **Riemannian metric**
- ★ (∇, ∇^*) : torsion-free **affine connections**
 $R^\nabla = 0, R^{\nabla^*} = 0 \rightarrow$ **dually flat**
- ★ $\nabla^{(\alpha)} = \frac{1+\alpha}{2}\nabla + \frac{1-\alpha}{2}\nabla^* : \alpha$ -**connections**

5

Doubly autoparallel submanifolds

[Uohashi&O 04] [OIT23]

- Def. Let (S, g, ∇, ∇^*) be a statistical manifold and M be its submanifold. We call M a **doubly autoparallel** submanifold in S when the followings hold:

- $\forall X, Y \in \mathcal{X}(M), \nabla_X Y \in \mathcal{X}(M)$
i.e. $H_M(X, Y) = 0$
- $\forall X, Y \in \mathcal{X}(M), \nabla_X^* Y \in \mathcal{X}(M)$
i.e. $H_M^*(X, Y) = 0$

6

Important Properties

Proposition 1 The following statements are equivalent:

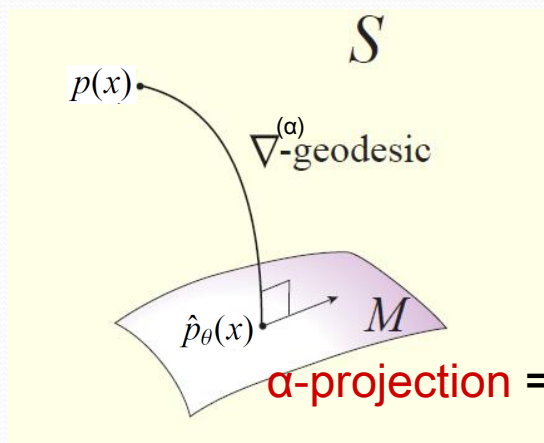
- 1) A submanifold M is doubly autoparallel (DA)
- 2) M is autoparallel w.r.t. the α -connections

$$\nabla^{(\alpha)} = \{(1 + \alpha)\nabla + (1 - \alpha)\nabla^*\}/2$$
 for **two different** α 's.
- 3) M is autoparallel w.r.t. **all** the α -connections.
- 4) **all** the α -geodesics connecting two points on M lay in M (if it is simply connected).
- 5) M is affinely constrained in both ∇ - and ∇^* -affine coordinates **if S is dually flat**.

7

Furthermore, for a parametric statistical model S

- If M is DA in S , then α -projections from p to M are unique for **all** α .



M is simultaneously an
exponential and mixture
family

$$\text{\textcolor{red}{\alpha-projection}} = \arg \min D^{(\alpha)}(p, M)$$

8

Related topics and applications

Symmetric cones

- MLE for structured covariance matrices is tractable (cast to convex program: **inversely linear structure**) [Anderson 70, Malley 94]
- **Explicitly solvable Semi-Definite Programs** [O 99]
- Structure of α -power means on symmetric cones [O 04]

9

Related topics and applications

Probability simplex

- Statistical models Markov-isomorphic to the probability simplex [Nagaoka 17]
- Characterization and classification of DA submfds in the probability simplex via Hadamard algebra [O&Ishi 18]
- Learning theory [Mutus&Ay 03]

Miscellaneous

- The self-similar (*Barenblatt–Pattle*) solution for the porous medium equation [O&Wada 10]

General statistical manifolds

- Purely geometric study [Satoh *et al.* 21]

10

Preliminaries

[Faraut&Korani 94]

- Euclidian space E with an inner product $(\cdot|\cdot)$
 - Symmetric cone Ω : open and convex in E
 - homogeneous
 $G(\Omega) = \{\tau \in GL(E) \mid \tau(\Omega) = \Omega\}$ acts transitively
 - self-dual w.r.t. an inner product of E
 $\Omega = \Omega^*, \quad \Omega^* = \{y \in E \mid (x|y) > 0, \forall x \in \overline{\Omega} \setminus \{0\}\}$
 - Euclidean Jordan algebra $(V, *)$
 - commutative
 - $x^2 * (x * y) = x * (x^2 * y)$, where $x^2 = x * x$
 - \exists associative inner-product $(x * y|z) = (y|x * z)$ on V
- Prop. $\Omega = \text{int}\{x^2 \mid x \in V\}$ is a symmetric cone in $V(=E)$.

11

- $L(x) : V \rightarrow V \quad L(x)y = x * y$
- $P(x, y) := L(x)L(y) + L(y)L(x) - L(x * y)$
- **Mutation:** $x \perp_a y := P(x, y)a$

isomorphic to $*$, the unit element: a^{-1}

Ex. the set of real symmetric pos. def. matrices $\text{PD}(n, \mathbf{R})$

$$V = \text{Sym}(n; \mathbf{R}), \quad X * Y = (XY + YX)/2$$

$$\tau_G(X) = GXG^T, \quad G \in GL(n, \mathbf{R})$$

$$(X|Y) = \text{tr}(XY), \quad \text{the unit: } I, \text{ the inverse: } X^{-1}$$

$$X \perp_A Y = (XAY + YAX)/2$$

12

Dually flat structure on Ω

- Logarithmic characteristic function on Ω

$$\psi(x) := \log \int_{\Omega^*} e^{-\langle s, x \rangle} ds,$$

- positive definite Hessian on Ω
- $x^{-1} = -\text{grad } \psi(x)$, $(\text{grad } f(x)|u) = D_u f(x)$

- a coordinate system (x^i) : $x = \sum_{i=1}^n x^i e_i$, $\{e_i\}_{i=1}^n$: a basis of E

- a dual coordinate system (s_i) :

$$x^{-1} = \sum_{i=1}^n s_i e^i, \quad \{e^i\}_{i=1}^n: \text{a basis of } E \text{ with } (e^i|e_j) = -\delta_j^i$$

13

- D : the canonical **flat** affine connection on E

- $\{x^1, \dots, x^n\}$: affine coordinate system, i.e., $D_{\frac{\partial}{\partial x^i}} \frac{\partial}{\partial x^j} = 0$

- g : Riemannian metric on Ω

$$g = Dd\psi = \sum_{i,j} \frac{\partial^2 \psi}{\partial x^i \partial x^j} dx^i dx^j.$$

- D' : the dual affine connection on Ω

$$Xg(Y, Z) = g(D_X Y, Z) + g(Y, D'_X Z)$$

$(g, D, D') : \text{dually flat structure on } \Omega$

14

Expression via Jordan algebra

Dually flat structure on Ω [Uohashi&O 04]

- Potential= log. char. func. $\psi(x) = -\log \det x$,
- Riemanian metric: $g_x \left(\frac{\partial}{\partial x^i}, \frac{\partial}{\partial x^j} \right) = (P(x)^{-1} e_i | e_j)$, $P(x) := P(x, x)$
- α -connections: $\left(\nabla_{\frac{\partial}{\partial x^i}}^{(\alpha)} \frac{\partial}{\partial x^j} \right)_x = (\alpha - 1)(e_i \perp_{x^{-1}} e_j)$

Ex. On $\text{PD}(n, \mathbf{R})$

- $\psi(x) = -\log \det X$, ($X = \sum_{i=1}^N x^i E_i$)
- $g_X \left(\frac{\partial}{\partial x^i}, \frac{\partial}{\partial x^j} \right) = \text{tr}(X^{-1} E_i X^{-1} E_j)$
- $\left(\nabla_{\frac{\partial}{\partial x^i}}^{(\alpha)} \frac{\partial}{\partial x^j} \right)_X = \frac{\alpha - 1}{2} (E_i X^{-1} E_j + E_j X^{-1} E_i)$

15

Characterization of DA submfds in Ω

Let W be a linear subspace in Jordan algebra $(V, *)$ and $p = q * q$ in a symmetric cone Ω .

Thm. [OIT23] The following 1)-3) are equivalent:

1) A submanifold $M = (W + p) \cap \Omega$ is DA, where

$$W + p = \{w + p \mid w \in W\}$$

2) For **all** x in M , $u \perp_{x^{-1}} v \in W$, ($u, v \in W$)

3) The subspace $P(q)^{-1}W$ is a Jordan subalgebra.

Rem. (a) 3) is able to be checked at the **single** point p

(b) $M = \{(W' + p^{-1}) \cap \Omega\}^{-1}$ with $W' = P(p)^{-1}W$

The proof is based on the property 5) in Prop. 1 (p.6)

16

(c) **Implication:** Classification of DA submflds in Ω reduces to that of Jordan subalgs of $(V, *)$.
(For $V = \text{Sym}(n, \mathbf{R})$ cf. [Jacobson 87], [Malley 87])

- Ex. - Jordan subalgebras in $\text{Sym}(n, \mathbf{R})$
 - 1) fixed eigen vectors, 2) doubly symmetric, etc.
- Two bases $\{E_i\}_{i=1}^m$ and $\{F^i\}_{i=1}^m$ of $\text{Sym}(n, \mathbf{R})$

$$\mathcal{M} = \{P \mid P = E_0 + \sum_{i=1}^m x^i E_i, \exists x = (x^i) \in \mathbf{R}^m\} \cap \text{PD}(n)$$

$$\mathcal{M} = \{P \mid P^{-1} = F^0 + \sum_{i=1}^m s_i F^i, \exists s = (s_i) \in \mathbf{R}^m\} \cap \text{PD}(n).$$

17

Application(1) Means on Positive Operators

[Kubo & Ando 80]

- Def. (Axioms of means)
 σ is a **mean** on self-conjugate positive operators

- i) $A \leq C, B \leq D \Rightarrow A\sigma B \leq C\sigma D$
- ii) $C(A\sigma B)C = (CAC)\sigma(CBC)$
- iii) $A_n \downarrow A, B_n \downarrow B \Rightarrow A_n\sigma B_n \downarrow A\sigma B$

where $A_n \downarrow A \stackrel{\text{def}}{\Leftrightarrow} (A_i \geq A_{i+1}, \forall i) \ \& \ (A_n \rightarrow A)$

- iv) $I\sigma I = I$

18

α -geodesics on $PD(n)$

- α -geodesic $P(s)$ boundary conds. : $P(0)=A, P(1)=B$

$$P^{(\alpha)}(s) = A^{1/2} \left\{ [(A^{-1/2} B A^{-1/2})^\alpha - I] s + I \right\}^{1/\alpha} A^{1/2}$$

$$\alpha = 1 \quad P(s) = A + s(B - A)$$

$$\alpha = 0 \quad \hat{P}(s) = A^{1/2} \exp(s \log A^{-1/2} B A^{-1/2}) A^{1/2}$$

$$\alpha = -1 \quad P^*(s) = \{A^{-1} + s(B^{-1} - A^{-1})\}^{-1}$$

$$AaB := P(1/2)$$

$$AgB := \hat{P}(1/2)$$

$$AhB := P^*(1/2)$$

$P^{(\alpha)}(1/2)$: a power mean

19

Means and α -geodesics on $PD(n)$ [O 04]

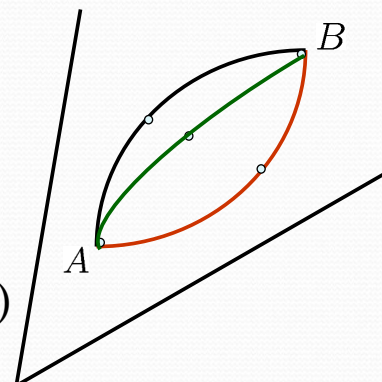
Thm. Points on α -geodesics for s in $[0,1]$ and α in $[-1,1]$ are 2-param. family of means, i.e.,

$$A\sigma_s^{(\alpha)} B = P^{(\alpha)}(s)$$

In particular, for fixed s in $[0,1]$

$$P^{(\alpha)}(s) > P^{(\beta)}(s),$$

$$1 \geq \alpha > \beta \geq -1 \quad \text{AGH ineq. (s=1/2)}$$



Cor. A and B are in a DA submanifold M

$$\Rightarrow A\sigma_s^{(\alpha)} B \in M, s \in [0,1], \alpha \in [-1,1]$$

20

App.(2) MLE for structured covariance matrices

- a sample covariance S in $\text{PD}(n, \mathbf{R})$
- a zero-mean Gaussian p.d.f. with a covariance mtx. Σ

$$p(x) = (2\pi)^{-n/2} (\det \Sigma)^{-1/2} \exp\left\{-\frac{1}{2}x^T \Sigma^{-1}x\right\}$$

- structured covariance mtx. (with linear constraints)

$$\Sigma \in \mathcal{M} = (E_0 + \mathcal{W}) \cap \text{PD}(n, \mathbf{R})$$

• Ex.

- Toeplitz matrices: $\{T = (t_{ij}) | t_{ij} = t_{ji} = y_{|i-j|}\}$
- zero-patterns: $\{\Sigma = (\sigma_{ij}) | \sigma_{ij} = \sigma_{ji} = 0, (i, j) \in \mathcal{E}\}$
- etc...

21

MLE for structured covariance matrices

- Negative logarithmic likelihood func (up to const.):

$$\ell(\Sigma) := -\log \det \Sigma^{-1} + \text{tr}(\Sigma^{-1}S) \rightarrow \min$$

- Rem Note that $-\log \det$ is a convex function.

- If \mathcal{M} is DA (inversely linear structure), then the **minimization** problem of $\ell(\Sigma)$ (**MLE**) s.t. $\Sigma \in \mathcal{M}$ is a **strictly convex program**.



Unique solution if it exists,
Numerically tractable (optimality eq. is linear)

22

App.(3) Convex program

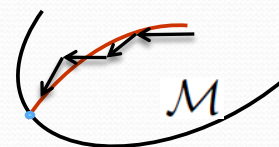
Affine-scaling method and IG

- General convex program: Convex set $\mathcal{M} \subset \mathbb{R}^n$ $c \in \mathbb{R}^n$

$$\text{minimize } c^T x, \quad \text{s.t. } x \in \overline{\mathcal{M}}$$
- Ψ : a **good** convex barrier function for \mathcal{M} ,
 1) $\Psi(x) \rightarrow +\infty$ ($x \rightarrow \text{bd}\mathcal{M}$), 2) h : p.d. Hessian, 3) **self-concordance**
- Gradient flow for Riemannian mfd (\mathcal{M}, h)

$$\dot{x} = \frac{dx}{dt} = -h(x)^{-1}c, \quad x(0) \in \mathcal{M}$$

 $x(t)$: **affine-scaling trajectory**
 (numerically traced)
- Legendre transform \Rightarrow **linearization**



$$\dot{s} = -c, \quad s_i = \frac{\partial \Psi}{\partial x^i}, \quad i = 1, \dots, n, \quad \hat{s} := -\lim_{t \rightarrow +\infty} ct + s(0)_{13}$$

- **➡** Opt. sol. is $\hat{x} = \text{grad}\Psi^*(\hat{s})$ (inverse Legendre trans.) ★
- **However**, we need to solve the nonlinear eq. $\hat{s} = \text{grad}\Psi(\hat{x})$ because getting **the explicit form of Ψ^*** from Ψ is difficult.

If we know the explicit form of $\Psi^* \Rightarrow$ a formula ★ for \hat{x}

Idea: 1) Ω : sym. cones $\Rightarrow \psi(x) = -\log \det x$, $\psi^*(s) = -\log \det s$,
 Legendre transform: $x \mapsto s = x^{-1}$

- 2) \mathcal{M} is realized as a DA submfd $\mathcal{M} = (a+W) \cap \Omega$ in Ω
 1) & 2) \Rightarrow linearized AS trajectory stays in \mathcal{M} and $\hat{x} = \hat{s}^{-1}$

- Typical Example: SemiDefinite Program (SDP)

$$\underset{P}{\text{minimize}} (C|P), \text{ s.t. } P = E_0 + \sum_{i=1}^m x^i E_i \in \overline{\mathcal{M}} = \overline{(E_0 + \mathcal{W}) \cap \text{PD}(n)}$$

- If \mathcal{M} is DA in $\text{PD}(n)$ and $P \in \mathcal{M}$

- 1. Set $F^0 = P^{-1}$, $F^i = -P^{-1} E_i F^{-1}$, then

$$\mathcal{M} = \{P \mid P^{-1} = F^0 + \sum_{i=1}^m s_i F^i, \exists s = (s_i) \in \mathbf{R}^m\} \cap \text{PD}(n)$$

- 2. Solve $\tilde{C} \in \text{span}\{F^i\}_{i=1}^m$ meeting

$$\forall P \in \mathcal{M}, \quad (C|P) = (\tilde{C}|P) + \text{const.}$$

25

- 3. Spectral decomposition

$$\tilde{C} = \begin{pmatrix} V_1 & V_2 \end{pmatrix} \begin{pmatrix} \Sigma_1 & O \\ O & O \end{pmatrix} \begin{pmatrix} V_1^T \\ V_2^T \end{pmatrix} = V_1 \Sigma_1 V_1^T$$

- 4. For $\forall P_0 \in \mathcal{M}$ with $S_0 = P_0^{-1}$, the opt. sol. is

$$\hat{P} = \lim_{t \rightarrow \infty} S(t)^{-1} = \lim_{t \rightarrow \infty} (-\tilde{C}t + S_0)^{-1} = P_0 - P_0 V_1 (V_1^T P_0 V_1)^{-1} V_1^T P_0$$

Rem. **Independent** of the objective function $(C|P)$ and an initial value P_0

26

Interior point method (IP) for Conic linear program

Notation :

- Vector space E of dimension n
- the dual vector space E^*
- $\langle s, x \rangle$: Paring $E^* \ni s \quad E \ni x$
- Ω : proper open convex cone in E
- Ω^* : the dual cone of Ω

$$\Omega^* := \{s \in E^* \mid \langle s, x \rangle > 0, \forall x \in \overline{\Omega} \setminus \{0\}\}$$
- T^* : (Orthogonal) dual subspace of $T \subset E$

$$T^* = \{s \in E^* \mid \langle s, x \rangle = 0, \forall x \in T\}$$

27

Conic Linear Program

Given

$$c \in E^*, f \in E \text{ and } T \subset E$$

- **Primal problem**

$$(P) \quad \begin{aligned} &\text{minimize } \langle c, x \rangle, \text{ s.t. } x \in \overline{\mathcal{P}}, \\ &\text{where } \mathcal{P} := (f + T) \cap \Omega, \end{aligned}$$

- **Dual problem**

$$(D) \quad \begin{aligned} &\text{maximize } \langle s, f \rangle, \text{ s.t. } s \in \overline{\mathcal{D}}, \\ &\text{where } \mathcal{D} := (c + T^*) \cap \Omega^*. \end{aligned}$$

28

Typical Examples

- Linear program (**LP**):

$$E = E^* = \mathbf{R}^n, \quad \Omega = \Omega^* = \mathbf{R}_{++}^n$$

- Semidefinite program (**SDP**):

$E = E^*$: the set of real symmetric matrices

$\Omega = \Omega^*$: the set of positive definite matrices

- Second order cone (Lorentz cone) program (**SOCP**)
- Mixture of the aboves

29

ϑ -normal barrier on an open convex cone Ω

[Nesterov & Nemirovski 94]

- Def. θ -normal barrier ψ on Ω

- A (smooth) convex function ψ satisfying, at each x in Ω ,

1) $\psi(tx) = \psi(x) - \vartheta \log t,$

2) $|(D^2\psi)_x(X, X, X)| \leq 2((D\psi)_x(X, X))^{3/2}$ self-concordance

for $\vartheta \geq 1, \forall t > 0$ and $\forall X \in T_x\Omega \cong E$

3) $\psi(x) \rightarrow +\infty$ ($x \rightarrow \text{bd } \Omega$),

Rem. (a) Existence for all Ω , but not with explicit forms

(symmetric cones \rightarrow **Yes, log. char. func.**)

(b) the Hessian is positive definite

(c) **Self-concordance** \Rightarrow the Newton method is efficient

- Ex. $\psi(x) = -\sum_{i=1}^n \log x^i$ (LP), $\psi(x) = -\log \det X$ (SDP)

Dually flat structure on Ω (revisited)

- D : the canonical **flat** affine connection on E

- $\{x^1, \dots, x^n\}$: affine coordinate system, i.e.,

- g : Riemannian metric on Ω

$$D_{\frac{\partial}{\partial x^i}} \frac{\partial}{\partial x^j} = 0$$

$$g = Dd\psi = \sum_{i,j} \frac{\partial^2 \psi}{\partial x^i \partial x^j} dx^i dx^j.$$

- D' : the dual affine connection on Ω

$$Xg(Y, Z) = g(D_X Y, Z) + g(Y, D'_X Z)$$

$(g, D, D') : \text{dually flat structure on } \Omega$

31

Remark

- $\{s_1, \dots, s_n\}$: dual coordinate system on E^* , s.t.

$$\langle s, x \rangle = \sum_i s_i(s) x^i(x)$$

- **Gradient map** $\iota : \Omega \rightarrow \Omega^*$ defined by

$$s_i \circ \iota = -\frac{\partial \psi}{\partial x^i} \quad \text{Legendre transform}$$

induces **dually flat structure on Ω^*** from (g, D, D')

- (1) D^* : the canonical **flat** affine connection on E^*

$$D_{\iota_*(X)}^* \iota_*(Y) = \iota_*(D'_X Y) \quad (\iota^* D^* = D')$$

D^* -autoparallel in Ω^* \longleftrightarrow D' -autoparallel in Ω

32

Remark

(2) Riemannian metric $g^* := D^*d\psi^*$ on Ω^*

$$g = \iota^* g^*$$

(3) $\langle \iota_*(X), Y \rangle = -g_x(X, Y)$

Hessian norm : We denote the length of X in $T_x\Omega \cong E$ by

$$\|X\|_x := \|Z\|_s := \sqrt{g_x(X, X)} = \sqrt{g_s^*(Z, Z)},$$

where $s = \iota(x)$ and $Z = \iota_*(X)$.

33

Curvature integral and iteration-complexity of IP

One of important computational performance indices for optimization algorithms is the **iteration-complexity**.

- Ω : **sym. cone** and $\mathcal{P} := (f + T) \cap \Omega$ is **DA**

\Rightarrow iteration-complexity=0 for (P)

- General case? Iter.-comp. is characterized by

- Curvature integrals along the **central trajectory** $\gamma_{\mathcal{P}}$

$$\int_{t_1}^{t_2} \|H_{\mathcal{P}}^*(\dot{\gamma}_{\mathcal{P}}(t), \dot{\gamma}_{\mathcal{P}}(t))\|_{\gamma_{\mathcal{P}}(t)}^{1/2} dt$$

- Similarly, for (D) curvature integrals along the dual c. t. $\gamma_{\mathcal{D}}$

$$\int_{t_1}^{t_2} \|H_{\mathcal{D}}(\dot{\gamma}_{\mathcal{D}}(t), \dot{\gamma}_{\mathcal{D}}(t))\|_{\gamma_{\mathcal{D}}(t)}^{1/2} dt$$

34

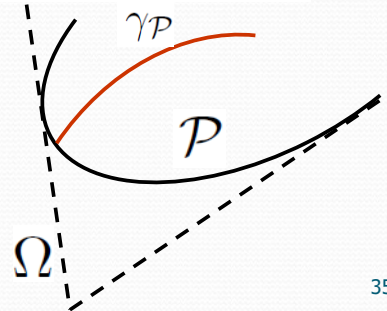
Central trajectory (= the special A-S trajectory)

- **Primal problem:** minimize $\langle c, x \rangle$, s.t. $x \in \overline{\mathcal{P}}$,
where $\mathcal{P} := (f + T) \cap \Omega$,

- $x(t) := \gamma_{\mathcal{P}}(t)$: the unique minimizer of
minimize $\underline{t}\langle c, x \rangle + \psi(x)$, s.t. $x \in \overline{\mathcal{P}}$.

for each $t > 0$

- $\gamma_{\mathcal{P}} := \{\gamma_{\mathcal{P}}(t) | t > 0\}$:
(Primal) **central trajectory**



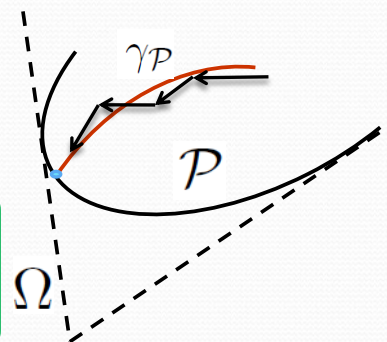
35

Central trajectory $\gamma_{\mathcal{P}}$

- Homotopy path to the **opt. sol. of the primal problem**, i.e., $x(t)$ converges when $t \rightarrow \infty$.
- Numerically tracing $\gamma_{\mathcal{P}}$ is the standard and efficient way to solve the primal problem.

Path-following method

Idea: consider the problem in the **dual cone Ω^*** in order to relate the complexity with the curvature



36

(1) Representation of feasible region

- a linear surj. operator $A : E \rightarrow \mathbf{R}^m$ s.t. $\text{Ker } A = T$

$$\mathcal{P} = \{x \in \Omega \mid Ax = b\},$$

$$\mathcal{D} = \{s \in \Omega^* \mid s = c - A^*y, y \in \mathbf{R}^m\}$$

where $A^* : \mathbf{R}^m \rightarrow E^*$ satisfying $y^T(Ax) = \langle A^*y, x \rangle$,

$$b := Af \in \mathbf{R}^m$$

- $\dim \mathcal{P} = n-m$, $\dim \mathcal{D} = m$

\mathcal{P} is D -autoparallel and \mathcal{D} is D^* -autoparallel

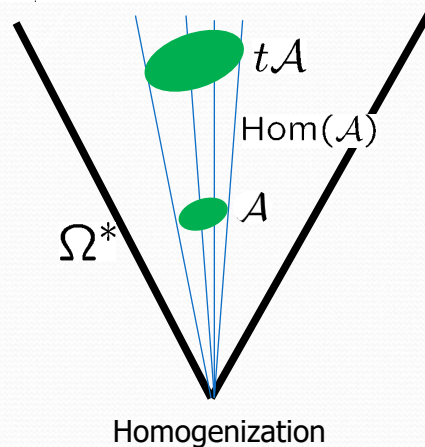
37

(2) Homogenization (conic hull)

- homogenization of \mathcal{D} in Ω^*

$$\text{Hom}(\mathcal{D}) := \bigcup_{t>0} t\mathcal{D}, \quad t\mathcal{D} := \{s \in \Omega^* \mid s = t\tilde{s}, \tilde{s} \in \mathcal{D}\}$$

- D^* -autoparallel
(because \mathcal{D} is.)
- $\dim \text{Hom}(\mathcal{D}) = m+1$



38

Lemma

The following relations hold in Ω^* :

$$\iota(\gamma_{\mathcal{P}}) = \iota(\mathcal{P}) \cap \text{Hom}(\mathcal{D})$$

$$s(t) := \iota(x(t)) = \iota(\mathcal{P}) \cap t\mathcal{D}$$

proved using the Lagrange function

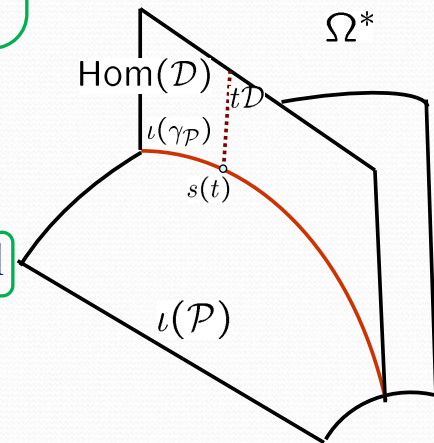
$$L(x, y) := t\langle c, x \rangle + \psi(x) + y^T(b - Ax)$$

$$\partial L / \partial x = 0 \rightarrow s \in t\mathcal{D}$$

\mathcal{P} is DA $\rightarrow \iota(\gamma_{\mathcal{P}})$ is D^* -autoparallel

Remark

$\iota(\mathcal{P})$ and $t\mathcal{D}$ are orthogonal
w.r.t. g^* at $s(t)$ by definition.



39

3. Geometric predictor-corrector algorithm (tracing $\gamma_{\mathcal{P}}$ in $\text{Hom}(\mathcal{D})$)

Ideal case

• Predictor

From $s(t) \in \iota(\gamma_{\mathcal{P}})$

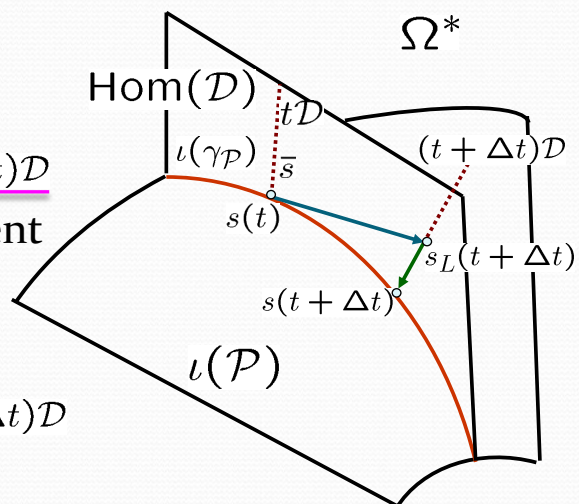
to $s_L(t + \Delta t) \in (t + \Delta t)\mathcal{D}$

with the direction tangent
to $\iota(\gamma_{\mathcal{P}})$

• Corrector

From $s_L(t + \Delta t) \in (t + \Delta t)\mathcal{D}$

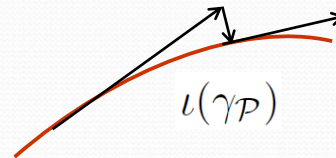
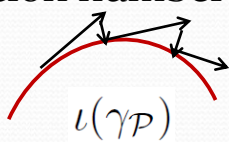
to $s(t + \Delta t) \in \iota(\gamma_{\mathcal{P}})$



40

Intuitive observation

- $H_{\mathcal{P}}^*(\dot{\gamma}_{\mathcal{P}}(t), \dot{\gamma}_{\mathcal{P}}(t))$: the **Euler-Schouten embedding curvature (second fundamental form)** of $\iota(\gamma_{\mathcal{P}})$ with respect to D^*
- If $H_{\mathcal{P}}^*(\dot{\gamma}_{\mathcal{P}}(t), \dot{\gamma}_{\mathcal{P}}(t))$ is small at t , so is expected the iteration number !?



Actually,

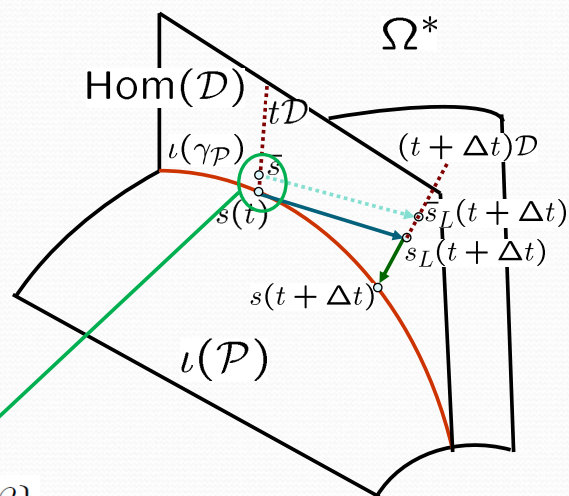
$$\ddot{s} = D_{\dot{s}}^* \dot{s} = \iota_*(H_{\mathcal{P}}^*(\dot{\gamma}_{\mathcal{P}}, \dot{\gamma}_{\mathcal{P}}))$$

41

Remark: practical case

- Cannot expect that the corrector returns precisely on $\iota(\gamma_{\mathcal{P}})$

- Consider the point \bar{s} in the **neighborhood** of $s(t) \in \iota(\gamma_{\mathcal{P}})$ in the sense of Riemannian metric



$$\mathcal{N}_t(\beta) := \{s \in t\mathcal{D} \mid \delta(s) \leq \beta\}$$

42

Predictor

- The differential equation expressing $\iota(\gamma_{\mathcal{P}})$:

$$\dot{s} = (\text{id} - \Pi_s^\perp)c = \frac{1}{t}(\text{id} - \Pi_s^\perp)s$$

where Π_s^\perp is the orthogonal projection w.r.t. g^* from E^* to $T^* = \text{Range} A^*$ at s .

- Hence, the predictor is defined by

$$\bar{s}_L(t + \Delta t) := \bar{s} + \Delta t(I - \Pi_{\bar{s}}^\perp)c \in (t + \Delta t)\mathcal{D}$$

43

Corrector

- Reduces to the following convex optimization on $t\bar{\mathcal{D}}$:

$$\text{minimize } F(s) := \langle s, f \rangle + \psi^*(s), \text{ s.t. } s \in t\bar{\mathcal{D}}$$

- Newton direction N** for this opt. problem:

$$D^*dF(X, N) = -dF(X), \forall X \in \mathcal{X}(t\bar{\mathcal{D}})$$

- Newton decrement:** measure of approximation of s to the optimal sol. s^*

$$\delta(s) := \|N\|_s$$

- We define the corrector with a single **Newton step** by:

$$\bar{s}_L^+(t + \Delta t) := \bar{s}_L(t + \Delta t) + N_{\bar{s}_L(t + \Delta t)}$$

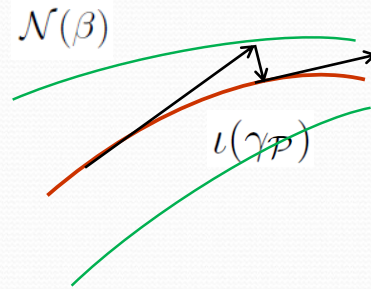
44

Tubular neighborhood

- The standard analysis technique in IP ensures the polynomiality of the complexity for this path-following strategy if all the generated points are near to $\iota(\gamma_{\mathcal{P}})$.
- Introduce the tubular neighborhood $\mathcal{N}(\beta)$ of $\iota(\gamma_{\mathcal{P}})$

$$\mathcal{N}(\beta) := \bigcup_{t \in (0, \infty)} \mathcal{N}_t(\beta),$$

where $\mathcal{N}_t(\beta) := \{s \in t\mathcal{D} \mid \delta(s) \leq \beta\}$.



45

4. Curvature integral and asymptotic iteration-complexity (Main result)

- **Assumption:** $\iota(\gamma_{\mathcal{P}})$ is **not D^* -autoparallel** (i.e. \mathcal{P} is not DA)

In this case, $\beta \rightarrow 0$ implies that $\Delta t \rightarrow 0$

- **Theorem**

For $0 < t_1 < t_2$ and $s_1 \in \mathcal{N}(\beta) \cap t_1\mathcal{D}$, let $\sharp(s_1, t_2, \beta)$ be **the iteration number** to find $s_2 \in \mathcal{N}(\beta) \cap t_2\mathcal{D}$. Then,

$$\lim_{\beta \rightarrow 0} \frac{\sqrt{\beta} \times \sharp(s_1, t_2, \beta)}{I_{\mathcal{P}}(t_1, t_2)} = 1,$$

where

$$I_{\mathcal{P}}(t_1, t_2) := \frac{1}{\sqrt{2}} \int_{t_1}^{t_2} \|H_{\mathcal{P}}^*(\dot{\gamma}_{\mathcal{P}}(t), \dot{\gamma}_{\mathcal{P}}(t))\|_{\dot{\gamma}_{\mathcal{P}}(t)}^{1/2} dt.$$

46

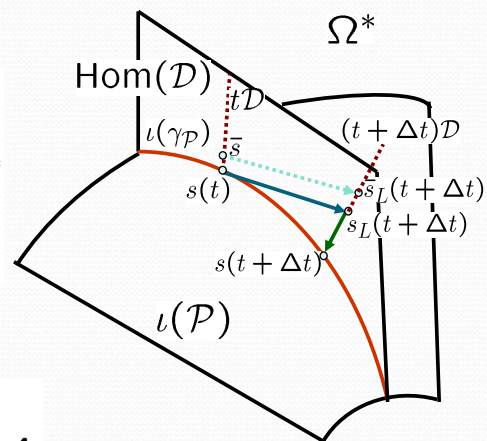
Outline of the proof

- Evaluate the Newton dec. of the predictor $\bar{s}_L(t + \Delta t)$ by $\|\ddot{s}(t)\|_{s(t)}$ (For each iteration)

$$\begin{aligned} & \delta(\bar{s}_L(t + \Delta t)) \\ &= \|s(t + \Delta t) - \bar{s}_L(t + \Delta t)\|_{\bar{s}_L(t + \Delta t)} + r_4 \\ &= \frac{(\Delta t)^2}{2} \|\ddot{s}(t)\|_{s(t)} + \delta(\bar{s}) + r_1 + r_2 + r_3, \end{aligned}$$

Rem. By the assumption

When $\beta \rightarrow 0$,

$$\delta(\bar{s}) \rightarrow 0, r_i \rightarrow 0 \quad \text{for } i = 1, \dots, 4,$$


47

Outline of the proof

- **Self-concordance** implies the following two inequalities:
Intermediate two relations for sufficiently small Δt and β .
(For each iteration)

- $$\begin{aligned} \sqrt{(1-\eta)\beta}(1-O(\sqrt{\beta})) &\leq \sqrt{w} - \sqrt{M_3}\delta(\bar{s}) \\ &\leq \frac{\Delta t}{\sqrt{2}} \|\ddot{s}(t)\|_{s(t)}^{1/2} + \sqrt{|r_1|} + \sqrt{M_3}(\Delta t)^2, \end{aligned}$$
- $$\begin{aligned} \frac{\Delta t}{\sqrt{2}} \|\ddot{s}(t)\|_{s(t)}^{1/2} - \sqrt{|r_1|} - \sqrt{M_3}(\Delta t)^2 \\ \leq \sqrt{w} + \sqrt{M_3}\delta(\bar{s}) \leq \sqrt{\beta}(1+O(\sqrt{\beta})) \end{aligned}$$

48

Outline of the proof

- Take summations of iterations

$$\begin{aligned}
 & \bullet \quad \sqrt{(1-\eta)\beta} \sum_{k=1}^{\sharp(s_1, t_2, \beta)} (1 - O(\sqrt{\beta})) \\
 & \qquad \leq \frac{1}{\sqrt{2}} \int_{t_1}^{t_2} \|\ddot{s}(t)\|_{s(t)}^{1/2} dt + M' \sqrt{\Delta t_{\max}}, \\
 & \bullet \quad \frac{1}{\sqrt{2}} \int_{t_1}^{t_2} \|\ddot{s}(t)\|_{s(t)}^{1/2} dt - M' \sqrt{\Delta t_{\max}} \\
 & \qquad \leq \sqrt{\beta} \sum_{k=1}^{\sharp(s_1, t_2, \beta)} (1 + O(\sqrt{\beta}))
 \end{aligned}$$

- Recall $\ddot{s} = D_{\dot{s}}^* \dot{s} = \iota_*(H_{\mathcal{P}}^*(\dot{\gamma}_{\mathcal{P}}, \dot{\gamma}_{\mathcal{P}}))$

49

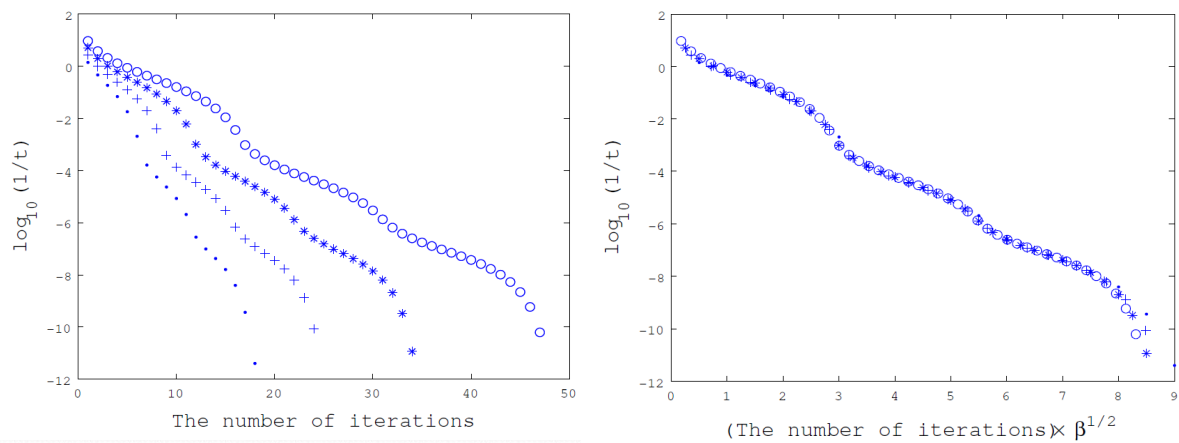
Remark

- An asymptotic result for $\beta \rightarrow 0$ (and hence, $\Delta t \rightarrow 0$)
 - \mathcal{P} is DA $\Rightarrow \iota(\gamma_{\mathcal{P}})$ is DA (D^* -autoparallel) $\Rightarrow \Delta t \rightarrow \infty \Rightarrow$ explicit sol.
- The same argument holds for the **dual** problem.
- The results are valid for general convex cones

50

Numerical experiment

- Curvature structure of CT for a certain LP



($\cdot : \beta = 1/4$, $+$: $\beta = 1/8$, $*$: $\beta = 1/16$, $\circ : \beta = 1/32$)

51

Curved part is Straight and
Straight part is Curved?(1)

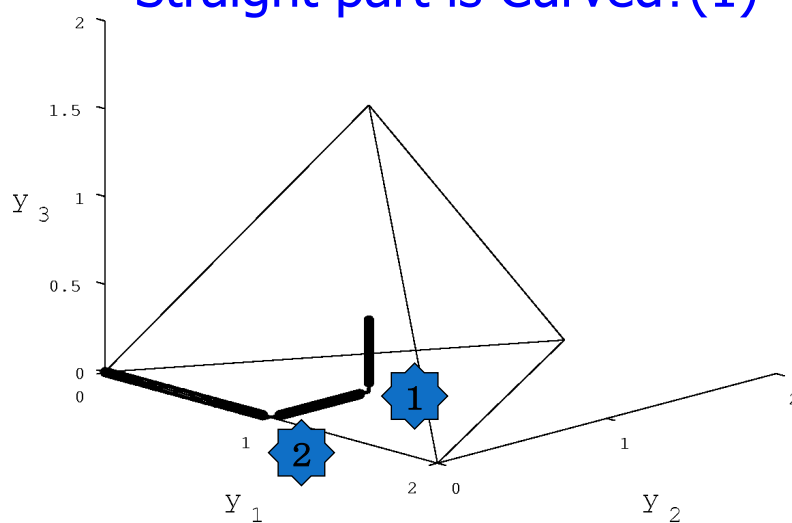


Figure 7: Exaple Figure

52

Proposition

It holds that $\|H_{\mathcal{P}}^*(\dot{\gamma}_{\mathcal{P}}, \dot{\gamma}_{\mathcal{P}})\|_{\gamma_{\mathcal{P}}(t)}^{1/2} \leq \frac{\sqrt{2\vartheta}}{t}$

ϑ : a constant determined by $\psi(x)$

- Remark

The above proposition gives the upper bound:

$$I_{\mathcal{P}}(t_1, t_2) \leq \sqrt{\vartheta} \log(t_2/t_1)$$

53

Further study for LP case

- Primal and Dual **Linear Program**:

$$\min c^T x$$

$$\text{s.t. } Ax = b, \quad x \geq 0, \quad A \in R^{m \times n}, b \in R^n$$

$$\max b^T y$$

$$\text{s.t. } s = c - A^T y, \quad s \geq 0,$$

54

Application to Primal-dual path-following (PDPF) method

- current **main-stream** IP (cheap in each iteration)
- The following quantity has been known to play an important and similar role in complexity analysis of **PDPF** method:

$$I_{PD}(t_1, t_2) = \int_{t_1}^{t_2} h_{PD}(t)^{1/2} dt$$

where $h_{PD}(t)$ is given by

$$h_{PD}(t) := \frac{1}{t^2} ((I_n - Q(t))e) \underline{*} (Q(t)e).$$

e : the unit element of Jordan product *

$Q(t)$: a certain projection matrix

55

Proposition

It holds that

$$h_{PD}(t)^2 = \left(\frac{1}{2} \|H_{\mathcal{P}}^*(\dot{\gamma}_{\mathcal{P}}(t), \dot{\gamma}_{\mathcal{P}}(t))\|_{\gamma_{\mathcal{P}}(t)} \right)^2 + \left(\frac{1}{2} \|H_{\mathcal{D}}(\dot{\gamma}_{\mathcal{D}}(t), \dot{\gamma}_{\mathcal{D}}(t))\|_{\gamma_{\mathcal{D}}(t)} \right)^2$$

Remark :

- **geometric implication** of the quantity of $I_{PD}(t_1, t_2)$

- inequalities $\max\{I_{\mathcal{P}}(t_1, t_2), I_{\mathcal{D}}(t_1, t_2)\} \leq I_{PD}(t_1, t_2) \leq I_{\mathcal{P}}(t_1, t_2) + I_{\mathcal{D}}(t_1, t_2).$

56

Concluding Remark

- Tractable characterization of **DA** submfdns in symmetric cones Ω
- Application to conic linear programs
 - Explicit sol. when the feasible region M is DA in Ω .
 - M is DA \Rightarrow **AS (CT) traj. is DA** (D^* -autoparallel) $\Rightarrow \Delta t \rightarrow \infty \Rightarrow$ explicit sol.
- **Extension**: # of iterations and **curvature integral** of CT
 - Asymptotic analysis ($\beta \rightarrow 0$)
 - Complemented by numerical experiment for finite β
 - Geometric structure of CT has an influence on complexity of the IP algorithm

57

- Relation among iteration-complexities of **primal-, dual- and primal-dual algorithms**.
- DA submanifolds in the set of invertible elements of Jordan algebras [OIT23]
- Future work: **Geometrical** study for general stat. mfd.
 - Various geometrical concepts for mutually dual connections and their characterizations (Furuhata *et al.*)
 - Classifications
 - Families of continuous probability densities
 - Applications (Ex. Study of ODE's on manifolds?)

58

Thank you for your attention

References:

[OIT₂₃] A. Ohara, H. Ishi and T. Tsuchiya,
Doubly autoparallel structure and curvature integrals:
Applications to iteration complexity for solving convex
programs,
Information Geometry (2023),
<https://doi.org/10.1007/s41884-023-00116-x> (Open Access).

Uncovering Data Symmetries: Estimating Covariance Matrix in High-Dimensional Setting With 'gips' R Package

Adam Przemysław Chojecki¹
Hideyuki Ishi²

¹Warsaw University of Technology (Poland)

²Osaka Metropolitan University

In high-dimensional settings, where the number of variables exceeds the number of observations, accurately estimating the covariance matrix poses a significant challenge. This talk presents a novel approach that leverages the identification of symmetries within the data to improve covariance matrix estimation. In the 'gips' R package [2], we implement the Bayesian model selection procedure within Gaussian vectors (invariant under the permutation group) introduced in [1].

Our method aims to capture the underlying low-dimensional structure by exploring the permutation symmetries within the data. Identifying symmetries enables us to interpret relationships in data in a new and natural way. The 'gips' package provides a comprehensive set of functions that facilitate identifying and utilizing symmetries, making it a valuable resource for researchers working with high-dimensional data.

We demonstrate the effectiveness of our approach through simulations and real-world data examples. Our results show that incorporating data symmetries leads to more reliable covariance matrix estimates, enabling better inference and decision-making. More results can be found in [3].

The presented novel approach contributes to the growing field of statistical methods for $p > n$, offering promising avenues for future research and practical applications.

References

- [1] P. Graczyk, H. Ishi, B. Kołodziejek and H. Massam., *Model selection in the space of Gaussian models invariant by symmetry*, The Annals of Statistics, 2022
- [2] A. Chojecki, P. Morgen, B. Kołodziejek, *gips: Gaussian Model Invariant by Permutation Symmetry*, CRAN, 2022, <https://CRAN.R-project.org/package=gips>, <https://przechoj.github.io/gips/>
- [3] A. Chojecki, P. Morgen and B. Kołodziejek, *Learning permutation symmetries with gips in R*, <https://arxiv.org/abs/2307.00790>

Undirected graphs - independence
Colored graphs - equality; gips
Conclusions
References

Uncovering Data Symmetries: Estimating Covariance Matrix in High-Dimensional Setting With 'gips' R Package

Adam Przemysław Chojecki

MiNI PW

14.12.2023

Adam Przemysław Chojecki

Uncovering Data Symmetries with 'gips'

Undirected graphs - independence
Colored graphs - equality; gips
Conclusions
References

What is gips?

gips

lifecycle stable CRAN 1.2.1 R-CMD-check passing codecov 96%



gips - Gaussian model Invariant by Permutation Symmetry

gips is an R package that looks for permutation symmetries in the multivariate Gaussian sample. Such symmetries reduce the free parameters in the unknown covariance matrix. This is especially useful when the number of variables is substantially larger than the number of observations.

gips will help you with two things:

1. Finding hidden symmetries between the variables. **gips** can be used as an exploratory tool for searching the space of permutation symmetries of the Gaussian vector. Useful in the Exploratory Data Analysis (EDA).
2. Covariance estimation. The Maximum Likelihood Estimator (MLE) for the covariance matrix is known to exist if and only if the number of variables is less or equal to the number of observations. Additional knowledge of symmetries significantly weakens this requirement. Moreover, the reduction of model dimension brings the advantage in terms of precision of covariance estimation.

Installation

From [CRAN](https://cran.r-project.org/web/packages/gips/index.html):

```
# Install the released version from CRAN:  
install.packages("gips")
```



Figure: Documentation of our R package, **gips**, <https://przechoj.github.io/gips/>

Adam Przemysław Chojecki

Uncovering Data Symmetries with 'gips'

Undirected graphs - independence
Colored graphs - equality; gips
Conclusions
References

Inspiring citation

*Life is the art of drawing
sufficient conclusions from
insufficient premises.*

— Samuel Butler

Adam Przemysław Chojewski

Uncovering Data Symmetries with 'gips'

Undirected graphs - independence
Colored graphs - equality; gips
Conclusions
References

Presentation plan

- 1 Undirected graphs - independence
 - Motivation
 - Example of the gain
 - How to estimate it?
 - Empirical data example
- 2 Colored graphs - equality; gips
 - Empirical data example - Additional Equalities
 - Permutational symmetry example
 - Comparison with similar methods
- 3 Conclusions

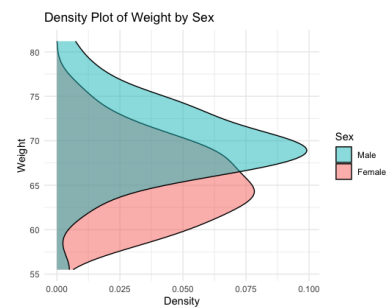
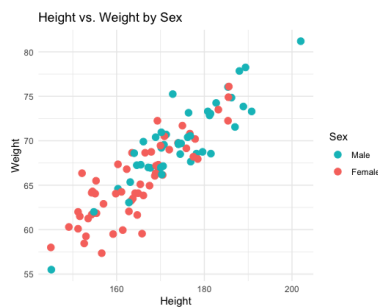
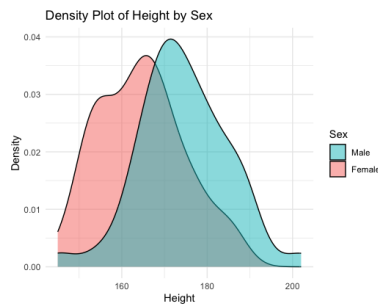
Adam Przemysław Chojewski

Uncovering Data Symmetries with 'gips'

Undirected graphs - independence
Colored graphs - equality; gips
Conclusions
References

Motivation
Example of the gain
How to estimate it?
Empirical data example

Example of conditional Independence



Adam Przemysław Chojewski

Uncovering Data Symmetries with 'gips'

Undirected graphs - independence
Colored graphs - equality; gips
Conclusions
References

Motivation
Example of the gain
How to estimate it?
Empirical data example

Example assumption

$$(W \perp\!\!\!\perp S) | H$$

$$p(W, H, S) = \frac{p(W, H)p(H, S)}{p(H)}$$

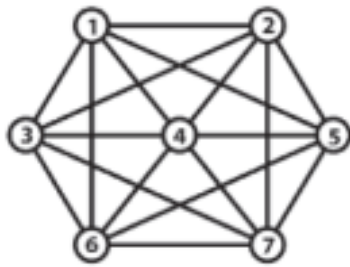
Adam Przemysław Chojewski

Uncovering Data Symmetries with 'gips'

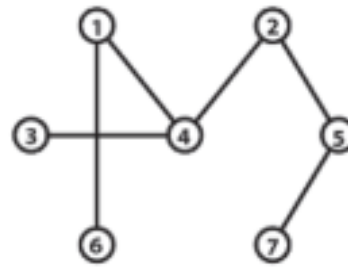
Undirected graphs - independence
Colored graphs - equality; gips
Conclusions
References

Motivation
Example of the gain
How to estimate it?
Empirical data example

Sparse Graph Models



Dense



Sparse

Adam Przemysław Chojewski

Uncovering Data Symmetries with 'gips'

Undirected graphs - independence
Colored graphs - equality; gips
Conclusions
References

Motivation
Example of the gain
How to estimate it?
Empirical data example

Covariance matrix estimation

$$X_1, X_2, \dots, X_n \sim \mathcal{N}_p(0, \Sigma)$$

When $n > p$:

- $\hat{\Sigma}_{MLE} = \frac{1}{n} \sum_{i=1}^n (X_i)(X_i)^T$; $\hat{K}_{MLE} = \hat{\Sigma}_{MLE}^{-1}$
- If we know the graph G & $n > \text{maximum clique}$, then $\hat{K}_{G,MLE} = \dots \hat{K}_{MLE}$ with zeros putted in places where G has no edge,
 $\hat{\Sigma}_{G,MLE} = \dots \hat{K}_{G,MLE}^{-1}$

More can be found in article [5] or in the book [4].

Authors of LASSO [6] also made the paper [1] where they introduce a GLASSO method for graph estimation.

Adam Przemysław Chojewski

Uncovering Data Symmetries with 'gips'

Undirected graphs - independence
Colored graphs - equality; gips
Conclusions
References

Motivation
Example of the gain
How to estimate it?
Empirical data example

Empirical data example

Table 1. Empirical concentrations $\times 1000$ (on or above the diagonal) and partial correlations (below the diagonal) for the examination marks in five mathematical subjects.

	Mechanics	Vectors	Algebra	Analysis	Statistics
Mechanics	5.24	-2.44	-2.74	0.01	-0.14
Vectors	0.33	10.43	-4.71	-0.79	-0.17
Algebra	0.23	0.28	26.95	-7.05	-4.70
Analysis	-0.00	0.08	0.43	9.88	-2.02
Statistics	0.02	0.02	0.36	0.25	6.45

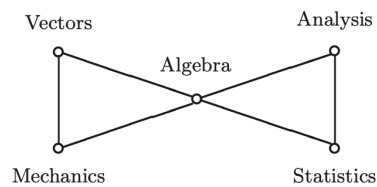


Fig. 1. Conditional independence structure of examination marks for 88 students.

Figure: Table and Fig., from [2]

Adam Przemysław Chojewski

Uncovering Data Symmetries with 'gips'

Undirected graphs - independence
Colored graphs - equality; gips
Conclusions
References

Empirical data example - Additional Equalities
Permutational symmetry example
Comparison with similar methods

Empirical data example

Table 1. Empirical concentrations $\times 1000$ (on or above the diagonal) and partial correlations (below the diagonal) for the examination marks in five mathematical subjects.

	Mechanics	Vectors	Algebra	Analysis	Statistics
Mechanics	5.24	-2.44	-2.74	0.01	-0.14
Vectors	0.33	10.43	-4.71	-0.79	-0.17
Algebra	0.23	0.28	26.95	-7.05	-4.70
Analysis	-0.00	0.08	0.43	9.88	-2.02
Statistics	0.02	0.02	0.36	0.25	6.45

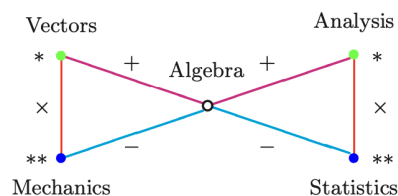


Fig. 8. Coloured graph of an RCOP symmetry model for the examination marks of 88 students. The distribution of the marks is unchanged if we simultaneously replace Vectors with Analysis and Mechanics with Statistics.

Figure: Table and Fig., from [2]

Adam Przemysław Chojewski

Uncovering Data Symmetries with 'gips'

Undirected graphs - independence
Colored graphs - equality; gips
Conclusions
References

Empirical data example - Additional Equalities
Permutational symmetry example
Comparison with similar methods

Permutational symmetry example

This is the real covariance matrix
We want to estimate it based on $n = 4$ observations

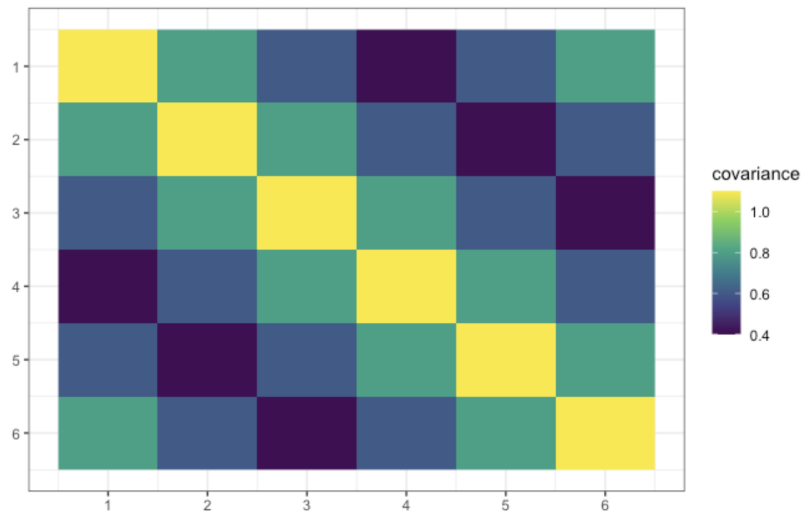


Figure: Fig., from [3]; Symmetry $\Gamma = \langle (1, 2, 3, 4, 5, 6) \rangle$

Adam Przemysław Chojecki

Uncovering Data Symmetries with 'gips'

Undirected graphs - independence
Colored graphs - equality; gips
Conclusions
References

Empirical data example - Additional Equalities
Permutational symmetry example
Comparison with similar methods

Permutational symmetry example

Covariance estimated in standard way

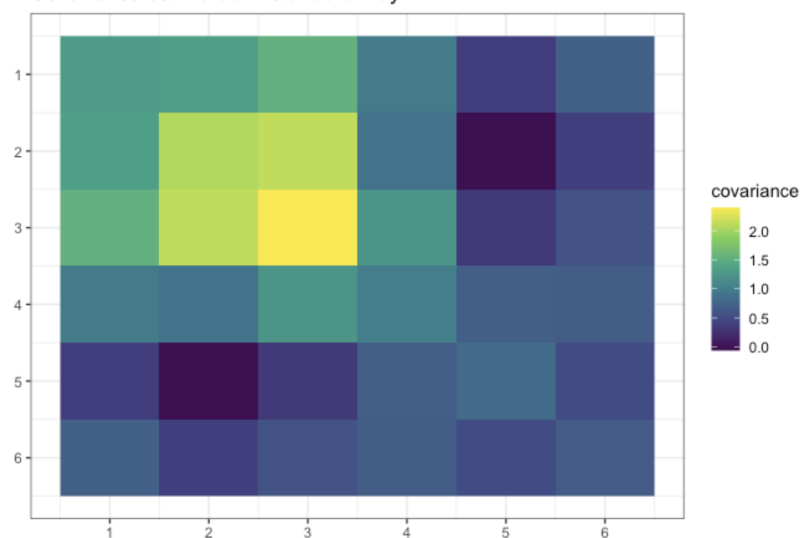


Figure: Fig., from [3]; no symmetry assumed

Adam Przemysław Chojecki

Uncovering Data Symmetries with 'gips'

Undirected graphs - independence
Colored graphs - equality; gips
Conclusions
References

Empirical data example - Additional Equalities
Permutational symmetry example
Comparison with similar methods

Permutational symmetry example

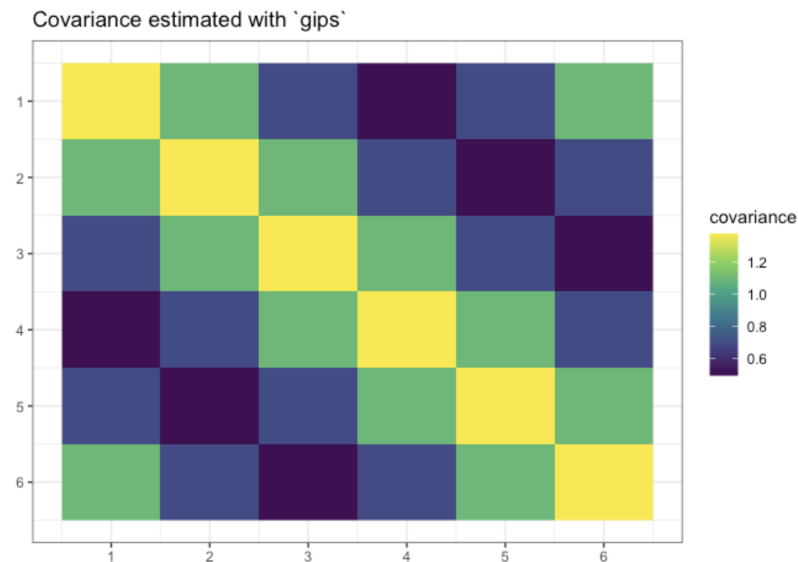


Figure: Fig., from [3]; gips automatically found symmetry $\Gamma = \langle (1, 2, 3, 4, 5, 6) \rangle$

Adam Przemysław Chojewski

Uncovering Data Symmetries with 'gips'

Undirected graphs - independence
Colored graphs - equality; gips
Conclusions
References

Empirical data example - Additional Equalities
Permutational symmetry example
Comparison with similar methods

Recognition on big matrix

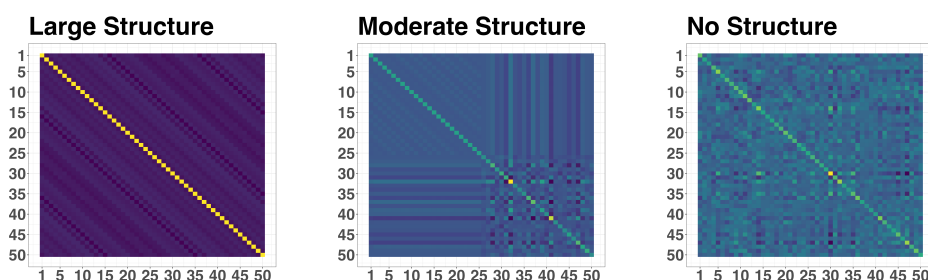


Figure: True covariance matrices corresponding to the three scenarios:
left panel - large structure $\pi_{\langle(1,2,\dots,50)\rangle}(S)$,
middle panel - moderate structure $\pi_{\langle(1,2,3,\dots,25)\rangle}(S)$,
right panel - no structure S .

Adam Przemysław Chojewski

Uncovering Data Symmetries with 'gips'

Undirected graphs - independence
Colored graphs - equality; gips
Conclusions
References

Empirical data example - Additional Equalities
Permutational symmetry example
Comparison with similar methods

Performance of matrix restoration

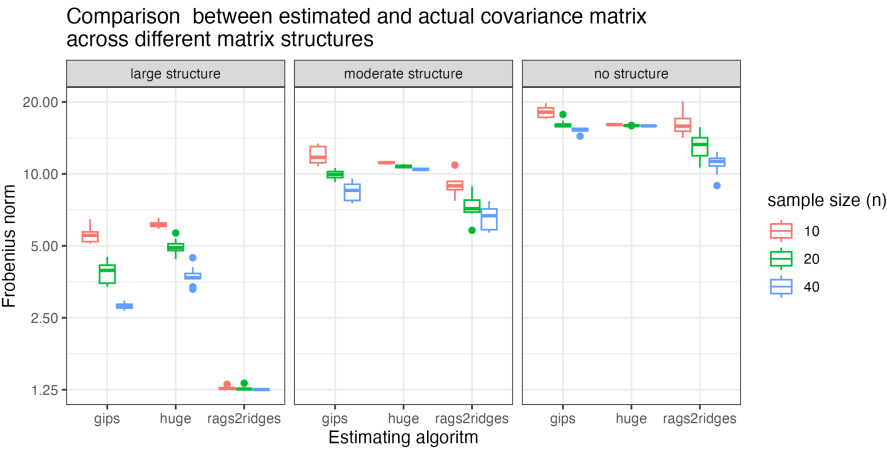


Figure: Frobenius norm (on a logarithmic scale) of the difference of the estimate and the true covariance matrix. 10 runs for each configuration.

Undirected graphs - independence
Colored graphs - equality; gips
Conclusions
References

Inspiring citation

If you torture the data long enough, it will confess

— Ronald Coase

Undirected graphs - independence
 Colored graphs - equality; gips
Conclusions
 References

Conclusions

- ① Graphs are a convenient language for recording relationships in data.
- ② Expert knowledge is an integral part of statistical analysis.
- ③ Expert knowledge can be captured with a graph and used in modeling.
- ④ There are methods for learning graphs from data, i.e., automatic learning of expert knowledge.
- ⑤ There are many ways to put constraints on parameters. The appropriate choice of these depends on the nature of the collected data.
- ⑥ It is possible, when one tries, to draw false conclusions from the data.

Adam Przemysław Chojewski

Uncovering Data Symmetries with 'gips'

Undirected graphs - independence
 Colored graphs - equality; gips
 Conclusions
References

Bibliography I

- [1] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. "Sparse inverse covariance estimation with the graphical lasso". In: *Biostatistics* 9.3 (Dec. 2007), pp. 432–441.
- [2] Søren Højsgaard and Steffen L. Lauritzen. "Graphical Gaussian Models with Edge and Vertex Symmetries". In: *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* 70.5 (2008), pp. 1005–1027. ISSN: 13697412, 14679868. URL: <http://www.jstor.org/stable/20203867> (visited on 10/30/2023).
- [3] *R package: gips - Gaussian model Invariant by Permutation Symmetry*. przechoj.github.io/gips. [Online; 14-12-2023].

Adam Przemysław Chojewski

Uncovering Data Symmetries with 'gips'

Undirected graphs - independence
 Colored graphs - equality; gips
 Conclusions
 References

Bibliography II

- [4] Alberto Roverato. *Graphical Models for Categorical Data*. SemStat Elements. Cambridge University Press, 2017. DOI: 10.1017/9781108277495.
- [5] KAYVAN SADEGHI and STEFFEN LAURITZEN. “Markov properties for mixed graphs”. In: *Bernoulli* 20.2 (2014), pp. 676–696. ISSN: 13507265. URL: <http://www.jstor.org/stable/42919409>.
- [6] Robert Tibshirani. “Regression Shrinkage and Selection Via the Lasso”. In: *Journal of the Royal Statistical Society: Series B (Methodological)* 58.1 (1996), pp. 267–288.

Adam Przemysław Chojewski

Uncovering Data Symmetries with 'gips'

Undirected graphs - independence
 Colored graphs - equality; gips
 Conclusions
 References

Thanks for attention
 Happy modelling!

Adam Przemysław Chojewski

Uncovering Data Symmetries with 'gips'

Mathematical complement

$$\sigma \in G_p \quad R(\sigma) := \sum_{i=1}^p E_{\sigma(i), i} \in GL(p, \mathbb{R}) \quad \left(\begin{array}{l} \text{e.g. } \sigma = (1 \ 2) \in G_3 \\ \Rightarrow R(\sigma) = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \end{array} \right)$$

permutation matrix

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix} \sim N(0, \Sigma) \quad (\Sigma \in S_{\text{ym}}^+(p, \mathbb{R}))$$

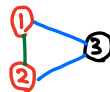
$$R(\sigma)X \stackrel{d}{=} X \Leftrightarrow {}^t R(\sigma) \Sigma R(\sigma) = \Sigma$$

identically distributed

$$\Gamma \subset G_p \quad \mathcal{P}_\Gamma := \left\{ \Sigma \in S_{\text{ym}}^+(p, \mathbb{R}) ; \forall \sigma \in \Gamma \quad {}^t R(\sigma) \Sigma R(\sigma) = \Sigma \right\}$$

subgr.

$$\text{Ex. } p=3, \Gamma = \langle (1 \ 2) \rangle$$



$$\mathcal{P}_\Gamma = \left\{ \Sigma = \begin{pmatrix} a & b & c \\ b & a & c \\ c & c & d \end{pmatrix} \in S_{\text{ym}}^+(3, \mathbb{R}) ; a, b, c, d \in \mathbb{R} \right\}$$

$$\mathcal{M}_\Gamma := \{ N(0, \Sigma) ; \Sigma \in \mathcal{P}_\Gamma \} : \text{statistical model of Gaussian distr. with } \Gamma\text{-invariance}$$

Problem For given samples $x^{(1)}, x^{(2)}, \dots, x^{(n)} \in \mathbb{R}^p$,
 find the most reasonable subgroup $\Gamma \subset G_p$.
 (model selection)

Bayesian approach:

$\Gamma \subset G_p$ and $K := \Sigma^{-1}$ are RANDOM VARIABLES

$$p(\Gamma) := \frac{1}{\#(\text{subgroups of } G_p)} \quad \text{uniform distribution}$$

Fix $\delta > 0$ and $D \in \text{Sym}^+(p, \mathbb{R})$ (hyperparameters)

$$p(K|\Gamma) := \frac{1}{I_\Gamma(\delta, D)} e^{-\text{tr} KD/2} (\det K)^{\frac{\delta-2}{2}} \mathbb{1}_{\mathcal{P}_\Gamma}(K)$$

Diaconis-Ylvisaka conjugate prior

where

$$I_\Gamma(\delta, D) := \int_{\mathcal{P}_\Gamma} e^{-\text{tr} KD/2} (\det K)^{\frac{\delta-2}{2}} dK$$

$$p(\underbrace{x^{(1)}, \dots, x^{(n)}}_{\Sigma^{-1}} | K) = (2\pi)^{-np/2} (\det K)^{\frac{n}{2}} e^{-\text{tr}(Ky)/2}$$

($y := \sum_{k=1}^n x^{(k)} x^{(k)\top} \in \text{Sym}(p, \mathbb{R})$)

Then

$$p(x^{(1)}, \dots, x^{(n)}, \Sigma, \Gamma) = p(x^{(1)}, \dots, x^{(n)} | \Sigma) p(\Sigma | \Gamma) p(\Gamma)$$

$$\text{We want to know } p(\Gamma | x^{(1)}, \dots, x^{(n)}) = \frac{p(x^{(1)}, \dots, x^{(n)} | \Gamma) p(\Gamma)}{p(x^{(1)}, \dots, x^{(n)})} \\ \propto p(x^{(1)}, \dots, x^{(n)} | \Gamma)$$

$$p(x^{(1)}, \dots, x^{(n)} | \Gamma) = \int_{\mathcal{P}_\Gamma} p(x^{(1)}, \dots, x^{(n)} | K) p(K | \Gamma) dK \\ = (2\pi)^{-np/2} \frac{I_\Gamma(\delta+n, D+y)}{I_\Gamma(\delta, D)}$$

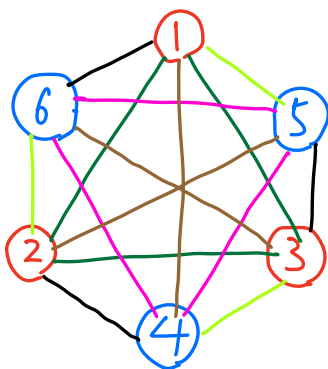
GIKM 2022 : EXACT FORMULA for

$$I_{\Gamma}(S, D) = \int_{\mathcal{P}_{\Gamma}} e^{-\text{tr}(KD)} (\det K)^{\frac{S-2}{2}} dK$$

\leadsto implemented in GIPS.

Point : Group representation of Γ

Ex. The case $\Gamma = \langle \sigma \rangle \subset \mathbb{S}_6$ with $\sigma = (1\ 2\ 3)(4\ 5\ 6)$



$$K = \begin{pmatrix} a & b & b & e & f & g \\ b & a & b & g & e & f \\ b & b & a & f & g & e \\ e & g & f & c & d & d \\ f & e & g & d & c & d \\ g & f & e & d & d & c \end{pmatrix} \in \mathcal{P}_{\Gamma}$$

Put

$$U_{\Gamma} := \begin{pmatrix} 1/\sqrt{3} & 0 & 2/\sqrt{6} & 0 & 0 & 0 \\ 1/\sqrt{3} & 0 & -1/\sqrt{6} & 1/\sqrt{2} & 0 & 0 \\ 1/\sqrt{3} & 0 & -1/\sqrt{6} & -1/\sqrt{2} & 0 & 0 \\ 0 & 1/\sqrt{3} & 0 & 0 & 2 & 0 \\ 0 & 1/\sqrt{3} & 0 & 0 & -1/\sqrt{6} & 1/\sqrt{2} \\ 0 & 1/\sqrt{3} & 0 & 0 & -1/\sqrt{6} & 1/\sqrt{2} \end{pmatrix}$$

Then

$${}^t U_{\Gamma} K U_{\Gamma} = \begin{pmatrix} A & B & 0 & 0 & 0 & 0 \\ B & C & 0 & 0 & 0 & 0 \\ 0 & 0 & D & 0 & F & G \\ 0 & 0 & 0 & D & -G & F \\ 0 & 0 & F & -G & E & 0 \\ 0 & 0 & G & F & 0 & E \end{pmatrix} \sim \begin{pmatrix} A & B \\ B & C \end{pmatrix} \oplus \begin{pmatrix} D & F-iG \\ F+iG & E \end{pmatrix}$$

U_Γ comes from an irreducible decomposition of \mathbb{R}^6 as Γ -m.

$$\mathbb{R}^6 = \mathbb{R} \begin{pmatrix} 1/\sqrt{3} \\ 1/\sqrt{3} \\ 1/\sqrt{3} \\ 0 \\ 0 \\ 0 \end{pmatrix} \oplus \mathbb{R} \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1/\sqrt{3} \\ 1/\sqrt{3} \\ 1/\sqrt{3} \end{pmatrix} \oplus \left\langle \begin{pmatrix} 2/\sqrt{6} \\ -1/\sqrt{6} \\ -1/\sqrt{6} \\ 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1/\sqrt{5} \\ 1/\sqrt{5} \\ 0 \\ 0 \\ 0 \end{pmatrix} \right\rangle_{\mathbb{R}} \oplus \left\langle \begin{pmatrix} 0 \\ 0 \\ 0 \\ 2/\sqrt{6} \\ -1/\sqrt{6} \\ -1/\sqrt{6} \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1/\sqrt{3} \\ 1/\sqrt{3} \end{pmatrix} \right\rangle_{\mathbb{R}}$$

$$\mathcal{P}_\Gamma \simeq \text{Sym}^+(2, \mathbb{R}) \oplus \text{Herm}^+(2, \mathbb{C})$$

$$\begin{array}{ccc} \mathbb{C} & & \\ \downarrow & & \\ K & \longmapsto & \begin{pmatrix} A & B \\ B & C \end{pmatrix} \oplus \begin{pmatrix} D & F-iG \\ F+iG & E \end{pmatrix} \end{array}$$

$$K = \pi_\Gamma(X) \in \mathcal{P}_\Gamma \quad (X \in \text{Sym}^+(6, \mathbb{R}))$$

ortho. proj. w.r.t. the trace inner product

$$\Rightarrow I_\Gamma(\delta, X) = (2\pi)^{3/2} 2^{-\delta} \Gamma\left(\frac{\delta+1}{2}\right) \Gamma\left(\frac{\delta}{2}\right) \Gamma(\delta) \Gamma(\delta-1) \\ \times (AC-B^2)^{-(\delta+1)/2} (DE-F^2-G^2)^{-\delta}$$

$$\Gamma \subset G_p$$

$$\mathbb{R}^p = \bigoplus_{i=1}^L V_i^{\oplus r_i}$$

find U_Γ (only for cyclic group in GIKM)

$$\Rightarrow \mathcal{P}_\Gamma \simeq \bigoplus_{i=1}^L \text{Herm}^+(r_i, \mathbb{K}_i)$$

$$\Rightarrow \text{formula for } I_\Gamma(\delta, D)$$

Maximum likelihood estimation for discrete exponential families, its geometry and combinatorics

Tomasz Skalski

Wrocław University of Science and Technology, Poland

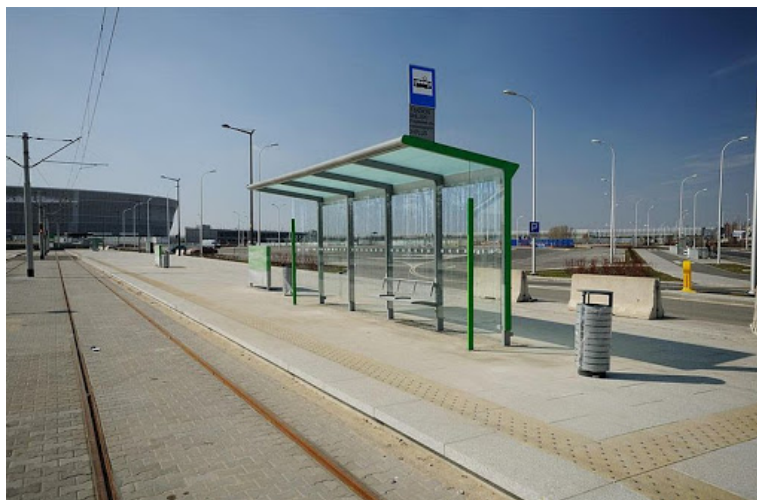
We discuss the existence of the maximum likelihood estimator for discrete exponential family on finite set. Using the newly introduced notion of sets of uniqueness, we present new criterion for the existence of MLE. We show how this criterion can be applied in various discrete settings, in which it can be easily solved using the tools from random graph theory and discrete geometry. Most notably, we discuss the MLE existence for exponential models of random graphs and for linear spaces spanned by Rademacher functions. Additionally, we give a few remarks concerning the existence of MLE for spaces spanned by products of Rademacher functions. As an application, we discuss the asymptotics of the size of independent identically distributed samples for which the maximum likelihood estimator exists with high probability.

References

- [1] K. Bogdan, M. Bosy, T. Skalski, *Maximum likelihood estimation for discrete exponential families and random graphs*, ALEA Lat. Am. J. Probab. Math. Stat. 19 (2022), no. 1, 1045–1070.

Osaka & on-line
2023/12/14

Estimation



Estimation



Estimation

- Waiting time τ – random variable
- Number of waiting passengers N – observed data
- $\hat{\tau}(N)$ – estimator of τ

When the estimator is "good"?

Popular choice – to maximise the likelihood function
(Maximum Likelihood Estimation, MLE)

In other words, we are looking for the parameter, under which the given situation is most likely to be present.

Notation

- $\mathcal{X} = \{x_1, \dots, x_K\}$ – finite state space, $K = |\mathcal{X}|$
- $\mu : \mathcal{X} \rightarrow (0, \infty)$ – weight function
- $\mathcal{B} \subset \mathbb{R}^{\mathcal{X}}$ – linear space of functions ($\phi = \mathbb{1} \in \mathcal{B}$)
- $\mathcal{B}_+ = \{\phi \in \mathcal{B} : \phi \geq 0\}$ – subclass (cone) of non-negative functions
- $N(\phi) = \sum_{x \in \mathcal{X}} e^{\phi(x)} \mu(x)$ – normalising constant (partition function)
- $p = e(\phi) = \frac{e^\phi}{N(\phi)}$ – exponential density
- $e(\mathcal{B}) = \{p = e(\phi) : \phi \in \mathcal{B}\}$ – exponential family

MLE

Definition

Let x_1, \dots, x_n be a sample from the finite set \mathcal{X} and let $\phi \in \mathcal{B}$. The likelihood function of $p = e(\phi)$ is defined as:

$$L_p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i).$$

Joint density: function of (x_1, \dots, x_n)

Likelihood: function of p

Definition

The $\hat{p} \in e(\mathcal{B})$ is called the maximum likelihood estimator (MLE), if

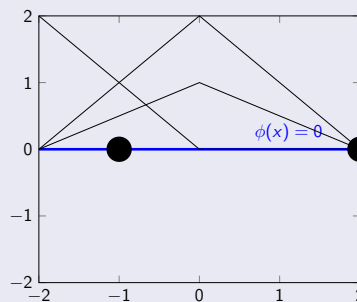
$$L_{\hat{p}}(x_1, \dots, x_n) = \sup_{p \in e(\mathcal{B})} L_p(x_1, \dots, x_n).$$

Definition

U is a set of uniqueness for \mathcal{B}_+ , if $[\phi \in \mathcal{B}_+, \phi(U) = 0] \Rightarrow [\phi \equiv 0]$.

Example

Again, let $\mathcal{X} = \{-2, -1, 0, 1, 2\}$ and let \mathcal{B} be the class of all the real functions on \mathcal{X} that are linear (affine) on $\{-2, -1, 0\}$ and on $\{0, 1, 2\}$.



Then the set $\{-1, 2\}$ is of uniqueness for \mathcal{B}_+ .

Existence of MLE – main criterion

Theorem (K. Bogdan, M. Bosy, TS (2022))

The maximum likelihood estimator for $e(\mathcal{B})$ and $x_1, \dots, x_n \in \mathcal{X}$ exists if and only if $\{x_1, \dots, x_n\}$ is a set of uniqueness for \mathcal{B}_+ .

Proof.

(\Rightarrow) If $\{x_1, \dots, x_n\}$ is not of uniqueness for \mathcal{B}_+ , we may subtract from every candidate for MLE ϕ a non-negative function ψ vanishing on $\{x_1, \dots, x_n\}$, so $\psi - \phi = \psi$ on $\{x_1, \dots, x_n\}$. Thus $N(\psi - \phi) < N(\psi)$ and the resulting likelihood is increased. \square

Existence of MLE – main criterion

Theorem (K. Bogdan, M. Bosy, TS (2022))

The maximum likelihood estimator for $e(\mathcal{B})$ and $x_1, \dots, x_n \in \mathcal{X}$ exists if and only if $\{x_1, \dots, x_n\}$ is a set of uniqueness for \mathcal{B}_+ .

Proof.

(\Leftarrow) We introduce a seminorm

$$\lambda_U(\phi) = \max_{\mathcal{X}}(\phi) - \min_U(\phi)$$

related to given set $U \subset \mathcal{X}$ and compare it with an oscillation seminorm $\lambda_{\mathcal{X}}(\phi)$. Both seminorms are comparable since U is a set of uniqueness for \mathcal{B}_+ . □

Applications

Two types of proposed applications:

- Conditions for the existence of MLE for specific exponential families
- Probability bounds for MLE for i.i.d. samples

For the i.i.d. random variables X_1, X_2, \dots valued in \mathcal{X} we define the following (random) time:

$$\nu_{\text{uniq}} = \inf\{n \geq 1 : \{X_1, \dots, X_n\} \text{ is a set of uniqueness for } \mathcal{B}_+\}.$$

Probabilistic tool — Threshold functions

Definition (Threshold)

A function $n^* = n^*(K)$ is a threshold of the size of the sample $\mathbb{X} = (X_1, \dots, X_n)$ for a given (monotone) property \mathcal{P} if

$$\lim_{K \rightarrow \infty} \mathbb{P}(\mathbb{X} \in \mathcal{P}) = \begin{cases} 0 & \text{if } n(K)/n^*(K) \rightarrow 0, \\ 1 & \text{if } n(K)/n^*(K) \rightarrow \infty, \end{cases} \quad K \rightarrow \infty.$$

Definition (Sharp threshold)

A function $n^* = n^*(K)$ is a sharp threshold of the size of the sample $\mathbb{X} = (X_1, \dots, X_n)$ for a given (monotone) property \mathcal{P} if for every $\varepsilon > 0$

$$\lim_{K \rightarrow \infty} \mathbb{P}(\mathbb{X} \in \mathcal{P}) = \begin{cases} 0 & \text{if } n(K)/n^*(K) < 1 - \varepsilon, \\ 1 & \text{if } n(K)/n^*(K) > 1 + \varepsilon. \end{cases}$$

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ≡ ≡ ≡

13 / 25

Applications – $\mathbb{R}^{\mathcal{X}}$

Let $\mathcal{B} = \mathbb{R}^{\mathcal{X}}$. As \mathcal{X} is the only set of uniqueness for \mathcal{B}_+ , we observe that

Lemma

MLE for $e(\mathbb{R}^{\mathcal{X}})$ and x_1, \dots, x_n exists if and only if $\{x_1, \dots, x_n\} = \mathcal{X}$,

i.e. each element of \mathcal{X} has to be reached by the sample
(Coupon Collector Problem).

Corollary

Let $\mathcal{B} = \mathbb{R}^{\mathcal{X}}$ and $K = |\mathcal{X}|$. Let X_1, X_2, \dots be independent random variables, each with uniform distribution on \mathcal{X} . Then, for every $c \in \mathbb{R}$,

$$\lim_{K \rightarrow \infty} \mathbb{P}(\nu_{\text{uniq}} < K \log K + Kc) = e^{-e^{-c}}.$$

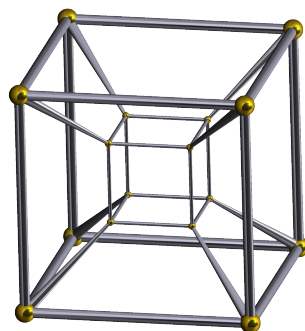
In particular, $n^*(K) = K \log K$ is a sharp threshold of the sample size for the existence of MLE for $e(\mathcal{X})$.

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ≡ ≡ ≡

14 / 25

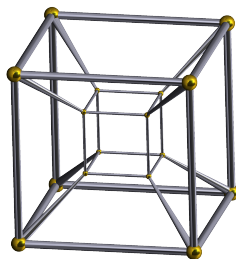
Examples – Rademacher functions

Vertices of a discrete hypercube $\{-1, +1\}^k \sim$ sequences of (-1) and $(+1)$ with length k



Examples – Discrete hypercube

Rademacher function \sim indicator of a half-cube in $\{-1, +1\}^k$



Product of Rademacher functions \sim indicator of a subcube of fixed size

Applications – Rademacher functions

For $k \in \mathbb{N}$ consider the discrete hypercube $\mathcal{X} = Q_k = \{-1, 1\}^k$.

For $j = 1, \dots, k$ we define Rademacher functions:

$$r_j(\chi) = \chi_j, \quad \chi = (\chi_1, \dots, \chi_k) \in Q_k$$

and denote $r_0(\chi) = 1$.

Here $K = |\mathcal{X}| = 2^k$.

Theorem (K. Bogdan, M. Bosy, TS (2022))

Let $\mathcal{B}^k = \text{Lin}\{r_0, r_1, \dots, r_k\}$. MLE for $e(\mathcal{B}^k)$ and $x_1, \dots, x_n \in Q_k$ exists if and only if for all $j = 1, \dots, k$ we have $\{r_j(x_1), \dots, r_j(x_n)\} = \{-1, 1\}$.

Satisfied if and only if $\{x_1, \dots, x_n\}$ intersects with every half-cube of Q_k , e.g. $\{x_1, -x_1\}$.

Applications – Rademacher functions

Theorem (K. Bogdan, M. Bosy, TS (2022))

Let $k \in \mathbb{N}$, $n(k) = \log_2 k + b + o(1)$. Let $X_1, \dots, X_{n(k)}$ be independent random variables, each with uniform distribution on Q_k . Then

$$\lim_{k \rightarrow \infty} \mathbb{P}(\{X_1, \dots, X_{n(k)}\} \text{ is of uniqueness for } \mathcal{B}_+) = \exp\{-2^{1-b}\}.$$

and $n^(K) = \log_2 k = \log_2 \log_2 K$ is a sharp threshold of the sample size for the existence of MLE for $e(\mathcal{B}^k)$ and i.i.d. uniform samples on Q_k .*

Proof idea: asymptotics of the maximum of i.i.d. geometric variables.

Applications – ERGM

We consider simple undirected graphs containing no loops or multiple edges. Let N and m denote the number of vertices and edges of the graph. By \mathcal{G}_N we denote the family of all graphs with N vertices. For graphs $G = (V, E_1)$, $H = (V, E_2)$ we let, as usual,

$$G \cup H := (V, E_1 \cup E_2), \quad G \cap H := (V, E_1 \cap E_2),$$

$$G \subset H \quad \equiv \quad E_1 \subset E_2.$$

We define $\chi_{r,s}(G) = 1 - 2\mathbb{1}_G(r, s)$ and consider the following linear space

$$\mathcal{B}^{\mathcal{G}_N} = \text{Lin} \left\{ 1, \chi_{r,s}(G) : 1 \leq r < s \leq N \right\}.$$

19 / 25

Consider coefficients $c \in \mathbb{R}^{\binom{V}{2}}$, indexed by the edges of the complete graph K_N , and the following exponential family:

$$e(\mathcal{B}^{\mathcal{G}_N}) = \left\{ p_c := e^{\phi_c - \psi(\phi_c)} : c \in \mathbb{R}^{\binom{V}{2}} \right\},$$

where

$$\phi_c(G) = \sum_{(r,s) \in \binom{V}{2}} c_{r,s} \chi_{r,s}(G), \quad \psi(\phi_c) = \log \sum_{G \in \mathcal{G}_N} e^{\phi_c(G)},$$

and $G \in \mathcal{G}_N$.

Observation

Fix $c \in \mathbb{R}^{\binom{V}{2}}$. In the random graph \mathbb{G} sampled from $p_c \in e(\mathcal{B}^{\mathcal{G}_N})$, each edge (r, s) appears independently with probability

$$p_{r,s} = \frac{e^{c_{r,s}}}{1 + e^{c_{r,s}}}.$$

20 / 25

Applications – ERGM

Theorem (K. Bogdan, M. Bosy, TS (2022))

MLE for $e(\mathcal{B}^{\mathcal{G}_N})$ and $G_1, \dots, G_n \in \mathcal{G}_N$ exists if and only if

$$\bigcup_{i=1}^n G_i = K_N \quad \text{and} \quad \bigcap_{i=1}^n G_i = \overline{K_N}.$$

Lemma (K. Bogdan, M. Bosy, TS (2022))

Let $\{G_1, \dots, G_n\}$ be independent random graphs from $p_c \in e(\mathcal{B}^{\mathcal{G}_N})$. Then the probability of the existence of MLE for $e(\mathcal{B}^{\mathcal{G}_N})$ equals

$$\prod_{1 \leq r < s \leq N} (1 - p_{r,s}^n - (1 - p_{r,s})^n).$$

In particular, $n^*(N) = \log N$ is a threshold of the sample size n for the existence of MLE for $e(\mathcal{B}^{\mathcal{G}_N})$.

21 / 25

Applications – Products of Rademacher functions

Let $k \in \mathbb{N}$, $1 \leq q \leq k$, and $\mathcal{B}_q^k = \text{Lin}\{w_S : S \subset \{1, \dots, k\} \text{ and } |S| \leq q\}$, where $w_S(x) = \prod_{i \in S} r_i(x)$, $x \in Q_k$, $S \subset \{1, \dots, k\}$, are the Walsh functions.

Observation

\mathcal{B}_q^k is the linear space spanned by indicator functions of the sub-cubes of Q_k , obtained by fixing q out of k coordinates.

- $q = 1$: Rademacher functions (already discussed)
- $q = 2$: Ising model (open problem)
- \dots : open problems
- $q = k - 1$: see the next slide

Applications – Products of $(k - 1)$ Rademacher functions

\mathcal{B}_{k-1}^k corresponds to indicators of edges of Q_k . Consider the following partition: $Q_k = \mathcal{E} \cup \mathcal{O}$:

Definition









- $\mathcal{E} := \{\chi \in Q_k : \chi \text{ has even number of positive coordinates}\}$
- $\mathcal{O} := \{\chi \in Q_k : \chi \text{ has odd number of positive coordinates}\}$

Theorem (K. Bogdan, M. Bosy, TS (2022))

MLE exists for $e(\mathcal{B}_{k-1}^k)$ and $x_1, \dots, x_n \in Q_k$ if and only if $\mathcal{E} \subset \{x_1, \dots, x_n\}$ or $\mathcal{O} \subset \{x_1, \dots, x_n\}$.

Thank you for your attention!

References

-  [Barndorff-Nielsen, O. \(1978\)](#)
"Information and exponential families in statistical theory", John Wiley & Sons Ltd., Chichester. Wiley Series in Probability and Mathematical Statistics.
-  [Bogdan, K., Bogdan, M. \(2000\)](#)
"On existence of maximum likelihood estimators in exponential families", Statistics, 34(2):137–149.
-  [Bogdan, K., Bosy, M., Skalski, T. \(2022\)](#)
"Maximum likelihood estimation for discrete exponential families and random graphs", ALEA Lat. Am. J. Probab. Math. Stat. 19 (2022), no. 1, 1045–1070.
-  [Crain, B. R. \(1976\)](#)
"Exponential models, maximum likelihood estimation, and the Haar condition", J. Amer. Statist. Assoc., 71 (355), 737–740.
-  [Eriksson, N., Fienberg, S. E., Rinaldo, A., Sullivant, S. \(2006\)](#)
"Polyhedral conditions for the nonexistence of the MLE for hierarchical log-linear models", J. Symbolic Comput., 41(2):222–233.
-  [Haberman, S. J. \(1974\)](#)
"The analysis of frequency data", The University of Chicago Press, Chicago, Ill.-London. Statistical Research Monographs, Vol. IV.
-  [Jacobsen, M. \(1989\)](#)
"Existence and unicity of MLEs in discrete exponential family distributions", Scand. J. Statist., 16 (4), 335–349.
-  [Rinaldo, A., Fienberg, S. E., Zhou, Y. \(2009\)](#)
"On the geometry of discrete exponential families with application to exponential random graph models", Electron. J. Stat., 3:446–484

Some open problems on minimum information dependence models

Tomonari Sei (The University of Tokyo)*

This is joint work with Keisuke Yano (The Institute of Statistical Mathematics).

1. Minimum information dependence model

Consider σ -finite measure spaces $(\mathcal{X}_i, \mathcal{F}_i, dx_i)$ for $i \in [d] = \{1, \dots, d\}$. Denote their product space by $(\mathcal{X}, \mathcal{F}, dx)$. Let r_i be a probability density function on \mathcal{X}_i for each i , which may have a nuisance parameter.

Suppose that we have a measurable map $h : \mathcal{X} \rightarrow \mathbb{R}^K$, which describes dependence among the d variables $(x_1, \dots, x_d) \in \mathcal{X}$. A minimum information dependence model [2] is defined by a set of probability density functions

$$p(x; \theta) = \exp \left(\theta^\top h(x) - \sum_{i=1}^d a_i(x_i; \theta) - \psi(\theta) \right) \prod_{i=1}^d r_i(x_i), \quad \theta \in \mathbb{R}^K,$$

with respect to dx , where $a_i(x_i; \theta)$ and $\psi(\theta)$ are determined by conditions

$$\int_{\mathcal{X}_{-i}} p(x; \theta) dx_{-i} = r_i(x_i), \quad i \in [d], \quad x_i \in \mathcal{X}_i,$$

and $\int_{\mathcal{X}} \sum_{i=1}^d a_i(x_i; \theta) p(x; \theta) dx = 0$. Here, $-i$ denotes the removal of the i -th coordinate. We call θ the canonical parameter, $h(x)$ the canonical statistic, $a_i(x_i; \theta)$ the adjusting function and $\psi(\theta)$ the potential function. An existence and uniqueness theorem under mild conditions is established by [2]. If the marginal distributions are uniform on the unit interval $[0, 1] \subset \mathbb{R}$, the model is called the minimum information copula model [1].

The potential function ψ is convex because it is characterized by

$$\psi(\theta) = \sup_{p \in \mathcal{M}} \left\{ \theta^\top \int h(x) p(x) dx - \int p(x) \log \frac{p(x)}{\prod_i r_i(x_i)} dx \right\},$$

where \mathcal{M} is the set of probability density functions with fixed marginals $r_i(x_i)$. The pair of the parameters θ and $\eta = \int h(x) p(x; \theta) dx$ induces a dually flat structure.

2. Open problems

There are a couple of open questions about the model. See [2] for details.

1. Is $\psi(\theta)$ analytic?
2. Find examples of $p(x; \theta)$ with closed expressions. Known examples are the Gaussian, multinomial and a class of circular distributions.
3. Construct an asymptotically efficient estimator of θ (without knowing r_i 's).

References

- [1] T. Bedford and K. J. Wilson (2014). On the construction of minimum information bivariate copula families, *Annals of the Institute of Statistical Mathematics*, **66**, 703–723.
- [2] T. Sei and K. Yano (2023). Minimum information dependence modeling, *Bernoulli*, accepted. (arXiv:2206.06792)

* e-mail: sei@mist.i.u-tokyo.ac.jp

Some open problems on minimum information dependence models

Tomonari Sei Keisuke Yano

The University of Tokyo The Institute of Statistical Mathematics

Dec 14–15, 2022 OCAMI Workshop
“Statistical Theories and Machine Learning Using Geometric Methods”

1 / 29

Introduction

- In this talk, we introduce a novel statistical model called the [minimum information dependence model](#).
- The model can describe dependence among variables separately from marginal distributions. (cf. copula theory)
- We discuss some open problems.

Reference: T. Sei and K. Yano (2023). Minimum information dependence modeling, *Bernoulli*, accepted. (arXiv:2206.06792)

2 / 29

Table of contents

1 Definition and examples

2 Properties

3 Statistical inference

3 / 29

Setup

- Consider σ -finite measure spaces $(\mathcal{X}_i, \mathcal{F}_i, dx_i)$ for $i \in [d] = \{1, \dots, d\}$.
- Denote their product as $(\mathcal{X}, \mathcal{F}, dx) = (\prod_i \mathcal{X}_i, \prod_i \mathcal{F}_i, \prod_i dx_i)$.
- How to construct a statistical model on the product space \mathcal{X} ?

Example

Fisher's iris data consists of 4 continuous and 1 categorical variables:

Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
5.1	3.5	1.4	0.2	setosa
4.9	3.0	1.4	0.2	setosa
4.7	3.2	1.3	0.2	setosa
\vdots	\vdots	\vdots	\vdots	\vdots
5.9	3.0	5.1	1.8	virginica

A natural choice is

$$\mathcal{X}_1 = \dots = \mathcal{X}_4 = \mathbb{R}, \quad \mathcal{X}_5 = \{\text{setosa}, \text{versicolor}, \text{virginica}\}.$$

equipped with the Lebesgue and counting measures, respectively.

4 / 29

Gaussian model

- Let us begin with the Euclidean case $\mathcal{X}_1 = \dots = \mathcal{X}_d = \mathbb{R}$.
- The d -dimensional Gaussian model on $\mathcal{X} = \mathbb{R}^d$ is

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right),$$

where μ and Σ are parameters.

- Properties:

- 1 The density function is written as

$$p(x; \mu, \Sigma) = \exp \left(\underbrace{-\sum_{i < j} \sigma^{ij} x_i x_j}_{\text{dependence}} \right) \underbrace{\prod_{i=1}^d A_i(x_i; \mu, \Sigma)}_{\text{independence}}$$

- 2 The marginal distributions are Gaussian.
- These properties characterize the Gaussian model. We use this idea for constructing more general statistical models.

5 / 29

The minimum information dependence model

- Now consider the general space $\mathcal{X} = \prod_{i=1}^d \mathcal{X}_i$.
- Suppose that we have
 - 1 a measurable map $h : \mathcal{X} \rightarrow \mathbb{R}^K$ and
 - 2 marginal density functions $r_i : \mathcal{X}_i \rightarrow \mathbb{R}_{>0}$ for $i \in [d]$.

Definition

A **minimum information dependence model** is defined by

$$p(x; \theta) = \exp \left(\theta^\top h(x) - \sum_{i=1}^d a_i(x_i; \theta) - \psi(\theta) \right) \prod_{i=1}^d r_i(x_i), \quad \theta \in \mathbb{R}^K,$$

where a_i and ψ are determined by conditions $\int_{\mathcal{X}_{-i}} p(x; \theta) dx_{-i} = r_i(x_i)$ and $\int_{\mathcal{X}} \sum_{i=1}^d a_i(x_i; \theta) p(x; \theta) dx = 0$.

- We call θ the canonical parameter, h the canonical statistic, a_i the adjusting functions and ψ the potential function.
- $\theta = 0$ corresponds to the independence model $\prod_{i=1}^d r_i(x_i)$.

6 / 29

Comparison

exp family	min info dep model
$p(x; \theta) = e^{\theta^\top h(x) - \psi(\theta)} p_0(x)$	$p(x; \theta) = e^{\theta^\top h(x) - \sum_i a_i(x_i; \theta) - \psi(\theta)} \prod_i r_i(x_i)$
θ	θ
h	h
ψ	ψ
—	a_i
$\int p dx = 1$	$\int p dx_{-i} = r_i$

7 / 29

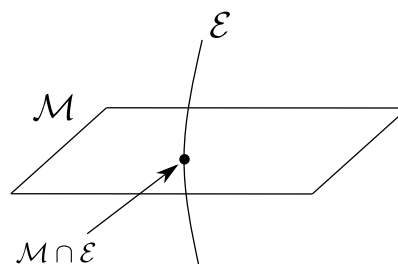
Picture

- Let \mathcal{P} be the set of probability density functions on \mathcal{X} .
- The density $p(x; \theta)$ is the unique intersection point of two manifolds

$$\mathcal{M} = \mathcal{M}(r_1, \dots, r_d) = \{p \in \mathcal{P} \mid \int p(x) dx_{-i} = r_i(x_i)\}$$

and

$$\mathcal{E} = \mathcal{E}(\theta) = \{e^{\theta^\top h(x) - \sum_i b_i(x_i)} \in \mathcal{P} \mid b_i : \mathcal{X}_i \rightarrow \mathbb{R}\}.$$



- Uniqueness follows from the Pythagorean theorem (discussed later).

8 / 29

Notation

Before giving examples, we define additional notations.

- The marginal $r_i(x_i)$ often has parameters. We denote it as $r_i(x_i; \nu)$, where ν is the (nuisance) parameter.
- The model is written as

$$p(x; \theta, \nu) = \exp \left(\theta^\top h(x) - \sum_{i=1}^d a_i(x_i; \theta, \nu) - \psi(\theta, \nu) \right) \prod_{i=1}^d r_i(x_i, \nu).$$

- We also use an abbreviation

$$p(x; \theta, \nu) = e^{\theta^\top h(x)} \prod_{i=1}^d A_i(x_i; \theta, \nu).$$

9 / 29

Example (1/3)

Multivariate Gaussian

- The d -dimensional normal distribution $N(\mu, \Sigma)$ is written as

$$p(x; \theta, \nu) = e^{\sum_{i < j} \theta_{ij} x_i x_j} \prod_{i=1}^d A_i(x_i; \theta, \nu),$$

$$\int p(x; \theta, \nu) dx_{-i} = \phi(x_i; \mu_i, \sigma_i^2), \quad i \in [d],$$

where the nuisance parameter is $\nu = (\mu_1, \dots, \mu_d, \sigma_1^2, \dots, \sigma_d^2)$.

- We obtain a bijection

$$(\mu, \Sigma) \mapsto (\theta, \nu)$$

$$\mathbb{R}^d \times \text{Sym}^+(d, \mathbb{R}) \rightarrow \mathbb{R}^{d(d-1)/2} \times \mathbb{R}^d \times \mathbb{R}_{>0}^d.$$

Lemma (Dempster 1972)

If A and B are positive definite matrices, then there exists a unique positive definite matrix C such that $C_{ii} = A_{ii}$ and $(C^{-1})_{ij} = (B^{-1})_{ij}$ ($i \neq j$).

10 / 29

Example (2/3)

Contingency table model (log-linear model)

- Consider finite spaces $\mathcal{X}_1 = [I]$ and $\mathcal{X}_2 = [J]$.
- A statistical model on $I \times J$ contingency tables is written as

$$p_{ij} = e^{\theta_{ij}} A_{1i}(\theta, \nu) A_{2j}(\theta, \nu), \quad (i, j) \in [I] \times [J],$$

$$\sum_j p_{ij} = \nu_{1i}, \quad \sum_i p_{ij} = \nu_{2j}.$$

- We obtain a bijection

$$(p_{ij}) \mapsto (\theta, \nu)$$

$$\Delta_{IJ-1} \rightarrow \mathbb{R}^{(I-1)(J-1)} \times \Delta_{(I-1)} \times \Delta_{(J-1)}$$

- Relevant studies:
 - Sinkhorn and Knopp (1967): an algorithm for finding A_{1i}, A_{2j}
 - Amari (2001): orthogonal foliation structure
 - Piantadosi et al. (2012): maximum entropy checkerboard copulas
 - Geenens (2020): similar construction of discrete families

11 / 29

Example (3/3)

Directional statistics

Circula

- Let $\mathcal{X}_1 = \mathcal{X}_2 = S^1$ (circle), identified with $[0, 2\pi)$.
- A distribution satisfying conditions

$$\int_0^{2\pi} p(x_1, x_2) dx_2 = \int_0^{2\pi} p(x_1, x_2) dx_1 = \frac{1}{2\pi}$$

is called a **circula** (Jones et al. 2015).

- The following density is circula for any function $h_0 : S^1 \rightarrow \mathbb{R}$.

$$e^{\theta h_0(x_1 - x_2) - \psi(\theta)}, \quad \psi(\theta) = \log \int_0^{2\pi} \exp(\theta h_0(z)) dz,$$

This is a minimum information dependence model with a **shift-invariant canonical statistic** $h(x_1, x_2) = h_0(x_1 - x_2)$.

- In this case, the adjusting functions are zero.

12 / 29

Existence and uniqueness theorem

- The examples so far have explicit density functions. But this is not the case in general. We need an existence and uniqueness theorem.
- Denote $H(x) = \theta^\top h(x)$. Suppose that $H \in L_1(\prod_i r_i dx)$.

Theorem 1 (SY 2023)

If there are integrable functions $b_i \in L_1(r_i dx_i)$ satisfying

$$\int e^{H(x) - \sum_i b_i(x_i)} \prod_i r_i(x_i) dx < \infty,$$

then there exist measurable $a_i : \mathcal{X}_i \rightarrow \mathbb{R}$ and $\psi \in \mathbb{R}$ such that

$$p(x) = e^{H(x) - \sum_i a_i(x_i)} \prod_i r_i(x_i) \in \mathcal{M}(r_1, \dots, r_d)$$

and $\int p \sum_i a_i dx = 0$. The function p is unique.

- The proof relies on Csiszár (1975) and Borwein et al. (1994).
- Remark: each a_i may not be integrable. This point was missed in Csiszár's paper.

13 / 29

Application of Theorem 1

3-factor interaction on \mathbb{R}^3

- Let $\mathcal{X}_1 = \mathcal{X}_2 = \mathcal{X}_3 = \mathbb{R}$ equipped with the Lebesgue measure.
- Suppose that $\int |x_i|^3 r_i(x_i) dx_i < \infty$ for each i .
- The minimum information dependence model

$$p(x_1, x_2, x_3) = e^{\theta x_1 x_2 x_3 - a_1(x_1) - a_2(x_2) - a_3(x_3) - \psi} r_1(x_1) r_2(x_2) r_3(x_3)$$

is well defined for arbitrary $\theta \in \mathbb{R}$.

- Indeed, the inequality $|x_1 x_2 x_3| \leq \frac{1}{3}(|x_1|^3 + |x_2|^3 + |x_3|^3)$ implies

$$\iiint e^{\theta x_1 x_2 x_3 - |\theta|(|x_1|^3 + |x_2|^3 + |x_3|^3)/3} r_1(x_1) r_2(x_2) r_3(x_3) dx_1 dx_2 dx_3 \leq 1 < \infty.$$

So $b_i(x_i) = |\theta| |x_i|^3 / 3$ satisfy the sufficient condition of the theorem.

In a similar manner, we can take any polynomial $h(x)$ as canonical statistics whenever the marginal distributions have all finite moments.

14 / 29

Pros and cons

Advantages of our model:

- Various type of dependence (such as conditional independence) can be incorporated.
- The domain and marginal distributions are almost arbitrary.
- Multivariate Gaussian and multinomial models are contained as particular instances.

Drawbacks:

- Computation of the adjusting and potential functions are difficult.
→ We can avoid them in estimation of θ as described later.

From practical side:

- How to choose h ? → polynomial, eigenfunctions etc.
- How to interpret θ ? → regression (not discussed today)

15 / 29

The first open problem

Problem

Find examples that have explicit density functions.

Observation

- In the circular case, we used the shift-invariant canonical statistic $h_0(x_1 - x_2)$. This may be generalized to other compact groups.

16 / 29

Table of contents

1 Definition and examples

2 Properties

3 Statistical inference

17 / 29

Pythagorean relationship

- Recall that $\mathcal{M} = \{p \in \mathcal{P} \mid \int p(x) dx_{-i} = r_i(x_i)\}$
- Let $\mathcal{E} = \mathcal{E}(\theta) = \{e^{\theta^\top h(x) - \sum_i b_i(x_i)} \prod_i r_i \in \mathcal{P} \mid \mathbf{b}_i \in L_1(r_i dx_i)\}$
- Let $D(p|q) = \int p \log(p/q) dx$ (Kullback–Leibler divergence).

Generalized Pythagorean theorem (Csiszár 1975)

If $p \in \mathcal{M}, q \in \mathcal{M} \cap \mathcal{E}, s \in \mathcal{E}$, then

$$D(p|s) = D(p|q) + D(q|s)$$

Proof:

$$\begin{aligned} D(p|s) - D(p|q) - D(q|s) &= \int (p - q) \log(q/s) dx \\ &= \int (p - q) \left(- \sum_i b_i^q + \sum_i b_i^s \right) dx \\ &= \sum_i \int (r_i - r_i) \left(- \sum_i b_i^q + \sum_i b_i^s \right) dx_i = 0. \end{aligned}$$

18 / 29

Consequence of Pythagorean

- For $p \in \mathcal{M}, q \in \mathcal{M} \cap \mathcal{E}, s \in \mathcal{E}$,

$$D(p|s) = D(p|q) + D(q|s).$$

In particular,

$$D(p|s) \geq D(q|s).$$

The equality is attained iff $p = q$.

- By substituting $q = e^{\theta^\top h - \sum_i a_i - \psi} \prod_i r_i$ and $s = e^{\theta^\top h - \sum_i b_i} \prod_i r_i$, we obtain

$$\int p \left(\log p - \theta^\top h + \sum_i b_i - \sum_i \log r_i \right) \geq \int q \left(- \sum_i a_i - \psi + \sum_i b_i \right).$$

- Since $\int p b_i = \int q b_i$ and $\int q \sum_i a_i = 0$, we have

$$\psi \geq \theta^\top \int p h - \int p \log \frac{p}{\prod_i r_i}$$

19 / 29

Potential function

- Now we obtained the following characterization of ψ :

$$\psi(\theta) = \sup_{p \in \mathcal{M}} \left\{ \theta^\top \int h p - \int p \log \frac{p}{\prod_i r_i} \right\},$$

which implies ψ is convex.

Open problem 2

Are $a_i(x_i; \theta)$ and $\psi(\theta)$ analytic with respect to θ ?

- Observation: For exponential families $e^{\theta^\top h(x) - \psi(\theta)} p_0(x)$, the potential function ψ is analytic on $\text{int}(\text{dom}(\psi))$. This is directly proved by Taylor expansion

$$\int e^{(\theta+\delta)^\top h(x)} p_0(x) dx = \sum_k \frac{\delta^k}{k!} \int h(x)^k e^{\theta^\top h(x)} p_0(x) dx.$$

20 / 29

Derivatives of the potential

- The model

$$\log p(x|\theta) = \theta^\top h(x) - \sum_i a_i(x_i; \theta) - \psi(\theta)$$

- Suppose that a_i and ψ are smooth, and integrals and derivatives are exchangeable.
- The derivative of ψ is the expectation of h :

$$\partial_\alpha \psi(\theta) = \int h_\alpha(x) p(x) dx$$

- The Hessian matrix is equal to the Fisher information

$$\begin{aligned} \partial_\alpha \partial_\beta \psi(\theta) = \int \{ & h_\alpha(x) - \sum_i \partial_\alpha a_i(x_i; \theta) - \partial_\alpha \psi(\theta) \} \\ & \cdot \{ h_\beta(x) - \sum_i \partial_\beta a_i(x_i; \theta) - \partial_\beta \psi(\theta) \} p(x) dx, \end{aligned}$$

which is **not equal to** the covariance of h .

- The quantity $\sum_i \partial_\alpha a_i + \partial_\alpha \psi$ is the orthogonal projection of $h_\alpha(x)$ to the space of additive functions $\{\sum_i b_i(x_i)\}$.

21 / 29

Comparison

exp. family	"mindemo"
θ	θ
h	h
ψ	ψ
—	a_i
$\nabla \psi = E(h)$	$\nabla \psi = E(h)$
$\nabla \nabla^\top \psi = \text{Cov}(h)$	$\nabla \nabla^\top \psi \neq \text{Cov}(h)$

22 / 29

Table of contents

1 Definition and examples

2 Properties

3 Statistical inference

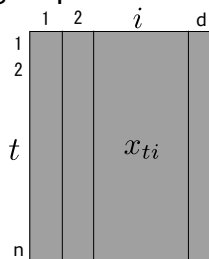
23 / 29

Conditional inference

- Consider a random sample $x(t) = (x_i(t))_{i=1}^d$ for $1 \leq t \leq n$.
- Decompose the data into “marginal” and rank statistics:

$$M = (\{x_{1i}, \dots, x_{ni}\})_{1 \leq i \leq d}, \quad \pi = (\pi_1, \dots, \pi_d) \in S_n^d$$

where S_n is the permutation group.



- Decomposition of the likelihood

$$\prod_{t=1}^n p(x(t); \theta, \nu) = f(\pi|M; \theta) \times g(M; \theta, \nu)$$

- We use the conditional likelihood $f(\pi|M; \theta)$ for inference of θ .
- M must have almost no information about θ ! (ancillary statistic)

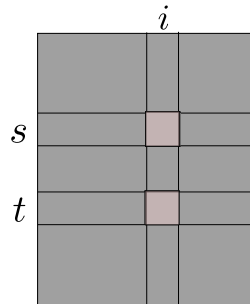
24 / 29

Parameter estimation

- The conditional likelihood is an exponential family on S_n^d :

$$f(\pi|M;\theta) = \frac{e^{\sum_{t=1}^n \theta^\top h((M \circ \pi)(t))}}{\sum_{\bar{\pi} \in S_n^d} e^{\sum_{t=1}^n \theta^\top h((M \circ \bar{\pi})(t))}}, \quad \pi \in S_n^d.$$

- The maximizer of $f(\pi|M;\theta)$ is called the conditional maximum likelihood estimator.
- We can sample π by MCMC. The implementation is not difficult.



- MCMC + conditional MLE (Geyer and Thompson 1992).
- Other option: Besag's pseudo likelihood estimator (Mukherjee 2016).

25 / 29

Consistency of conditional MLE

- The conditional MLE has consistency ($\hat{\theta} \rightarrow \theta$ as $n \rightarrow \infty$).
- Assumptions
 - ① The domain \mathcal{X} satisfies some metric entropy condition. The function h is Lipschitz. The parameter space Θ is bounded.
 - ② a_i is smooth in θ .

Theorem 1 (SY 2023)

Under Assumption 1, the likelihood ratio is close to the conditional likelihood ratio. (a rough statement)

Corollary (SY 2023)

Under Assumptions 1 and 2, we have $\hat{\theta} \rightarrow \theta$ in probability.

Open problem 3

Is the conditional MLE asymptotically normal and efficient?

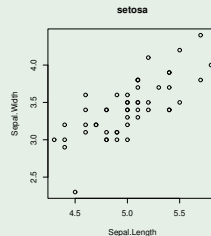
- For finite-set cases, Haberman (1977) proved this is correct.

26 / 29

Conditional inference works

Example (Fisher's iris data)

- Consider sepal length and sepal width of Setosa. $\mathcal{X}_1 = \mathcal{X}_2 = \mathbb{R}$.



- We compared the MLE of Gaussian model and the conditional MLE of the proposed model with $h(x_1, x_2) = x_1 x_2$. Only the difference is whether the marginal model is specified or not.
- Result: $\hat{\theta}_{\text{MLE}} = 1.66$ and $\hat{\theta}_{\text{CLE}} = 1.72$ (standard errors are 0.4).
- Almost the same!

Other applications: penguin data and earthquake data. See SY (2023).

27 / 29

Conclusion

- We proposed the minimum information dependence model specified by canonical statistic h and the marginal densities r_i .
- Estimation is performed via conditional likelihood and permutation.
- The model is applicable to any domain of data.

Open problems:

- Further examples
- Properties of adjusting and potential functions
- Proof of asymptotic efficiency of the conditional MLE (or other estimators)

Thank you for your attention.

28 / 29

References I

- Amari, S. (2001). Information geometry on hierarchy of probability distributions, *IEEE Trans. Info. Theo.*, **47** (5), 1701–1711.
- Arnold, R. and Jupp, P. (2013). Statistics of orthogonal axial frames, *Biometrika*, **100** (3), 571–586.
- Bedford, T. (2002). Interactive expert assignment of minimally-informative copulae, In: Lindqvist B, editor. Proceedings of mathematical methods in reliability.
- Bedford, T. and Wilson, K. J. (2014). On the construction of minimum information bivariate copula families, *Ann. Inst. Stat. Math.*, **66**, 703–723.
- Bedford, T., Daneshkhah, A. and Wilson, K. J. (2016). Approximate uncertainty modeling in risk analysis with vine copulas, *Risk Analysis*, **36** (4), 792–815.
- Borwein, J. M., Lewis, A. S. and Nussbaum, R. D. (1994). Entropy minimization, DAD problems, and doubly stochastic kernels, *J. Funct. Anal.*, **123**, 264–307.
- Chen, Y. and Sei, T. (2022). A proper scoring rule for minimum information copulas, 2022, preprint. arxiv:2204.03118.
- Choi, L., Blume, J. D. and Dupont, W. D. (2015). Elucidating the foundations of statistical inference with 2×2 tables, *PLoS ONE*, **10** (4), e0121263, 1–22.
- Csiszár, I. (1975). I-divergence geometry of probability distributions and minimization problems, *Ann. Probab.*, **3** (1), 146–158.
- Dempster, A. P. (1972). Covariance selection, *Biometrics*, **28**, 157–175.

27 / 29

References II

- Eckstein, S. and Nutz, M. (2021). Quantitative stability of regularized optimal transport and convergence of Sinkhorn's algorithm, arXiv:2110.06798.
- Geenens, G. (2020). Copula modeling for discrete random vectors, *Dependence Modeling*, **8** (1), 417–440.
- Geyer, C. J. and Thompson, E. A. (1992). Constrained Monte Carlo maximum likelihood for dependent data, *J. Roy. Statist. Soc. Ser. B*, **54** (3), 657–683.
- Haberman, S. J. (1977), *The Analysis of Frequency Data*, Univ. of Chicago Press.
- Hayakawa, J. and Takemura, A. (2016). Estimation of exponential-polynomial distribution by holonomic gradient descent, *Comm. Statist. Theo. Meth*, **45** (23), 6860–6882.
- Jansen, M. J. W. (1997). Maximum entropy distributions with prescribed marginals and normal score correlations, In: Beneš, V. and Štěpán, J. (ed.), *Distributions with given Marginals and Moment Problems*, 87–92, Springer.
- Jaynes, E. T. (1957). Information theory and statistical mechanics, *Phys. Rev.*, **106** (4), 620–630.
- Jones, M.C., Pewsey, A. and Kato S. (2015). On a class of circulas: copulas for circular distributions, *Ann. Inst. Stat. Math*, **67**, 843–862.
- Meeuwissen, A. M. H. and Bedford, T. (1997). Minimally informative distributions with given rank correlation for use in uncertainty analysis, *J. Statist. Comput. Simul.*, **57**, 143–174.

28 / 29

References III

- Mukherjee, S. (2016). Estimation in exponential families on permutations, *Ann. Statist.*, **44** (2), 853–875.
- Peyré, G. and Cuturi, M. (2019). Computational Optimal Transport: With Applications to Data Science, *Foundations and Trends in Machine Learning*, **11** (5-6), 355–607. (Preprint arxiv:1803.00567)
- Piantadosi, J., Howlett, P. and Borwein, J. (2012). Copulas with maximum entropy, *Optim. Lett.*, **6**, 99–125.
- Sei, T. and Yano, K. (2022). Minimum information dependence modeling, preprint, arxiv:2206.06792.
- Sinkhorn, R. and Knopp, P. (1967). Concerning nonnegative matrices and doubly stochastic matrices, *Pacific J. Math.*, **21** (2), 343–348.
- Takayama, N., Kuriki, S., Takemura, A. (2018) A-hypergeometric distributions and Newton polytopes, *Advances in Applied Mathematics*, **99**, 109–133.
- Yang, E., Ravikumar, P., Allen, G. I., and Liu, Z. (2015). Graphical models via univariate exponential family distributions. *J. Mach. Learn. Res.*, **16** (1), 3813–3847.
- Zhang, M. and Bedford, T. (2018). Vine copula approximation: a generic method for coping with conditional dependence, *Stat. Comput.*, **28**, 219–237.
- Zhu, Y. and Reid, N. (1994). Information, ancillary, and sufficiency in the presence of nuisance parameters, *Canad. J. Stat.*, **22**, 111–123.

Stein identity, Poincaré inequality and exponential integrability on a metric measure space

Tomonari Sei (The University of Tokyo)*¹

Ushio Tanaka (Osaka Metropolitan University)*²

Stein identity due to Stein [3] characterises a Gaussian random variable. The identity further implies Poincaré inequality in terms of the Gaussian random variable:

Theorem (Chernoff [1]). Let g be absolutely continuous function and X a random variable following standard normal distribution such that $g(X)$ has finite variance. Then it follows that

$$\text{Var}(g(X)) \leq \mathbb{E}[g'(X)^2]$$

with equality if and only if $g(X)$ is linear.

We showed a discrete Stein identity based on the idea of Sei [2] and that it implies discrete Poincaré inequality. We study them from the point of view of geometric analysis; in particular, we restrict our attention to implications with regard to discrete Poincaré inequality to show exponential integrability of a function on a discrete metric measure space.

References

- [1] H. Chernoff, A note on an inequality involving the normal distribution, *Ann. Probab.* **9** (1981), 533–535.
- [2] T. Sei, Coordinate-wise transformation of probability distributions to achieve a Stein-type identity, *Inf. Geom.* **5** (2022), 325–354.
- [3] C.M. Stein, Estimation of the mean of a multivariate normal distribution, *Ann. Statist.* **9** (1981), 1135–1151.

*¹ e-mail: tomonari-sei@g.ecc.u-tokyo.ac.jp

*² e-mail: utanaka@omu.ac.jp

Infinite Dimensional Parameterized Measure Models

Hikaru Watanabe

The University of Tokyo

Shunichi Amari pointed out that generalization of finite dimensional Information Geometry to infinite dimensional Information Geometry is an important problem in his book “Methods of Information Geometry” (1993) [1]. After that, in 2015, Nihat Ay, Jürgen Jost, Hông Vân Lê and Lorenz Schwachhöfer proposed the notion of (infinite dimensional) parameterized measure models and statistical models, which is a natural generalization of finite dimensional statistical manifolds ([2]). Moreover, they proved that Exponential statistical manifold, which is built by Giovanni Pistone and Carlo Sempì ([4]), can be regarded as an example of infinite dimensional statistical models.

In this talk, we prove that Exponential manifold by reproducing kernel Hilbert spaces, which is constructed by Kenji Fukumizu ([3]), is also an example of infinite dimensional statistical models. Furthermore, we replace reproducing kernel Hilbert spaces with reproducing kernel Banach spaces (see [5] for reproducing kernel Banach spaces). The aim of this replacement is to show that existing finite dimensional open exponential family is obtained when a certain reproducing kernel Banach space is set.

Moreover, if a parameterized measure model or a statistical model satisfies the condition called n -integrability, we can obtain covariant n -tensor. Especially, covariant 2-tensor is Fisher metric if it is positive definite, and covariant 3-tensor is Amari-Chentsov tensor. However we must be careful about tensor fields on infinite dimensional manifolds because there are two kinds of definition of tensor fields on infinite dimensional manifolds, unlike tensor fields on finite dimensional manifolds. Nihat Ay, Jürgen Jost, Hông Vân Lê Lorenz Schwachhöfer introduced covariant n -tensor as C^0 “weak tensor fields”. However in this talk we show that one can introduce covariant n -tensors as C^∞ “strong tensor fields” when parameterized measure model or statistical model is C^∞ .

References

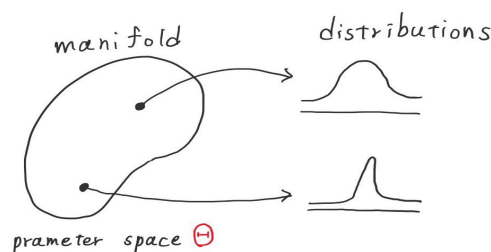
- [1] Amari, S and Nagaoka, H.: Methods of Information Geometry . Iwanami (1993)
- [2] Ay, N., Jost, J., Lê, H.V., Schwachhöfer, L.: Information Geometry, Ergebnisse der Mathematik und ihrer Grenzgebiete. Springer, Berlin (2017)
- [3] Fukumizu, K.: Exponential Manifold by Reproducing Kernel Hilbert Spaces, pp. 291–306. Cambridge University Press, Cambridge (2010).
<https://doi.org/10.1017/CBO9780511642401.019>
- [4] Pistone, G., Sempì, C.: An infinite-dimensional geometric structure on the space of all the probability measures equivalent to a given one. Ann. Stat. 23(5), 1543–1561 (1995)
- [5] Lin, R., Zhang, H., Zhang, J.: On reproducing kernel Banach spaces: Generic definitions and unified framework of constructions, (2019). arXiv:1901.01002.

previous works
questions
result

previous works about finite dimensional information geometry
previous works about infinite dimensional information geometry

Rao(1945)

Rao: Regarding statistical model as manifold
Fisher metric: Riemannian metric



Navigation icons: back, forward, search, etc.

Hikaru Watanabe

Infinite dimensional parameterized measure model

previous works
questions
result

previous works about finite dimensional information geometry
previous works about infinite dimensional information geometry

exponential family and Fisher metric

Example 1 (exponential family)

Ω : m'ble set

$\mu \in P(\Omega)$

$\Theta \subset \mathbb{R}^n$: parameter space

f_1, \dots, f_n : m'ble function on Ω

Some assumptions

Exponential family is a family of distributions

$$P := \left\{ \frac{\exp(\sum_{k=1}^n \theta_k f_k)}{\int_{\Omega} \exp(\sum_{k=1}^n \theta_k f_k) d\mu} \mu \mid \theta \in \Theta \right\}$$

Example 2

examples of exponential family:

Normal distributions, Bernoulli distributions, etc.

Proposition 3 (Fisher metric)

$$g_{\theta} \left(\frac{\partial}{\partial \theta_i}, \frac{\partial}{\partial \theta_j} \right) = E_{\theta}[(f_i - E_{\theta}[f_i])(f_j - E_{\theta}[f_j])]$$

\hookrightarrow variance form

Navigation icons: back, forward, search, etc.

Hikaru Watanabe

Infinite dimensional parameterized measure model

- previous works
- questions
- result

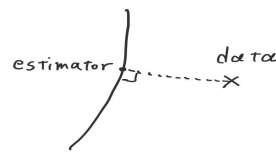
- previous works about finite dimensional information geometry
- previous works about infinite dimensional information geometry

Efron(1975)

Efron: connecting geometry and statistics

statistical curvature: geometrical

statistical estimation by statistical curvature



◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

Hikaru Watanabe

Infinite dimensional parameterized measure model

- previous works
- questions
- result

- previous works about finite dimensional information geometry
- previous works about infinite dimensional information geometry

Amari(1980)

Amari: formalizing connecting geometry and statistics

 α, e, m -connection

statistical estimation by e, m -connection

dually flat structure

AC-tensor

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ 🔍 ↺

Hikaru Watanabe

Infinite dimensional parameterized measure model

- previous works about finite dimensional information geometry
- previous works about infinite dimensional information geometry

Infinite dimensional parameterized measure model

previous works questions result	previous works about finite dimensional information geometry previous works about infinite dimensional information geometry
problems based on the pointout	

problems based on the pointout

- 1 How to formalize infinite dimensional manifold?
- 2 How to regard infinite dimensional distributions as infinite dimensional manifold?
- 3 How to define geometrical notions such as Fisher metric and AC-tensor on infinite dimensional information geometry?

Hikaru Watanabe	Infinite dimensional parameterized measure model
-----------------	--

previous works questions result	previous works about finite dimensional information geometry previous works about infinite dimensional information geometry
Lang(1962)	

(1)How to formalize infinite dimensional manifold?
 \hookrightarrow Lang: Banach manifold
 Banach manifold: manifold modeled on Banach space

Hikaru Watanabe	Infinite dimensional parameterized measure model
-----------------	--

previous works questions result	previous works about finite dimensional information geometry previous works about infinite dimensional information geometry
---------------------------------------	--

three studies

(2) How to regard infinite dimensional distributions as infinite dimensional manifolds?

↔ There are three studies.

- (2-1) Pistone and Sempi(1995):
exponential statistical manifold by Orlicz space
- (2-2) Fukumizu(2010):
exponential manifold by reproducing kernel Hilbert space
- (2-3) Ay, Jost, Lê and Schwachhöfer(2017):
parameterized measure model and statistical model

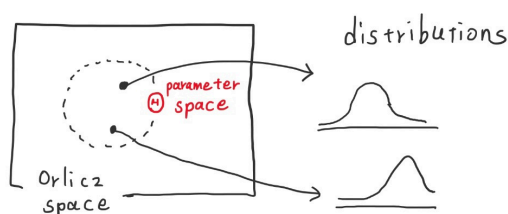
Hikaru Watanabe	Infinite dimensional parameterized measure model
-----------------	--

previous works questions result	previous works about finite dimensional information geometry previous works about infinite dimensional information geometry
---------------------------------------	--

Pistone and Sempi(1995)

(2) How to regard infinite dimensional distributions as infinite dimensional manifolds?

↔ Pistone and Sempi: exponential statistical manifold by Orlicz space



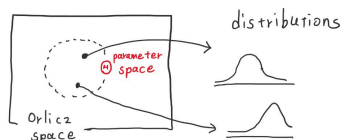
Hikaru Watanabe	Infinite dimensional parameterized measure model
-----------------	--

previous works
questions
result

previous works about finite dimensional information geometry
previous works about infinite dimensional information geometry

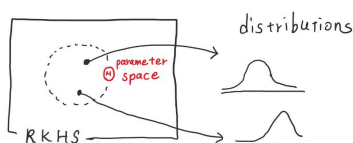
examples of PMM and SM

Exponential statistical model by Pistone and Sempi



is a SM.

Exponential manifold by Fukumizu



is a SM.

Navigation icons: back, forward, search, etc.

Hikaru Watanabe

Infinite dimensional parameterized measure model

previous works
questions
result

previous works about finite dimensional information geometry
previous works about infinite dimensional information geometry

Ay, Jost, Lê and Schwachhöfer(2017)

(3) How to define geometrical notions such as Fisher metric and AC-tensor on infinite dimensional information geometry?

↔ Ay, Jost, Lê and Schwachhöfer: covariant n -tensor

covariant n -tensor: a **generalization** of Fisher metric and AC-tensor.

covariant 2-tensor: Fisher metric in PMM and SM

covariant 3-tensor: AC-tensor in PMM and SM

Navigation icons: back, forward, search, etc.

Hikaru Watanabe

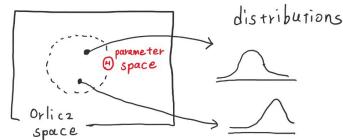
Infinite dimensional parameterized measure model

previous works
questions
result

previous works about finite dimensional information geometry
previous works about infinite dimensional information geometry

Ay, Jost, Lê and Schwachhöfer(2017)

In the case of the PMM before normalized to exponential statistical model



Proposition 11 (Ay, Jost, Lê and Schwachhöfer(2017))

τ^n : covariant n -tensor

$$\tau_f^n(v_1, \dots, v_n) = E_f[v_1 \cdots v_n]$$

where $f \in \Theta, v_1, \dots, v_n \in$ (this Orlicz space).

Especially covariant 2, 3-tensor is as follows.

$$\tau_f^2(v_1, v_2) = E_f[v_1 v_2]$$

$$\tau_f^3(v_1, v_2, v_3) = E_f[v_1 v_2 v_3]$$

Navigation icons: back, forward, search, etc.

Hikaru Watanabe

Infinite dimensional parameterized measure model

previous works
questions
result

questions
other good examples
covariant 2, 3-tensor

questions

We want to understand

- I Are there any other good examples of parameterized measure models and statistical models?
- II Give explicit formula of covariant n -tensor on normalized statistical models.

Navigation icons: back, forward, search, etc.

Hikaru Watanabe

Infinite dimensional parameterized measure model

I Are there any other good examples of parameterized measure models and statistical models?
(I think) we want more examples.

In finite dimensional exponential family

$$g_{\theta}(\frac{\partial}{\partial \theta_i}, \frac{\partial}{\partial \theta_j}) = E_{\theta}[(f_i - E_{\theta}[f_i])(f_j - E_{\theta}[f_j])]$$

$$\Gamma_f(\frac{\partial}{\partial \theta_i}, \frac{\partial}{\partial \theta_j}, \frac{\partial}{\partial \theta_k}) = E_\theta[(f_i - E_\theta[f_i])(f_j - E_\theta[f_j])(f_k - E_\theta[f_k])]$$

In the PMM before normalized to exponential statistical manifold

$$\tau_f^2(v_1, v_2) = E_f[v_1 v_2]$$

$$\tau_f^3(v_1, v_2, v_3) = E_f[v_1 v_2 v_3]$$

previous works
questions
result

other good examples
covariant n -tensor

other good examples

I Are there any other good examples of parameterized measure models and statistical models?

↔ We can construct exponential manifolds by reproducing kernel Banach space.

reproducing kernel Banach space(RKBS): Banach spaces of functions, a **generalization** of RKHS to Banach space

Hikaru Watanabe

Infinite dimensional parameterized measure model

previous works
questions
result

other good examples
covariant n -tensor

RKBS

Definition 12

Ω : set
 B : Banach space of functions on Ω .
 $B \rightarrow \mathbb{R}, f \mapsto f(x)$ is continuous.

Is RKBS really generalization of RKHS?

Definition 13

Ω : set
 $(H, \langle \cdot, \cdot \rangle)$: Hilbert space of functions on Ω
 $\exists K: \Omega \times \Omega \rightarrow \mathbb{R}$: positive definite (kernel function)

- 1 $K(x, \cdot) \in H \quad (x \in \Omega)$
- 2 $f(x) = \langle f, K(x, \cdot) \rangle \quad (x \in \Omega, f \in H)$

 $\Leftrightarrow H \rightarrow \mathbb{R}, f \mapsto f(x)$ is continuous.

Hikaru Watanabe

Infinite dimensional parameterized measure model

What is the difficulty of $RKHS \rightarrow RKBS$?
 \hookrightarrow non existing of kernel function K

$$\sqrt{K(x, x)} = \|e_x\|_{B^*}$$

We can obtain exponential manifold by RKBS when we replace RKHS with RKBS and $\sqrt{K(x, x)}$ with $\|e_x\|_{B^*}$ in the discussion of Fukumizu.

What is a benefit of RKBS against RKHS?
 \hookrightarrow examples of RKBS not RKHS Ω : top sp

$$C_b(\Omega) := \{f: \Omega \rightarrow \mathbb{R} \mid f \text{ is continuous and bounded.}\}$$

$$\|f\|_{C_b(\Omega)} := \sup|f(x)|$$

$$C_0(\Omega) := \{f: \Omega \rightarrow \mathbb{R} \mid f \text{ is continuous and vanishing at } \infty\}$$

$$\|f\|_{C_0(\Omega)} := \sup |f(x)|$$

Hikaru Watanabe

II Give explicit formula of covariant n -tensor on normalized statistical models.

- X : Banach manifold
- $p: X \rightarrow P(\Omega)$
- (X, p) : a parameterized measure model
- (Y, q) : the normalization of (X, p) (statistical model)
- Under some assumption

where $y \in Y, \mu := q(y), f_i \mu := dp_y(v_i), \tau^n$: covariant n -tensor on (Y, q) .

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ 🔍 ↺






 τ^n : covariant n -tensor on exponential statistical manifold

Especially covariant 2, 3-tensor are

$$\tau_f^3[v_1, v_2, v_3] = E_f[(v_1 - E_f[v_1])(v_2 - E_f[v_2])(v_3 - E_f[v_3])]$$

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ 🔍 ↺

reference

-  Amari, S and Nagaoka, H.: Methods of Information Geometry . Iwanami (1993)
-  Ay, N., Jost, J., Lê, H.V., Schwachhöfer, L.: Information Geometry, Ergebnisse der Mathematikund ihrer Grenzgebiete. Springer, Berlin (2017)
-  Fukumizu, K.: Exponential Manifold by Reproducing Kernel Hilbert Spaces, pp. 291–306. Cambridge University Press, Cambridge (2010).
<https://doi.org/10.1017/CBO9780511642401.019>
-  Pistone, G., Sempi, C.: An infinite-dimensional geometric structure on the space of all the probability measures equivalent to a given one. Ann. Stat. 23(5), 1543–1561 (1995)
-  Lin, R., Zhang, H., Zhang, J.: On reproducing kernel Banach spaces: Generic definitions and unified framework of constructions, (2019). arXiv:1901.01002.

Neural-Kernel Conditional Mean Embeddings

Eiki Shimizu (SOKENDAI)

With a positive definite kernel, conditional distributions can be embedded into an associated reproducing kernel Hilbert space (RKHS). Such approaches are referred to as Kernel Conditional Mean Embeddings (CMEs), and have been applied to various applications such as causal inference and kernelized Bayes rule.

Although several appealing theoretical properties have been shown, CMEs are yet to be suitable for modern ML tasks: the empirical estimates involve inversion of a Gram matrix, and the hyperparameter selection is not straightforward. To address these issues, we propose a method that combines neural networks (NNs) with CMEs.

Our approach is simple, and can be interpreted as a NN trained with a RKHS loss. This allows us to replace the matrix inversion with a NN model, while taking advantage of NN's ability to learn useful features. We also provide a strategy to efficiently optimize the hyperparameter of the kernel, without relying on median heuristics or cross validation. We demonstrate the effectiveness of our method with ML related tasks, where the estimation of conditional distribution plays an important role.

Neural-Kernel Conditional Mean Embeddings

Eiki Shimizu, SOKENDAI

Intro

- PhD student working with Prof. Kenji Fukumizu
- Interested in : Kernel methods, Bayesian Inference
and their applications to Deep Learning models
- This presentation is about
Kernel method + DeepNN = useful conditional density estimator

Outline

1. Brief review of kernel methods and kernel mean embeddings
2. Propose a new conditional density estimator,
Share experimental results
3. Applications to more complicated ML tasks: Reinforcement Learning

Section 1

RKHS and Kernel Mean Embeddings

RKHS and Notations

RKHS

A symmetric function $k_{\mathcal{X}}: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a reproducing kernel of $\mathcal{H}_{\mathcal{X}}$ if and only if

- $\forall x \in \mathcal{X}, k_{\mathcal{X}}(x, \cdot) \in \mathcal{H}_{\mathcal{X}}$
- $\forall x \in \mathcal{X} \text{ and } \forall f \in \mathcal{H}_{\mathcal{X}}, f(x) = \langle f, k_{\mathcal{X}}(x, \cdot) \rangle_{\mathcal{H}_{\mathcal{X}}}$

A space $\mathcal{H}_{\mathcal{X}}$ is called reproducing kernel Hilbert space (RKHS)

Notations

Let (X, Y) be a random variable on $\mathcal{X} \times \mathcal{Y}$ with distribution P and density function $p(x, y)$, and $k_{\mathcal{X}}$ and $k_{\mathcal{Y}}$ be positive definite kernel corresponding to $\mathcal{H}_{\mathcal{X}}$ and $\mathcal{H}_{\mathcal{Y}}$ respectively.

We denote feature maps as $\psi(x) = k_{\mathcal{X}}(x, \cdot)$ and $\phi(y) = k_{\mathcal{Y}}(y, \cdot)$

Kernel Mean Embeddings

$$m_{P(X)} = \mathbb{E}_P[\psi(X)] \in \mathcal{H}_{\mathcal{X}}, \quad \langle f, m_{P(X)} \rangle_{\mathcal{H}_{\mathcal{X}}} = \mathbb{E}_P[f(X)]$$

The embedding uniquely defines the probability distribution (the mapping is injective) if the kernel $k_{\mathcal{X}}$ is *characteristic*. Popular kernels like Gaussian kernel has this property.

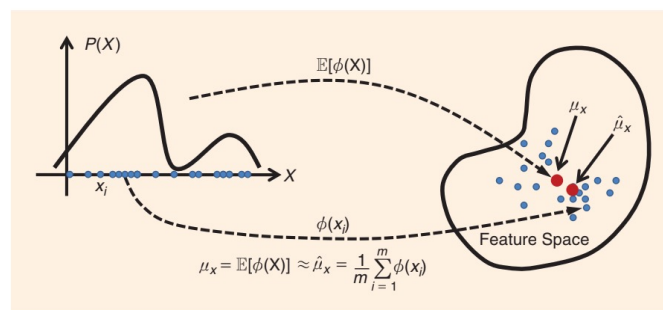


Figure taken from [1]

Kernel Covariance Operators

$$C_{XX} = E[\psi(X) \otimes \psi(X)], C_{XY} = E[\psi(X) \otimes \phi(Y)]$$

Generalizes finite-dimensional covariance matrices to the case of infinite feature spaces. Always exist for bounded kernels.

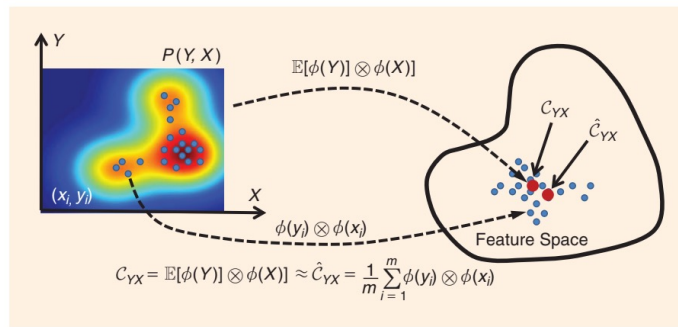


Figure taken from [1]

Kernel Conditional Mean Embeddings (CME)

$$m_{P(Y|X)}(x) = E_P[\phi(Y)|X = x] \in \mathcal{H}_Y$$

Under condition $E_P[g(Y)|X = x] \in \mathcal{H}_Y$ for all $g \in \mathcal{H}_Y$, there exists an operator $C_{Y|X}$ such that $m_{P(Y|X)}(x) = C_{Y|X}\psi(x)$.

This is a minimizer of the RKHS loss l

$$l(C_{Y|X}) = E_P \left[\|\phi(Y) - C_{Y|X}\psi(X)\|_{\mathcal{H}_Y}^2 \right]$$

The closed form solution is

$$C_{Y|X} = C_{YX}(C_{XX})^{-1}$$

Empirical Estimate

$$\hat{l}(C_{Y|X}) = \frac{1}{n} \sum_{i=1}^n \|\phi(y_i) - C_{Y|X} \psi(x_i)\|_{\mathcal{H}_Y}^2 + \lambda \|C_{Y|X}\|_{HS}$$

The solution is:

$$\hat{C}_{Y|X} = \Phi(\mathbf{K}_X + \lambda I)^{-1} \Psi^\top$$

where, $\Phi = (\phi(y_1), \dots, \phi(y_n))$, $\Psi = (\psi(x_1), \dots, \psi(x_n))$, $\mathbf{K}_X = \Psi^\top \Psi$

Thus,

$$\hat{m}_{P(Y|X)}(x) = \sum_{i=1}^n \beta_i(x) \phi(y_i) = \Phi \beta(x)$$

where, $\beta(x) = (\mathbf{K}_X + \lambda I)^{-1} \mathbf{k}_X$

Interpretations

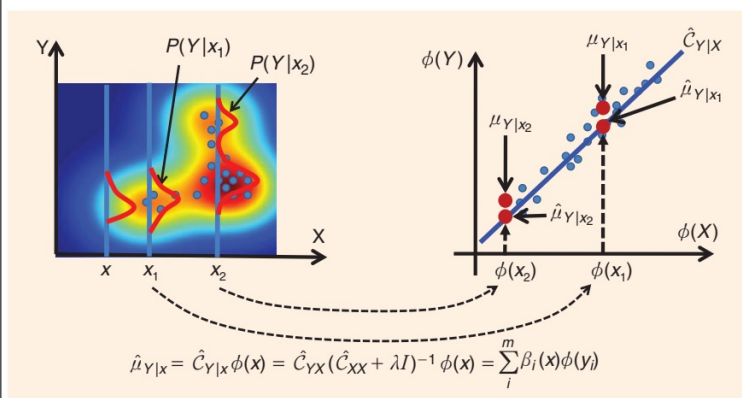


Figure taken from [1]

- Regression in the function space
- If k_y is a linear kernel, this is just a Kernel Ridge Regression

- Weighted particle view

$$\hat{m}_{P(Y|X)}(x) = \sum_{i=1}^n \beta_i(x) \phi(y_i)$$

$\beta_i(x)$ weights particles $\phi(y_i)$

Note that the weights does not necessarily have to be positive nor sum up to one

Kernel herding

Can we “sample” from the KMEs?

$$\hat{m}_P = \sum_{i=1}^N w_i k(x_i, \cdot) \rightarrow \hat{x}_1, \dots, \hat{x}_m \text{ (samples)}$$

Kernel herding is a deterministic sampling approach used to obtain super samples \hat{x}_m . Selects \tilde{x} greedily and iterates the following:

$$\hat{x}_{l+1} = \operatorname{argmin}_{\tilde{x}} \left\| \frac{1}{l+1} \left\{ \sum_{j=1}^l k(\hat{x}_j, \cdot) + k(\tilde{x}, \cdot) \right\} - \hat{m}_P \right\|_{\mathcal{H}_k}^2$$

Practical limitations

- The computational cost of $(\mathbf{K}_X + \lambda I)^{-1}$ is $O(n^3)$, and does not scale to large dataset
- RKHS features are pre-specified feature maps. This may lead to poor performance when input variables are high-dimensional, or possess highly non-linear structure
- Hyperparameter selection for k_X , k_Y and λ is not straightforward. While the choice significantly affects the performance, particularly for k_Y , standard procedures like cross-validation can not be applied

Section 2

Proposal: Integrate CME with DNNs

Proposal: Big idea

CME Recap:

$$\hat{m}_{P(Y|X)}(x) = \sum_{i=1}^n \beta_i(x) \phi(y_i) = \boldsymbol{\Phi} \boldsymbol{\beta}(x)$$

where, $\boldsymbol{\beta}(x) = (\mathbf{K}_X + \lambda I)^{-1} \mathbf{k}_X$

Why don't we just replace $\boldsymbol{\beta}(x)$ with a DNN model?

$$\hat{m}_{P(Y|X)}(x) = \sum_{a=1}^M \phi(\eta_a) f_a(x; \theta)$$

where, $f(x; \theta): \mathcal{X} \rightarrow \mathbb{R}^M$ is DNN, and $\eta \in \mathcal{Y}$ corresponds to M atoms/particles

Objective function (for CDE)

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \left\| \phi(y_i) - \sum_{a=1}^M \phi(\eta_a) f_a(x_i; \theta) \right\|_{\mathcal{H}_y}^2$$

$$\Leftrightarrow \min_{\theta} \frac{1}{n} \sum_i \left\{ k_y(y_i, y_i) - 2 \sum_a k_y(y_i, \eta_a) f_a(x_i; \theta) + \sum_{a,b} k_y(\eta_a, \eta_b) f_a(x_i, \theta) f_b(x_i, \theta) \right\}$$

Simply use this loss function, and everything else (e.g. implementation, training procedures) is the same as the standard DNN!

“What’s there to be happy about? Job’s not finished.” Kobe Bryant

Positives (solves first two issues):

- No matrix inversion!
- Can use mini-batch optimization for efficient training
- DNN models learn features: $k_{\mathcal{X}}$ implicitly tuned!

Negatives:

- We still have k_y left to be tuned, and this is not easy
- The objective function is defined in terms of the RKHS norm of \mathcal{H}_y . If we change the kernel parameter, the definition of the objective function also changes

What should we do with k_y ?

Empirically, we find the following strategy to work quite well

1. Use a positive definite kernel that also has density interpretation

We use the Gaussian-Density Kernel

$$k_\sigma(y, y') = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\|y - y'\|^2}{2\sigma^2}\right)$$

2. Use the following objective function (will just call it the RKHS loss):

$$\min_{\theta, \sigma} \frac{1}{n} \sum_i \left\{ -2 \sum_a k_\sigma(y_i, \eta_a) f_a(x_i; \theta) + \sum_{a,b} k_\sigma(\eta_a, \eta_b) f_a(x_i, \theta) f_b(x_i, \theta) \right\}$$

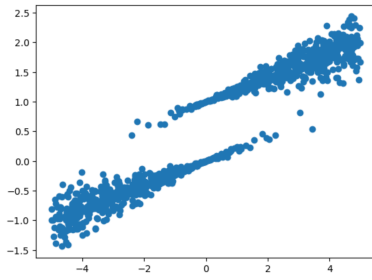
Experiments (Only to be shared during the presentation)

- Note that this work is still in progress
- Experiments on 1-dimensional conditional density estimation
(output variables are 1-dimensional, but input variables can be multi-dimensional)
- Show preliminary results on:
 1. Toy data simulations
 2. More realistic setting with the UCI dataset

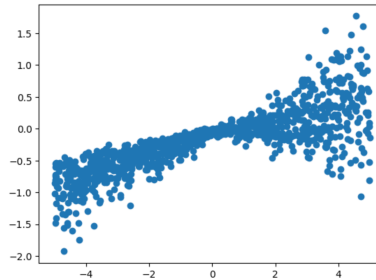
Toy data settings

- 3 settings

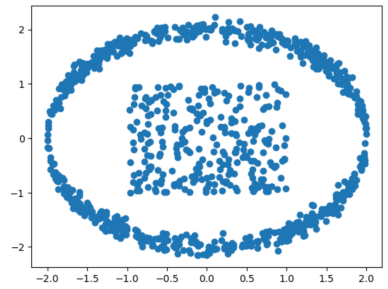
Bimodal (MoG with hetero noise)



Skewed normal



Ring & Box



Evaluation

1. Train each model with 5000 training data
2. Generate enough samples from $\hat{p}(y|X = x)$
3. Calculate the Wasserstein distance between samples and the true conditional distribution for each x , and average them in the end

Wasserstein distance:

$$l_1(u, v) = \inf_{\pi \in \Gamma(u, v)} \int_{\mathbb{R} \times \mathbb{R}} |x - y| d\pi(x, y)$$

where, $\Gamma(u, v)$ is the set of distribution whose marginals are u and v

Competitors Part1

Deep Feature approach (DF)

Learn adaptive feature $\psi_\theta(x)$ represented by DNN: $\mathcal{X} \rightarrow \mathbb{R}^d$

$$\beta(x) = \Psi_\theta^\top (\Psi_\theta \Psi_\theta^\top + \lambda I)^{-1} \psi_\theta(x)$$

Several successful applications such as in causal inference (IV regression) [3]

- Positives

Explicitly learn feature, The matrix inverse can be done with $O(d^3)$, Compatible with mini-batch optimization

- Negatives

Still need to tune λ and the kernel hyperparameter of k_y

For k_y , typically rely on the median heuristic

$$\sigma = \text{median}\{\|X_i - X_j\| \mid i, j = 1, \dots, n\}$$

Competitors Part2

Mixture Density Networks (MDN)

DNN learns the weight, the mean and the variance of a GMM

$$\hat{p}(y|x) = \sum_{k=1}^K w_k(x; \theta) \mathcal{N}(y | \mu_k(x, \theta), \sigma_k^2(x; \theta))$$

Several successful applications such as in likelihood-free inference [4]

- Positives

Easy to implement, sample and evaluate the log-likelihood

- Negatives

Optimizing 3 parameters of the GMM jointly may lead to numerical instability and over-fitting

UCI datasets

- “Real world” dataset

Dataset	Boston	Concrete	Energy	Kin8nm	Naval	Power	Protein
(N, P)	(506, 13)	(1030, 8)	(768, 8)	(8192, 8)	(11, 934, 16)	(9568, 4)	(45, 730, 9)

- Both input and output variables are normalized before training

Evaluation

Quantile Interval Coverage Error (QICE) [5]:

1. Generate enough samples from $\hat{p}(y|X = x)$
2. Divide them into $M=10$ bins and get 10 quantile intervals with the boundary $\hat{y}_n^{\text{low}_m}$ and $\hat{y}_n^{\text{high}_m}$
3. Calculate

$$\text{QICE} = \frac{1}{M} \sum_{m=1}^M \left| r_m - \frac{1}{M} \right|, \text{ where } r_m = \frac{1}{N} \sum_{n=1}^N \mathbf{1}_{y_n \geq \hat{y}_n^{\text{low}_m}} \cdot \mathbf{1}_{y_n \geq \hat{y}_n^{\text{high}_m}}$$

In the optimal scenario, about 10% of true data shall fall into each of the 10 quartile intervals, and QICE reaches 0

New competitor: The Big Boss

Diffusion model

Used in “Generative AI”, super strong performance on image generation tasks

[5] proposed conditional version of this model, enabling flexible conditional density modelling

Though this model demonstrates SOTA level performance on several tasks, it is computationally costly (requires two NNs/optimizations), and the sampling may be slow

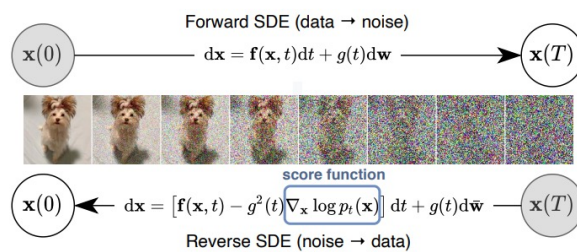


Figure 1: **Solving a reverse-time SDE yields a score-based generative model.** Transforming data to a simple noise distribution can be accomplished with a continuous-time SDE. This SDE can be reversed if we know the score of the distribution at each intermediate time step, $\nabla_x \log p_t(x)$.

Figure taken from [6]

Section 3

Application to Distributional RL

Reinforcement Learning framework

1. An agent on S_t interacts with the environment (takes an action A_t)
2. Moves to the next state S_{t+1} , and gets a reward R_{t+1}
3. Repeat until the agent learns a good policy

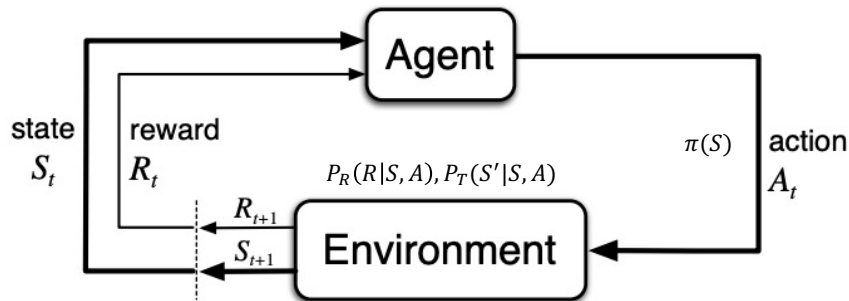


Figure taken from [7]

Q-learning: basic idea

Learn state-action value function Q^π of a policy π , which corresponds to expected discounted sum of rewards

$$Q^\pi(s, a) = E_{s,r} \left[\sum_t \gamma^t r_t \right], a_t = \pi(s_t)$$

where, $\gamma \in [0,1)$ is the discount factor.

This satisfies the Bellman equation/backup:

$$Q^\pi(s, a) \leftarrow E_{s',r} [r + \gamma Q^\pi(s', \pi(s'))]$$

Similarly, define Bellman operator \mathcal{T}^π :

$$\mathcal{T}^\pi Q^\pi(s, a) = E_{s',r} [r + \gamma Q^\pi(s', \pi(s'))]$$

Q-learning: algorithm

Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
 Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$
 Loop for each episode:
 Initialize S
 Loop for each step of episode:
 Choose A from S using policy derived from Q (e.g., ε -greedy)
 Take action A , observe R, S'
 $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$
 $S \leftarrow S'$
 until S is terminal

Figure taken from [7]

Notes

- $Q(s, a)$ can be approximated well by DNN
- In that case, we stop the gradient for $\max_a Q(s, a)$

Distributional RL (DRL)

Model distribution over state-action value instead of just expectation

$$Z^\pi(s, a) \stackrel{\text{D}}{=} \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)$$

Distributional Bellman equation [8]:

$$Z^\pi(s, a) \leftarrow R(s, a) + \gamma P^\pi Z^\pi(s, a)$$

Distributional Bellman operator \mathcal{T}^π :

$$\mathcal{T}^\pi Z^\pi(s, a) \stackrel{\text{D}}{=} R(s, a) + \gamma P^\pi Z^\pi(s, a)$$

Categorical DQN (CDQN)

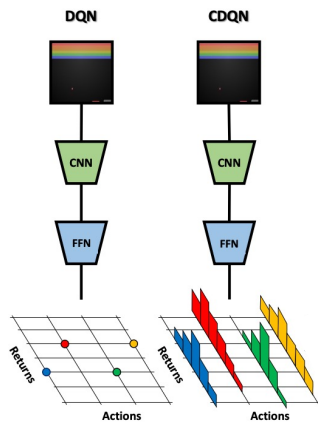


Figure taken from [9]

CDQN [8] basically models the distribution with histogram, and use cross-entropy loss for the loss function

A famous approach named "Rainbow" uses this CDQN structure and achieved SOTA performance on Atari games few years ago

Technically, there needs to be a heuristic step to enable cross-entropy loss to be used

Can we

- Represent the distribution better, and
- Compare two distributions in a more principled way?

Proposal: MCMD DRL

Maximum Conditional Mean Discrepancy (MCMD)[10]:

$$\widehat{\text{MCMD}}^2(P_{Y|X}, P_{Y'|X'}) = \left\| \hat{m}_{P(Y|X)}(x) - \hat{m}_{P(Y'|X')}(x) \right\|_{\mathcal{H}_Y}^2$$

Empirical estimate with our proposed model:

$$\sum_{a,b} k_Y(\eta_a, \eta_b) f_a(x_i, \theta) f_b(x_i, \theta) - 2 \sum_{a,b} k_Y(\eta_a, \eta'_b) f_a(x_i, \theta) f'_b(x_i, \theta) + \sum_{a,b} k_Y(\eta'_a, \eta'_b) f'_a(x_i, \theta) f'_b(x_i, \theta)$$

To make this Distributional Bellman equation, simply calculate

$$\widehat{\text{MCMD}}(P_{Y|X}, \mathcal{T}^\pi P_{Y|X})$$

Why DRL is the perfect application

In Q-learning framework,

- Evaluate the learned state-action value every step
- This is done for millions and billions of steps
- The accumulated “data”/experience can be as large as steps taken

In this case, we never want to see matrix inversion. Our approach can be applied to DQN-style learning in a straight manner, and offers principled way to represent and compare distributions.

Experimental setups/results only shared during the presentation

- Test on “Classic Control” provided by Gymnasium environment
- Though “classic”, SOTA DeepRL approaches can easily fail in these environments, with bad modelling or hyperparameters
- Run for 5,000,000 steps, and evaluate agents every 100 steps

Summary

- New conditional density estimator by combining KME and DNN
- Overcomes some of the KME limitations:
matrix inverse, hyperparameter tunings
- Good performances on density estimation tasks
- Promising performances on RL tasks, address bandwidth selection issues at the same time

Future work: make it Bayesian

- Large scale settings / multi-dimensional density estimation
- Application to kernelized Bayes Rule
- Combine with Gaussian Process, Bayesian Deep Learning?
- Can we do the Bayesian model selection for hyperparameter tuning?
- Could some “Geometric methods” incorporated into our work/KMEs

References

- [1] Song et.al. Kernel Embeddings of Conditional Distributions: A Unified Kernel Framework for Nonparametric Inference in Graphical Methods. *IEEE Signal Proceeding Magazine*, 2013
- [2] Hsu and Ramos. Bayesian Learning of Conditional Kernel Mean Embeddings for Automatic Likelihood-Free Inference. *Proceeding of the 22nd International Conference on Artificial Intelligence and Statistics*. 2019
- [3] Xu et.al. Learning Deep Features in Instrumental Variable Regression. *International Conference on Learning Representations*. 2020

References

- [4] Papamakorios and Murray. Fast ϵ -free inference of Simulation Models with Bayesian Conditional Density Estimation. *Advances in Neural Information Processing Systems*. 2016
- [5] Han et.al. CARD: Classification and Regression Diffusion Models. *Advances in Neural Information Processing Systems*. 2022
- [6] Song et.al. Score-Based Generative Modelling through Stochastic Differential Equations. *International Conference on Learning Representations*. 2021
- [7] Sutton and Barto. Reinforcement Learning: An Introduction
MIT Press. 2018

References

- [8] Bellemare et.al. A Distributional Perspective on Reinforcement Learning. *Proceedings of 34th International Conference on Machine Learning*. 2017
- [9] Theate et.al. Distributional Reinforcement Learning with Unconstrained Monotonic Neural Networks. *Neurocomputing*. 2023
- [10] Park and Muandet. A Measure-Theoretic Approach to Kernel Conditional Mean Embeddings. *Advances in Neural Information Processing Systems*. 2020
- [11] Biggs et.al. MMD-FUSE: Learning and Combining Kernels for Two-Sample Testing Without Data Splitting. *Advances in Neural Information Processing Systems*. 2023

Bonferroni method and tube method for heavy-tailed distributions

Satoshi Kuriki (Inst. Statist. Math.), Evgeny Spodarev (Ulm Univ.)

The Bonferroni method is the simplest method for approximating the maximum distribution of the statistics T_1, \dots, T_n :

$$\mathbb{P}\left(\max_i T_i > c\right) \lesssim \sum_i \mathbb{P}(T_i > c).$$

When the joint distribution of $(T_i)_{1 \leq i \leq n}$ is Gaussian, the relative approximation error

$$\Delta(c) = \frac{\sum_i \mathbb{P}(T_i > c) - \mathbb{P}(\max_i T_i > c)}{\sum_i \mathbb{P}(T_i > c)}$$

is exponentially small as $c \rightarrow \infty$. This is shown in the tube method that includes the Bonferroni method as a special case. However, when the statistics are studentized by a common standard deviation estimator $\hat{\sigma}$, the distribution of $(T_i)_{1 \leq i \leq n}$ becomes a heavy-tailed distribution such as the multivariate t -distribution. In this talk, we evaluate the relative error $\Delta(c)$ in such cases.

We first set the class of random vectors $(T_i)_{1 \leq i \leq n}$ with a correlation structure (ρ_{ij}) and a specified tail behavior. Let $x_1, \dots, x_n \in \mathbb{R}^n$ be unit vectors such that $\langle x_i, x_j \rangle = \rho_{ij}$. By using standard Gaussian random vector $\xi \sim \mathcal{N}_n(0, I_n)$ and $\hat{\sigma}$ independent of ξ , we let $T_i = \langle x_i, \xi \rangle / \hat{\sigma}$. The marginal distribution function of T_i^2 is denoted by F (identical for all i), and the tail distribution of F is parameterized with parameters (β, ℓ) ($\beta \leq 1$, ℓ is a slowly varying function):

$$1 - F(x) \sim C \exp\left(-\int_{x_0}^x q(t) dt\right), \quad q(t) = \frac{\ell(t)}{t^\beta},$$

where we assume that the limit $\lim_{x \rightarrow \infty} \ell(x) = \gamma \in (0, \infty) \cup \{\infty\}$ exists. This distribution family includes the light-tailed distribution (exponential, super-exponential), the long-tailed distribution $\text{RV}_{-\gamma}$ (regularly varying distribution with index $-\gamma$), and intermediate cases (subexponential distribution) \mathcal{S} .

β	$\beta = 0$	$\beta = 0$	$\beta \in (0, 1)$	$\beta = 1$	$\beta = 1$
γ	$\gamma = \infty$	$\gamma < \infty$	$\gamma \leq \infty$	$\gamma = \infty$	$\gamma < \infty$
	super-exponential	exponential	\mathcal{S}	\mathcal{S}	$\text{RV}_{-\gamma}$

Theorem 1. Suppose that $\beta < 1$ or $\gamma = \infty$. Then,

$$\log \Delta(c) \sim -c^{2(1-\beta)} \ell(c^2) g_\beta(\cos^2 \theta_{\text{cri}}), \quad c \rightarrow \infty,$$

where $g_\beta(y) = \frac{y^{\beta-1}-1}{1-\beta}$ ($\beta < 1$), $-\log y$ ($\beta = 1$), and $\theta_{\text{cri}} = \frac{1}{2} \min_{i < j} \cos^{-1} \rho_{ij}$ is the critical radius (reach). Moreover, $\Delta(c) \rightarrow 0$ ($c \rightarrow \infty$).

Example 1. If F is the chi-square distribution ($\beta = 0$, $\gamma = 1/2$), or the log-normal distribution ($\beta = 1$, $\ell(x) = \log x$, $\gamma = \infty$), then $\log \Delta(c) \sim -(1/2)c^2 \tan^2 \theta_{\text{cri}}$, or $-\log(c^2)(-\log \cos^2 \theta_{\text{cri}})$, respectively.

Theorem 2 (Regularly varying distributions). Suppose that $\beta = 1$ and $\gamma < \infty$. By using independent random variables $\tilde{B} \sim B_{\gamma+\frac{1}{2}, \frac{n-1}{2}}$ and $V_i \sim \text{Unif}(\{x \in \mathbb{R}^n \mid \|x\| = 1, \langle x, x_i \rangle = 0\})$, we have

$$\lim_{c \rightarrow \infty} \Delta(c) = \frac{1}{n} \sum_{i=1}^n \mathbb{P}\left(\tilde{B} < \cos^2 \theta_i(V_i)\right), \quad \theta_i(v) = \tan^{-1} \min_{j \neq i} \frac{1 - \rho_{ij}}{\langle x_j, v \rangle}.$$

Moreover, $0 \leq \lim_{c \rightarrow \infty} \Delta(c) \leq \bar{\Delta} := \mathbb{P}(\tilde{B} < \cos^2 \theta_{\text{cri}})$.

These evaluations for $\Delta(c)$ can be generalized to the tube formula.

Statistical Theories and Machine Learning Using Geometric Methods

Date : December 14-15, 2023 (Japan Standard Time)

Venue : Academic Extension Center (Osaka Metropolitan University)

Contents : Workshop (Hybrid: physical/virtual)

- This workshop is supported by Osaka Central Advanced Mathematical Institute (MEXT Promotion of Distinctive Joint Research Center Program JPMXP0723833165), Osaka Metropolitan University.

Organizers: Koichi Tojo (RIKEN AIP: koichi.tojo@riken.jp), Hideto Nakashima (ISM), Yoshihiko Konno (OMU), Hideyuki Ishi (OMU), Kenji Fukumizu (ISM)

Program

- December 14 (Thursday):

13:00–13:50 **Hiroto Inoue** (Nishinippon Institute of Technology)

Mean-variance joint statistic valued in a real Siegel domain

14:00–14:50 **Eren Mehmet Kiral** (Keio University)

Bayesian Learning with Lie Groups

15:00–15:50 **Hajime Fujita** (Japan Women's University)

The generalized Pythagorean theorem on the compactifications of certain dually flat spaces via toric geometry

16:10–17:00 **Atsumi Ohara** (University of Fukui)

Doubly autoparallel structure and curvature integrals: An application to iteration complexity analysis of convex optimization

17:10–18:00 **Adam Chojecki** (Warsaw University of Technology), **Hideyuki Ishi** (Osaka Metropolitan University)

Uncovering Data Symmetries: Estimating Covariance Matrix in High-Dimensional Setting With 'gips' R Package

18:10–19:00 **Tomasz Skalski** (Wroclaw University of Science and Technology)

Maximum likelihood estimation for discrete exponential families, its geometry and combinatorics

- December 15 (Friday):

10:00–10:50 **Tomonari Sei** (The University of Tokyo)

Some open problems on minimum information dependence models

11:00–11:50 **Tomonari Sei** (The University of Tokyo), **Ushio Tanaka** (Osaka Metropolitan University)

Stein identity, Poincaré inequality and exponential integrability on a metric measure space

13:50–14:40 **Hikaru Watanabe** (The University of Tokyo)

Infinite dimensional parameterized measure models

14:50–15:40 **Eiki Shimizu** (SOKENDAI)

Neural-Kernel Conditional Mean Embeddings

15:50–16:40 **Satoshi Kuriki** (The Institute of Statistical Mathematics)

Bonferroni method and tube method for heavy-tailed distributions