

Sampling From A Manifold *

Persi Diaconis[§], Susan Holmes[†] and Mehrdad Shahshahani[‡]

Susan Holmes and Persi Diaconis
Department of Statistics
Sequoia Hall
CA 94305 Stanford, USA.
e-mail: susan@stat.stanford.edu

Mehrdad Shahshahani
Mathematics Institute
Teheran, Iran.
e-mail: mshahshahani@gmail.com

Abstract: We develop algorithms for sampling from a probability distribution on a sub-manifold embedded in \mathbb{R}^n . Applications are given to the evaluation of algorithms in ‘Topological Statistics’; to goodness of fit tests in exponential families and to Neyman’s smooth test. This article is partially expository, giving an introduction to the tools of geometric measure theory.

AMS 2000 subject classifications: Primary 60K35, 60K35; secondary 60K35.

Keywords and phrases: manifold, conditional distribution, geometric measure theory, sampling.

1. Introduction

A variety of inferential tasks require drawing samples from a probability distribution on a manifold. This occurs in sampling from the posterior distribution on constrained parameter spaces (eg covariance matrices), in testing goodness of fit for exponential families conditional on sufficient statistics (eg the sum and product of the observations in a Gamma family), and in generating data to test algorithms in topological statistics.

In our applications, we found that examples involved domains with corners and non smooth functions (eg $\max(|x_1|, |x_2|, \dots, |x_n|)$). We found a useful set of tools in geometric measure theory. One of our goals is to explain and illustrate this in the usual language of probability and statistics.

To introduce the subject, consider the following two examples, used as illustrations throughout.

Example 1A: The Curved Torus Figure 2 shows a picture of 1000 points on the torus

$$\mathcal{M} = \{[(R + r \cos(\theta)) \cos(\psi), (R + r \cos(\theta)) \sin(\psi), r \sin(\theta)]\}, \quad (1.1)$$

*This work was part of a project funded by the French ANR under a Chaire d’Excellence at the University of Nice Sophia-Antipolis.

[†]Supported by a DARPA grant HR 0011-04-1-0025.

[‡]Supported by the NIH grant NIH-R01GM086884.

[§]Supported by NSF grant DMS 0804324

$0 \leq \theta, \psi < 2\pi$ for $R > r > 0$. The torus is formed by taking a circle of radius r in the (x, z) plane, centered at $x = r, z = 0$ and rotating it around the z axis.

Formula (1.1) gives the embedding of \mathcal{M} as a compact 2-dimensional manifold in \mathbb{R}^3 . As such, \mathcal{M} inherits a natural area measure: roughly, take a region on \mathcal{M} , thicken it out by ϵ to be fully 3-dimensional, compute the usual volume of the thickened region and take the limit of this area divided by ϵ as $\epsilon \rightarrow 0$. This area measure can be normalized to be a probability measure $\tilde{\mathcal{H}}^2(dx)$ on \mathcal{M} . The points shown are sampled from $\tilde{\mathcal{H}}^2(dx)$.

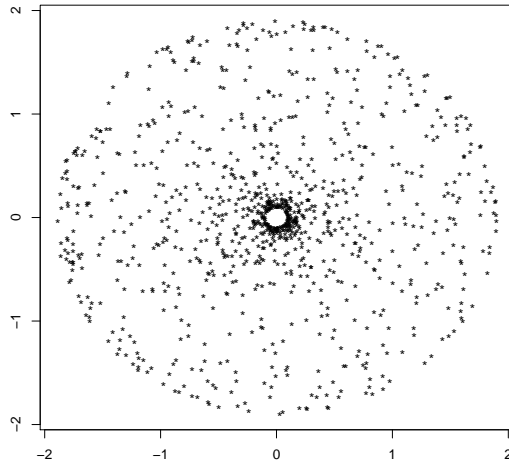


Figure 1: A sample of 1000 points from the naïve measure on a torus with $R=1, r=0.9$

Note that the sampled points are denser in regions with higher curvature such as the inside of the torus. This distribution is from the naïve choice: choose θ and ψ uniformly and map onto \mathcal{M} using (1.1). Figure (2.3) show both correctly and incorrectly generated points, see next section.

Such samples, with noise added, are used to calibrate topological algorithms for estimating dimension, number of components and homology in the emerging field of topological statistics. Examples such as two linked tori on the seven sphere and Klein bottles are shown to come up naturally in image analysis (Carlsson, Carlsson and de Silva, 2006).

Example 1B: Testing the Gamma Distribution For fixed $n \geq 3, S, P > 0$, let

$$\mathcal{M} = \left\{ (x_1, \dots, x_n); \quad x_i > 0, \quad \sum_{i=1}^n x_i = S, \quad \prod_{i=1}^n x_i = P \right\}. \quad (1.2)$$

This is a compact $(n - 2)$ -dimensional submanifold in \mathbb{R}^n . The need for samples from \mathcal{M} comes up in testing if random variables X_1, X_2, \dots, X_n are independently drawn from the Gamma density

$$\frac{e^{-x/\sigma} x^{a-1}}{\sigma^a \Gamma(a)} \quad 0 < x < \infty, \quad (1.3)$$

with $\sigma, a > 0$ unknown parameters. The sufficient statistics for σ, a are $S = \sum_{i=1}^n X_i, P = \prod_{i=1}^n X_i$. In numerous writings, R. A. Fisher suggested using the conditional distribution of X_1, \dots, X_n given S, P to give exact goodness of fit tests. These ideas are reviewed in section 3 below. The conditional distribution has a simple density with respect to $\bar{\mathcal{H}}^{n-2}(dx)$ leading to practical algorithms for random generation and testing. The proposed tests are different than the ones in Kallioras, Koutrouvelis and Canavos (2006) or Pettitt (1978). Goldman and Whelan (2000) and Yang (2006) explain interesting applications of these tests in modern evolutionary analyses of DNA.

Related Literature

There has been a steady interest in statistics on manifolds. The development of mean and variance estimators appears in Pennec (2006) and Bhattacharya and Patrangenaru (2003). The book by Bhattacharya and Bhattacharya (2012) about data on the shape space manifold contains several interesting results. Data on the sphere and the projective space are discussed in Beran (1979), Fisher, Lewis and Embleton (1993) and Watson (1983). Data on more general manifolds appear in Giné (1975). One widespread example occurs in physics and chemistry problems involving configurations of atoms with some inter-atomic distances or angles fixed; see Fixman (1974) or Ciccotti and Ryckaert (1986). Any of these settings give rise to the need for Monte Carlo sampling on manifolds.

There are well-known algorithms for sampling from the uniform distribution on compact groups and other homogeneous spaces. For instance, Eaton (1983) proves that if an $n \times n$ matrix is filled with iid standard normals and the QR decomposition is carried out, then the Q part is distributed as the uniform distribution on the orthogonal group (Haar measure). Mezzadri (2007); Diaconis and Shahshahani (1986) develop this. There are also elegant algorithms for sampling from the boundary of compact, convex sets in \mathbb{R}^n (Bélisle, Romeijn and Smith, 1993; Boender, Caron, McDonald, Kan, Romeijn, Smith, Telgen and Vorst, 1991). A different procedure, the Lalley and Robbins (1987) “princess-and monster” algorithm has been studied for sampling from the boundaries of more general sets (Comets, Popov, Schütz and Vachkovskaia, 2009; Narayanan and Niyogi, 2008). These algorithms are based on moving within the interior of the bounded set reflecting off the boundary. They are different from the present procedures and may be very effective when applicable. We do not know previous literature on sampling from more general manifolds.

Of course, conditional probability densities are standard fare, even with very general conditioning. However, explicit description of area measure and the use of the co-area formula is not so common. We only know of the foundational monograph by Tjur (1974). This contains a good history. The development is both more and less general. Tjur works with Riemannian manifolds and smooth functions. We work with embedded manifolds but allow Lipschitz functions such as max/min. Tjur gives a self-contained development based on Radon measures. We are able to use more classical foundations from standard sources. Tjur’s valuable monograph was written before the computer revolution. We emphasize techniques useful for sampling.

This paper studies the following problem of sampling from \mathcal{M} , an m -dimensional submanifold in \mathbb{R}^n . Consider $f(x) \geq 0$ such that $\int_{\mathcal{M}} f(x) \bar{\mathcal{H}}^m(dx) < \infty$ with $\bar{\mathcal{H}}^m(dx)$ the m -dimensional area measure on \mathcal{M} . Samples are to be drawn from the normalized version of f . Section 2 gives basic definitions for submanifolds, area measure, Jacobians and the co-area formula. These

notions are illustrated on examples 1A,1B.

Section 3 develops the theory for exponential families, Section 4 that of Neyman's smooth test.

The algorithms presented are reasonably standard Markov chain Monte Carlo methods supplemented by some geometrical tricks and the tools of geometric measure theory. We hope they will be useful to researchers who face similar problems.

The subject developed here may be considered as a continuous analog of algebraic statistics as initiated in Diaconis and Sturmfels (1998) and reviewed in Drton, Sturmfels and Sullivan (2009). That theory began by developing algorithms for sampling from the conditional distribution of discrete exponential families given their sufficient statistics. There, finding ways of moving around on the space of data sets with given sufficient statistics leaned on tools from computational algebra (Gröbner bases). Here, the same task is studied using direct geometric analysis and tools such as the curve selection lemma.

2. Definitions and Tools

The classical subject of calculus on manifolds has an enormous expository literature. We have found the elementary treatment of Hubbard and Hubbard (2007) readable and useful. In our applications, pieces of manifolds with corners occur naturally. For example, testing the three-parameter Gamma density gives rise to

$$\mathcal{M} = \left\{ (x_1, \dots, x_n); \ x_i > 0, \sum_{i=1}^n x_i = S, \prod_{i=1}^n x_i = P, \min x_i \geq m \right\}.$$

Geometric measure theory provides convenient tools. We use Federer (1996), denoted [F], as a principle reference. The introductory account by Morgan (2009) gives a useful taste of the subject matter. Recent references are Mattila (1999), Krantz and Parks (2008).

2.1. First Definitions

A function $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is *Lipschitz* if $|f(x) - f(y)| \leq c|x - y|$ for some finite, positive c . Euclidean distance is used for $|\cdot|$ on both sides. A set in \mathbb{R}^n is m -rectifiable [F, p. 251] if it is the Lipschitz image of a bounded subset in \mathbb{R}^m . This is a very rich class of sets, discussed at length in the references above. All of the sets that arise in our applications are rectifiable.

Use $\lambda^n(dx)$ for Lebesgue measure on the Lebesgue measurable sets of \mathbb{R}^n . Given any subset $A \subseteq \mathbb{R}^n$, define the m -dimensional *Hausdorff measure* $\bar{\mathcal{H}}^m(A)$ by

$$\bar{\mathcal{H}}^m(A) = \lim_{\delta \rightarrow 0} \inf_{\substack{A \subseteq \cup S_i \\ \text{diam}(S_i) \leq \delta}} \sum \alpha_m \left(\frac{\text{diam}(S_i)}{2} \right)^m$$

The infimum is taken over all countable coverings S_i of A with $\text{diam}(S_i) = \sup\{|x - y| : x, y \in S_i\}$ and $\alpha_m = \Gamma(\frac{1}{2})^m / \Gamma[(\frac{m}{2}) + 1]$, the volume of the unit ball in \mathbb{R}^m . Hausdorff measure is an outer measure which is countably additive on the Borel sets of \mathbb{R}^n . It serves as area measure for subsets. If the set A above is m -rectifiable, the coverings above can be restricted to balls or cubes [F,

Sect. 3.2.26]. For a closed set A , [F, Sect. 3.2.39] shows $\bar{\mathcal{H}}^m(A) = \lim_{\epsilon \rightarrow 0} \lambda^n \{x : \text{dist}(x, A) < \epsilon\} / \alpha_{(n-m)} \epsilon^{n-m}$, thus justifying the heuristic definition of area measure in Example A of Section 1.

To actually compute area measure, the Jacobian is an essential tool. Call $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ differentiable at $x \in \mathbb{R}^m$ if there exists a linear map $L : \mathbb{R}^m \rightarrow \mathbb{R}^n$ with

$$\lim_{h \rightarrow 0} |f(x+h) - f(x) - L(h)|/|h| = 0.$$

The linear map L is denoted $Df(x)$ when it exists. A celebrated theorem of Rademacher [F, Sect. 3.1.6] says that a Lipschitz function is differentiable at λ^m a.e. $x \in \mathbb{R}^m$. For a differentiable function, Df can be computed using partial derivatives $D_i(x) = \lim_{h \rightarrow 0} (f(x_1, \dots, x_i + h, \dots, x_m) - f(x))/h$. As usual, the derivative matrix is

$$(Df(x))_{i,j} = D_i f_j(x) \quad 1 \leq i \leq m, 1 \leq j \leq n$$

If $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is differentiable at x , the k -dimensional Jacobian $J_k f(x)$ may be defined as the norm of the derivative matrix [F, page 241]. Geometrically $J_k f(x)$ is defined as the maximum k -dimensional volume of the image under $Df(x)$ of a unit k -dimensional cube in \mathbb{R}^m (the maximum over all possible rotations of the cube under orthogonal rotations in \mathcal{O}_m (Morgan, 2009, p. 25)). As usual, if $\text{rank } Df(x) < k$, $J_k f(x) = 0$. If $\text{rank } Df(x) = k$, then $(J_k f(x))^2$ equals the sum of squares of the determinants of the $k \times k$ submatrices of $Df(x)$. Usually, $k = m$ or n . Then $(J_k f(x))^2$ equals the determinant of the $k \times k$ product of $Df(x)$ and its transpose. If $k = m = n$, $J_k f(x)$ is the absolute value of the determinant of $Df(x)$.

2.2. The Area Formula

The basic area formula [F, Sect. 3.2.5] is a useful extension of the classical change of variables formula of calculus.

Theorem: Area Formula *If $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is Lipschitz and $m \leq n$, then*

$$\int_A g(f(x)) J_m f(x) \lambda^m(dx) = \int_{\mathbb{R}^n} g(y) N(f|A, y) \bar{\mathcal{H}}^m(dy) \quad (2.1)$$

whenever A is λ^m measurable, $g : \mathbb{R}^n \rightarrow \mathbb{R}$ is Borel, and $N(f|A, y) = \#\{x \in A : f(x) = y\}$.

Remarks

1. In this paper, f is usually a parameterization of a submanifold \mathcal{M} , so f is 1 – 1 and the right-hand integral is the surface area integral of g over $f(A)$. The left side shows how to carry out this integral using Lebesgue measure on \mathbb{R}^m and the Jacobian. It shows that sampling from the density $J_m f(x)$ (normalized) on \mathbb{R}^m and then mapping onto \mathcal{M} via f gives a sample from the area measure.

2. There are many extensions and refinements of the area formula [F, Sect. 3.2]. In particular [F, Sect. 3.2.20] extends things to approximately differentiable functions and [F, Sect. 3.2.46] extends from Euclidean space to Riemannian manifolds.

Example 1A continued: The Curved Torus For the parameterization given in Example 1A, the curved torus is the Lipschitz image of $\{0 \leq \theta, \psi < 2\pi\}$, with

$$f(\theta, \psi) = (R + r \cos(\theta)) \cos(\psi), (R + r \cos(\theta)) \sin(\psi), r \sin(\theta) \quad (2.2)$$

$$Df(\theta, \psi) = \begin{bmatrix} -r \sin(\theta) \cos(\psi) & -(R + r \cos(\theta)) \sin(\psi) \\ -r \sin(\theta) \sin(\psi) & (R + r \cos(\theta)) \cos(\psi) \\ r \cos(\theta) & 0 \end{bmatrix} \quad (2.3)$$

$$J_2^2 f(\theta, \psi) = \det \begin{bmatrix} r^2 & 0 \\ 0 & (R + r \cos(\theta))^2 \end{bmatrix} = r^2 (R + r \cos(\theta))^2 \quad (2.4)$$

As explained in Section 2, \mathcal{M} is parametrized by $\mathcal{U} = \{\theta, \psi, 0 \leq \theta, \psi < 2\pi\}$ and the task reduces to sampling (θ, ψ) from the density $g(\theta, \psi) = (\frac{1}{4\pi^2})(1 + (r/R) \cos \theta)$. A random number generator outputs points that we assume are uniformly distributed on $[0, 1]$ and the task reduces to converting these into a sample from g . From the form of g , the measure factors into the uniform density for ψ on $[0, 2\pi)$ and the density

$$g_1(\theta) = \frac{1}{2\pi} \left(1 + \frac{r}{R} \cos \theta\right) \quad \text{on } 0 \leq \theta < 2\pi.$$

We may sample points from g_1 by rejection sampling (Hammersley and Handscomb, 1964). The function $(1 + (r/R) \cos \theta)$ is enclosed in the box $0 \leq \theta < 2\pi, [1 - (r/R) < \eta < 1 + (r/R)]$. Choose points (θ, η) uniformly in this box from two-dimensional Lebesgue measure. This uses two calls to the underlying uniform random number generator. If $\eta < 1 + (r/R) \cos \theta$, output θ . If not, choose again, continuing until the condition holds. The resulting θ is distributed as g_1 . Sample code for this is in algorithm 1.

Algorithm 1 Rejection Sampling yielding g_1 .

```

reject=function(n=100,r=0.5,R=1){
#Rejection sampler
xvec=runif(n,0,2*pi)
yvec=runif(n,0,1/pi)
fx=(1+(r/R)*cos(xvec))/(2*pi)
return(xvec[yvec<fx]) }

```

What we get is a density with support $[0, 2\pi]$. See Figures 1 and 2 below.

Example 1B continued: Sum and Product Fixed Here

$$\mathcal{M} = \left\{ (x_1, \dots, x_n); \quad x_i > 0, \quad \sum_{i=1}^n x_i = S, \quad \prod_{i=1}^n x_i = P \right\}.$$

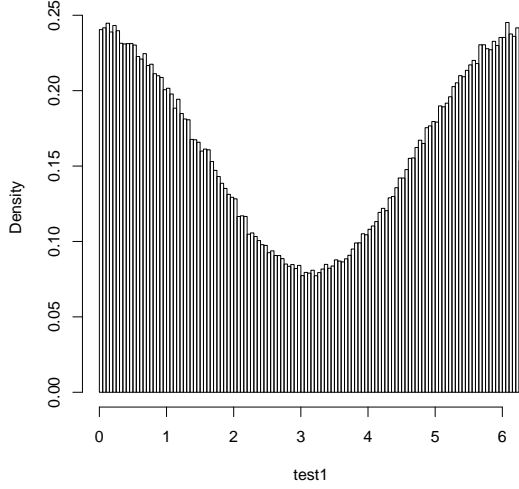


FIG 1. Rejection sampling density proportional to $1 + \frac{r}{R} \cos(\theta)$

The constraints S, P satisfy $0 < P^{1/n} \leq S/n$ because of the arithmetic-geometric mean inequality. Any such S, P can occur. To find a parameterization of \mathcal{M} consider the projection

$$\begin{aligned} \Pi : \mathcal{M} &\rightarrow \mathbb{R}^{n-2} \\ (x_1, \dots, x_n) &\rightarrow (x_3, \dots, x_n) \end{aligned}$$

Let $s = x_3 + \dots + x_n = S - t$ with $t \geq 0$ and $x_3 x_4 \dots x_n = p$. The equations $x_1 + x_2 = t, x_1 x_2 = P/p$ have a positive real solution if and only if $t^2 \geq 4P/p$. In this case the solution is the pair

$$\{x_1, x_2\} = \left(t \pm \sqrt{t^2 - \frac{4P}{p}} \right) / 2.$$

One way to parametrize \mathcal{M} is to define

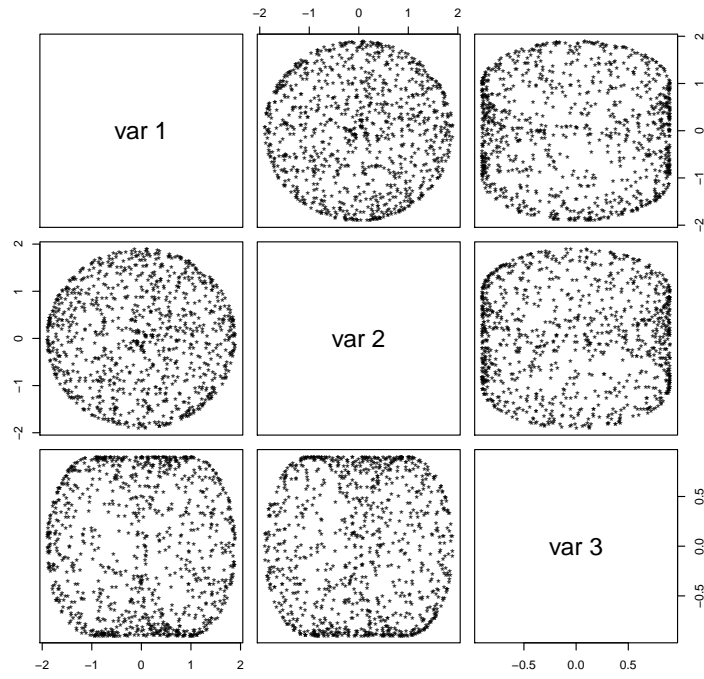
$$\mathcal{M}^+ = \{x \in \mathcal{M} : x_1 \geq x_2\}, \quad \mathcal{M}^- = \{x \in \mathcal{M} : x_1 < x_2\} \quad (2.5)$$

Define $\mathcal{U} = \{(x_3, \dots, x_n) : x_i > 0, s < S, p < 4P/(S-4)^2\}$ $s = \sum_{i=3}^n x_i, p = \prod_{i=3}^n x_i$;
 $f : \mathcal{U} \rightarrow \mathcal{M}^+$ is defined by

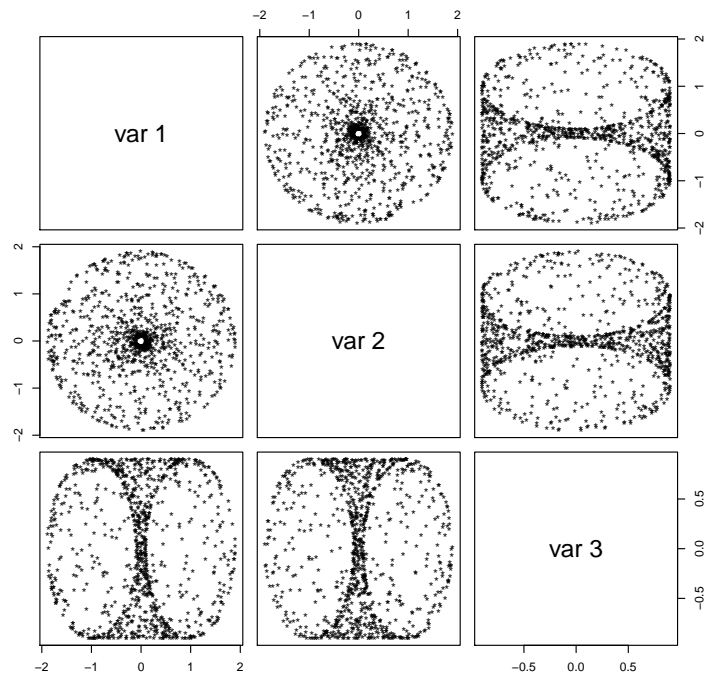
$$f(x_3, \dots, x_n) = (f_1(x_3, \dots, x_n), f_2(x_3, \dots, x_n), x_3, \dots, x_n) \quad (2.6)$$

$$\text{with } f_1(x_3, \dots, x_n) = \frac{(S - \sum_{i=3}^n x_i) + \sqrt{(S - \sum_{i=3}^n x_i)^2 - \frac{4P}{\prod_{i=3}^n x_i}}}{2}$$

$$\text{and } f_2(x_3, \dots, x_n) = \frac{(S - \sum_{i=3}^n x_i) - \sqrt{(S - \sum_{i=3}^n x_i)^2 - \frac{4P}{\prod_{i=3}^n x_i}}}{2}$$



Correctly generated points uniformly on the torus



Basic Uniform on Parameters.

FIG 2. Top figure shows a 3D representation of a sample of size 1000 with parameters $R = 1$, $r = 0.9$, the lower figure shows the incorrectly sampled points, although the difference is not obvious visually, standard tests pick up the difference between the two distributions.

The derivative is the $n \times (n - 2)$ matrix

$$Df = \begin{bmatrix} D_3 f_1 & D_4 f_1 & \cdots & D_n f_1 \\ D_3 f_2 & D_4 f_2 & \cdots & D_n f_2 \\ 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ & & \vdots & \\ 0 & 0 & \cdots & 1 \end{bmatrix} \quad (2.7)$$

$(J_{n-2}f(x))^2 = \det((Df)^T Df)$ is the determinant of a matrix of form $I_{n-2} + VV^T + WW^T$ with V^T, W^T the first and second rows of Df . A well-known determinant identity reduces this to a 2×2 determinant; if B is $p \times m$ and C is $m \times p$ then $\det(I_p + BC) = \det(I_m + CB)$. It follows that

$$((J_{n-2}f(x))^2 = \det\left(I_2 + \begin{pmatrix} V^T V & V^T W \\ V^T W & W^T W \end{pmatrix}\right). \quad (2.8)$$

To summarize:

Proposition 1. *The density of the $(n-2)$ -dimensional area measure $\bar{\mathcal{H}}^{n-2}$ on the submanifold \mathcal{M}^+ in (2.5) parametrized by $f : \mathcal{U} \rightarrow \mathcal{M}^+$ is $J_{n-2}f(x)$ of (2.8), above with V, W the first two rows of matrix (2.7).*

Remarks

1. Up to sets of $\bar{\mathcal{H}}^{n-2}$ measure 0, a similar result holds for \mathcal{M}^- . Since \mathcal{M}^- and \mathcal{M}^+ patently have the same area measure, it is easy to sample from \mathcal{M} using an additional coin flip to randomize the first two coordinates.

2. Of course, any $(n - 2)$ -tuple of coordinates can be used for the parameterization. In practical simulation, it might be wise to sample from \mathcal{M} as above and follow this by a random permutation of the coordinates.

3. The function f defined in (2.6), is only locally Lipschitz because of p in the denominator. However, \mathcal{U} may be decomposed into a countable union of pieces with f Lipschitz on each piece. Because the formula for $J_n f$ is local, the proposition is true as stated.

Given S, P , the manifold \mathcal{M}^+ is parametrized by \mathcal{U} of Example 1B. The task of sampling from area measure on \mathcal{M} is reduced to sampling from $J_{n-2}f(x)$ in \mathcal{U} . One problem here is that although $z = \int_{\mathcal{U}} J_{n-2}f(x) \lambda^{n-2}(dx) < \infty$ and $J_{n-2}f/z$ is a probability density on \mathcal{U} , the value of z is unknown. This standard problem may be approached by the Metropolis algorithm, the Gibbs sampler, importance sampling, or by the hit-and-run algorithm in many variations (see Liu (2001), Andersen and Diaconis (2008) for background). Here, we briefly explain the Metropolis algorithm for sampling from $J_{n-2}f$. This generates a Markov chain X_0, X_1, X_2, \dots starting from $X_0 = x_0$, a fixed point in \mathcal{U} . From $X_n = x \in \mathcal{U}$, we propose $y \in \mathbb{R}^M$, choosing $y = x + \epsilon$ with ϵ chosen (say) uniformly from a unit cube centered at x . Then,

$$X_{n+1} = \begin{cases} y & \text{with probability } \min\left(\frac{J_{n-2}f(y)}{J_{n-2}f(x)}, 1\right) \\ x & \text{otherwise.} \end{cases}$$

Since $J_{n-2}f(y)$ is taken as 0 outside \mathcal{U} , note that $X_{n+1} \in \mathcal{U}$. Standard theory shows that for n large, $P(X_n \in A) \sim \int_A \frac{J_{n-2}f(x)}{z} \lambda^M(dx)$. Careful evaluation of how large n must be to make this approximation valid is an open research problem both here and in most real applications of the Metropolis algorithm (see Diaconis and Saloff-Coste (1998) and Diaconis, Lebeau and Michel (2010b)). A host of heuristics are available for monitoring convergence; for adapting the choice of the proposal for ϵ and for efficient use of the output. We will not discuss these further here.

Several further examples admitting an explicit parameterization, with computations of Jf , are in Hubbard and Hubbard (2007, Chap. 5) which is enthusiastically recommended to newcomers.

2.3. Conditional Densities and the Co-Area Formula

Federer's co-area formula gives an explicit density for the conditional distribution. The main tool is:

Theorem: Co-Area Formula [F, Sect. 3.2.12] *Suppose that $\Phi : \mathbb{R}^M \rightarrow \mathbb{R}^N$ is Lipschitz with $M > N$. Then*

$$\int_{\mathbb{R}^M} g(x) J_N \Phi(x) \lambda^M(dx) = \int_{\mathbb{R}^N} \int_{\Phi^{-1}(y)} g(x) \bar{\mathcal{H}}^{M-N}(dx) \lambda^N(dy). \quad (2.9)$$

In (2.9), g is Lebesgue measurable from $\mathbb{R}^M \rightarrow \mathbb{R}$ and $J_N \Phi$ is defined in Section 2.1.

Recall next the definition of a regular conditional probability. Let (Ω, \mathcal{F}, P) be a probability space and $\mathcal{C} \subseteq \mathcal{F}$ a sub-sigma algebra. A function $P(w, dw)$ from $(\Omega \times \mathcal{F})$ into $[0, 1]$ is a regular conditional probability for P given \mathcal{C} if

$$\text{For each } w \in \Omega, P(w, \cdot) \text{ is a probability measure on } \mathcal{F}. \quad (2.10a)$$

$$\text{For each } F \in \mathcal{F}, \text{ the function } w \mapsto P(w, F) \text{ is } \mathcal{C} \text{ measurable}. \quad (2.10b)$$

$$\text{For } C \in \mathcal{C}, F \in \mathcal{F}, P(C \cap F) = \int_C P(w, F) P(dw). \quad (2.10c)$$

Let $p(x)$ be a probability density on \mathbb{R}^M with respect to $\lambda^M(dx)$. Let $\Phi : \mathbb{R}^M \rightarrow \mathbb{R}^N$ be Lipschitz with $M > N$. From Rademacher's Theorem, Φ is differentiable at almost every x , and $J_N \Phi(x)$ can be computed by the usual rules.

Proposition 2. *Suppose that $J_N \Phi(x)$ exists and is strictly positive for all x where $p(x) > 0$. Then*

(a) *The marginal density of Φ is absolutely continuous with density*

$$m(y) = \int_{\Phi^{-1}(y)} \frac{p(x)}{J_N \Phi(x)} \bar{\mathcal{H}}^{M-N}(dx) \text{ with respect to } \lambda^M(dy).$$

(b) *If $m(y) \in \{0, \infty\}$, set $\mathcal{Q}(y, F) = \delta_{x^*}(F)$ for some fixed $x^* \in \mathbb{R}^M$. Else set*

$$\mathcal{Q}(y, F) = \frac{1}{m(y)} \int_{\Phi^{-1}(y) \cap F} \frac{p(x)}{J_N \Phi(x)} \bar{\mathcal{H}}^{M-N}(dx).$$

Set $P(x, F) = \mathcal{Q}(\Phi(x), F)$. Then P is a regular conditional probability for $P(dx) = p(x)\lambda^M(dx)$ given $\mathcal{C} = \Phi^{-1}(\mathcal{B})$ with \mathcal{B} the Lebesgue measurable sets in \mathbb{R}^N .

Proof. Clearly (2.10a) and (2.10b) are satisfied. To show (2.10c), fix $C \in \mathcal{C}$ and F a Lebesgue measurable set in \mathbb{R}^M . Take g in (2.1) to be

$$\delta_{C \cap F}(x) \frac{p(x)}{J_N \Phi(x)} \quad \text{with } g(x) \text{ defined as 0 if } p(x) = 0.$$

Where $\delta_{C \cap F}$ denotes the indicator function of the intersection $C \cap F$.

The co-area formula shows

$$\begin{aligned} P(C \cap F) &= \int_{C \cap F} p(x)\lambda^M(dx) = \int_{\mathbb{R}^N} \int_{\Phi^{-1}(y)} \delta_C(x) \delta_F(x) \frac{p(x)}{J_N \Phi(x)} \bar{\mathcal{H}}^{M-N}(dx) \lambda^N(dy) \\ &= \int_C \int_{\Phi^{-1}(y) \cap F} \frac{p(x)}{J_N \Phi(x)} \bar{\mathcal{H}}^{M-N}(dx) \lambda^N(dy). \end{aligned}$$

Let $C_0 = \{y : m(y) = 0\}$, $C_\infty = \{y : m(y) = \infty\}$, $C^+ = (C_0 \cup C_\infty)^C$. Taking $C = F = \mathbb{R}^M$, we see $\lambda^N(C_\infty) = 0$. For $y \in C_0$, $\int_{\Phi^{-1}(y) \cap F} \frac{p(x)}{J_N \Phi(x)} \bar{\mathcal{H}}^{M-N}(dx) = 0$. Hence, the integrals equal

$$\begin{aligned} &\int_{C \cap C^+} \int_{\Phi^{-1}(y) \cap F} \frac{p(x)}{J_N \Phi(x)} \bar{\mathcal{H}}^{M-N}(dx) \lambda^N(dy) \\ &= \int_{C \cap C^+} \frac{m(y)}{m(y)} \int_{\Phi^{-1}(y) \cap F} \frac{p(x)}{J_N \Phi(x)} \bar{\mathcal{H}}^{M-N}(dx) \lambda^N(dy) \\ &= \int_C m(y) \mathcal{Q}(y, F) \lambda^N(dy) \\ &= \int_C P(x, F) P(dx) \end{aligned}$$

□

Remark Of course, $m(y)$ can be 0, if $\Phi^{-1}(y)$ is empty or p vanishes there. Similarly, $m(y)$ can be infinite: consider (following Tjur [1972, Sect. 30]) a set of finite area in \mathbb{R}^2 of the shape shown in Figure 3. Let $p(x)$ be the normalized indicator of this set. Let $\Phi(x, y) = x$, so $J_N \Phi(x) = 1$. Then $m(0) = \infty$.

Example 1A (continued): From (1.1) the torus is $\{(x, y, z) \in \mathbb{R}^3\}$

$$x = (R + r \cos(\theta)) \cos(\psi), \quad y = (R + r \cos(\theta)) \sin(\psi), \quad z = r \sin(\theta)$$

$0 \leq \theta, \psi < 2\pi$ for $R > r > 0$. What is the conditional distribution in (θ, ψ) space given that $x = 0$? In the notation of Proposition 2,

$$p(\theta, \psi) = \begin{cases} \frac{1}{2\pi} (1 + \frac{r}{R} \cos(\theta)) & 0 \leq \theta, \psi < 2\pi \\ 0 & \text{elsewhere} \end{cases}$$

The function $\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}$ is

$$\Phi(\theta, \psi) = (R + r \cos(\theta)) \cos(\psi)$$

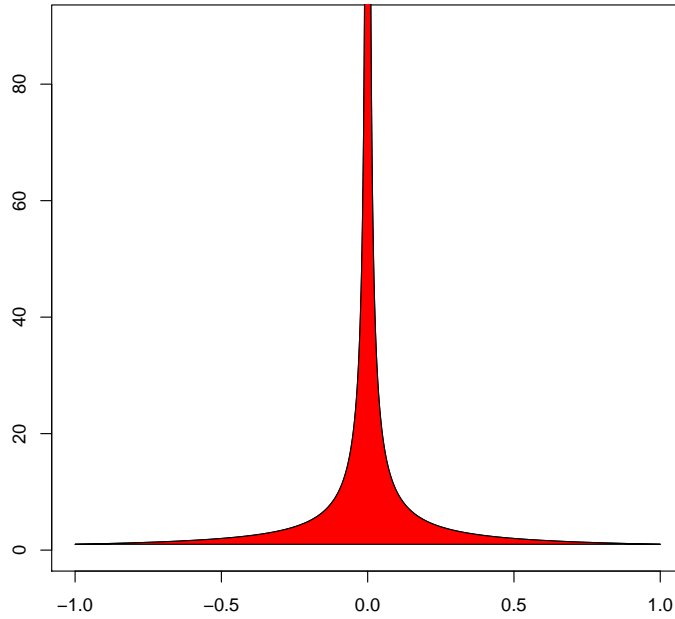


FIG 3. Instance of Infinite region.

Thus

$$(J\Phi)^2 = (r \sin(\theta) \cos(\psi))^2 + ((R + r \cos(\theta)) \sin(\psi))^2$$

$$\Phi^{-1}(0) = \{(\theta, \psi), 0 \leq \theta < 2\pi, \psi \in \{\frac{\pi}{2}, \frac{3\pi}{2}\}\}$$

It follows that $J\Phi(\theta, \frac{\pi}{2}) = J\Phi(\theta, \frac{3\pi}{2}) = R + r \cos(\theta)$. This is proportional to $p(\theta, \psi)$ and Proposition 2b says that the conditional distribution is uniform on the two line segments that make up $\Phi^{-1}(0)$ and assigns equal mass to each segment.

Example 1B (continued) Consider the area measure on \mathcal{M}^+ of (2.5). Proposition 1 above shows that \mathcal{M}^+ is parametrized by a map f from the set $U \subset \mathbb{R}^{n-2}$ and gives an explicit expression for the corresponding probability density. One standard method for sampling from this density is to use the Gibbs sampler. This entails sampling from the conditional distribution given the values at some of the coordinates. One simple implementation which uses Proposition 2 is this: \mathcal{M}^+ is given as an embedded manifold in \mathbb{R}^n . From $(x_1, x_2, \dots, x_n) \in \mathcal{M}^+$, choose three coordinates uniformly at random, fix the remaining $(n-3)$ coordinates. The map f of Proposition 1 composed with the projection onto the corresponding $(n-3)$ space gives a map $\Phi: \mathcal{U} \rightarrow \mathbb{R}^{n-3}$. The conditional density given $\Phi = y$ is explicitly given by Proposition 2. Here $\Phi^{-1}(0)$ is a one dimensional curve and the sampling problem reduces to a standard task. We omit further details.

Example 3C: How Not To Sample Here is a mistake to avoid. Let \mathcal{M} be a compact embedded manifold. To sample from the area measure, the following scheme presents itself. Suppose that for each point $x \in \mathcal{M}$ a neighborhood $\mathcal{N}_x \subseteq \mathcal{M}$ is specified (e.g., a ball of specified radius on \mathcal{M}). Suppose it is possible to sample from the area measure restricted to \mathcal{N}_x . It seems plausible that this drives a Markov chain with area measure globally. This is an error. Perhaps the easiest way to see through the problem is to consider the discrete case:

Consider a finite connected undirected graph with vertex set \mathcal{X} and edge set \mathcal{E} . Let $\pi(x) > 0$, $\sum_{x \in \mathcal{X}} \pi(x) = 1$ be a probability distribution \mathcal{X} . Suppose for each point $x \in \mathcal{X}$, a neighborhood \mathcal{N}_x is defined. These may be arbitrary finite sets; we do not need $x \in \mathcal{N}_x$, but will assume $y \in \mathcal{N}_x \leftrightarrow x \in \mathcal{N}_y$. For example, we may take $\mathcal{N}_x = B_r(x)$, the r -ball using graph distance. A Markov chain on \mathcal{X} is defined as follows:

From x , choose $y \in \mathcal{N}_x$ with probability π restricted to \mathcal{N}_x . Thus

$$K(x, y) = \begin{cases} \frac{\pi(y)}{\pi(\mathcal{N}_x)} & \text{if } y \in \mathcal{N}_x \\ 0 & \text{otherwise} \end{cases} \quad (2.11)$$

Lemma 1. *The chain (2.11) is reversible with reversing measure*

$$\sigma(x) = \frac{\pi(\mathcal{N}_x)\pi(x)}{z}, \text{ with } z \text{ a normalizing constant.} \quad (2.12)$$

Proof. If $K(x, y) = 0$, then $K(y, x) = 0$, so reversibility holds. Otherwise

$$\sigma(x)K(x, y) = \frac{\pi(\mathcal{N}_x)\pi(x)}{z} \frac{\pi(y)}{\pi(\mathcal{N}_x)} = \frac{\pi(x)\pi(y)}{z} = \sigma(y)K(y, x).$$

□

Remarks

1. Thus, unless $\pi(\mathcal{N}_x) = \text{constant}$, $\sigma(x) \neq \pi(x)$.
2. In the continuous setting, sampling locally from area measure $\bar{\mathcal{H}}$, this chain has stationary density proportional to $\bar{\mathcal{H}}(\mathcal{N}_x)$. An analysis of rates of convergence for this walk on compact Riemannian manifolds in Lebeau and Michel (2010).
3. On a curve, with distance measured by arc length, $\bar{\mathcal{H}}(B_r(x))$ is constant for r suitably small because of the volume of tubes theorem. However, this is no longer true for higher-dimensional manifolds with non-constant Gaussian curvature.
4. We may use the Metropolis algorithm to change the stationary distribution from σ in (2.12) to π . The chain is $\mathcal{M}(x, y) = \pi(y) \min(\frac{1}{\pi(\mathcal{N}_x)}, \frac{1}{\pi(\mathcal{N}_y)})$ for $x \neq y \in \mathcal{N}_x$. Note that this requires knowledge of $\pi(\mathcal{N}_x)$, $\pi(\mathcal{N}_y)$.

3. Exponential Families, Conditional Densities and the Co-Area Formula

One motivation for the current work is conditional testing in statistical problems. This is a central topic of classical statistics beginning with R. A. Fisher's exact test for independence in

contingency tables and the Neyman–Pearson theory of uniformly most powerful unbiased tests for exponential families. The best general reference for these topics is (Lehmann and Romano, 2005, Chap. 4, 5, 10) See also the survey in Diaconis and Sturmfels (1998) and the techniques and references in Lindqvist and Taraldsen (2005, 2006).

The problems addressed in the present paper are a continuous analog. Section 3.1 below presents exponential families in a version convenient for applications. Section 3.2 briefly discusses conditional densities and sufficiency. Section 3.3 uses the co-area formula to give a useful expression for the conditional density, given a sufficient statistic, with respect to the area measure. These formulae are applied in Section 4.

3.1. Exponential Families

Many widely-used families of probability measures, such as the Gamma family of Example 1B, have a common exponential form. Theorems and properties can be derived generally and then applied in specific cases. A good first reference for this material is (Lehmann and Romano, 2005, Sect. 2.7). The specialist monographs of Barndorff-Nielsen (1978), Brown (1986) and Letac (1992) may be supplemented by the references in Diaconis, Khare and Saloff-Coste (2010a) to give an overview of this basic subject.

Let $T : \mathbb{R}^a \rightarrow \mathbb{R}^b$ be a measurable function. Let $\Theta \subseteq \mathbb{R}^b$ be a non-empty open set and $\psi : \Theta \rightarrow \mathbb{R}^b$ a measurable function. Let $f(x) : \mathbb{R}^a \rightarrow \mathbb{R}_+$ be measurable and suppose

$$0 < z(\theta) = \int_{\mathbb{R}^a} f(x) e^{\psi(\theta) \bullet T(x)} \lambda^a(dx) < \infty \text{ for each } \theta \in \Theta.$$

Definition The family of probability densities

$$P_\theta(x) = z^{-1}(\theta) f(x) e^{\psi(\theta) \bullet T(x)} \quad \theta \in \Theta \quad (3.1)$$

is called the exponential family generated by (f, Θ, ψ, T) .

For the Gamma family in Example 1B, $a = 1, b = 2, T(x) = \begin{cases} (x, \log x) & x > 0 \\ 0 & \text{otherwise} \end{cases}$

$$\Theta = \mathbb{R}_+^2 = \{(\sigma, a) : \sigma, a > 0\}, \quad \psi(\sigma, a) = \left(-\frac{1}{\sigma}, a-1\right), \quad z(\theta) = \sigma^a \Gamma(a) \quad f(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}.$$

The exponential families here are a subclass, in having absolutely continuous densities whose support does not depend on θ .

3.2. Sufficiency

The product measure on $(\mathbb{R}^a)^n$ generated by P_θ of (3.1) has density

$$z(\theta)^{-n} \prod_{i=1}^n f(x_i) e^{\psi(\theta) \bullet \sum_{i=1}^n T(x_i)}.$$

The function $\bar{T} = \sum_{i=1}^n T(x_i)$ is called a sufficient statistic for the family. The references above show that the distribution of the product measure conditional on \bar{T} does not depend on θ . Conversely, the Koopman–Pitman–Darmois theorem says if P_θ is a family of measures on \mathbb{R}^a with T locally Lipschitz and for some $n \geq 2$, the distribution of P_θ^n , conditional on T , does not depend on θ , then P_θ is an exponential family. See Hipp (1974) for a careful statement; see Diaconis (1988) for background and further references on sufficiency.

For the Gamma family, \bar{T} is equivalent to $S = \sum_{i=1}^n x_i$, $P = \prod_{i=1}^n x_i$ as used throughout.

3.3. Conditional Densities and the Co-Area Formula

This dual to the area formula is explained in Section 2.3 above. We may use it directly to compute an expression for the conditional density of an exponential family given a sufficient statistic.

Theorem 1. *With notation as above, for $na > b$ consider an exponential family (3.1) based on a Lipschitz $T : \mathbb{R}^a \rightarrow \mathbb{R}^b$. Let $\bar{T} : \mathbb{R}^{na} \rightarrow \mathbb{R}^b (= \sum_{i=1}^n T(x_i))$ and suppose $J_b \bar{T}(x) \neq 0$ for $\prod f(x_i) \neq 0$. Define $\mathcal{M}_t = \{\mathbf{x} \in (\mathbb{R}^a)^n : \bar{T}(\mathbf{x}) = t\}$. Then, the conditional density on \mathcal{M}_t with respect to the area measure is*

$$W^{-1} \prod_{i=1}^n f(\mathbf{x}_i) / J_b \bar{T}(\mathbf{x}). \quad (3.2)$$

with the normalizing constant $W = W_t$ taken to be $\int \prod_{i=1}^n f(\mathbf{x}_i) / J_b \bar{T}(\mathbf{x}) \bar{\mathcal{H}}^{M-N}(d\mathbf{x})$ if this integral is in $(0, \infty)$.

Proof. In the co-area formula take $\Psi = \bar{T} : (\mathbb{R}^a)^n \rightarrow \mathbb{R}^b$. Thus $M = na$, $N = b$. For $h : \mathbb{R}^M \rightarrow \mathbb{R}^N$ bounded continuous, set

$$g(\mathbf{x}) = \begin{cases} \frac{h(\mathbf{x})}{J_N \Psi(\mathbf{x})} \prod_{i=1}^n P_\theta(x_i) & \text{if } J_N \Psi(\mathbf{x}) \neq 0 \\ 0 & \text{otherwise} \end{cases}.$$

Then $\Psi^{-1}(t) = \mathcal{M}_t$ and the co-area formula shows that \mathcal{M}_t has positive, finite total area measure for λ^N a.e.t. Further

$$\int h(\mathbf{x}) \prod_{i=1}^n P_\theta(\mathbf{x}_i) \lambda^M(d\mathbf{x}) = \int_{\mathbb{R}^b} e^{\Psi(\theta) \cdot t} z(\theta)^{-n} \int_{\mathcal{M}_t} \frac{h(\mathbf{x}) \prod_{i=1}^n f(\mathbf{x}_i)}{J_N \Psi(\mathbf{x})} \bar{\mathcal{H}}^{M-N}(d\mathbf{x}) \lambda^N(dt).$$

This formula says that (3.2) is a regular conditional probability for the product measure $\prod_{i=1}^n P_\theta(\mathbf{x})$ given $\bar{T} = t$. \square

Remarks

1. Since the conditional density (3.2) does not depend on θ , \bar{T} is a sufficient statistic.
2. The calculation shows that the marginal density of \bar{T} is $e^{\Psi(\theta) \cdot t} / z(\theta)^n W$ with respect to $\lambda^b(dt)$. Thus the induced measures of \bar{T} form an exponential family.

Example: Gamma Family With $\bar{T} : \mathbb{R}^n \rightarrow \mathbb{R}^2$ given by $\bar{T}(\mathbf{x}) = (\sum_{i=1}^n x_i, \sum_{i=1}^n \log x_i)$, for $n > 2$,

$$D\bar{T}(\mathbf{x}) = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ \frac{1}{x_1} & \frac{1}{x_2} & \cdots & \frac{1}{x_n} \end{bmatrix}, \quad J_2^2 \bar{T}(\mathbf{x}) = \sum_{i < j} \left(\frac{1}{x_i} - \frac{1}{x_j} \right)^2.$$

From Theorem 1, we may sample from the conditional distribution of the Gamma family given $\bar{T} = t$ on \mathcal{M}_t by sampling from the probability density (w.r.t. λ^{n-2}) proportional to

$$\frac{J_{n-2} f(x_3, \dots, x_n)}{J_2 \bar{T}(f(x_3, \dots, x_n))}$$

on \mathcal{U} and f defined in (2.6), followed by randomizing the first two coordinates by a fair coin toss.

4. Neyman's Smooth test and the Gibb's sampler

This section illustrates a useful general procedure (the Gibbs sampler) in a natural example: Neyman's smooth test for goodness of fit. The problem reduces to sampling from an explicit density $f(x_1, x_2, \dots, x_n)$ on the following submanifold: fix m and $p_1 \geq p_2 \geq \dots \geq p_m$. Let

$$\mathcal{M}_{\mathbf{p}} = \{x_1, x_2, \dots, x_n, 0 \leq x_i \leq 1, \sum_{j=1}^n x_j^i = p_i, 1 \leq i \leq m\}.$$

In Neyman's case, $m = 4$, assume this for now. The idea underlying our algorithm, developed below, is to pick a uniformly chosen subset of $m + 1 = 5$ coordinates with probability $1/\binom{n}{5}$, say the first five. Set $\bar{p}_i = \sum_{j=1}^5 x_j^i$. The submanifold

$$\mathcal{M}_{\bar{\mathbf{p}}} = \{x_1, x_2, \dots, x_5, 0 \leq x_i \leq 1, \sum_{j=1}^5 x_j^i = \bar{p}_i, 1 \leq i \leq 4\} \quad (4.1)$$

is a curve which lies both on the submanifold $\mathcal{M}_{\mathbf{p}}$ and in \mathbb{R}^5 . We may sample from the conditional density on the curve and replacing (x_1, x_2, \dots, x_5) by the sampled values gives a new point on $\mathcal{M}_{\mathbf{p}}$.

Repeatedly choosing fresh five-tuples gives a connected reversible Markov chain on $\mathcal{M}_{\mathbf{p}}$ with f as its stationary density. In the present section we find it convenient to work directly with the density of f with respect to the area measure, avoiding the extra step of local coordinates. Neyman's smooth test is developed in 4.1, the relevant conditional densities are derived in 4.2 and 4.3 contains a study of the ergodicity of this chain. Section 4.4 develops an idea of Besag and Clifford (1989) for valid testing with non ergodic chains.

4.1. Neyman's Smooth test

Consider investigating the following null hypothesis; fix F a distribution function of a continuous random variable. Let

$$H_0 : X_1, X_2, X_3, \dots, X_n \sim iid \quad F \quad (4.2)$$

Then

$$Y_i = F(X_i)$$

are iid uniform on $[0, 1]$. Neyman (1937) developed a test of H_0 based on testing $\theta = 0$ in the model

$$f_\theta(y) = z^{-1} e^{\theta_1 y + \theta_2 y^2 + \theta_3 y^3 + \theta_4 y^4}, \quad 0 \leq y \leq 1 \quad (4.3)$$

This test (and its modifications by David (1939); Barton (1953, 1956)) has been shown to have a good empirical track record and comes in for repeated favorable mention in Lehman and Romano's survey of testing goodness of fit (Lehmann and Romano, 2005, chapter 9). That chapter also explains the difficulty of such omnibus testing problems. One justification for this test is that if the data are from a *smooth* distribution F , using a simple χ^2 test loses information because it breaks the data into categorical bins, losing the actual ordering of the bins.

Any smooth positive probability density $h(y)$ on $[0, 1]$ can be expanded as

$$h(y) = e^{\log h(y)} = e^{\sum_{i=0}^{\infty} \theta_i y^i}$$

The four parameter exponential family is a commonsense truncation of this non-parametric model. Fan (1996) has developed tests based on m term approximations with m chosen adaptively from the data.

In the rest of this section we investigate the adequacy of the truncation (4.3) (with $m = 4$) by testing if the model (4.3) fits the data. Thus given data Y_1, Y_2, \dots, Y_n in $[0, 1]$, we develop conditional tests of the model (4.3). These ideas work for every m and could be used as input to Fan's adaptive procedure. The four dimensional sufficient statistics for the family (4.3) is

$$\mathbf{p} = (p_1, p_2, p_3, p_4), \quad p_i = \sum_{j=1}^n Y_j^i$$

The conditional procedures explained in section 4.2 are based on the conditional distribution of the model f_θ given \mathbf{p} . This is supported on

$$\mathcal{M}_{\mathbf{p}} = \{(x_1, x_2, \dots, x_n), 0 \leq x_i \leq 1, \sum_{j=1}^n x_j^i = p_i, 1 \leq i \leq 4\} \quad (4.4)$$

This is a compact $n - 4$ dimensional submanifold of $[0, 1]^n$. To actually construct a test, a test statistic must be chosen. Neyman's test of section 4.1 was based on the L^2 norm of the averages of the first four orthogonal polynomials for the uniform distribution on $[0, 1]$. Under (4.3) the sum of these norms should have been an approximate chi-square (4) distribution. We may follow Neyman, using a further orthogonal polynomial as the test statistic but calibrating it with the exact conditional distribution.

4.2. The Gibbs Sampler

The Gibbs sampler is well developed in Liu (2001). As usually explained, to sample from a probability density $g(z_1, z_2, \dots, z_n)$ on \mathbb{R}^n one begins at a sample point $z_0 = (z_1^0, z_2^0, \dots, z_n^0)$ and changes coordinates sequentially: first to $(z_1^1, z_2^0, \dots, z_n^0)$ then to $(z_1^1, z_2^1, \dots, z_n^0)$...then $z_1 =$

$(z_1^1, z_2^1, \dots, z_n^1)$. The i th change is made by sampling from the conditional distribution of the i th coordinate given all the rest. The one dimensional problem is supposed to be easy to do. The transition from z^0 to z^1 is one step of the Gibbs sampler. Proceeding as above to z^2, z^3, \dots gives a Markov chain with g as stationary distribution. In the present problem

- (a) It is not possible to change just one coordinate and stay on the surface (4.1). The minimal change is in five coordinates resulting in the curve

$$\{(x_1, x_2, \dots, x_5) : 0 \leq x_i \leq 1 \sum_i^5 = p_j\}. \quad (4.5)$$

- (b) Instead of random sampling, one can systematically run through all sets of five coordinates using for instance a Gray code approach as in Diaconis and Holmes (1994).
- (c) Sampling from a conditional distribution on the curve in (a) is not so simple and instead a single Metropolis step is proposed. This is sometimes called ‘Metropolis on Gibbs’ in the literature, for notational clarity we suppose that the five chosen coordinates are the first five. Let P be the conditional distribution for the model (4.3) on the submanifold (4.1). Let Q be the conditional measure on the curve (4.5). The following proposition determines the density of Q with respect to arc-length.

Proposition 3. *The measure Q on the curve (4.5) has density with respect to arc-length*

$$q(x_1, x_2, x_3, x_4, x_5) = z^{-1} \sqrt{J_4^{-1}} \quad z^{-1} \text{ a normalizing constant}$$

$$J_4 = \det \begin{pmatrix} 5 & 2\bar{p}_1 & 3\bar{p}_2 & 4\bar{p}_3 \\ 2\bar{p}_1 & 4\bar{p}_2 & 6\bar{p}_3 & 8\bar{p}_4 \\ 3\bar{p}_2 & 6\bar{p}_3 & 9\bar{p}_4 & 12\bar{p}_5 \\ 4\bar{p}_3 & 8\bar{p}_4 & 12\bar{p}_5 & 16\bar{p}_5 \end{pmatrix} \quad \bar{p}_i = \sum_{j=1}^5 x_j^i, 1 \leq i \leq 5 \quad (4.6)$$

Proof. By the usual calculus of double conditioning, Q is the conditional distribution of the product measure f_θ^5 on $[0, 1]^5$ given $\bar{p}_1, \bar{p}_2, \bar{p}_3, \bar{p}_4$. Now use Theorem 1 of section 3.3. The mapping $\bar{T} : [0, 1]^5 \rightarrow \mathbb{R}^4$ takes $T(y_1, y_2, y_3, y_4, y_5) = (\bar{p}_1, \bar{p}_2, \bar{p}_3, \bar{p}_4)$. Clearly the 5×4 derivative $D\bar{T}$ is

$$D\bar{T} = \begin{pmatrix} 1 & 2y_1 & 3y_1^2 & 4y_1^3 \\ 1 & 2y_2 & 3y_2^2 & 4y_2^3 \\ 1 & 2y_3 & 3y_3^2 & 4y_3^3 \\ 1 & 2y_4 & 3y_4^2 & 4y_4^3 \\ 1 & 2y_5 & 3y_5^2 & 4y_5^3 \end{pmatrix}$$

so that J_4 is given by (4.6) as claimed.

Remark:

For general m , the density is proportional to $J_m^{-\frac{1}{2}}$ with J_m having i, j entry $i \cdot j p_{i+j-2}$, $1 \leq i, j \leq m$. The following algorithm combines the ideas above to give a reversible Markov chain for sampling from the conditional distribution of the model 4.3 on the manifold \mathcal{M}_p . From $\mathbf{x} \in \mathcal{M}_p$

- (a) Choose five coordinates uniformly at random. Without loss, suppose these are the first five, calculate $\bar{p}_i = \sum_{j=1}^5 x_j^i$, $1 \leq i \leq 4$.

- (b) Pick a small parameter ϵ , then choose one of the five coordinates uniformly at random without loss, suppose the first coordinate has been chosen. Change x_1 to $y_1 = x_1 + \epsilon_1$ with ϵ_1 chosen uniformly in $[x_1 - \epsilon, x_1 + \epsilon]$. Solve for y_2, y_3, y_4, y_5 so that $\mathbf{y} = (y_1, y_2, y_3, y_4, y_5) \in \mathcal{M}_{\bar{p}}$ as in (4.1).
- (c) Calculate $J_4(\mathbf{x}), J_4(\mathbf{y})$ from Proposition 3 above. If $J_4(\mathbf{x}) \geq J_4(\mathbf{y})$ the algorithm moves to \mathbf{y} . If $J_4(\mathbf{x}) < J_4(\mathbf{y})$ flip a coin with success probability

$$\sqrt{\frac{J_4(\mathbf{x})}{J_4(\mathbf{y})}}$$

If success move to \mathbf{y} , otherwise stay at \mathbf{x}

Remarks:

- 1 For $m \leq 4$, calculations for solving the \mathbf{y} can be done in closed form as they involve at most quartic equations. For higher m a variety of numerical procedures are available.
- 2 Of course, if \mathbf{y} in step (b) is outside $[0, 1]^5$, the algorithm stays at \mathbf{x} .
- 3 We began studying the problem hoping to parametrize the curve (4.1) and sample directly from the arc length measure. This proved impractical. The technique we have developed seems easier and is applicable to general continuous exponential families.

4.3. Ergodicity

Let $P_j(x) = x_1^j + \dots + x_n^j$ and S be the set defined by

$$0 < x_1 < \dots < x_n < 1, \quad P_1(x) = c_1, \dots, P_4(x) = c_4. \quad (4.7)$$

The closure of S will be denoted by \bar{S} . We also assume that $1 \geq c_1 > c_2 > c_3 > c_4 > 0$ which is a necessary condition for the existence of a solution to (4.7). Assume that the system (4.7) has a solution.

Lemma 2. *Let $y \in S$ be a solution to (4.7). Then there is a submanifold of dimension $n - 4$ passing through $y \in S$. Furthermore the orthogonal projection of S near y on any coordinate line x_j contains a neighborhood of y_j .*

Proof. We have $dP_j(x) = jx_1^{j-1}dx_1 + \dots + jx_n^{j-1}dx_n$. Therefore to show the first assertion it suffices to show that the matrix

$$\begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \\ x_1^2 & x_2^2 & \dots & x_n^2 \\ x_1^3 & x_2^3 & \dots & x_n^3 \end{pmatrix}$$

has rank 4 which is immediate. The second assertion follows from the fact that the locally the system can be solved near y as function of any $n - 4$ coordinates. \square

Lemma 3. *Let $k \geq 6$ and $y \in S$ be a solution to (4.7). Consider the solution of the system (4.7) subject to the additional requirements*

$$x_j = y_j, \quad \text{for } j \geq k.$$

Then there is a submanifold of dimension $k - 5$ of solutions passing through y . For $k = 6$ the solution is a curve and its projection on the coordinate line x_i , $1 \leq i \leq 5$ contains a neighborhood of y_i .

Proof. We look at the differentials dP_j , $j = 1, 2, 3, 4$ and dx_j , $j \geq k$. To prove the first assertion it suffices to show that the $(n - k + 5) \times n$ matrix

$$\begin{pmatrix} 1 & 1 & \dots & \dots & 1 & & \\ x_1 & x_2 & \dots & \dots & \dots & x_n & \\ x_1^2 & x_2^2 & \dots & \dots & \dots & x_n^2 & \\ x_1^3 & x_2^3 & \dots & \dots & \dots & x_n^3 & \\ 0 & 0 & \dots & 1 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 0 & \dots & 1 \end{pmatrix}$$

has rank $n - k + 5$ which is obvious. The second assertion follows from the fact we can solve for $n - 1$ coordinates in terms of any one of x_i 's for $i = 1, 2, 3, 4$. \square

Let $M \subset S$ be a connected component of S . We consider the following process in M . Given that the process is at $y = (y_1, \dots, y_n) \in M$, one chooses five coordinates i_1, \dots, i_5 and the process can move to any point along the curve defined by

$$x_j = y_j, \quad \text{for } j \neq i_1, \dots, i_5.$$

The question we want to answer is whether any two points in M *communicate* in the sense that one can move from one to the other in a finite number of iterations. More technically, we say two points, y, z are sufficiently close if, given $y \in M$ there is $\delta > 0$ such that if z is within δ of y , then one can move from y to z in finite number of iterations. The positive number δ may depend on y .

Lemma 4. *If two points $y, z \in M$ are sufficiently close then they communicate.*

Proof. We do induction on n . The case $n = 5$ is clear. Let $i_1, \dots, i_5 = 1, \dots, 5$ and $k = 6$ in the notation of Lemma 3. Then the determinant of the matrix in the proof of Lemma 3 is

$$\pm \prod_{i < j < 6} (x_i - x_j).$$

Therefore if y and z sufficiently close then one can move from (y_1, \dots, y_n) to a point $(z_1, y'_2, \dots, y'_5, y_6, \dots, y_n)$ by the second assertion of Lemma 3. Now the induction hypothesis applies to complete the proof. \square

Proposition 4. *Any pair of points in M communicate.*

Proof. Starting at $y \in M$ we show that the set of points in M that can be reached in a finite number of steps from y is both open and closed in M . The required result then follows from connectedness of M . From Lemma 4 it follows that the set of points that can be reached from y in finitely many iterations is open. To show closed-ness let $y = y^{(1)}, y^{(2)}, \dots$ be a sequence of

points each of which can be reached in finitely many steps from y and assume $y^{(m)} \rightarrow z \in M$. Then for all m sufficiently large the point $y^{(m)}$ lies in a sufficiently small neighborhood of z and Lemma 4 is applicable to show that z can be reached in finitely many steps from such $y^{(m)}$ proving the required ‘closed-ness’. \square

Let S' be the set defined by

$$0 \leq x_1, \dots, x_n \leq 1, \quad P_1(x) = c_1, \dots, P_4(x) = c_4, \quad (4.8)$$

and M' be a connected component of S' . We consider the process in M' where in addition we allow any permutation of the coordinates as well as evolution described in M .

Proposition 5. *Any pair of points in M' communicate.*

Proof. For points away from the set V consisting of the boundary hyperplanes of the unit cube in \mathbb{R}^n and the generalized diagonal $\bigcup_{i \neq j} \{x_i = x_j\}$ the assertion follows from Proposition 4. Applying the Curve Selection Lemma (see for example Milnor (1968)) we move away from V in one step, and then Proposition 4 is applicable. \square

4.4. Valid tests and connectedness

For many applications of the present techniques, it is only a conjecture that the algorithms are ergodic. Consider the manifold $\mathcal{M}_{\mathbf{p}}$ above based on the first four sample moments. Choosing 5 coordinates and sampling from the correct conditional distribution on the resulting curve gives a way of moving around on $\mathcal{M}_{\mathbf{p}}$. However it has not been proved that this algorithm is connected; Indeed Proposition 5 of section 4.3 only shows that the algorithm goes between points in the same connected component (in the topological sense) in finitely many steps.

Bormeshenko (2009) gave a difficult proof that the analogous problem based on changing 3 coordinates on the manifold determined by the sum and the sum of squares is connected and we certainly conjecture this for any number of moments.

If these samples are used for goodness of fit test, there is a valid test procedure available, even in the absence of connectedness, by adapting an idea of Besag and Clifford (1989).

The idea is simple. Let \mathcal{X} be the original data. This gives rise to a point x_0^* on $\mathcal{M}_{\mathbf{p}}$. Suppose $K(x, dy)$ is a Markov chain with the correct stationary distribution on the connected component containing x_0^* . Fix a number of steps T^* and run this chain T^* steps to get y^* say. Then run the time reversed chain, starting at y^* for T^* steps and independently repeat this B^* times (starting at y^* each time). This results in $(x_1^*, x_1^*, \dots, x_{B^*}^*) \in \mathcal{M}_{\mathbf{p}}$. The $B^* + 1$ values $(x_0^*, x_1^*, x_1^*, \dots, x_{B^*}^*)$ are exchangeable, so the relative position of any test statistic $s(x_0^*)$ among $s(x_i^*)$ is uniform under the null hypothesis. If $s(x_0^*)$ is among the extreme values of these statistics then a valid rejection is possible.

Acknowledgements

We thank Hans Andersen, Olena Bormishenko, Greg Brumfiel, Brian White, Leonid Pekelis and an anonymous referee for help with this paper.

References

- ANDERSEN, H. and DIACONIS, P. (2008). Hit and run as a unifying device. *Journal de la SSF* 5–28.
- BARNDORFF-NIELSEN, O. (1978). *Information and exponential families in statistical theory*. Wiley, NY.
- BARTON, D. (1953). On Neyman’s test of goodness of fit and its power with respect to a particular system of alternatives. *Skand. Aktuar.*, **36** 24–63.
- BARTON, D. (1956). Neyman’s ψ_k^2 test of goodness of fit when the null hypothesis is composite. *Skand. Aktuar.*, **39** 216–46.
- BÉLISLE, C., ROMELJN, H. and SMITH, R. (1993). Hit-and-run algorithms for generating multivariate distributions. *Mathematics of Operations Research* 255–266.
- BERAN, R. (1979). Exponential models for directional data. *The Annals of Statistics* 1162–1178.
- BESAG, J. and CLIFFORD, P. (1989). Generalized Monte Carlo significance tests. *Biometrika*, **76** 633–642.
- BHATTACHARYA, A. and BHATTACHARYA, R. (2012). *Nonparametric Inference On Manifolds With Applications To Shape Spaces*. IMS, Cambridge University Press, Cambridge, UK.
- BHATTACHARYA, R. and PATRANGENARU, V. (2003). Large sample theory of intrinsic and extrinsic sample means on manifolds. I. *Annals of Statistics* 1–29.
- BOENDER, C., CARON, R., McDONALD, J., KAN, A., ROMELJN, H., SMITH, R., TELGEN, J. and VORST, A. (1991). Shake-and-bake algorithms for generating uniform points on the boundary of bounded polyhedra. *Operations research* 945–954.
- BORMESHENKO, O. (2009). Walking around by three flipping. Unpublished manuscript.
- BROWN, L. D. (1986). *Fundamentals of statistical exponential families: with applications in statistical decision theory*. Institute of Mathematical Statistics, Hayworth, CA, USA.
- CARLSSON, E., CARLSSON, G. and DE SILVA, V. (2006). An algebraic topological method for feature identification. *Internat. J. Comput. Geom. Appl.*, **16** 291–314.
- CICCOTTI, G. and RYCKAERT, J. P. (1986). Molecular dynamics simulation of rigid molecules. Tech. rep.
- COMETS, F., POPOV, S., SCHÜTZ, G. and VACHKOVSKAIA, M. (2009). Billiards in a general domain with random reflections. *Archive for rational mechanics and analysis*, **191** 497–537.
- DAVID, F. (1939). On Neyman’s ”smooth” test for goodness of fit i. distribution of the criterion ψ^2 when the hypothesis tested is true. *Biometrika*, **31** 191–199.
- DIACONIS, P. (1988). Sufficiency as statistical symmetry. In *Proceedings of the AMS Centennial Symposium*. Amer. Math. Soc., Providence, RI, 15–26.
- DIACONIS, P. and HOLMES, S. (1994). Gray codes for randomization procedures. *Statistics and Computing* 287–302.
- DIACONIS, P., KHARE, K. and SALOFF-COSTE, L. (2010a). Gibbs sampling, conjugate priors and coupling. *Sankhya*, **72** 136–169.
- DIACONIS, P., LEBEAU, G. and MICHEL, L. (2010b). Geometric analysis for the metropolis algorithm on lipschitz domains. *Inventiones Mathematicae* 1–43.
- DIACONIS, P. and SALOFF-COSTE, L. (1998). What do we know about the metropolis algorithm? *Journal of Computer and System Sciences*, **57** 20–36.
- DIACONIS, P. and SHAHSHAHANI, M. (1986). On square roots of the uniform distribution on compact groups. *Proc. Amer. Math. Soc.*, **98** 341–348.
- DIACONIS, P. and STURMFELS, B. (1998). Algebraic algorithms for sampling from conditional

- distributions. *Ann. Statist.*, **26** 363–397.
- DRTON, M., STURMFELS, B. and SULLIVANT, S. (2009). *Lectures on algebraic statistics*. Birkhauser.
- EATON, M. (1983). *Multivariate statistics: a vector space approach*. Wiley, New York.
- FAN, J. (1996). Test of significance based on wavelet thresholding and neyman’s truncation. *Journal of the American Statistical Association* 674–688.
- FEDERER, H. (1996). *Geometric Measure Theory*. Springer, Berlin.
- FISHER, N., LEWIS, T. and EMBLETON, B. (1993). *Statistical analysis of spherical data*. Cambridge Univ Pr.
- FIXMAN, M. (1974). Classical statistical mechanics of constraints: A theorem and application to polymers. *Proc Natl Acad Sci U S A*, **71** 3050–3053.
- GINÉ, E. (1975). Invariant tests for uniformity on compact Riemannian manifolds based on Sobolev norms. *Ann. Statist.*, **3** 1243–1266.
- GOLDMAN, N. and WHELAN, S. (2000). Statistical tests of gamma-distributed rate heterogeneity in models of sequence evolution in phylogenetics. *Mol Biol Evol*, **17** 975–8.
- HAMMERSLEY, J. M. and HANDSCOMB, D. C. (1964). *Monte Carlo methods*. Methuen, London.
- HIPP, C. (1974). Sufficient statistics and exponential families. *The Annals of Statistics* 1283–1292.
- HUBBARD, J. H. and HUBBARD, B. B. (2007). *Vector calculus, linear algebra, and differential forms : a unified approach*. 3rd ed. Matrix Editions, Ithaca, NY.
- KALLIORAS, A., KOUTROUVELIS, I. and CANAVOS, G. (2006). Testing the fit of gamma distributions using the empirical moment generating function. *Communications in Statistics—Theory and Methods*, **35** 527–540.
- KRANTZ, S. and PARKS, H. (2008). *Geometric Integration Theory*. Birkhauser.
- LALLEY, S. and ROBBINS, H. (1987). Asymptotically minimax stochastic search strategies in the plane. *Proceedings of the National Academy of Sciences*, **84** 2111–2112.
- LEBEAU, G. and MICHEL, L. (2010). Semi-classical analysis of a random walk on a manifold. *The Annals of Probability*, **38** 277–315.
- LEHMANN, E. and ROMANO, J. (2005). *Testing statistical hypotheses*. Springer Verlag.
- LETAC, G. (1992). *Lectures on natural exponential families and their variance functions*. Conselho Nacional de Desenvolvimento Científico e Tecnológico, Instituto de Matemática Pura e Aplicada.
- LINDQVIST, B. and TARALDSEN, G. (2005). Monte carlo conditioning on a sufficient statistic. *Biometrika*, **92** 451–464.
- LINDQVIST, B. and TARALDSEN, G. (2006). Conditional monte carlo based on sufficient statistics with applications. *Festschrift Doksum*.
- LIU, J. S. (2001). *Monte Carlo strategies in scientific computing*. Springer, New York.
- MATTILA, P. (1999). *Geometry of Sets and Measures in Euclidean Spaces: Fractals and Rectifiability*. Cambridge studies in advanced mathematics, Cambridge University Press.
- MEZZADRI, F. (2007). How to generate random matrices from the classical compact groups. *Notices Amer. Math. Soc.*, **54** 592–604.
- MILNOR, J. (1968). *Singular points of complex hypersurfaces*. 61, Princeton Univ Press.
- MORGAN, F. (2009). *Geometric measure theory: a beginner’s guide*. 3rd Ed, Academic Press.
- NARAYANAN, H. and NIYOGI, P. (2008). Sampling hypersurfaces through diffusion. *Approximation, Randomization and Combinatorial Optimization. Algorithms and Techniques* 535–548.
- NEYMAN, J. (1937). “smooth” test for goodness of fit. *Skand Aktuariotioskr*, **20** 149–199.

- PENNEC, X. (2006). Intrinsic statistics on Riemannian manifolds: Basic tools for geometric measurements. *J. Math. Imaging Vis.*, **25** 127–154.
- PETTITT, A. (1978). Generalized Cramer-von Mises statistics for the gamma distribution. *Biometrika*, **65** 232–5.
- TJUR, T. (1974). *Conditional probability distributions*. Lecture notes - Institute of Mathematical Statistics, University of Copenhagen ; 2, Institute of Mathematical Statistics, University of Copenhagen, Copenhagen.
- WATSON, G. (1983). *Statistics on Spheres*, vol. 6. Wiley-Interscience.
- YANG, Z. (2006). *Computational molecular evolution*. Oxford University Press, Oxford.