

Dual Geometries in Statistics

Paul Vos
East Carolina University
vosp@mail.ecu.edu

29 July 2004

[Home Page](#)

[Title Page](#)

[Contents](#)



Page 1 of 46

[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

Goals

My Primary Goal: By describing some important statistical results from theoretical and applied statistics using ideas and concepts from dual geometries you will have some understanding of the important role these geometries play, and perhaps, motivate some of you to study these geometries further.

Another Goal: Introduce some basic terminology of (dual) differential geometry. [Terms]

Home Page

Title Page

Contents



Page 2 of 46

Go Back

Full Screen

Close

Quit

Outline

- I. Relative Information Loss
- II. MLE and Sufficiency
- III. Generalized Linear Models
- Summary

Home Page

Title Page

Contents



Page 3 of 46

Go Back

Full Screen

Close

Quit

Curved Exponential Family (CEF)

Consider a regular exponential family of order k ($y \in \mathcal{Y} \subset \mathbb{R}^k$) having density

$$p(y | \eta) = \exp\{y'\eta - \psi(\eta)\}h(y)$$

A subset

$$\mathcal{M} = \{p : p(y | \eta(\theta)) \mid \theta \in \Theta\} \subset \mathcal{S} = \{p(y | \eta)\}$$

is a (one-parameter) curved exponential family (CEF) if-f the map $\theta \mapsto \eta(\theta)$ having domain $\Theta \subset \mathbb{R}^1$ is an imbedding. [Manifold]

Home Page

Title Page

Contents



Page 4 of 46

Go Back

Full Screen

Close

Quit

Relative Information Loss

Limiting Relative Info Loss = $\lim I(\theta)^{-1}[nI(\theta) - I^T(\theta)]$

Result: Let $\bar{y} = \frac{1}{n}(y_1 + \dots + y_n) \in \mathcal{Y} \subset R^k$ have density in a CEF \mathcal{M} and let $T = T(\bar{y})$ be a statistic. **Then**

$$\text{Rel. Info. Loss } (T \text{ for } \mathcal{M}) = \gamma^2 + \frac{1}{2}\beta^2$$

where

γ^2 is the statistical curvature of \mathcal{M}
 β^2 is the curvature **for** T

Home Page

Title Page

Contents

◀◀ ▶▶

◀ ▶

Page 5 of 46

Go Back

Full Screen

Close

Quit

Comments

$$\text{Rel. Info. Loss } (T \text{ for } \mathcal{M}) = \gamma^2 + \frac{1}{2}\beta^2$$

- Duality of γ^2 and β^2
- Parameter Invariance of γ^2 and β^2
- Efficient Notation (especially multivariate generalization)
- Geometric Intuition

[Home Page](#)

[Title Page](#)

[Contents](#)



Page 6 of 46

[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

Bib. Notes: Relative Information Loss

- Statistical geometry has been called Fisher-Efron-Amari theory.
- Fisher (1922, 1925)
- Efron (1975, 1978)
- Amari (1990)

[Home Page](#)

[Title Page](#)

[Contents](#)



Page 7 of 46

[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

Log Likelihood Function

Defn: Let \mathcal{M} be a set of densities with common support $\mathcal{Y} \subset R^k$ and let $y \in \mathcal{Y}$. The **log likelihood** is the function $\ell_y : \mathcal{M} \mapsto R^1$ defined by $\ell_y(p) = \log(p(y))$.

Exp Fam: $\ell_y(p) = y'\eta - \psi(\eta) + \log(h(y))$

Statistical Importance:

- Fundamental inference procedure
- Generates superior inference procedures.

Home Page

Title Page

Contents

◀◀ ▶▶

◀ ▶

Page 8 of 46

Go Back

Full Screen

Close

Quit

Maximum Likelihood

Defn: If $\hat{p} \in \mathcal{M}$ such that $\ell_y(\hat{p}) = \max_{p \in \mathcal{M}} \ell_y(p)$ then \hat{p} is a **maximum likelihood estimate** (mle) of \mathcal{M} (based on y).

Note 1: The mle is usually expressed in terms of some parameterization (co-ordinate system). Eg, $\hat{\theta} = \theta(\hat{p})$ or $\hat{\beta} = \beta(\hat{p})$.

Properties of $\hat{\theta}$ are expressed by considering $\hat{\theta}$ as a random variable. (Bias, variance, and exact or asymptotic distribution.)

Note 2: Many times the parameterization is not arbitrary (θ could be a mean or a slope).

Home Page

Title Page

Contents

◀◀ ▶▶

◀ ▶

Page 9 of 46

Go Back

Full Screen

Close

Quit

Maximum Likelihood (cont)

Geometry: For CEF \mathcal{M} ,

$$\arg \max_{p \in \mathcal{M}} \ell_y(p) = \arg \min_{p \in \mathcal{M}} D(p_y, p)$$

Maximum Likelihood = Minimum Divergence

Likelihood \longleftrightarrow Divergence

Statistical Concepts \longleftrightarrow Geometric Quantities

[Home Page](#)

[Title Page](#)

[Contents](#)



Page 10 of 46

[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

Divergence

Defn: Let Ξ be an open subset of R^k and let $\xi_1, \xi_2 \in \Xi$. A **divergence on Ξ** , denoted by $D_{\Xi}(\xi_1, \xi_2)$, is a smooth function taking values in $[0, \infty]$ such that for any $\xi_1, \xi_2 \in \Xi$ the following hold:

1. $D_{\Xi}(\xi_1, \xi_2) \geq 0$ with equality holding if-f $\xi_1 = \xi_2$.
2. $\frac{\partial}{\partial \xi_1^r} D_{\Xi}(\xi_1, \xi_2) \big|_{\xi_1 = \xi_2} = \frac{\partial}{\partial \xi_2^r} D_{\Xi}(\xi_1, \xi_2) \big|_{\xi_1 = \xi_2} = 0$
3. $g_{rs} = \frac{\partial^2}{\partial \xi_1^r \partial \xi_1^s} D_{\Xi}(\xi_1, \xi_2)$ is positive definite and a function of ξ_1 alone.

Home Page

Title Page

Contents



Page 11 of 46

Go Back

Full Screen

Close

Quit

Divergence (cont)

Defn: For a **space** of points \mathcal{S} , $D : \mathcal{S} \times \mathcal{S} \mapsto R^1$ is a **divergence on \mathcal{S}** if there exists a parameterization $f : \mathcal{S} \mapsto \Xi$ such that $D(f^{-1}(\xi_1), f^{-1}(\xi_2))$ is a divergence on Ξ .

Note 1: The parameterization exists when \mathcal{S} is globally **flat** with respect to a **nonmetric connection**.

connection \leftrightarrow geometry [**Connection**] [**Metric**]

metric connection \leftrightarrow “straight lines” minimize distance

flat \leftrightarrow “nice geometry” (R^n)

Note 2: \mathcal{S} globally **flat** in another **nonmetric connection**.

Note 3: **Dual** connections: 1-connection – 1-connection.

Home Page

Title Page

Contents

◀◀ ▶▶

◀ ▶

Page 12 of 46

Go Back

Full Screen

Close

Quit

Divergence (Dual) Geometry

- $D(p_1, p_2)$ behaves like a squared distance
- D is not symmetric
- Two Geometries used to describe D
- Pythagorean Relationship:

$$D(p_y, p) = D(p_y, \hat{p}) + D(\hat{p}, p)$$

Home Page

Title Page

Contents



Page 13 of 46

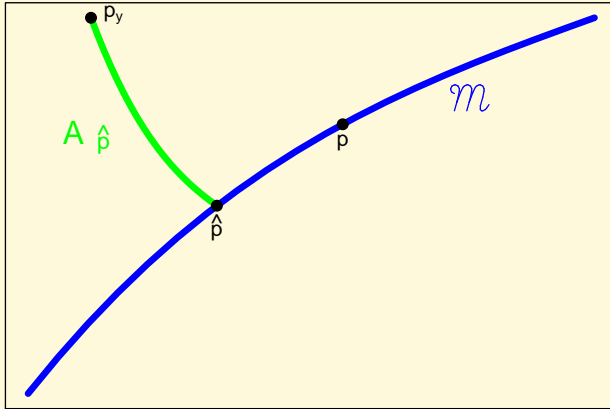
Go Back

Full Screen

Close

Quit

Divergence Geometry (cont)



Home Page

Title Page

Contents



Page 14 of 46

Go Back

Full Screen

Close

Quit

Law of Cosines

Triangle: having sides of length a, b, c and angle C opposite side of length c .

$$c^2 = a^2 + b^2 - 2ab \cos C$$

Squared Distances in R^k : $p_1, p_2, p_3 \in R^k$, so

$$c^2 = \| p_1 - p_3 \|^2$$

$$a^2 = \| p_1 - p_2 \|^2$$

$$b^2 = \| p_2 - p_3 \|^2$$

$$C = \text{angle between } p_3 - p_2 \text{ and } p_1 - p_2$$

Home Page

Title Page

Contents

◀◀ ▶▶

◀ ▶

Page 15 of 46

Go Back

Full Screen

Close

Quit

$$\begin{aligned}
\| p_1 - p_3 \|^2 &= \| p_1 - p_2 \|^2 + \| p_2 - p_3 \|^2 \\
&\quad - 2 \| p_1 - p_2 \| \cdot \| p_2 - p_3 \| \cos \theta \\
&= \| p_1 - p_2 \|^2 + \| p_2 - p_3 \|^2 - 2 \langle p_1 - p_2, p_2 - p_3 \rangle
\end{aligned}$$

Setting $\| p - q \|^2 = 2D(p, q)$ and $p - q = v(q, p)$, we have

$$\begin{aligned}
2D(p_1, p_3) &= 2D(p_1, p_2) + 2D(p_2, p_3) - 2 \langle v(p_2, p_1), v(p_2, p_3) \rangle \\
\text{iff } D(p_1, p_3) &= D(p_1, p_2) + D(p_2, p_3) - \langle v(p_2, p_1), v(p_2, p_3) \rangle
\end{aligned}$$

Divergences: For points $p_1, p_2, p_3 \in \mathcal{S}$,

$$D(p_1, p_3) = D(p_1, p_2) + D(p_2, p_3) - \langle v(p_2, p_1), v^*(p_2, p_3) \rangle_{p_2}$$

Home Page

Title Page

Contents



Page 16 of 46

Go Back

Full Screen

Close

Quit

Divergence in Statistics

- Log Likelihood Function
(properties of estimators)
- Quasi-likelihood Function
(weaken distributional assumptions)
- Kullback-Leibler Information (Distance) (Divergence for Exp Fam, approximating Exp Fam)

[Home Page](#)

[Title Page](#)

[Contents](#)



Page 17 of 46

[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

Sufficiency

Defn: Let \mathcal{M} be a family densities and let Y be an observation from $p \in \mathcal{M}$. A statistic $T(Y)$ is sufficient for \mathcal{M} if the conditional distribution of Y given $T = t$ is the same for all $p \in \mathcal{M}$.

Factorization Theorem: T is sufficient for \mathcal{M} if-f there exists functions u_1 and u_2 such that

$$\log(p(y)) = u_1(y) + u_2(p, T(y))$$

Note 1: For CEF \mathcal{M} , $\log(p(y)) = C(y) - D(p_y, p)$.

Note 2: When the Pythagorean Relationship holds $\log(p(y)) = C(y) - D(p_y, \hat{p}) - D(\hat{p}, p)$.

Home Page

Title Page

Contents

◀◀ ▶▶

◀ ▶

Page 18 of 46

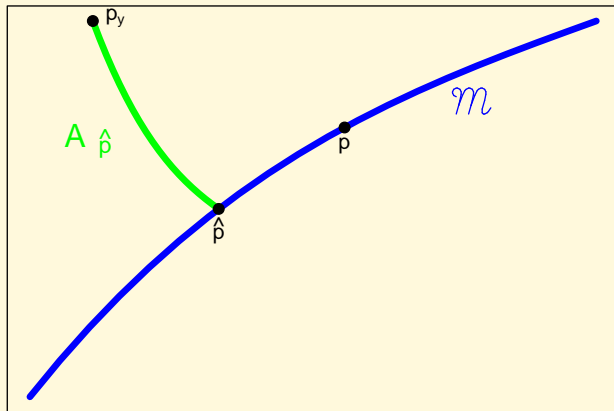
Go Back

Full Screen

Close

Quit

Sufficiency (cont)



$$\begin{aligned}\log(p(y)) &= C(y) - D(p_y, \hat{p}) - D(\hat{p}, p) \\ &= u_1(y) + u_2(p, T(y)), T(y) = \hat{p}\end{aligned}$$

Home Page

Title Page

Contents

◀◀

▶▶

◀

▶

Page 19 of 46

Go Back

Full Screen

Close

Quit

Sufficiency (cont)

- Pythagorean Relationship requires 3 conditions
 - 1-geodesic connecting \hat{p} to p
 - -1 -geodesic connecting p_y to \hat{p}
 - Right angle between these geodesics
- Geometry of Maximum Likelihood Estimation
- Geometry of Exponential Families
- Information Loss and Recovery

Home Page

Title Page

Contents

◀◀ ▶▶

◀ ▶

Page 20 of 46

Go Back

Full Screen

Close

Quit

Bib. Notes: MLE and Sufficiency

- Ideas related to MLE (divergence): contrast functions (Eguchi, 1983), Csiszar's (1967) f -divergence; minimum divergence estimation (Vos, 1997; Kass and Vos, 1997).
- Ideas Related to sufficiency: information recovery (Efron and Hinkley, 1978; Amari, 1990; Skovgaard, 1985) and higher order efficiency (Eguchi, 1983, 1985; Amari, 1990).
- Global Information Recovery: Marriott and Vos, 2004.

Home Page

Title Page

Contents



Page 21 of 46

Go Back

Full Screen

Close

Quit

Global Information Recovery

Helix Model: $X \sim N(\mu, I_{3 \times 3})$ where $\mu = (r \cos \theta, r \sin \theta, \theta d)'$, r and d are fixed.

Information in a statistic T : Ability to approximate the likelihood function using a function that depends on X only through T . **Information Recovery:** $T = (\hat{\theta}, A)$, find A .

Principles of (Local) Information Recovery

- **Minimize Expected Info Loss:** minimize $I^X - I^T$
- **Observed Fisher Information:** $A = I_{obs}$
- **Ancillarity:** Choose A having distn not a function of θ

Home Page

Title Page

Contents



Page 22 of 46

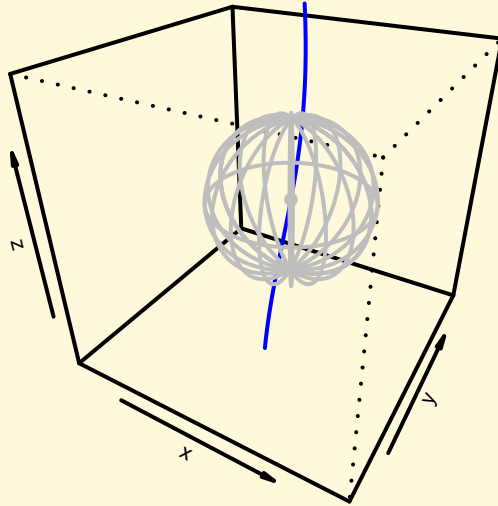
Go Back

Full Screen

Close

Quit

Helix Model: $r = 0.5$, $d = 3$



[Home Page](#)

[Title Page](#)

[Contents](#)



Page 23 of 46

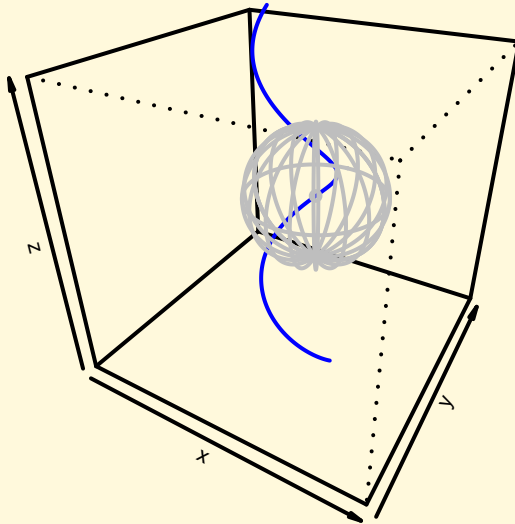
[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

Helix Model: $r = 1, d = 1$



Home Page

Title Page

Contents



Page 24 of 46

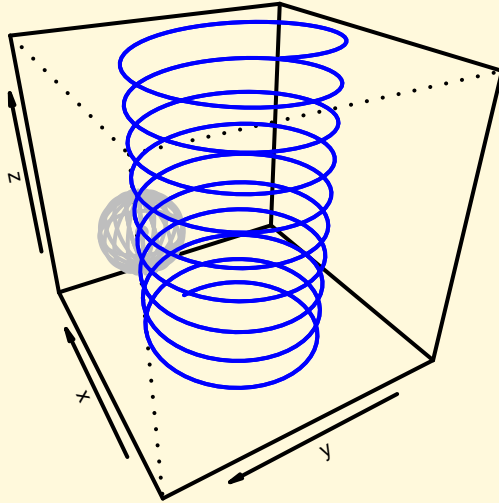
Go Back

Full Screen

Close

Quit

Helix Model: $r = 5$, $d = 0.3$



Home Page

Title Page

Contents



Page 25 of 46

Go Back

Full Screen

Close

Quit

Remarks on Helix Models/Info Recovery

- Local methods of information recovery can fail.
- Global method (spectral decomposition) does not fail: indicates number of required “information directions”.
- Explains why local methods (asymptotics) work.

[Home Page](#)

[Title Page](#)

[Contents](#)



Page 26 of 46

[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

Generalized Linear Models

- Model the relationship between a response variable Y and one or more explanatory variables x_1, \dots, x_k .
- Example Y is college freshman GPA and x is SAT score.
- Model $Y | x$; see graph.
 - Mean $(Y | x) = \mu(x) = \beta_1 x + \beta_0$
 - Var $(Y | x) = \sigma^2(\mu(x))$
 - Shape $(Y | x)$: exponential family or not specified

Home Page

Title Page

Contents



Page 27 of 46

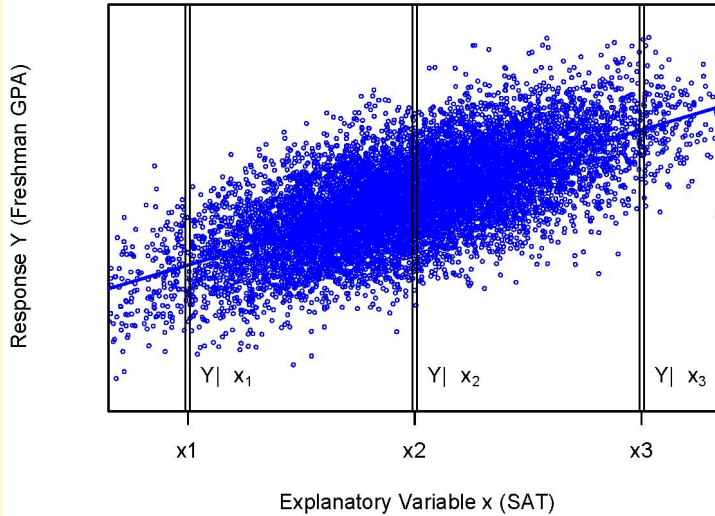
Go Back

Full Screen

Close

Quit

Population View of Linear Regression



[Home Page](#)

[Title Page](#)

[Contents](#)



Page 28 of 46

[Go Back](#)

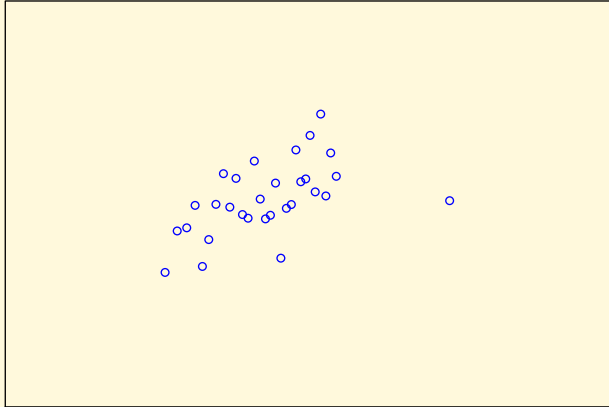
[Full Screen](#)

[Close](#)

[Quit](#)

Sample View of Linear Regression (n=30)

Response Y (Freshman GPA)



Explanatory Variable x (SAT)

[Home Page](#)

[Title Page](#)

[Contents](#)



Page 29 of 46

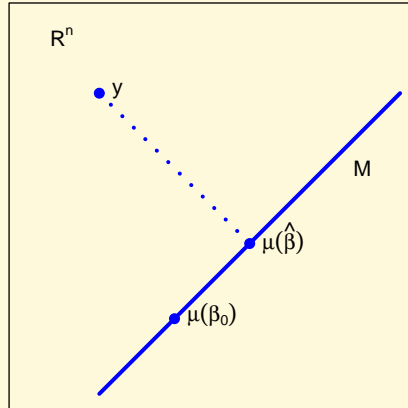
[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

Model View of Linear Regression



[Home Page](#)

[Title Page](#)

[Contents](#)



Page 30 of 46

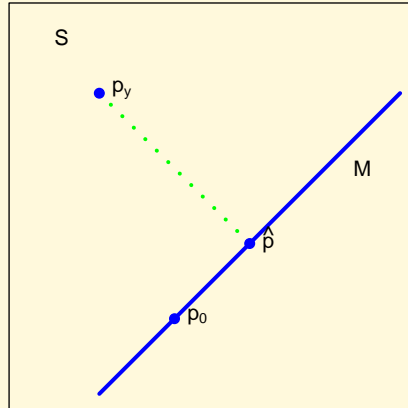
[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

Model View of a GLM



Home Page

Title Page

Contents



Page 31 of 46

Go Back

Full Screen

Close

Quit

Geometry of GLMs

- Structure
 - Dual Geometries
 - \mathcal{S} is of dimension n
 - $\mathcal{M} \subset \mathcal{S}$ has dimension $k + 1$
 - \mathcal{M} is often flat in +1 connection
 - Riemannian metric often diagonal in divergence parameter
- Role
 - Estimation Algorithm
 - Model checking

Home Page

Title Page

Contents



Page 32 of 46

Go Back

Full Screen

Close

Quit

Summary

- Parameter Invariance
- Economical Notation
- Intuition
- Importance of Pythagorean Relationship

Home Page

Title Page

Contents



Page 33 of 46

Go Back

Full Screen

Close

Quit

Geometry Concepts Summary

Terminology	Concept
manifold	set endowed with geometric structure
metric	orthogonality
angle	departures from ortho.
α -connection	straight lines
α -curvature	departures from strght. lines
divergence	“distance” of interest
Pythagorean Relationship (Law of Cosines)	Relationship among divergence, strght lines, and ortho. (angle)

Home Page

Title Page

Contents



Page 34 of 46

Go Back

Full Screen

Close

Quit

References

- Amari, S.I.** (1990) *Differential-Geometrical Methods in Statistics*, Springer, New York.
- Csiszar, I.** (1967) On topological properties of f -divergence. *Studia Sci. Math. Hungar.*, **2**, 329-339.
- Efron, B.** (1975) Defining the curvature of a statistical problem, *Ann. Statist.*, **3**, 1189-1217.
- Efron, B.** (1978) The geometry of exponential families, *Ann. Statist.*, **6**, 362-376.

[Home Page](#)

[Title Page](#)

[Contents](#)



Page 35 of 46

[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

Efron, B., and Hinkley, D.V. (1978) Assessing the accuracy of the maximum likelihood estimator: observed versus expected Fisher information. *Biometrika*, **65**, 457-481.

Eguchi, S. (1983) Second order efficiency of minimum contrast estimators in a curved exponential family. *Ann. Statist.*, **11**, 793-803.

Eguchi, S. (1985) A differential geometric approach to statistical inference on the basis of contrast functionals. *Hiroshima Math. J.*, **15**, 341-391.

Fisher, R.A. (1922) On the mathematical foundations of theoretical statistics. *Philos. Trans. Royal Soc. London, Ser. A*, **222** , 309-368.

[Home Page](#)

[Title Page](#)

[Contents](#)

[◀◀](#) [▶▶](#)

[◀](#) [▶](#)

Page 36 of 46

[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

Fisher, R.A. (1925) Theory of statistical estimation. *Proc. Camb. Phil Soc.*, **22**, 700-725.

Kass, R.E, and Vos, P.W. (1997) *Geometrical Foundations of Asymptotic Inference*, John Wiley & Sons, New York.

Marriott, P. and Vos, P.W. (2004) On the global geometry of parametric models and information recovery, *Bernoulli*, **10** (2), 1-11.

Skovgaard, I.M. (1985) A second-order investigation of asymptotic ancillarity. *Ann. Statist.*, **13**, 534-551.

Vos, P.W. (1997) Minimum divergence estimation. In *Encyclopedia of Statistical Sciences Update II* (ed. S. Kotz and C. Read) Wiley, New York.

[Home Page](#)

[Title Page](#)

[Contents](#)



Page **37** of **46**

[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

Manifolds

Consider a two binomial problem with n_1 and n_2 fixed. There is a **set** of densities:

$$\mathcal{S} = \left\{ f(x_1, x_2) : f = \binom{n_1}{x_1} p_1^{x_1} (1 - p_1)^{n_1 - x_1} \times \binom{n_2}{x_2} p_2^{x_2} (1 - p_2)^{n_2 - x_2}, p_1, p_2 \in (0, 1) \right\}$$

In fact, \mathcal{S} can be given additional structure: when (p_1, p_2) is close to (p'_1, p'_2) the models (densities) are not too different – statistically, a **family**; mathematically a **topological manifold**.

Home Page

Title Page

Contents



Page 38 of 46

Go Back

Full Screen

Close

Quit

NOTE: There are other parameterizations besides $p = (p_1, p_2)'$, but the properties of closeness still hold.

- The hypothesis $H_0 : p_1 - p_2 = .2$ defines a subfamily (imbedded submanifold)

$$\mathcal{M} = \left\{ f \in \mathcal{S} : p_1 - p_2 = .2 \right\}$$

- The parameterization $p : \mathcal{M} \mapsto R^2$ defines an image, $p(\mathcal{M})$, in the real plane. ($p(\mathcal{S})$ is the open unit square.)
- Another parameterization, say the log odds $\theta = \log(p/(1-p))$, defines another image, $\theta(\mathcal{M})$. ($\theta(\mathcal{S})$ is R^2 .)

Home Page

Title Page

Contents



Page 39 of 46

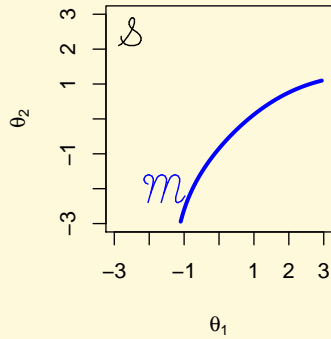
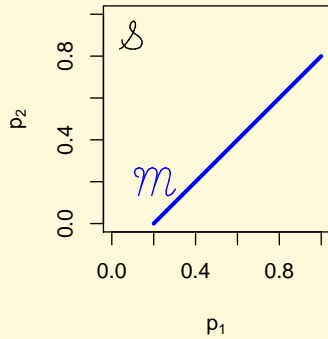
Go Back

Full Screen

Close

Quit

Two Visualizations



Home Page

Title Page

Contents



Page 40 of 46

Go Back

Full Screen

Close

Quit

Comments on Visualizations

- Different parameterizations provide different pictures – different geometries – different connections.

$p(\mathcal{M})$ is a straight line (geodesic); $\theta(\mathcal{M})$ is not.

- There is no one *correct* geometry.
(straight line \neq curve of shortest length)

- Geometries judged on their utility.

For studying sufficiency, should \mathcal{M} be considered a straight or curved line?

Home Page

Title Page

Contents



Page 41 of 46

Go Back

Full Screen

Close

Quit

- A statistical manifold is a topological manifold (parameterization into the reals) endowed with the geometric properties of
 - straight lines – departures measured by curvature
 - orthogonality – departures measured by angle
 - divergence (replaces the notion of length)

RETURN

Home Page

Title Page

Contents



Page 42 of 46

Go Back

Full Screen

Close

Quit

Connections

- **Concept:** straight lines are ones whose tangent vector does not change direction (or length).

In order compare tangent vectors along a curve, there needs to be a specified mapping or connection between tangent spaces – this is provided by the connection.

- **Definition:** based on derivatives of the log likelihood or divergence functions.
- **Notation:** $\nabla_{X_p}^\alpha Y, \Gamma_{jk}^{(\alpha)i}; \nabla_X^\alpha Y(p), \Gamma_{jk}^{(\alpha)i}(p)$

$$\nabla_{\partial_j}^\alpha \partial_k = \sum_i \Gamma_{jk}^{(\alpha)i} \partial_i = \Gamma_{jk}^{(\alpha)i} \partial_i \text{ where } \partial_i = \frac{\partial}{\partial \xi^i}$$

Home Page

Title Page

Contents



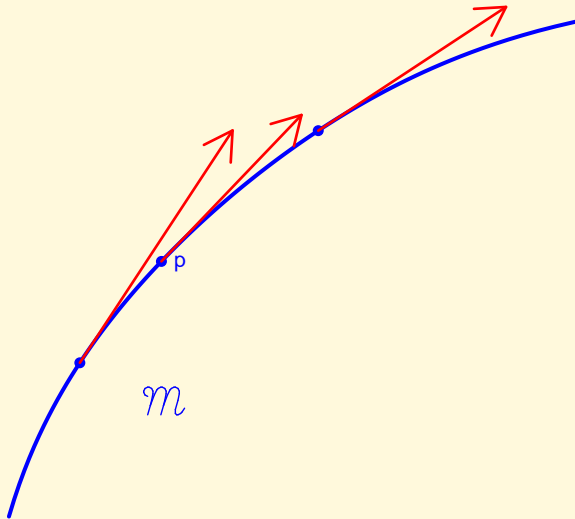
Page 43 of 46

Go Back

Full Screen

Close

Quit



Home Page

Title Page

Contents



Page 44 of 46

Go Back

Full Screen

Close

Quit

Metric

- **Concept:** Metric defines an inner product which in turn provides the idea of orthogonal vectors (and curves). Also gives lengths of vectors.

Metric is actually a (smooth) collection of inner products – depending on p – that gives lengths of curves.

However, we do not use an α family of metrics.

- **Definition:** based on derivatives of the log likelihood (Fisher information) or divergence functions.

[Home Page](#)

[Title Page](#)

[Contents](#)



Page 45 of 46

[Go Back](#)

[Full Screen](#)

[Close](#)

[Quit](#)

- **Notation:** $\langle X_p, Y_p \rangle_p, \langle X, Y \rangle, g_{ij}, g^{ij}; g_{ij}(p), g^{ij}(p)$

$$\langle \partial_i, \partial_j \rangle = g_{ij} \text{ where } \partial_i = \frac{\partial}{\partial \xi^i}$$

$$\langle \partial^i, \partial^j \rangle = g^{ij} \text{ where } \partial^i = \frac{\partial}{\partial \xi_i^*}$$

Or, matrix $[g^{ij}]$ is the matrix inverse of $[g_{ij}]$ – there is no ambiguity because of the properties of the dual parameter ξ^* if it exists.

Home Page

Title Page

Contents



Page 46 of 46

Go Back

Full Screen

Close

Quit