

2013年統計数理研究所夏期大学院
統計数理研究所 セミナー室 1
9月26日 14時40分～17時50分

ISM Summer School 2013

Information Geometry

Shinto Eguchi

The Institute of Statistical Mathematics, Japan

Email: eguchi@ism.ac.jp, komori@ism.ac.jp

Url: <http://www.ism.ac.jp/~eguchi/>



Outline



Information geometry

historical remarks

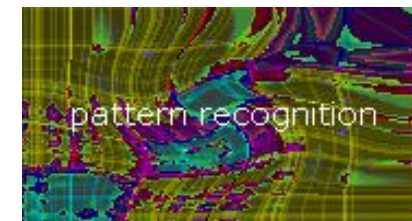
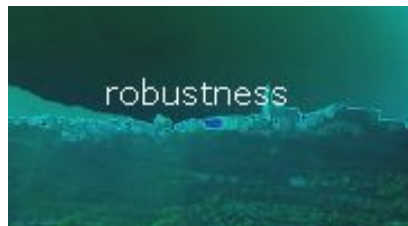
treasure example

Information divergence class

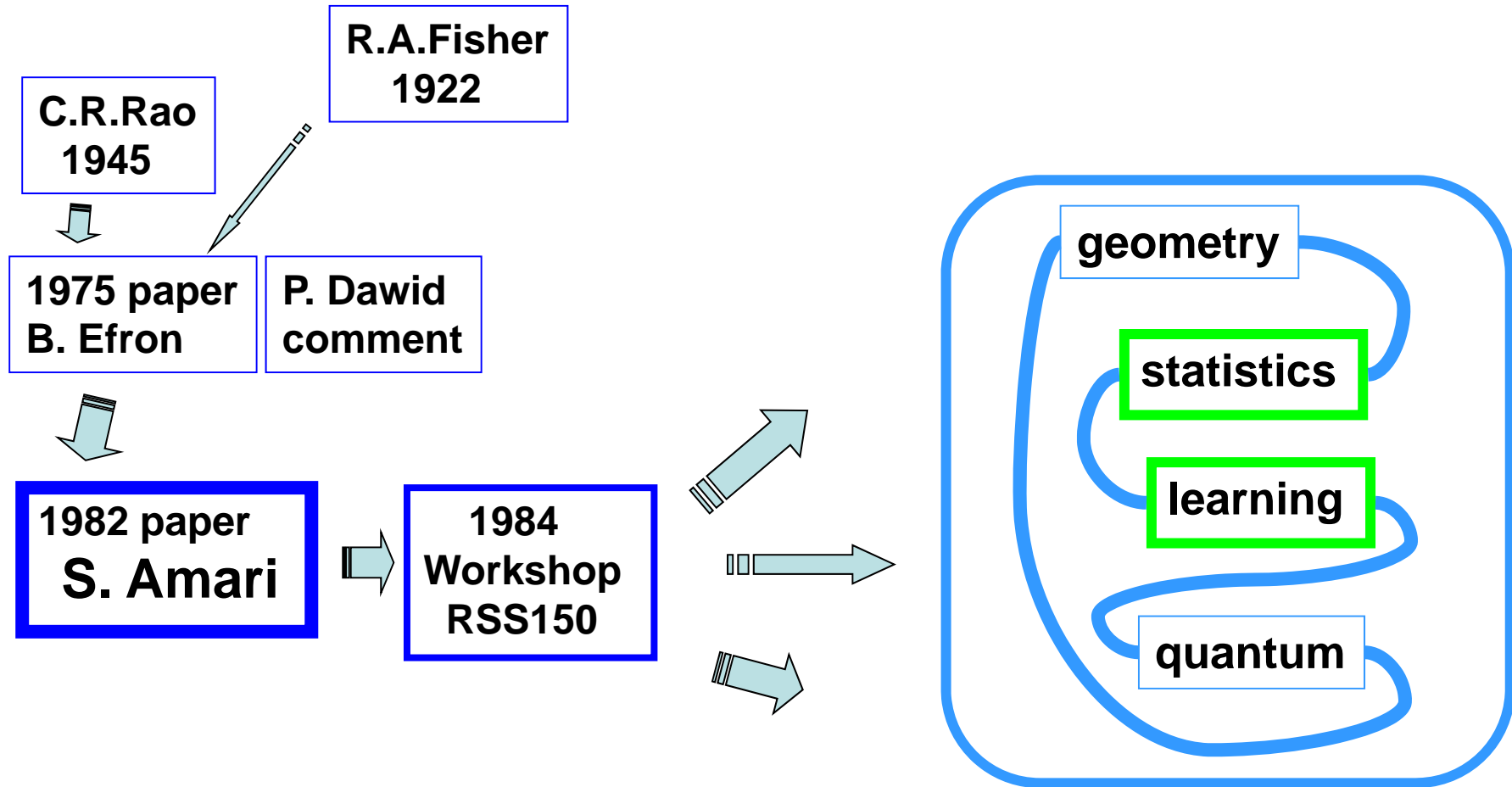
robust statistical methods

robust ICA

robust kernel PCA



Historical comment



Dual Riemannian Geometry

Dual Riemannian geometry gives reformulation for Fisher's foundation

Cf. Einstein field equations (1916)

Information Geometry aims to geometrize

A dualistic structure between modeling and estimation

Cf. Erlangen program (Klein, 1872)

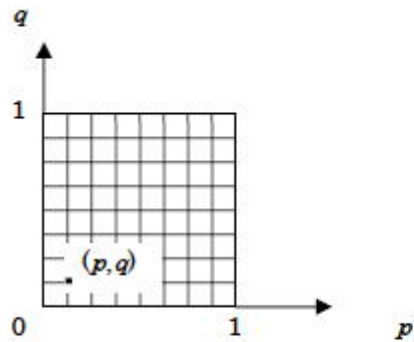
http://en.wikipedia.org/wiki/Erlangen_program

Estimation is projection of data onto model.

“Estimation is an action by projection and model is an object to be projected”

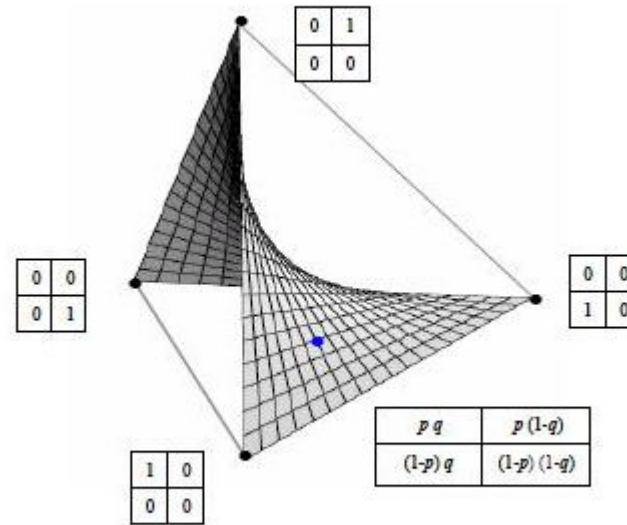
The interplay of action and object is elucidated by a geometric structure.

2 x 2 table



	$p q$	$p (1-q)$
	$(1-p)q$	$(1-p) (1-q)$

two-by-two table



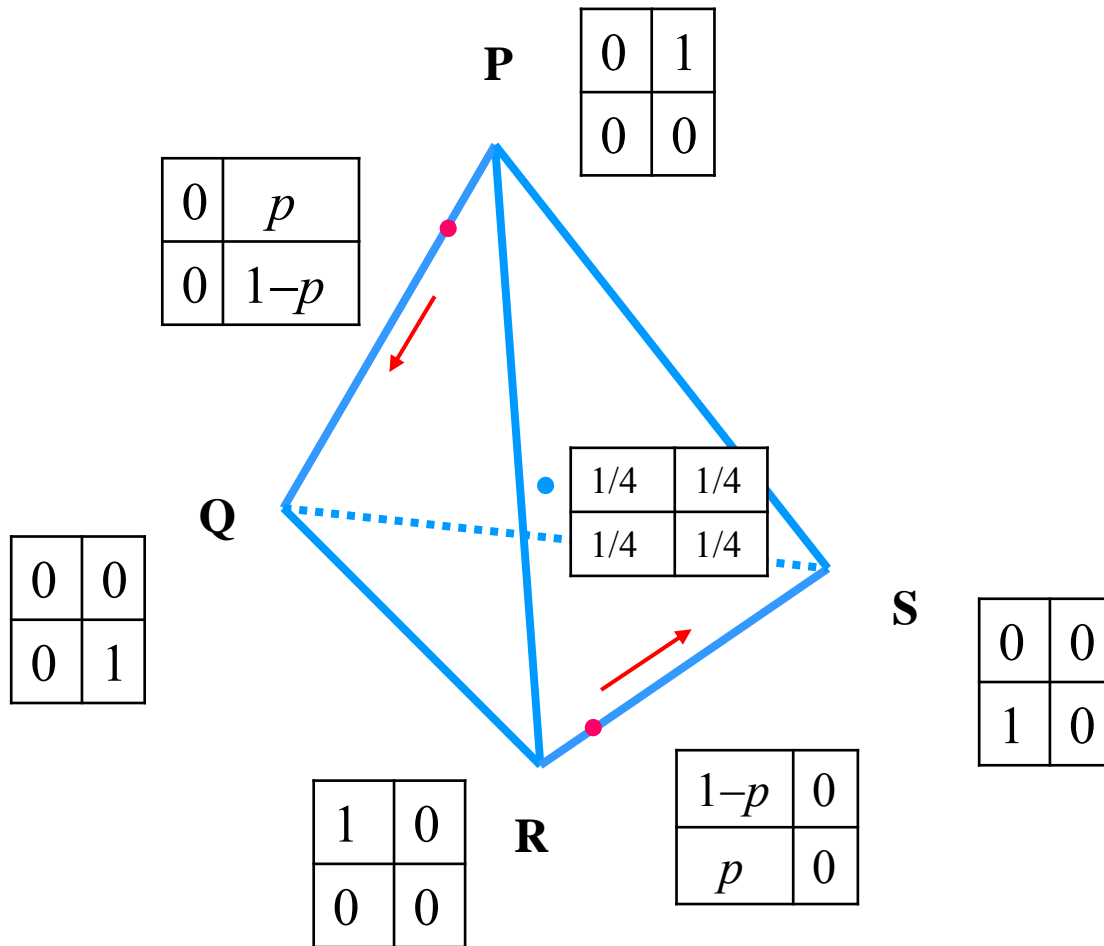
Independent 2x2 table

The space of all 2×2 tables associates with a regular tetrahedron

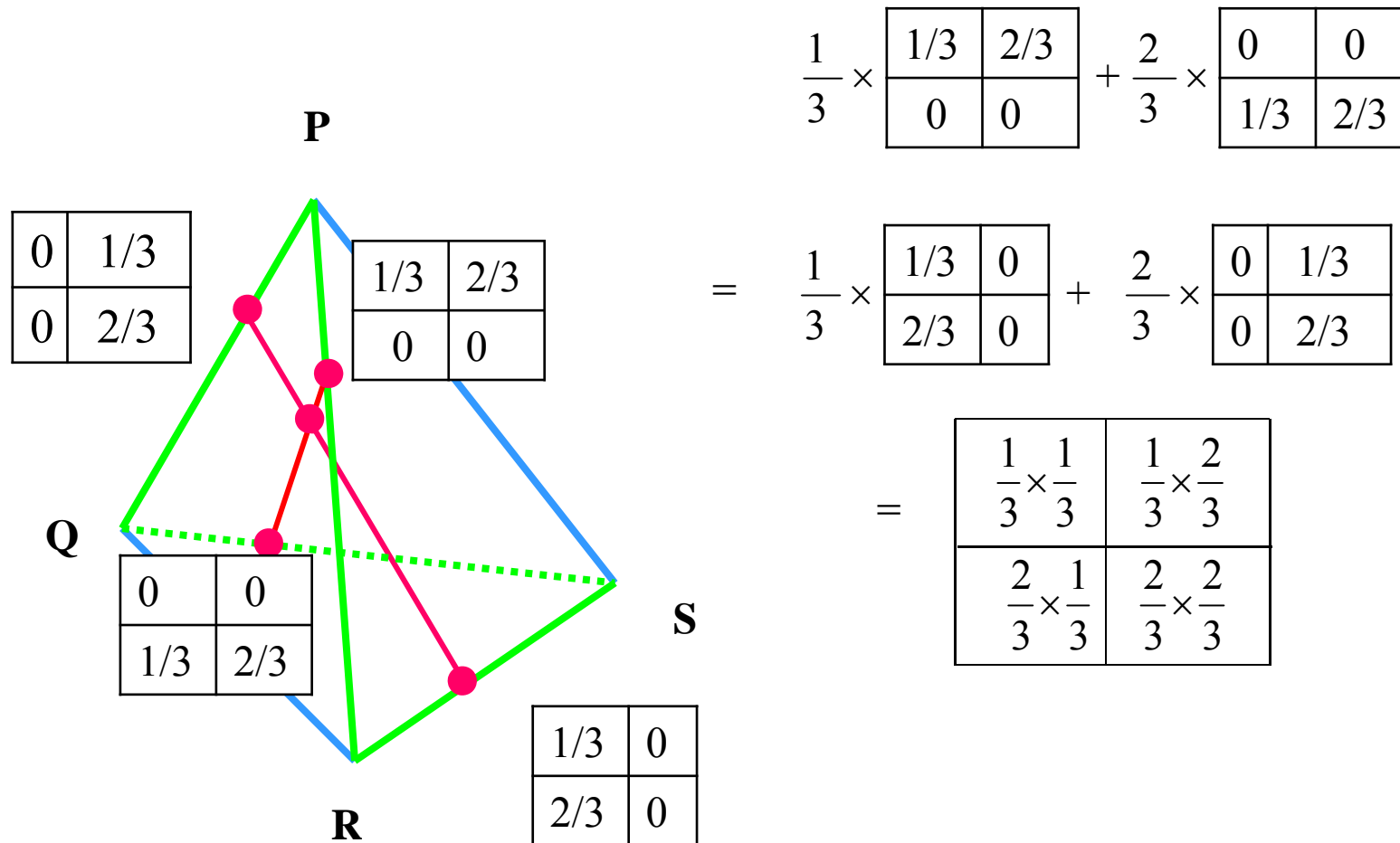
**The space of all independent 2×2 tables associates with a ruled surface
In the regular tetrahedron**

**We know the ruled surface is exponential geodesic, but
does anyone know the minimality of the ruled surface?**

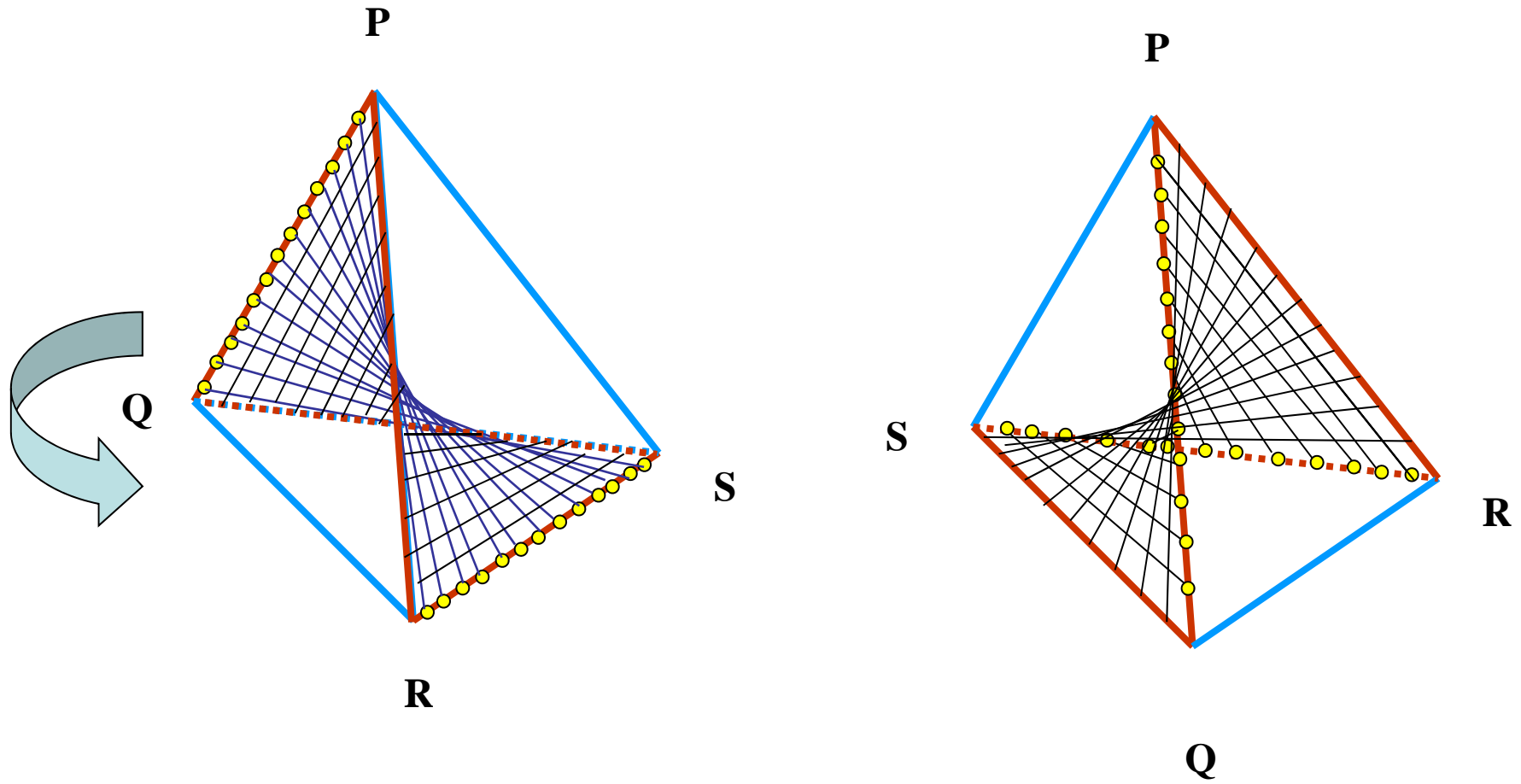
regular tetrahedron

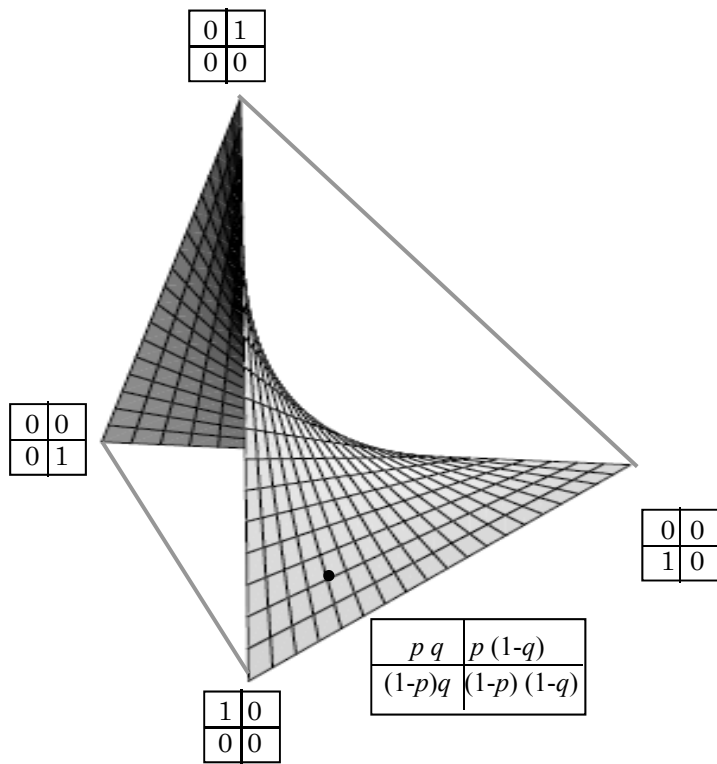


regular tetrahedron



Ruled surface





	A	B	
C	x	y	e
D	z	w	f
	g	h	n

$$e = x + y, \quad f = z + w,$$

$$g = x + z, \quad h = y + w$$

$$n = x + y + z + w$$

	A	B	
C	pq	$p(1-q)$	P
D	$(1-p)q$	$(1-p)(1-q)$	$1-p$
	q	$1-q$	1

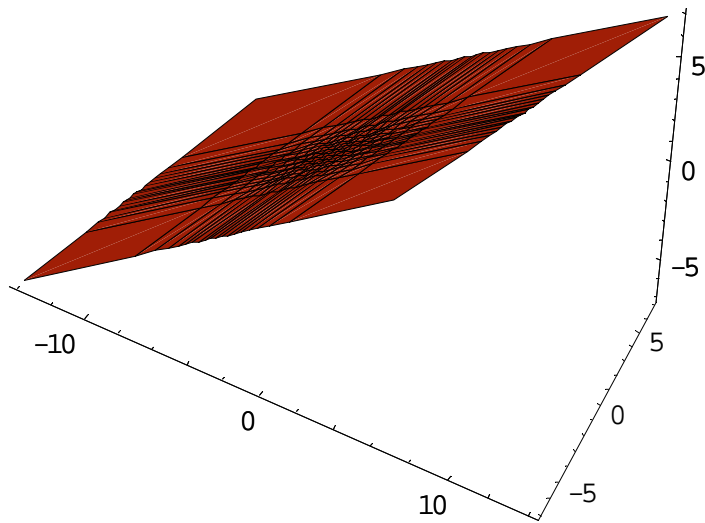
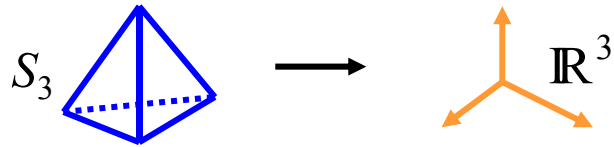
$$\chi^2 = n \frac{(xw - yz)^2}{efgh}$$

$$= n \frac{\{pq(1-p)(1-q) - p(1-q)(1-p)q\}^2}{p(1-p)q(1-q)}$$

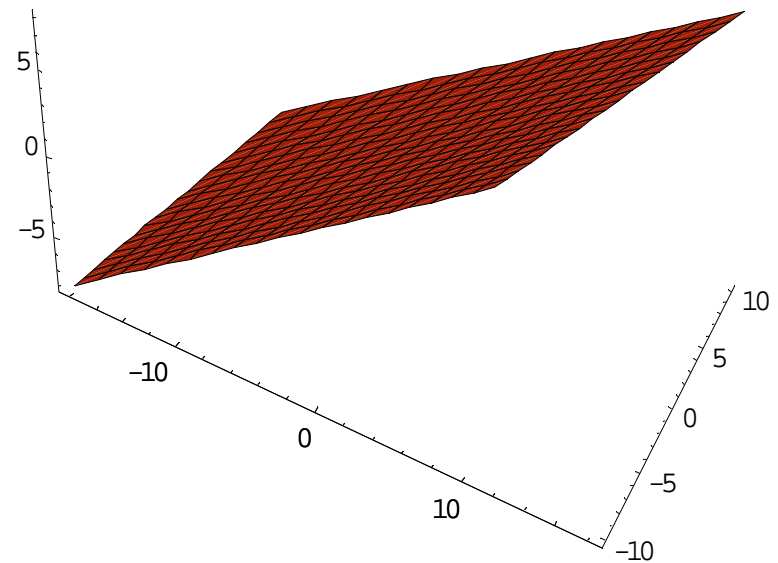
$$= 0$$

e-geodesical

$$(\pi_{11}, \pi_{12}, \pi_{21}) \longrightarrow \left(\log \frac{\pi_{11}}{\pi_{22}}, \log \frac{\pi_{12}}{\pi_{22}}, \log \frac{\pi_{21}}{\pi_{22}} \right) \text{ where } \pi_{22} = 1 - \pi_{11} - \pi_{12} - \pi_{21}$$

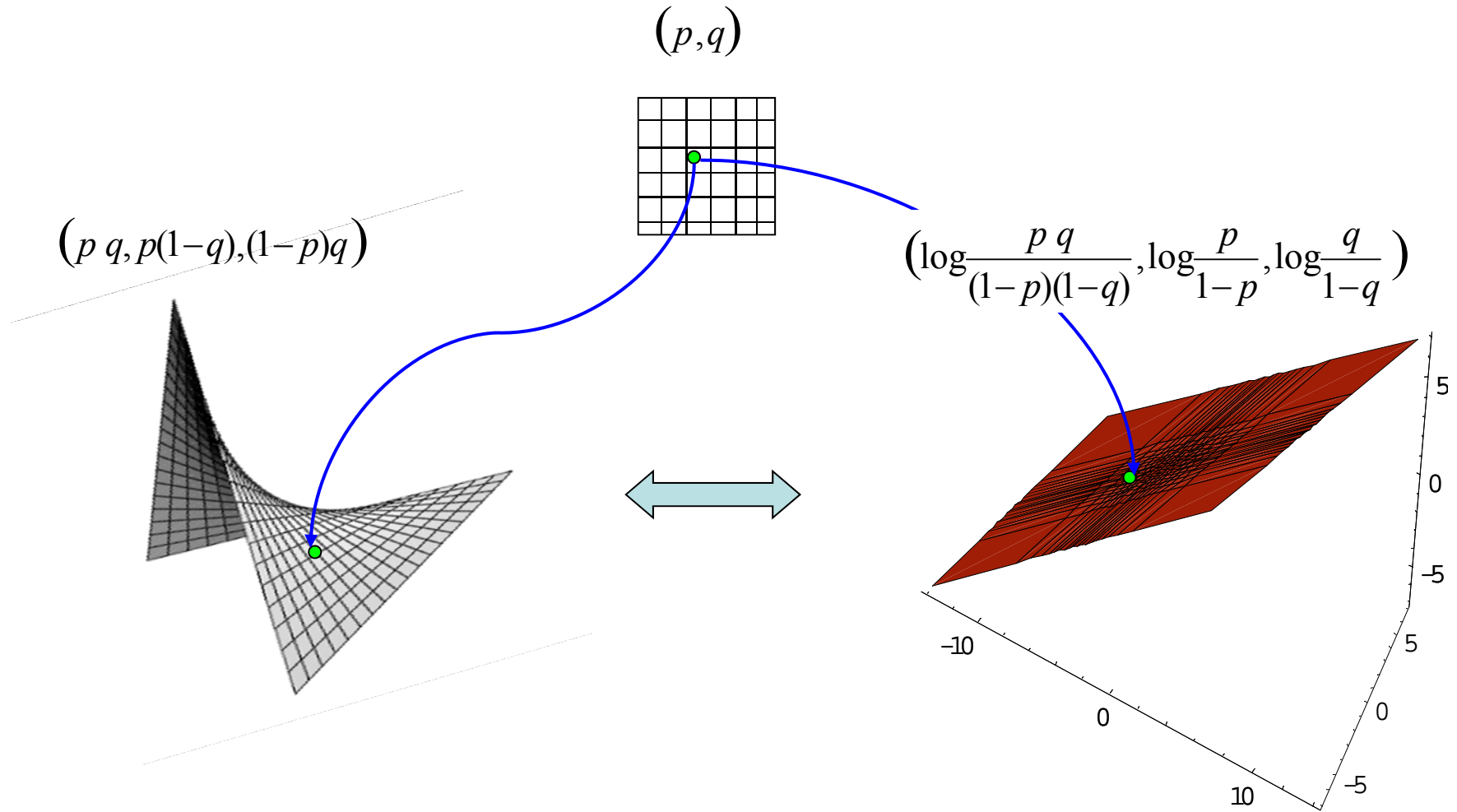


$$\left\{ \left(\log \frac{pq}{(1-p)(1-q)}, \log \frac{p}{1-p}, \log \frac{q}{1-q} \right), 0 < p < 1, 0 < q < 1 \right\}$$



$$\{(x+y, x, y), x \in \mathbb{R}, y \in \mathbb{R}\}$$

Two parametrization

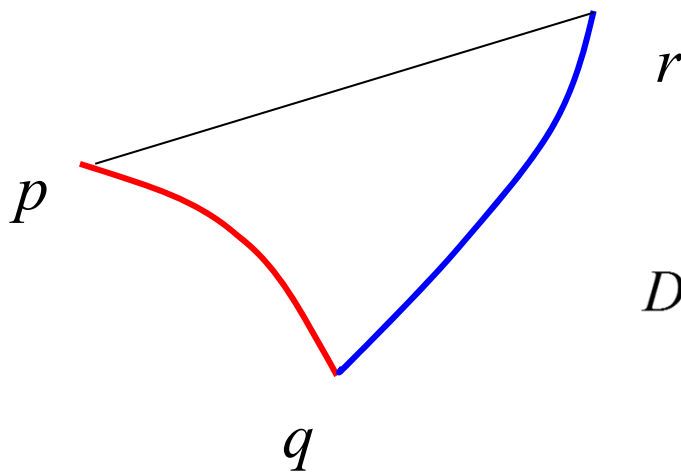


Kullback-Leibler Divergence

Let $p(\mathbf{x})$ and $q(\mathbf{x})$ be probability density functions.

Then **Kullback-Leibler Divergence** is defined by

$$D_{\text{KL}}(p, q) = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$$



$$D_{\text{KL}}(p, r) = D_{\text{KL}}(p, q) + D_{\text{KL}}(q, r)$$

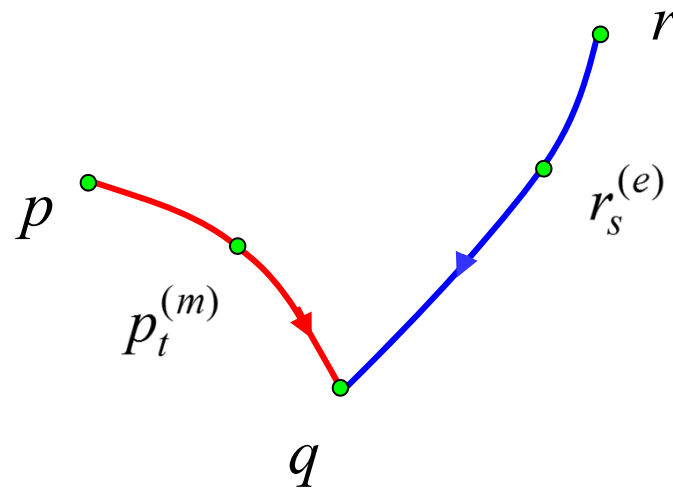
Two one-parameter families

Let \mathbb{P} be the space of all pdfs on a data space.

m-geodesic $p_t^{(m)}(\mathbf{x}) = (1-t)p(\mathbf{x}) + tq(\mathbf{x}), \quad (p, q \in \mathbb{P})$

e-geodesic $r_s^{(e)}(\mathbf{x}) = c_s \{r(\mathbf{x})\}^{1-s} \{q(\mathbf{x})\}^s, \quad (q, r \in \mathbb{P})$

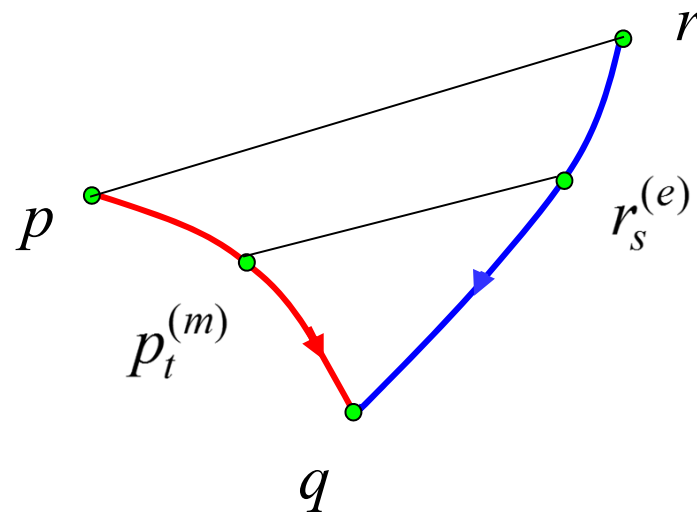
ここで $c_s = 1 / \int \{r(\mathbf{x})\}^{1-s} \{q(\mathbf{x})\}^s d\mathbf{x}$



Pythagoras theorem Amari-Nagaoka (1982)

$$D(p,r) = D(p,q) + D(q,r) \quad \Rightarrow$$

$$D(p_t^{(m)}, r_s^{(e)}) = D(p_t^{(m)}, q) + D(q, r_s^{(e)}) \quad (\forall (s,t) \in [0,1] \times [0,1])$$



Proof

$$\begin{aligned} & D_{\text{KL}}(p_t^{(m)}, r_s^{(e)}) - \{D_{\text{KL}}(p_t^{(m)}, q) + D_{\text{KL}}(q, r_s^{(e)})\} \\ &= \int (p_t^{(m)}(\mathbf{x}) - q(\mathbf{x}))(\log q(\mathbf{x}) - \log r_s^{(e)}(\mathbf{x}))d\mathbf{x} \\ &= \int (1-t)(p(\mathbf{x}) - q(\mathbf{x}))\{(1-s)(\log q(\mathbf{x}) - \log r(\mathbf{x}) - \log c_s)\}d\mathbf{x} \\ &= (1-t)(1-s) \int (p(\mathbf{x}) - q(\mathbf{x}))(\log q(\mathbf{x}) - \log r(\mathbf{x}))d\mathbf{x} \\ &\quad - (1-t)\log c_s \int (p(\mathbf{x}) - q(\mathbf{x}))d\mathbf{x} \\ &= (1-t)(1-s)\{D_{\text{KL}}(p, r) - D_{\text{KL}}(p, q) - D_{\text{KL}}(q, r)\} \\ &= 0 \end{aligned}$$

□

ABC in differential geometry

Riemannian metric defines an inner product on any tangent space of a manifold.

$$(M, g) \quad g : \mathbf{X}(M) \times \mathbf{X}(M) \rightarrow \mathbb{R}$$

Geodesic = a minimal arc between x and y

$$\gamma = \operatorname{argmin}_{c=\{x(t):x(0)=x,x(1)=y\}_{t=0}^1} \int_0^1 \sqrt{g(\dot{x}(t), \dot{x}(t))} dt$$

Linear connection defines parallelism along a vector field.

$$\nabla : \mathbf{X}(M) \times \mathbf{X}(M) \rightarrow \mathbf{X}(M)$$

$$(1) \quad \nabla_{fX} Y = f \nabla_X Y,$$

$$(2) \quad \nabla_X (fY) = f \nabla_X Y + (Xf)Y \quad (\forall f \in \mathbf{F}(M), \forall X, Y \in \mathbf{X}(M))$$

$$\text{Componetwise} \quad \nabla_{X_i} X_j = \sum_{k=1}^d \Gamma_{ij}^k X_k$$

Cf. slide 41

Geodesic

A one-parameter family (curve) $C = \{ \boldsymbol{\theta}(t) : -\varepsilon \leq t \leq \varepsilon \}$ is called geodesic with respect to a linear connection $\{ \Gamma_{jk}^i(\boldsymbol{\theta}) : 1 \leq i, j, k \leq p \}$

$$\ddot{\theta}^i(t) + \sum_{k=1}^p \sum_{j=1}^p \Gamma_{jk}^i(\boldsymbol{\theta}(t)) \dot{\theta}^k(t) \dot{\theta}^j(t) = 0 \quad (\forall i = 1, \dots, p) \quad \text{speed of acceleration}$$

Remark 2:

If $\Gamma_{jk}^i(\boldsymbol{\theta}) = 0$ ($1 \leq i, j, k \leq p$), any geodesic is a line. Newton's law of inertia

The property is not invariant with parametrization.

The geodesic C is expressed by a transform rule from parameter $\boldsymbol{\theta}$ to $\boldsymbol{\omega}$

$$\ddot{\omega}^a(t) + \sum_{c=1}^p \sum_{b=1}^p \tilde{\Gamma}_{cb}^a(\boldsymbol{\omega}(t)) \dot{\omega}^b(t) \dot{\omega}^c(t) = 0 \quad (\forall a = 1, \dots, p)$$

$$\tilde{\Gamma}_{cb}^a(\boldsymbol{\omega}) = \sum_{k=1}^p \sum_{j=1}^p \sum_{i=1}^p \bar{B}_c^k \bar{B}_b^j B_i^a \Gamma_{jk}^i(\boldsymbol{\theta}) + \sum_{a=1}^p B_i^a \frac{\partial B_b^i}{\partial \omega^c}, \quad B_i^a = \frac{\partial \omega^a}{\partial \theta^i}, \quad \bar{B}_a^i = \frac{\partial \theta^i}{\partial \omega^a} \text{ for } \omega^a = t^a(\boldsymbol{\theta})$$

Change rule on connections

Let us consider a transform ϕ from parameter θ to ω . Then a geodesic C is written by $\{\omega(t) = \phi(\theta(t)) : -\varepsilon < t < \varepsilon\}$. Hence we have the following:

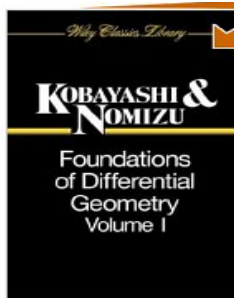
$$\dot{\omega}^a(t) = \sum_{i=1}^p B_i^a(\theta(t)) \dot{\theta}^i(t) \quad (\forall a = 1, \dots, p) \quad \left(B_i^a = \frac{\partial \phi^a}{\partial \theta^i} \right)$$

$$\ddot{\omega}^a(t) = \sum_{i=1}^p B_i^a(\theta(t)) \ddot{\theta}^i(t) + \sum_{j=1}^p \sum_{i=1}^p \frac{\partial B_i^a(\theta(t))}{\partial \theta^j} \dot{\theta}^i(t) \dot{\theta}^j(t) \quad (\forall a = 1, \dots, p)$$

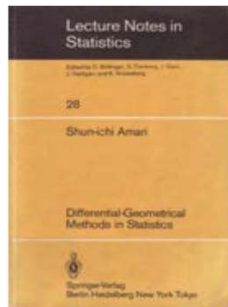
$$\ddot{\omega}^a(t) + \sum_{c=1}^p \sum_{b=1}^p \tilde{\Gamma}_{cb}^a(\omega(t)) \dot{\omega}^b(t) \dot{\omega}^c(t) = 0 \quad (\forall a = 1, \dots, p)$$

$$\tilde{\Gamma}_{cb}^a(\omega) = \sum_{k=1}^p \sum_{j=1}^p \sum_{i=1}^p \bar{B}_c^k \bar{B}_b^j B_i^a \Gamma_{jk}^i(\theta) + \sum_{a=1}^p B_i^a \frac{\partial B_b^i}{\partial \omega^c}, \quad \bar{B}_a^i = \frac{\partial (\phi^{-1})^i}{\partial \omega^a} \text{ for } \omega^a = \phi^a(\theta)$$

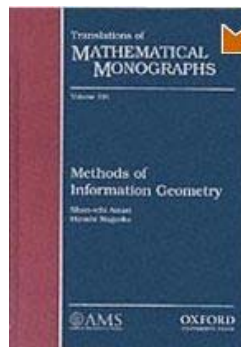
What are reference books?



Foundations of Differential Geometry (Wiley Classics Library)
Shoshichi Kobayashi, Katsumi Nomizu



Differential Geometrical Methods in Statistics
Shun-Ichi Amari 1985 年 Springer



Methods of Information Geometry
Shun-Ichi Amari , Hiroshi Nagaoka
Amer Mathematical Society (2001)

Statistical model and information

Statistical model $M = \{p_{\theta}(\mathbf{x}) = p(\mathbf{x}, \theta) : \theta \in \Theta\}$ ($\Theta \subset \mathbb{R}^p$)

Score vector $s(\mathbf{x}, \theta) = \frac{\partial}{\partial \theta} \log p(\mathbf{x}, \theta)$

Space of score vector $T_{\theta} = \{\alpha^{\top} s(\cdot, \theta) : \alpha \in \mathbb{R}^p\}$

Fisher information metric $g_{\theta}(u, v) = E_{\theta} \{uv\}$ ($\forall u, v \in T_{\theta}$)

Fisher information matrix $I_{\theta} = E_{\theta} \{s(\mathbf{x}, \theta)s(\mathbf{x}, \theta)^{\top}\}$

Score vector space

$$\begin{aligned}
 \text{Note 1: } u \in T_\theta &\Rightarrow E_\theta(u) = \int u(\mathbf{x}, \theta) p(\mathbf{x}, \theta) d\mathbf{x} = \int \boldsymbol{\alpha}^\top \mathbf{s}(\mathbf{x}, \theta) p(\mathbf{x}, \theta) d\mathbf{x} \\
 &= \boldsymbol{\alpha}^\top \int \mathbf{s}(\mathbf{x}, \theta) p(\mathbf{x}, \theta) d\mathbf{x} = 0 \\
 u, v \in T_\theta &\Rightarrow g_\theta(u, v) = \int u(\mathbf{x}, \theta) v(\mathbf{x}, \theta) p(\mathbf{x}, \theta) d\mathbf{x} \\
 &= \boldsymbol{\alpha}^\top \left\{ \int \mathbf{s}(\mathbf{x}, \theta) \mathbf{s}(\mathbf{x}, \theta)^\top p(\mathbf{x}, \theta) d\mathbf{x} \right\} \boldsymbol{\beta} = \boldsymbol{\alpha}^\top I_\theta \boldsymbol{\beta}
 \end{aligned}$$

The space of all random variables with mean 0 and finite variance

$$S_\theta = \{t(\mathbf{x}, \theta) : E_\theta(t(\mathbf{x}, \theta)) = 0, V_\theta(t(\mathbf{x}, \theta)) < \infty\}$$

Note 2: S_θ is an infinite dimensional vector space that includes T_θ

For $u(\mathbf{x}, \theta) = \boldsymbol{\alpha}^\top \mathbf{s}(\mathbf{x}, \theta)$

$$u(\mathbf{x}, \theta)^k - E_\theta\{u(\mathbf{x}, \theta)^k\}, \quad \frac{\partial^k}{\partial \theta_{i_1} \dots \partial \theta_{i_k}} u(\mathbf{x}, \theta) - E_\theta\left\{ \frac{\partial^k}{\partial \theta_{i_1} \dots \partial \theta_{i_k}} u(\mathbf{x}, \theta) \right\}, \quad k = 1, 2, \dots$$

belong to S_θ Cf. Tiatis (2006)

Linear connections

Score vector $\mathbf{s}(\mathbf{x}, \boldsymbol{\theta}) = (s_i(\mathbf{x}, \boldsymbol{\theta}))_{i=1}^p$

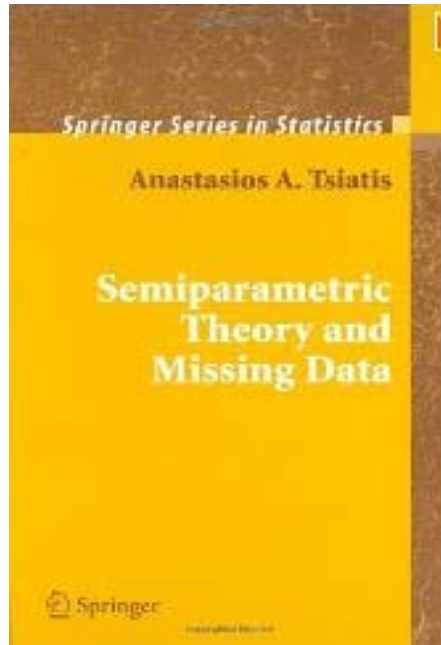
m-connection $\Gamma_{jk}^m(\boldsymbol{\theta}) = \sum_{i'=1}^p g^{ii'} \left[E_{\boldsymbol{\theta}} \left\{ \frac{\partial s_k}{\partial \theta^j} s_{i'} \right\} + E_{\boldsymbol{\theta}} \{ s_k s_j s_{i'} \} \right] \quad (1 \leq i, j, k \leq p)$

e-connection $\Gamma_{jk}^e(\boldsymbol{\theta}) = \sum_{i'=1}^p g^{ii'} E_{\boldsymbol{\theta}} \left\{ \frac{\partial s_k}{\partial \theta^j} s_{i'} \right\} \quad (1 \leq i, j, k \leq p)$

m-geodesic $p_t^{(m)}(\mathbf{x}) = (1-t)p(\mathbf{x}) + tq(\mathbf{x}), \quad (p, q \in \mathbb{P})$

e-geodesic $r_s^{(e)}(\mathbf{x}) = c_s \{r(\mathbf{x})\}^{1-s} \{q(\mathbf{x})\}^s, \quad (q, r \in \mathbb{P})$

Semiparametric theory



- 1 Introduction to Semiparametric Models
- 2 Hilbert Space for Random Vectors
- 3 The Geometry of Influence Functions
- 4 Semiparametric Models
- 5 Other Examples of Semiparametric Models
- 6 Models and Methods for Missing Data
- 7 Missing and Coarsening at Random for Semiparametric Models
- 8 The Nuisance Tangent Space and Its Orthogonal Complement
- 9 ...14.

Semiparametric Theory and Missing Data
(Springer Series in Statistics)
Anastasios Tsiatis

Geodesical models

Definition (i) Statistical model M is said to be **e-geodesical** if

$$p(\mathbf{x}), q(\mathbf{x}) \in M \Rightarrow c_t p(\mathbf{x})^{1-t} q(\mathbf{x})^t \in M \quad (\forall t \in (0,1))$$

$$\text{where } c_t = 1 / \int p(\mathbf{x})^{1-t} q(\mathbf{x})^t d\mathbf{x}$$

(ii) Statistical model M is said to be **m-geodesic** if

$$p(\mathbf{x}), q(\mathbf{x}) \in M \Rightarrow (1-t)p(\mathbf{x}) + tq(\mathbf{x}) \in M \quad (\forall t \in (0,1))$$

Note : Let \mathcal{P} be a space of all probability density functions.

By definition \mathcal{P} is e-geodesical and m-geodesical.

However the theoretical framework for \mathcal{P} is not perfectly complete.

Cf. Pistone and Sempi (1995, AS)

Two type of modeling

Let $p_0(\mathbf{x})$ be a pdf and $\mathbf{t}(\mathbf{x})$ a p -dimensional statistic.

Exponential model $M^{(e)} = \{ p(\mathbf{x}, \boldsymbol{\theta}) = p_0(\mathbf{x}) \exp \{ \boldsymbol{\theta}^T \mathbf{t}(\mathbf{x}) - \kappa(\boldsymbol{\theta}) \}; \boldsymbol{\theta} \in \Theta \}$

$$\Theta = \{ \boldsymbol{\theta} \in \mathbb{R}^p : \kappa(\boldsymbol{\theta}) < \infty \}$$

Matched mean model $M^{(m)} = \{ p(\mathbf{x}) : E_p \{ \mathbf{t}(\mathbf{x}) \} = E_{p_0} \{ \mathbf{t}(\mathbf{x}) \} \}$

Let $\{ p_i(\mathbf{x}) : i = 0, 1, \dots, I \}$ be a set of pdfs.

Exponential model $M^{(e)} = \{ p(\mathbf{x}, \boldsymbol{\theta}) = p_0(\mathbf{x}) \exp \{ \sum_{i=1}^I \theta_i \log \frac{p_i(\mathbf{x})}{p_0(\mathbf{x})} - \kappa(\boldsymbol{\theta}) \}; \boldsymbol{\theta} \in \Theta \}$

$$\Theta = \{ \boldsymbol{\theta} \in \mathbb{R}^p : \kappa(\boldsymbol{\theta}) < \infty \}$$

Mixture model $M^{(m)} = \{ p(\mathbf{x}, \boldsymbol{\theta}) = (1 - \sum_{i=1}^I \theta_i) p_0(\mathbf{x}) + \sum_{i=1}^I \theta_i p_i(\mathbf{x}) : \boldsymbol{\theta} \in \Theta \}$

$$\Theta = \{ \boldsymbol{\theta} = (\theta_1, \dots, \theta_I) : 0 < \sum_{i=1}^I \theta_i < 1, \theta_i > 0 (\forall i = 1, \dots, I) \}$$

Statistical functional

A functional $f(p)$ is called a statistical functional if $p(\mathbf{x})$ is a pdf. (Hampel, 1974)

A statistical functional $f(p)$ is said to be Fisher-consistent for a model

$M = \{p_{\theta}(\mathbf{x}) : \theta \in \Theta\}$ if $f(p)$ satisfies that

$$(1) \quad f(p) \in \Theta,$$

$$(2) \quad f(p_{\theta}) = \theta \quad (\forall \theta \in \Theta)$$

Example. For a normal model $M = \{p_{\theta}(\mathbf{x}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\{-\frac{(x-\mu)^2}{2\sigma^2}\}, \theta = (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}^+\}$

$$f(p) = \left(\int xp(x) dx, \int x^2 dx - \left(\int xp(x) dx \right)^2 \right)^T$$

is Fisher-consistent for $\theta = (\mu, \sigma^2)$.

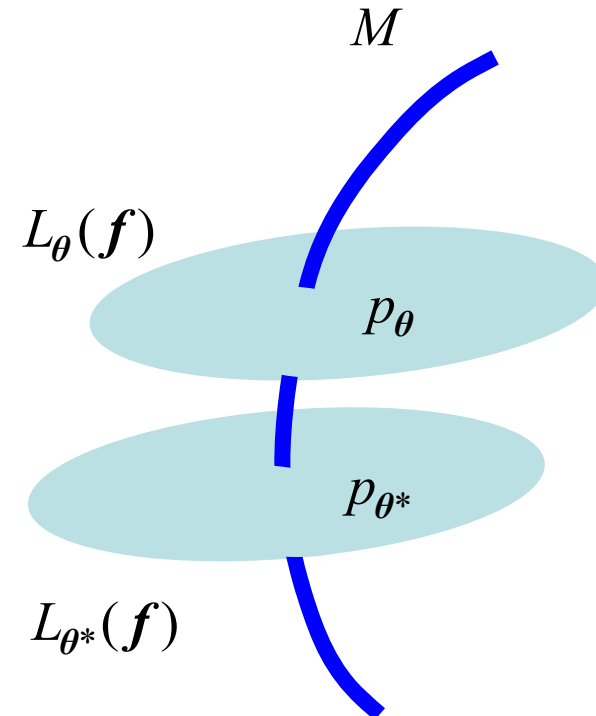
Transversality

Let $f(p)$ be Fisher consistent functional.

$$L_{\theta}(f) = \{p : f(p) = \theta\},$$

which is called a leaf transverse to M .

- (1) $L_{\theta}(f) \cap M = \{p_{\theta}\}$
- (2) $\bigoplus_{\theta \in \Theta} L_{\theta}(f)$ is a local neighborhood.



Foliation structure

Statistical model $M = \{p_\theta(\mathbf{x}) = p(\mathbf{x}, \theta) : \theta \in \Theta\}$ ($\Theta \subset \mathbb{R}^p$)

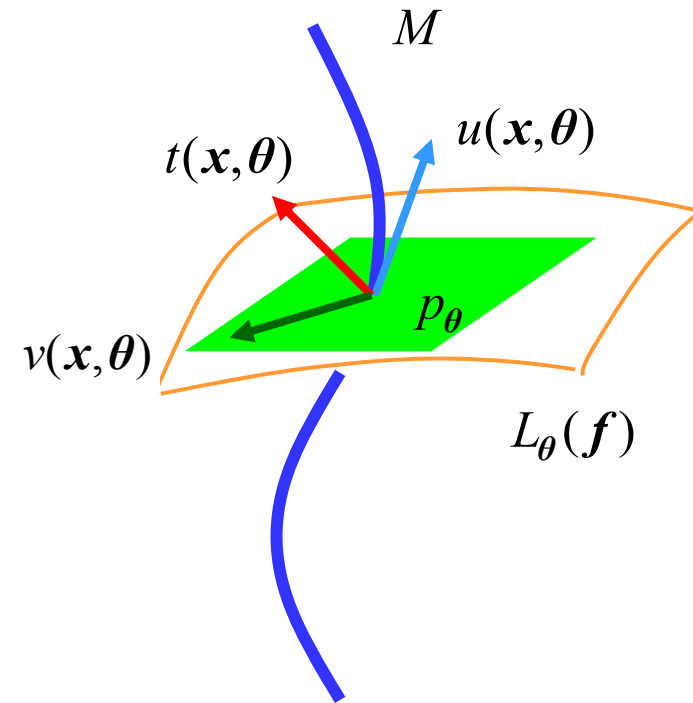
Foliation $\mathbf{P} = \bigoplus_{\theta \in \Theta} L_\theta(\mathbf{f})$

Decomposition of tangent spaces

$$T_{p_\theta}(\mathbf{P}) = T_{p_\theta}(M) \oplus T_{p_\theta}(L_\theta(\mathbf{f}))$$

$\forall t \in T_{p_\theta}(\mathbf{P}), \exists u \in T_{p_\theta}(M), \exists v \in T_{p_\theta}(L_\theta(\mathbf{f}))$

such that $t(\mathbf{x}, \theta) = u(\mathbf{x}, \theta) + v(\mathbf{x}, \theta)$



Transversality for MLE

Let $p_0(\mathbf{x})$ be a pdf and $\mathbf{t}(\mathbf{x})$ a p -dimensional statistic.

$$\text{Exponential model } M^{(e)} = \{ p(\mathbf{x}, \boldsymbol{\theta}) = p_0(\mathbf{x}) \exp \{ \boldsymbol{\theta}^T \mathbf{t}(\mathbf{x}) - \kappa(\boldsymbol{\theta}) \}; \boldsymbol{\theta} \in \Theta \}$$
$$\Theta = \{ \boldsymbol{\theta} \in \mathbb{R}^p : \kappa(\boldsymbol{\theta}) < \infty \}$$

For the exponential model $M^{(e)}$ the MLE functional

$$\mathbf{f}_{\text{ML}}(p) =: \arg \max_{\boldsymbol{\theta} \in \Theta} E_p \{ \log p(\mathbf{x}, \boldsymbol{\theta}) \}$$

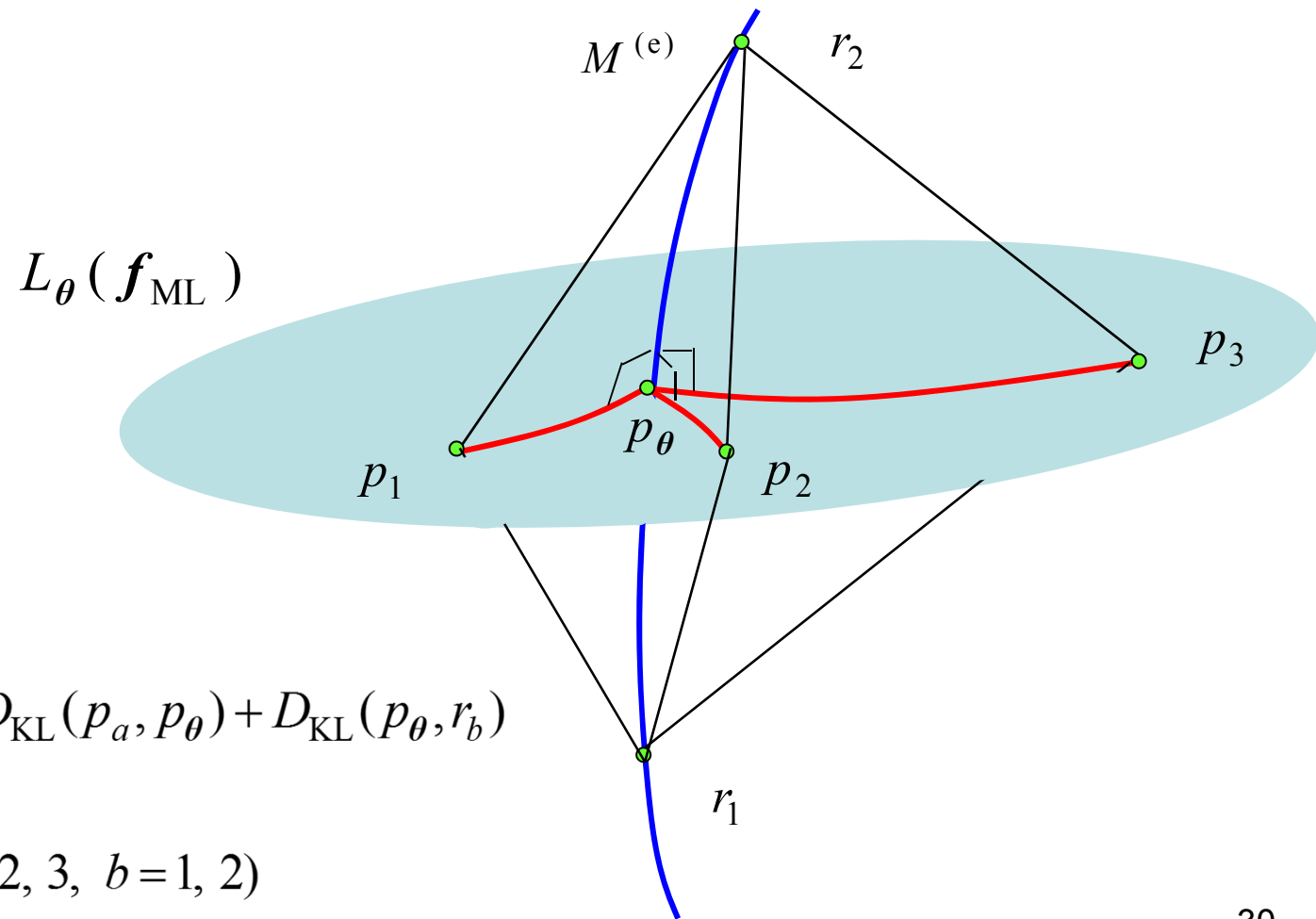
is written by

$$\mathbf{f}_{\text{ML}}(p) = \arg \left\{ E_p \{ \mathbf{t}(\mathbf{x}) \} = E_{p(\cdot, \boldsymbol{\theta})} \{ \mathbf{t}(\mathbf{x}) \} \right\}$$

Hence the foliation associated with \mathbf{f}_{ML} is a matched mean model.

$$L_{\boldsymbol{\theta}}(\mathbf{f}_{\text{ML}}) = \{ p : E_p \{ \mathbf{t}(\mathbf{x}) \} = E_{p(\cdot, \boldsymbol{\theta})} \{ \mathbf{t}(\mathbf{x}) \} \}$$

Maximum likelihood foliation



$$D_{\text{KL}}(p_a, r_b) = D_{\text{KL}}(p_a, p_\theta) + D_{\text{KL}}(p_\theta, r_b)$$

$$(a = 1, 2, 3, b = 1, 2)$$

Estimating function

Statistical model $M = \{p_\theta(\mathbf{x}) = p(\mathbf{x}, \theta) : \theta \in \Theta\}$ ($\Theta \subset \mathbb{R}^p$)

p -variable function $\mathbf{u}(\mathbf{x}, \theta)$ is unbiased

def
 \Leftrightarrow

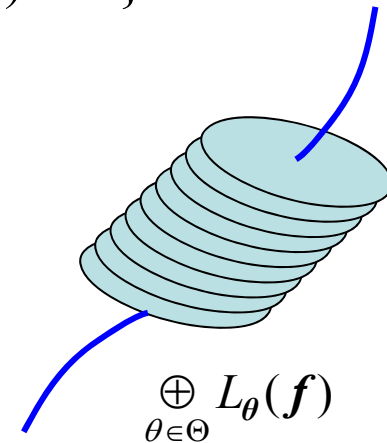
$$E_{p(\cdot, \theta)}\{\mathbf{u}(\mathbf{x}, \theta)\} = \mathbf{0}, \det\left(E_{p(\cdot, \theta)}\left\{\frac{\partial \mathbf{u}(\mathbf{x}, \theta)}{\partial \theta}\right\}\right) \neq 0 \quad (\forall \theta \in \Theta)$$

The statistical functional $\mathbf{f}(p) = \arg \operatorname{solve}_{\theta \in \Theta} \{E_p \mathbf{u}(\mathbf{x}, \theta) = \mathbf{0}\}$

is Fisher-consistent, and the leaf transverse to M

$$L_\theta(\mathbf{f}) = \{p : E_p \mathbf{u}(\mathbf{x}, \theta) = \mathbf{0}\}$$

is m -geodesical.



Minimum divergence geometry

$D : M \times M \rightarrow \mathbf{R}$ is an **information divergence** on a statistical model M

- \iff
- (i) $D(p, q) \geq 0$ with equality if and only if $p = q$
 - (ii) D is a differentialble on $M \times M$

Let
$$D(Z \cdots | XY \cdots)(p) = Z_p \cdots X_q Y_q \cdots D(p, q) |_{q=p}$$

Then we get a **Riemannian metric** and **dual connections** on M (Eguchi, 1983,1992)

$$g^{(D)}(X, Y) = -D(X | Y)$$

$$g^{(D)}(\nabla_X Y, Z) = -D(XY | Z) \quad (\forall Z \in X(M))$$

$$g^{(D)}(\nabla_X^* Y, Z) = -D(Z | XY) \quad (\forall Z \in X(M))$$

Remarks

$g^{(D)}$ is a Riemann metric

$$D(X | \cdot) = 0$$

$$g^{(D)}(X, Y) - g^{(D)}(Y, X) = D(XY | \cdot) - D(YX | \cdot) = D([X, Y] | \cdot) = 0$$

$$g^{(D)}(X, Y) = -D(X | Y) = D(XY | \cdot) = D(\cdot | XY)$$

∇_X and ∇_X^* are dual connections

$$g^{(D)}(\nabla_{fX} Y, Z) = -f D(XY | Z) = g^{(D)}(f \nabla_X Y, Z) \quad (\forall Z)$$

$$g^{(D)}(\nabla_X (fY), Z) = -D(X(fY) | Z) = g^{(D)}(f \nabla_X Y + (Xf)Y, Z) \quad (\forall Z)$$

$$\underline{Xg^{(D)}(Y, Z)} = g^{(D)}(\nabla_X Y, Z) + g^{(D)}(Y, \nabla_X^* Z) = \underline{g^{(D)}(\bar{\nabla}_X Y, Z) + g^{(D)}(Y, \bar{\nabla}_X Z)}$$

$$\text{where } \bar{\nabla}_X = \frac{1}{2}(\nabla_X + \nabla_X^*)$$

U divergence

Let a triple (U, u, ξ) such that U is a convex function, $u = U'$, $\xi = u^{-1}$

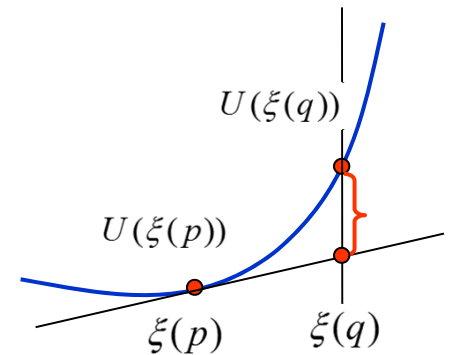
$$\xi(s) = \arg \max_{-\infty < t < \infty} \{ts - U(t)\}$$

$$u(t) = \arg \max_{-\infty < s < \infty} \{st - U^*(s)\}$$

$$\text{where } U^*(s) = s\xi(s) - U(\xi(s))$$

$$U \text{ cross-entropy } L_U(p, q) = \int \{-p\xi(q) + U(\xi(q))\}$$

$$U \text{ entropy } H_U(p) = L_U(p, p) = - \int U^*(p)$$



$$U \text{ divergence } D_U(p, q) = L_U(p, q) - H_U(p)$$

Example of U divergence

$$D_U(p, q) = \int U(\xi(q)) - U(\xi(p)) - p\{\xi(q) - \xi(p)\}$$

KL divergence

$$(U(t), u(t), \xi(s)) = (\exp(t), \exp(t), \log(s))$$

$$D_{\text{KL}}(p, q) = \int \{q - p - p(\log q - \log p)\} = \int p(\log p - \log q)$$

Beta (power) divergence

$$(U_\beta(t), u_\beta(t), \xi_\beta(s)) = \left(\frac{(1 + \beta t)^{(1+\beta)/\beta}}{1 + \beta}, (1 + \beta t)^{1/\beta}, \frac{s^\beta - 1}{\beta} \right)$$

$$D_\beta(p, q) = \int \left\{ \frac{q^{\beta+1} - p^{\beta+1}}{\beta+1} - p \frac{(q^\beta - p^\beta)}{\beta} \right\}$$

Note $\lim_{\beta \rightarrow 0} U_\beta(t) = \exp(t)$

$$\lim_{\beta \rightarrow 0} D_\beta(p, q) = D_{\text{KL}}(p, q)$$

Geometric formula with D_U

$$\begin{aligned}
 (g^{(U)}, \nabla^{(U)}, \nabla^{*(U)}) \text{ s.t. } & \quad g^{(U)}(X, Y) = -D_U(X | Y) \\
 & \quad g^{(U)}(\nabla_X Y, Z) = -D_U(XY | Z) \quad (\forall Z \in X(M)) \\
 & \quad g^{(D)}(\nabla_X^* Y, Z) = -D(Z | XY) \quad (\forall Z \in X(M))
 \end{aligned}$$

$$g_{ij}^{(U)}(\theta) = \int \frac{\partial}{\partial \theta^i} q(x, \theta) \frac{\partial}{\partial \theta^j} \xi(q(x, \theta)) dx$$

$$\Gamma_{ij,k}^{(U)}(\theta) = \int \frac{\partial^2}{\partial \theta^i \partial \theta^j} q(x, \theta) \frac{\partial}{\partial \theta^k} \xi(q(x, \theta)) dx$$

$$\Gamma_{ij,k}^{*(U)}(\theta) = \int \frac{\partial}{\partial \theta^k} q(x, \theta) \frac{\partial^2}{\partial \theta^i \partial \theta^j} \xi(q(x, \theta)) dx$$

$$g^{(U)} = g \quad \iff U = \text{exp} \quad (\text{KL divergence})$$

$$\underline{\nabla^{(U)}} = \overset{\text{m}}{\nabla} \quad (\forall U) \quad (\text{Eguchi, 2005})$$

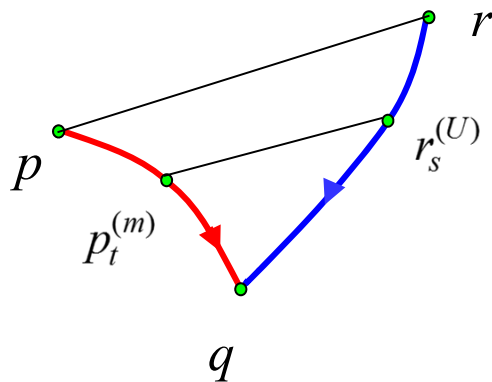
Triangle with D_U

mixture geodesic $\bar{p}_t^{(m)}(\mathbf{x}) = (1-t)p(\mathbf{x}) + tq(\mathbf{x})$

U geodesic $r_s^{(U)}(\mathbf{x}) = u((1-s)\xi(r(\mathbf{x})) + s\xi(q(\mathbf{x})) + \kappa_s)$

$$D_U(p_t^{(m)}, r_s^{(U)}) - D_U(p_t^{(m)}, q) - D_U(q, r_s^{(U)})$$

$$= (1-t)(1-s)\{D_U(p, r) - D_U(p, q) - D_U(q, r)\} \quad (\forall (s, t) \in [0, 1] \times [0, 1])$$



$$\{p_t^{(m)}\} \perp_q \{r_s^{(U)}\} \Rightarrow$$

$$D_U(p_t^{(m)}, r_s^{(U)}) = D_U(p_t^{(m)}, q) + D_U(q, r_s^{(U)})$$

**Information divergence class
and
robust statistical methods**

Light and shadow of MLE

<ol style="list-style-type: none">1. Invariance under data-transformations2. Asymptotic efficiency	<ol style="list-style-type: none">1. Non-robust2. Over-fitting
---	---

Log-likelihood on exponential family

Sufficiency and efficiency

Likelihood method \log \rightleftharpoons \exp

U -method ξ \rightleftharpoons u

Max U -entropy distribution

Let us fix a statistics $\mathbf{t}(\mathbf{x})$.

Equal mean space $\Gamma_\tau = \{p : E_p\{\mathbf{t}(X)\} = \boldsymbol{\tau}\}$

$$p^*(\mathbf{x}) = \arg \max_{p \in \Gamma_\tau} H_U(p)$$

$$p^*(\mathbf{x}) = u(\boldsymbol{\theta}^\top \mathbf{t}(\mathbf{x}) - \kappa(\boldsymbol{\theta})) \quad \mathbf{U}\text{-model}$$

Euler-Lagrange's 1st variation

$$\begin{aligned} & \frac{\partial}{\partial \varepsilon} \{H((1-\varepsilon)p^* + \varepsilon q) - \boldsymbol{\theta}^\top (\mathbf{t} - \boldsymbol{\tau})((1-\varepsilon)p^* + \varepsilon q) - \lambda((1-\varepsilon)p^* + \varepsilon q - 1)\} \Big|_{\varepsilon=0} \\ &= \left\langle q \{ \xi(p^*) - \boldsymbol{\theta}^\top \mathbf{t} + \kappa \} \right\rangle \quad (\forall q \in \Gamma_\tau) \end{aligned}$$

$$\text{Hence} \quad \xi(p^*(\mathbf{x})) = \boldsymbol{\theta}^\top \mathbf{t}(\mathbf{x}) - \kappa(\boldsymbol{\theta})$$

U -estimate

Let $p(\mathbf{x})$ be a data density function with statistical model $q_{\theta}(\mathbf{x})$

U -loss function $L_U(\boldsymbol{\theta}) = L_U(p, q_{\boldsymbol{\theta}}) = -E_p\{\xi(q_{\boldsymbol{\theta}})\} + \langle U(\xi(q_{\boldsymbol{\theta}})) \rangle$

U -empirical loss function $L_U^{\text{emp}}(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^n \xi(q_{\boldsymbol{\theta}}(\mathbf{x}_i)) + \langle U(\xi(q_{\boldsymbol{\theta}})) \rangle$

U -estimate $\hat{\boldsymbol{\theta}}_U = \arg \min_{\boldsymbol{\theta} \in \Theta} L_U^{\text{emp}}(\boldsymbol{\theta})$

U -estimating function $s_U(\mathbf{x}, \boldsymbol{\theta}) = w(\mathbf{x}, \boldsymbol{\theta}) s(\mathbf{x}, \boldsymbol{\theta}) - E_{q_{\boldsymbol{\theta}}}\{w(\mathbf{X}, \boldsymbol{\theta}) s(\mathbf{X}, \boldsymbol{\theta})\}$

where $w(\mathbf{x}, \boldsymbol{\theta}) = q_{\boldsymbol{\theta}}(\mathbf{x}) \xi'(q_{\boldsymbol{\theta}}(\mathbf{x}))$, $s(\mathbf{x}, \boldsymbol{\theta})$ is score function.

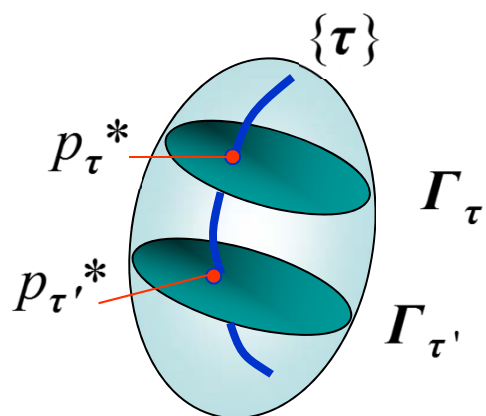
Γ -minimax

Let us fix a statistics $t(\mathbf{x})$.

Equal mean space $\Gamma_\tau = \{p : E_p\{t(\mathbf{X})\} = \tau\}$

$$\max_{p \in \Gamma_\tau} \min_{q \in \Gamma_\tau} L_U(p, q) = L_U(p_\tau^*, p_\tau^*) = \min_{q \in \Gamma_\tau} \max_{p \in \Gamma_\tau} L_U(p, q)$$

$$p_\tau^* = \arg \max_{p \in \Gamma_\tau} H_U(p)$$



where $p_\tau^*(\mathbf{x}) = u(\boldsymbol{\theta}^\top t(\mathbf{x}) - \kappa_\theta)$,

$$\int t(\mathbf{x}) u(\boldsymbol{\theta}^\top t(\mathbf{x}) - \kappa_\theta) = \tau$$

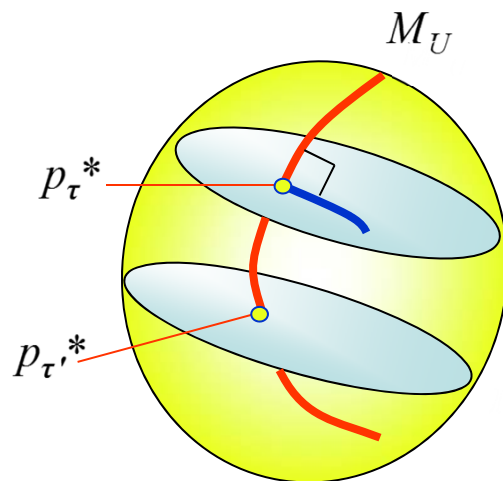
U -estimator for U -model

$$U\text{-model} \quad M_U = \{q_\theta(\mathbf{x}) = u(\boldsymbol{\theta}^T \mathbf{t}(\mathbf{x}) - \kappa_\theta) : \boldsymbol{\theta} \in \Theta\}$$

$$U\text{-estimator of mean parameter } \boldsymbol{\tau} = E_{q_\theta} \{\mathbf{t}(\mathbf{x})\} \text{ is } \hat{\boldsymbol{\tau}}_U = \frac{1}{n} \sum_{i=1}^n \mathbf{t}(\mathbf{x}_i)$$

We observe that

$$\begin{aligned} L_U^{\text{emp}}(\boldsymbol{\theta}) &= \frac{1}{n} \sum_{i=1}^n \xi(q_\theta(\mathbf{x}_i)) - \langle U(\xi(q_\theta)) \rangle \\ &= \frac{1}{n} \sum_{i=1}^n \{\boldsymbol{\theta}^T \mathbf{t}(\mathbf{x}_i) - \kappa_\theta\} - \langle U(\boldsymbol{\theta}^T \mathbf{t}(\mathbf{x}) - \kappa_\theta) \rangle \end{aligned}$$



Γ_τ which implies

$$\frac{\partial}{\partial \boldsymbol{\theta}} L_U^{\text{emp}}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \mathbf{t}(\mathbf{x}_i) - E_{q_\theta} \{\mathbf{t}(X)\} = 0$$

$\Gamma_{\tau'}$

Furthermore,

$$L_U^{\text{emp}}(\boldsymbol{\theta}) - L_U^{\text{emp}}(\hat{\boldsymbol{\theta}}_U) = D_U(q_{\hat{\boldsymbol{\theta}}_U}, q_\theta).$$

Influence function

Statistical functional

$$T_U(G) = \arg \min_{\theta \in \Theta} \left\{ - \int \xi(p(x, \theta)) dG(x) + \int U(\xi(p(x, \theta))) dx \right\}$$

Influence function

$$\text{IF}(T, x) \equiv \left[\frac{\partial}{\partial \varepsilon} T(G_\varepsilon) \right]_{\varepsilon=0} \quad G_\varepsilon = (1-\varepsilon)F_\theta + \varepsilon\delta_x$$

$$\text{IF}(T_U, x) = J_U^{-1}(\theta) [w(x, \theta)S(x, \theta) - E_\theta \{w(x, \theta)S(x, \theta)\}]$$

$$\text{GES}(T_U) = \sup_x \|\text{IF}(T_U, x)\|$$

Efficiency

Asymptotic variance

$$\sqrt{n}(\hat{\theta}_U - \theta) \underset{D}{\Rightarrow} N(0, J_U(\theta)^{-1} H_U(\theta) J_U(\theta)^{-1})$$

where

$$J_U(\theta) = E(w(X, \theta) S(X, \theta) S(X, \theta)^T),$$
$$H_U(\theta) = \text{Var}(w(X, \theta) S(X, \theta))$$

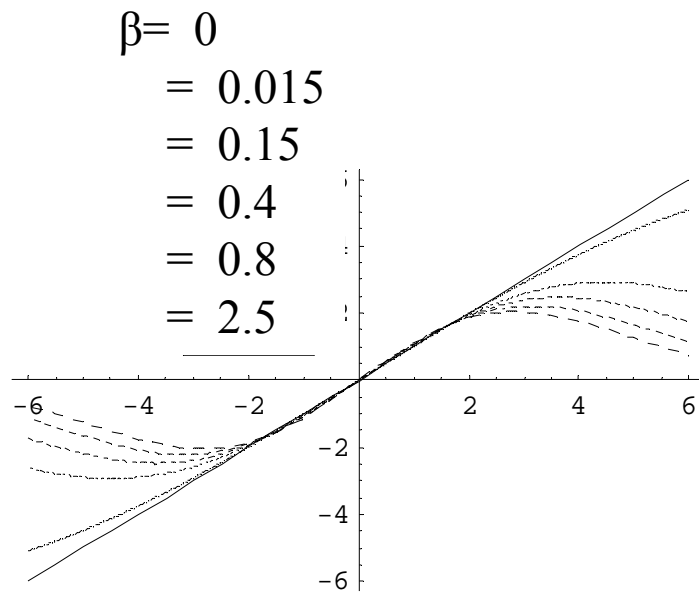
Information inequality

$$I(\theta)^{-1} \leq J_U(\theta)^{-1} H_U(\theta) J_U(\theta)^{-1}$$

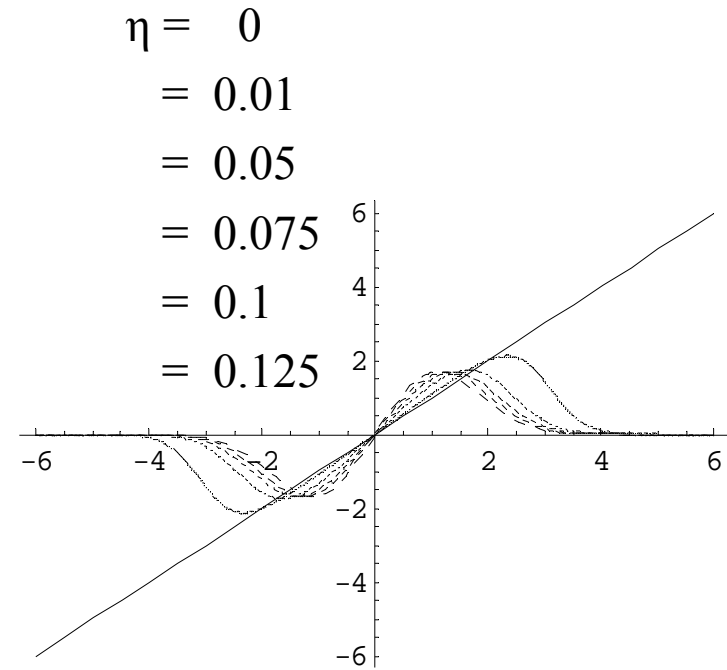
(equality holds iff $U(x) = \exp(x)$)

Normal mean

Influence function



β -power estimates



η -sigmoid estimates

$$U_{\eta}(t) = \exp(t) + \eta t$$

Gross Error Sensitivity

β -power estimates

η -sigmoid estimates

β	効率	GES	η	効率	GES
0	1	∞	0	1	∞
0.01	0.97	6.16	0.015	0.972	1.9
0.05	0.861	2.92	0.15	0.873	1.04
0.075	0.799	2.47	0.4	0.802	0.678
0.1	0.742	2.21	0.8	0.753	0.455
0.125	0.689	2.04	2.5	0.694	0.197

Multivariate normal

pdf is

$$f(y, \mu, \Sigma) = ((2\pi)^p \det \Sigma)^{-\frac{1}{2}} \exp\left\{-\frac{(y - \mu)^T \Sigma^{-1} (y - \mu)}{2}\right\}$$

Estimating equation

$$\left\{ \begin{array}{l} \frac{1}{n} \sum w(\mathbf{x}_i, \boldsymbol{\theta}) (\mathbf{x}_i - \boldsymbol{\mu}) = 0 \\ \frac{1}{n} \sum w(\mathbf{x}, \boldsymbol{\theta}) \{(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T - \Sigma\} = c_U(\boldsymbol{\theta}) \Sigma \end{array} \right.$$

$$\boldsymbol{\theta} = (\boldsymbol{\mu}, \Sigma)$$

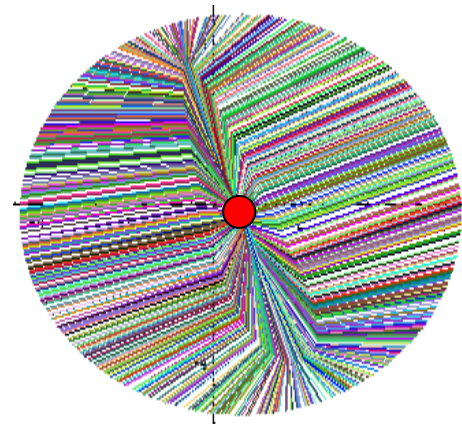
U -algorithm $\theta_k = (\mu_k, \Sigma_k) \rightarrow \theta_{k+1} = (\mu_{k+1}, \Sigma_{k+1})$

繰り返し重み付け平均と分散

$$\mu_{k+1} = \frac{\sum w(x_i, \theta_k) x_j}{\sum w(x_i, \theta_k)} \quad \Sigma_{k+1} = \frac{\sum w(x_i, \theta_k) (x_i - \mu_k)(x_i - \mu_k)^T}{\sum \{w(x_i, \theta_k) - c_U(\det \Sigma_k)\}},$$

Under a mild condition

$$L_U(\theta_{k+1}) \geq L_U(\theta_k) \quad (\forall k = 1, \dots)$$



Simulation

ε -contamination $G_\varepsilon^{(1)} = (1-\varepsilon)N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}\right) + \varepsilon N\left(\begin{pmatrix} 5 \\ -5 \end{pmatrix}, \begin{pmatrix} 4 & 0 \\ 0 & 1 \end{pmatrix}\right)$

$G_\varepsilon^{(2)} = (1-\varepsilon)N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}\right) + \varepsilon N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 9 & 9 \\ 9 & 9 \end{pmatrix}\right)$

	$\varepsilon = 0$	\rightarrow	$\varepsilon = 0.05$	
KL error $D_{\text{KL}}(\hat{\theta}, \theta_0)$	3.03	\rightarrow	39.24	under $G_\varepsilon^{(1)}$
with MLE $\hat{\theta}$	2.70	\rightarrow	16.65	under $G_\varepsilon^{(2)}$

β-power estimates v.s. η-sigmoid estimates

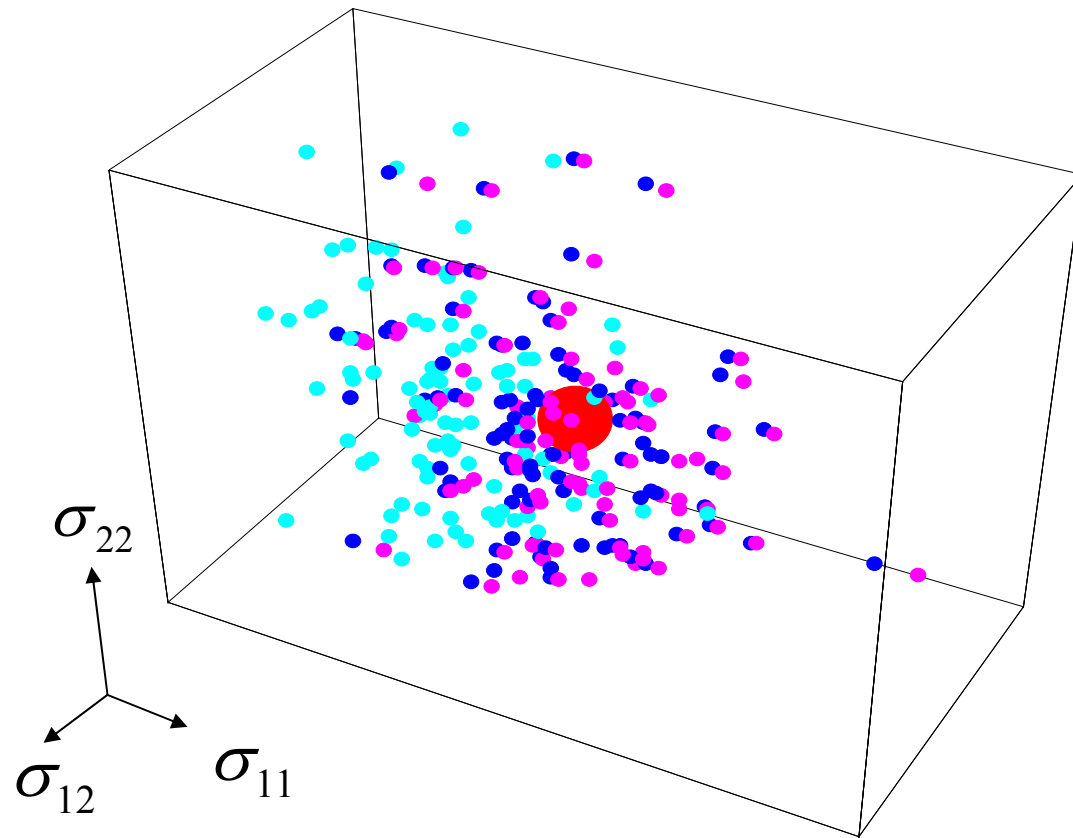
	β	KL error
$G_\varepsilon^{(1)}$	0	39.24
	0.01	35.30
	0.05	20.93
	0.10	8.91
	0.20	12.40
	0.30	31.64

	η	KL error
$G_\varepsilon^{(1)}$	0	39.24
	0.0001	23.04
	0.00025	6.70
	0.0005	4.46
	0.00075	4.64
	0.001	6.04

	β	KL error
$G_\varepsilon^{(2)}$	0	16.5
	0.01	13.91
	0.05	6.65
	0.10	5.14
	0.20	12.20
	0.30	29.68

	η	KL error
$G_\varepsilon^{(2)}$	0	16.5
	0.0005	3.36
	0.00075	3.19
	0.001	3.9
	0.002	3.04
	0.003	3.07

Plot of MLEs (no outliers)

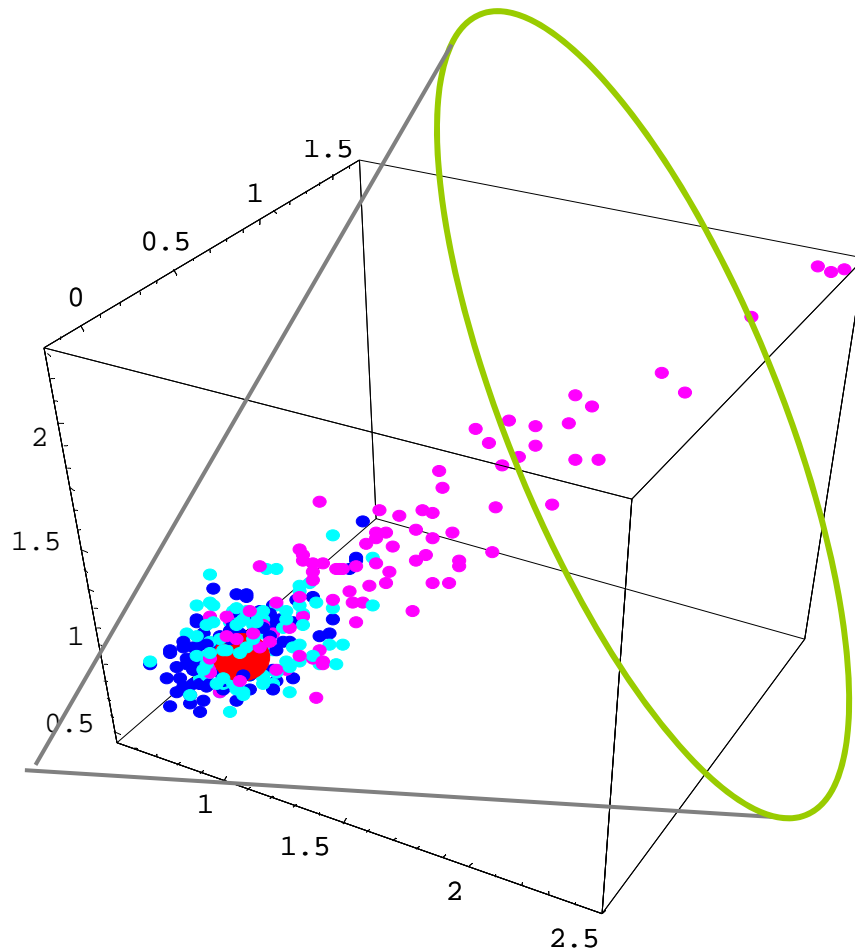


100 replications with
100 size of sample
under Normal dis.

$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{bmatrix}$$

- η - Est ($\eta=0.0025$)
- β - Est ($\beta=0.1$)
- MLE
- true $(1,0,1)$

Plot of MLEs with outliers



100 replications with
100 size of sample
under $G_{\varepsilon}^{(2)}$

- η - Est ($\eta=0.0025$)
- β - Est ($\beta=0.1$)
- MLE
- true (1,0,1)

Selection for tuning parameter

Squared loss function

$$\text{Loss}(\hat{\theta}) = \frac{1}{2} \int \{f(y, \hat{\theta}) - g(y)\}^2 dy$$

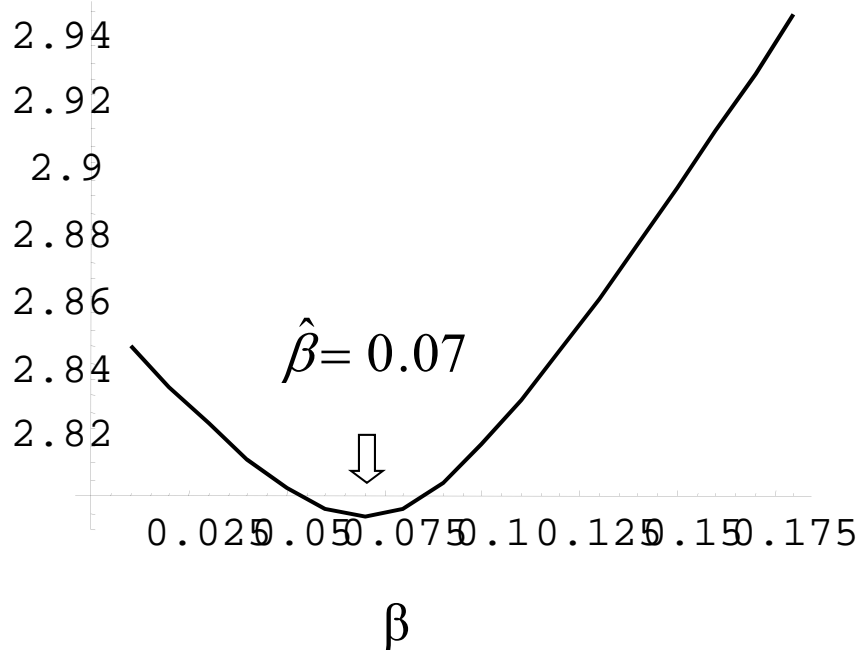
$$\text{CV}(\beta) = -\frac{1}{n} \sum_{i=1}^n f(x_i, \hat{\theta}_{\beta}^{(-i)}) + \frac{1}{2} \int f(y, \hat{\theta}_{\beta})^2 dy$$

$$\hat{\beta} = \arg \min_{\beta} \text{CV}(\beta)$$

Approximate $\hat{\theta}_{\beta}^{(-i)} \approx \hat{\theta}_{\beta} + \frac{1}{n-1} \text{IF}(x_i, \hat{\theta}_{\beta})$

Normal with mean $\begin{pmatrix} 0 \\ 0 \end{pmatrix}$ and variance $\begin{pmatrix} .26 & -.1 \\ -.1 & .26 \end{pmatrix}$

$CV(\beta)$



MLE $\begin{pmatrix} 0.054 \\ -0.081 \end{pmatrix}$ $\begin{pmatrix} .228 & -.126 \\ -.126 & .261 \end{pmatrix}$

MLE $\begin{pmatrix} 0.204 \\ -0.184 \end{pmatrix}$ $\begin{pmatrix} 1.059 & -.263 \\ -.263 & .383 \end{pmatrix}$

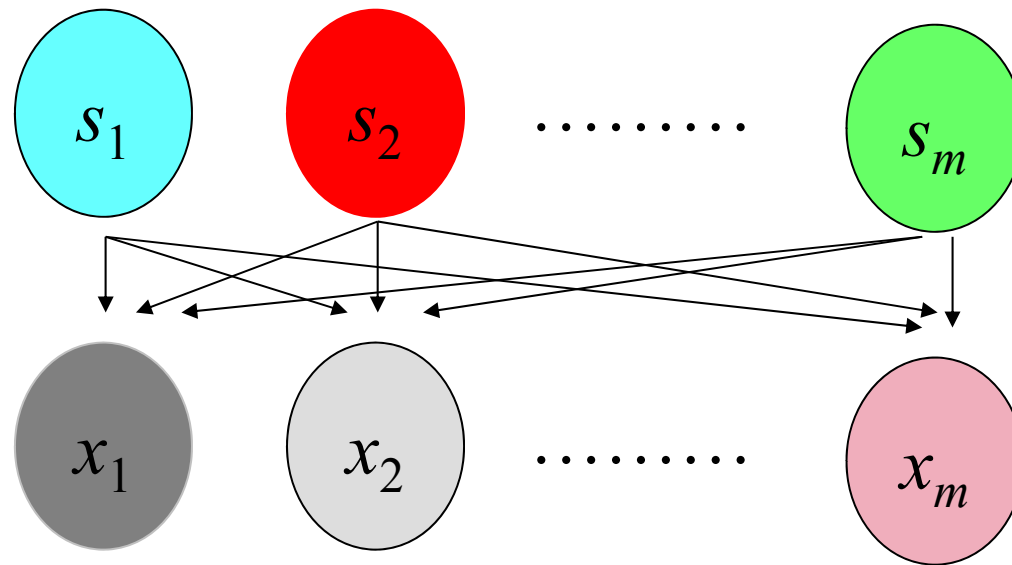
β - Est $\begin{pmatrix} 0.086 \\ -0.132 \end{pmatrix}$ $\begin{pmatrix} .293 & -.134 \\ -.134 & .286 \end{pmatrix}$



What is ICA?

Cocktail party effect

Blind source separation



ICA model

(Independent signals) $s = (s_1, \dots, s_m) \sim p(s) = p_1(s_1) \cdots p_m(s_m)$
 $E(S_1) = 0, \dots, E(S_m) = 0$

(Linear mixture of signals) $W \in \mathbf{R}^{m \times m}$, $\mu \in \mathbf{R}^m$ s.t. $x = W^{-1}s + \mu$

$$f(\mathbf{x}, W, p) = |\det(W)| p_1(\mathbf{w}_1(\mathbf{x} - \boldsymbol{\mu}_1)) \cdots p_m(\mathbf{w}_m(\mathbf{x} - \boldsymbol{\mu}_m))$$

Aim is to learn W from a dataset (x_1, \dots, x_n)

in which $p(s) = p_1(s_1) \cdots p_m(s_m)$ is unknown.

ICA likelihood

Log-likelihood function

$$\ell(W) = \sum_{i=1}^n \sum_{j=1}^m \log p_j(\mathbf{w}_j(\mathbf{x}_i - \boldsymbol{\mu}_j)) + \log |\det(W)|$$

Estimating equation

$$F(x, W, p) = \frac{\partial \ell(x, W, p)}{\partial W} = (I_m - h(Wx)(Wx)^T) W^{-T}$$

$$h(s) = \left(\frac{\partial \log p_1(s_1)}{\partial s_1}, \dots, \frac{\partial \log p_m(s_m)}{\partial s_m} \right)$$

Natural gradient algorithm, Amari et al. (1996)

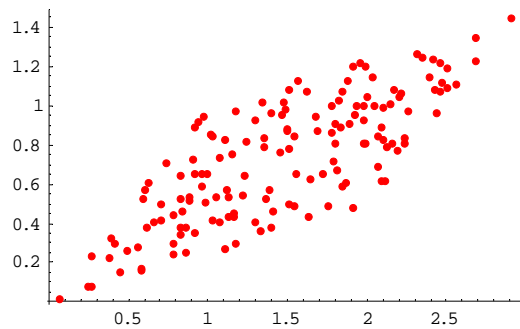
Beta-ICA

β power equation $\frac{1}{n} \sum_{i=1}^n f(x_i, W, \mu)^\beta F(x_i, W, \mu) = B_\beta(W, \mu)$

decomposability: $\forall s \neq t$

$$\begin{aligned} & \prod_{q \neq s, q \neq t} \mathbb{E}[\{p(w_q X - \mu_q)\}^\beta] \\ & \quad \times \mathbb{E}[\{p_s(w_s X - \mu_s)\}^\beta h_s(w_s X - \mu_s)] \\ & \quad \times \mathbb{E}[\{p_t(w_t X - \mu_t)\}^\beta w_t X] = 0 \end{aligned}$$

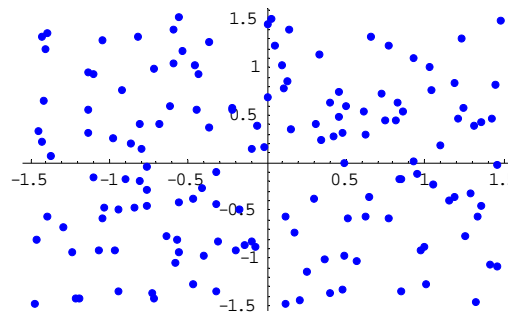
Likelihood ICA



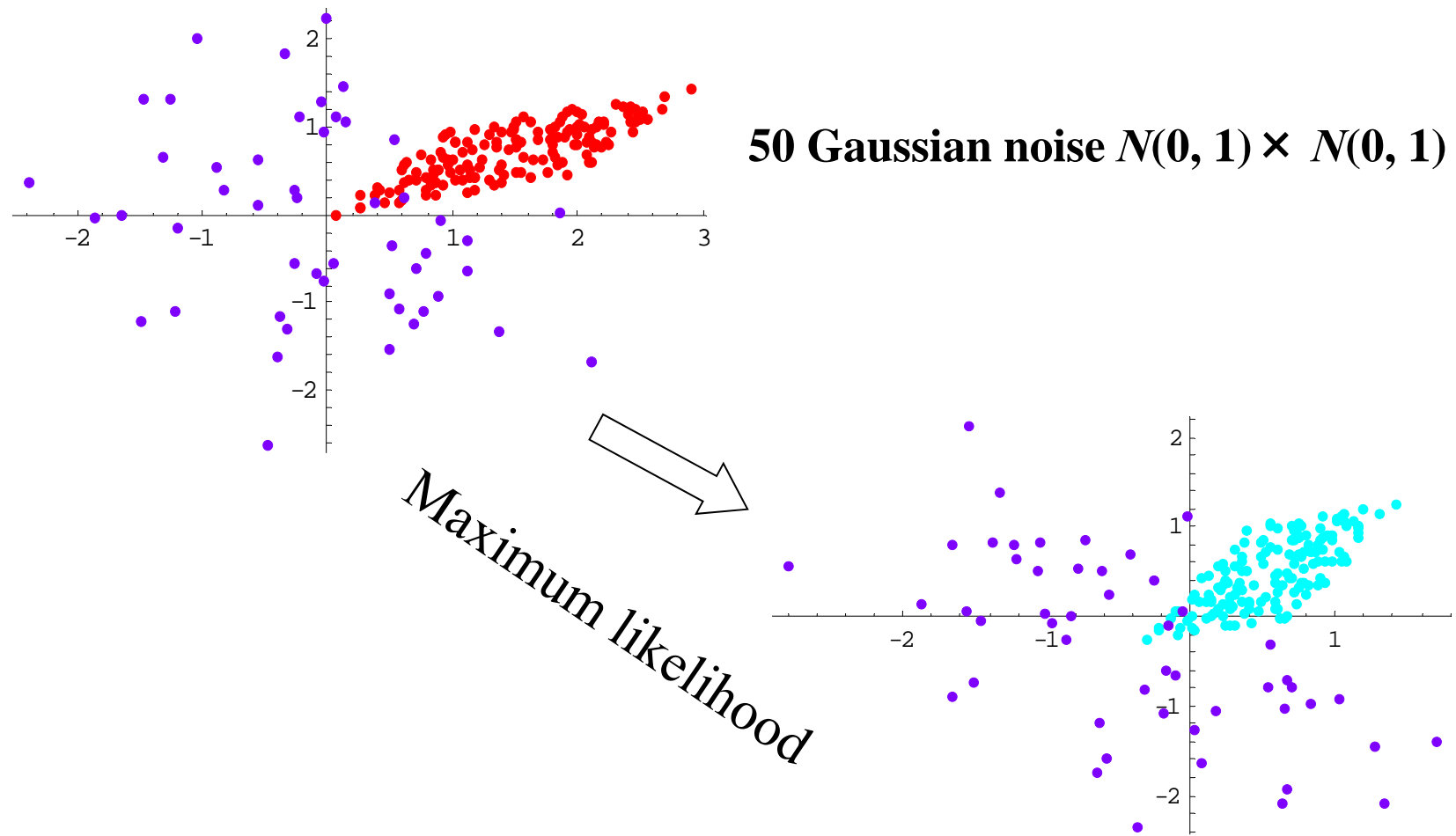
150 signals $U(0,1) \times U(0,1)$

Mixture matrix $W^{-1} = \begin{bmatrix} 1 & 2 \\ 1 & 0.5 \end{bmatrix}$

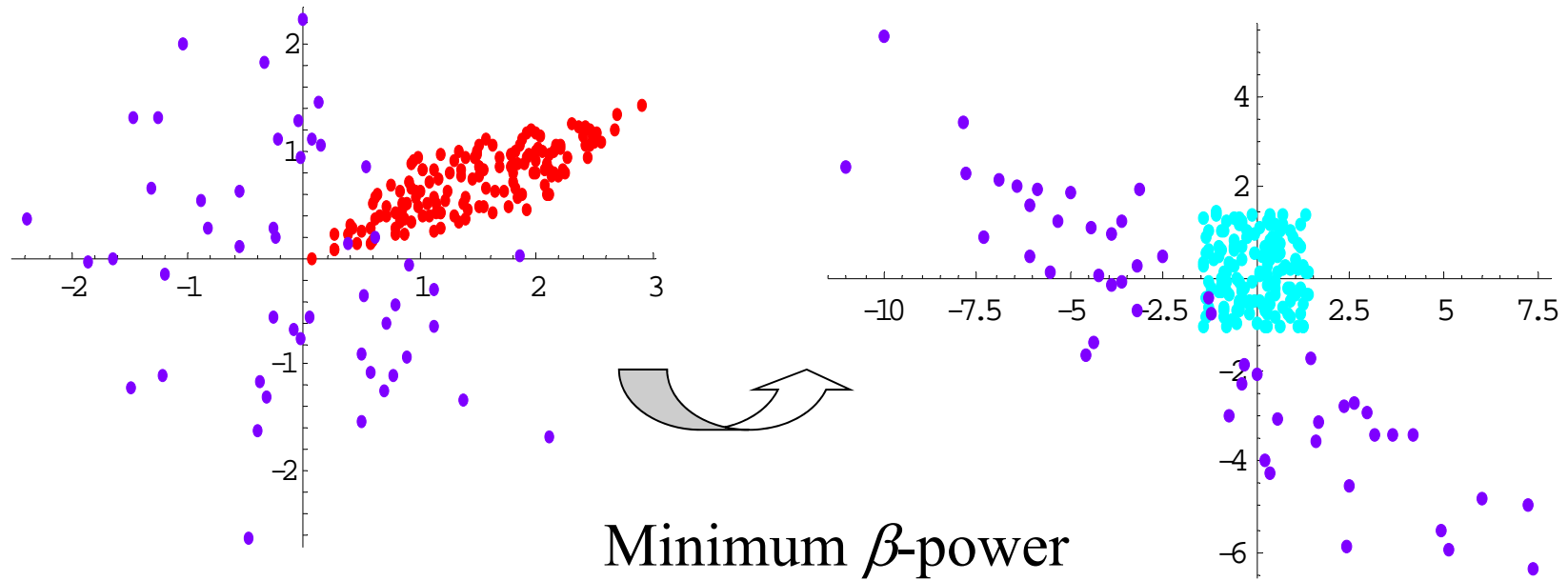
Maximum likelihood



Non-robustness



β power-ICA ($\beta=0.2$)



Article

Projective Power Entropy and Maximum Tsallis Entropy Distributions

Shinto Eguchi *, Osamu Komori and Shogo Kato

Tsallisべきエントロピーに射影不変性を課したエントロピー、ダイバージェンスを考察した。平均と分散を一定にする分布族で射影べきエントロピー最大分布モデルが正規分布, t -分布, ウェグナー分布を含む形で導出された。そのモデル上での平均と分散の素直な推定法を提案し, 推定量が標本平均と標本分散になることが示された。

Projective γ -power divergence

γ-cross entropy	$C_\gamma(g, f) = -\mathbb{E}_g \left\{ \frac{f(X)}{\ f\ _q} \right\}^\gamma$
	where $\ f\ _q = \left(\int f(x)^q dx \right)^{\frac{1}{q}}$ with $q = \gamma + 1$
γ-diagonal entropy	$H_\gamma(f) = C_\gamma(f, f)$
γ-divergence	$D_\gamma(g, f) = C_\gamma(g, f) - H_\gamma(g)$

Cf. Tsallis (1988), Basu et al (1998), Fujisawa-Eguchi (2008), Eguchi-Kato (2010)

Remark. $\gamma = 0$ reduces to Boltzmann-Shannon entropy and Kullback-Leibler divergence

Max γ -entropy

Equal moment space

$$F(\mu, \Sigma) = \{f(x) : E_f(X) = \mu, V_f(X) = \Sigma\}$$

γ -entropy

$$H_\gamma(f) = -\|f\|_q \quad (q = 1 + \gamma)$$

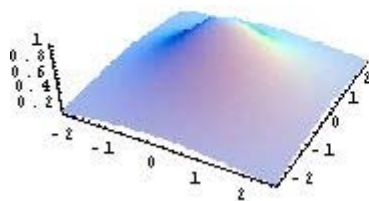
Max γ -entropy

$$H_\gamma(f_\gamma(\cdot, \mu, \Sigma)) = \max_{f \in \mathcal{F}_\gamma(\mu, \Sigma)} H_\gamma(f)$$

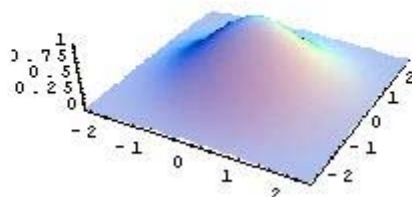
γ -model

$$f_\gamma(x, \mu, \Sigma) = c_\gamma \det(2\pi\Sigma)^{-\frac{1}{2}} \left\{1 - \frac{1}{2} \kappa_\gamma (x - \mu)^T \Sigma^{-1} (x - \mu)\right\}_+^{\frac{1}{\gamma}}$$

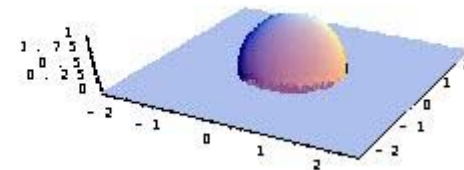
$$F_\gamma(\mu, \Sigma) = \{f \in F(\mu, \Sigma) : \text{Support}(f) \subseteq \{x : (x - \mu)^T \Sigma^{-1} (x - \mu) \leq \frac{1}{2} \kappa_\gamma\}\}$$



$\gamma = -0.3$ (t-distribution)



$\gamma = 0$ (Gaussian)



$\gamma = 2$ (Wigner)

γ -estimation

Parametric model $M = \{f(x, \theta) : \theta \in \Theta\}$

γ -Loss function $L_\gamma(\theta) = -\frac{1}{n} z_\gamma(\theta) \sum_{i=1}^n f(x_i, \theta)^\gamma$, $z_\gamma(\theta) = \left(\int f(x, \theta)^{1+\gamma} dx \right)^{-\frac{\gamma}{1+\gamma}}$

$L_\gamma(\theta) \approx C_\gamma(g, f(\cdot, \theta))$ if $x_1, \dots, x_n \sim g(x)$

γ -estimator $\hat{\theta}_\gamma = \arg \min_{\theta \in \Theta} L_\gamma(\theta)$

Gaussian $f_{\mu, \Sigma}(x) = \det(2\pi\Sigma)^{-\frac{1}{2}} \exp\{-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\}$

γ -estimator $\hat{\theta}_\gamma = (\hat{\mu}_\gamma, \hat{\Sigma}_\gamma)$

$$\hat{\mu}_\gamma = \frac{\sum \{f_{\hat{\mu}_\gamma, \hat{\Sigma}_\gamma}(x_i)\}^\gamma x_i}{\sum \{f_{\hat{\mu}_\gamma, \hat{\Sigma}_\gamma}(x_i)\}^\gamma}$$

$$\hat{\Sigma}_\gamma = \frac{\sum \{f_{\hat{\mu}_\gamma, \hat{\Sigma}_\gamma}(x_i)\}^\gamma (x_i - \hat{\mu}_\gamma)(x_i - \hat{\mu}_\gamma)^\top}{(1 + \gamma) \sum \{f_{\hat{\mu}_\gamma, \hat{\Sigma}_\gamma}(x_i)\}^\gamma}$$

Remark $(\hat{\mu}_\gamma, \hat{\Sigma}_\gamma)$ is super-robust. Cf. Fujisawa-Eguchi (2008)

γ -estimation on γ -model

Let (x_1, \dots, x_n) be a random sample from $f_\gamma(\cdot, \mu, \Sigma)$

γ -loss function

$$L_\gamma(\mu, \Sigma) = \frac{1}{n} \sum_{i=1}^n \det(\Sigma)^{-\frac{1}{2}\frac{\gamma}{\gamma+1}} \left\{ 1 - \frac{1}{2} \kappa_\gamma (x_i - \mu)^\top \Sigma^{-1} (x_i - \mu) \right\}$$

$$= \det(\Xi)^{\frac{\gamma}{p\gamma+2\gamma+2}} - \frac{1}{2} \kappa_\gamma \{ (\mu - \bar{x})^\top \Xi (\mu - \bar{x}) + \text{tr}(S \Xi) \}$$

where $\Xi = \det(\Sigma)^{-\frac{1}{2}\frac{\gamma}{2\gamma+1}} \Sigma^{-1}$

Theorem.

$$L_\gamma(\bar{x}, S) = \min_{(\mu, \Sigma)} L_\gamma(\mu, \Sigma)$$

(Pf). $L_\gamma(\mu, \Sigma) - L_\gamma(\bar{x}, S) = \det(\Xi)^{\frac{\gamma}{p\gamma+2\gamma+2}} - \frac{1}{2} \kappa_\gamma \{ (\mu - \bar{x})^\top \Xi (\mu - \bar{x}) + \text{tr}(S \Xi) \}$

$$- \det(\hat{\Xi})^{\frac{\gamma}{p\gamma+2\gamma+2}} + \frac{1}{2} \kappa_\gamma \text{tr}(S \hat{\Xi}) \}$$

$$= C_\gamma(f_{\bar{x}, S}, f_{\mu, \Sigma}) - C_\gamma(f_{\bar{x}, S}, f_{\bar{x}, S}) = D_\gamma(f_{\bar{x}, S}, f_{\mu, \Sigma}) \geq 0$$

$(\gamma$ -estimator , γ' -model)

model \ loss function	normal-model	γ -model
0-loss (- log-likelihood)	(\bar{x}, S) MLE	Robust model MLE
γ -loss	Robust γ-estimator	(\bar{x}, S) Moment estimator

Gamma-ICA

Mixture of independent signals $X = AS + \mu,$

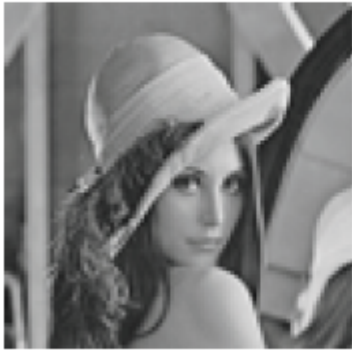
Whitening $Z = \Sigma^{-1/2}(X - \mu)$

Preprocessed form $Z = \tilde{A}S, \quad \tilde{A} = \Sigma^{-1/2}A$

Statistical model $f_Z(z; \mathbf{W}) = |\det(\mathbf{W})| \cdot \prod_{j=1}^p f_j(\mathbf{w}_j^\top z) = f(\mathbf{W}^\top z).$

$$\gamma\text{-ICA} \left\{ \begin{array}{l} (\hat{\mu}, \hat{\Sigma}) = \underset{\mu, \Sigma}{\operatorname{argmin}} \mathcal{D}_\gamma(\hat{g}_X, \xi_{\mu, \Sigma}) \\ \hat{\mathbf{W}} = \underset{\mathbf{W} \in \mathcal{SO}_p}{\operatorname{argmin}} \mathcal{D}_\gamma(\hat{g}_Z, f_Z(\cdot; \mathbf{W})) \end{array} \right.$$

original image



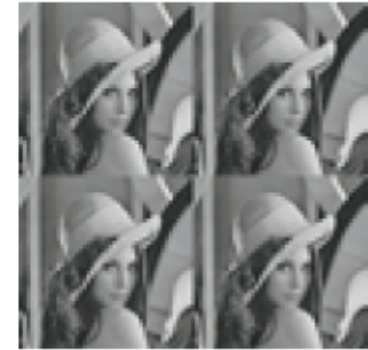
original image



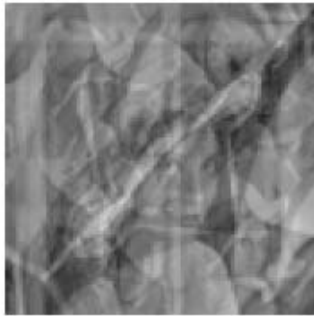
original image



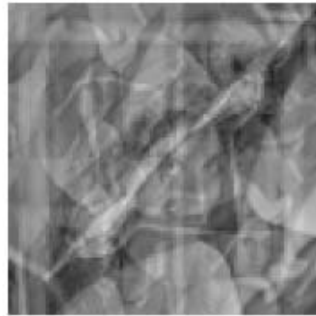
original image



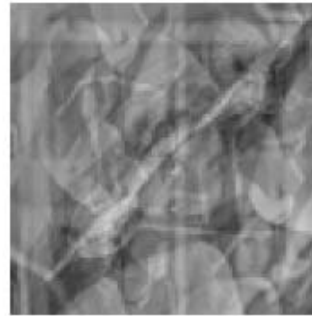
mixed image



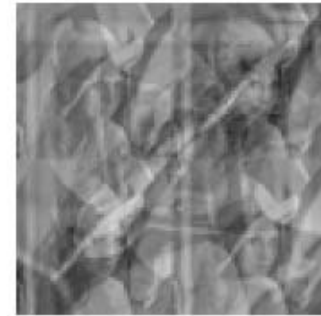
mixed image



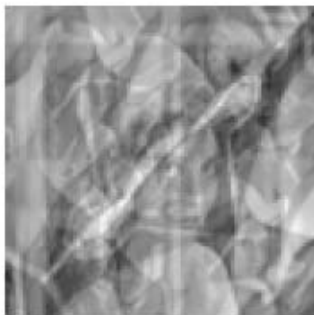
mixed image



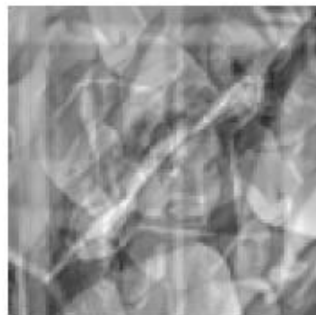
mixed image



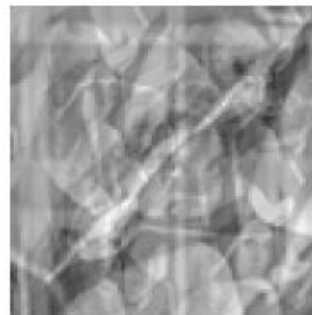
filtered image



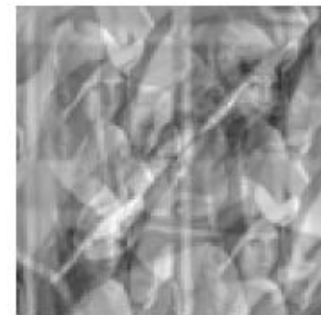
filtered image



filtered image



filtered image



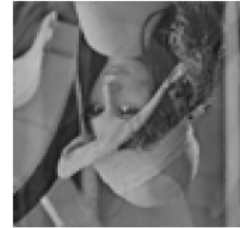
γ -ICA



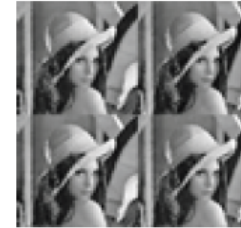
γ -ICA



γ -ICA



γ -ICA



γ -ICA (filtered)



γ -ICA (filtered)



γ -ICA (filtered)



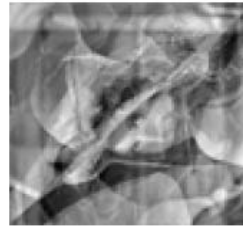
γ -ICA (filtered)



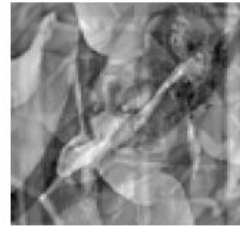
MLE-ICA



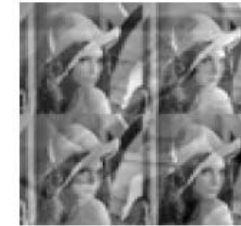
MLE-ICA



MLE-ICA



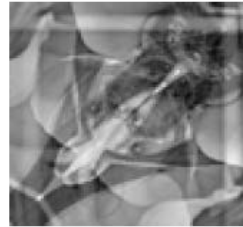
MLE-ICA



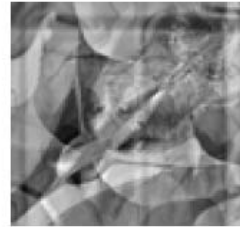
MLE-ICA (filtered)



MLE-ICA (filtered)

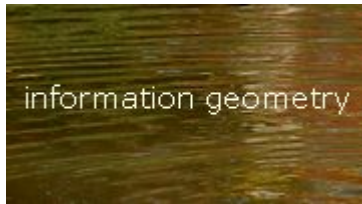


MLE-ICA (filtered)

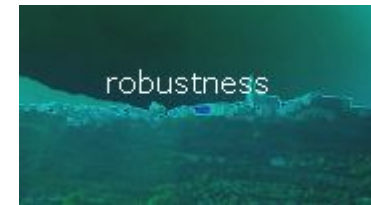


MLE-ICA (filtered)

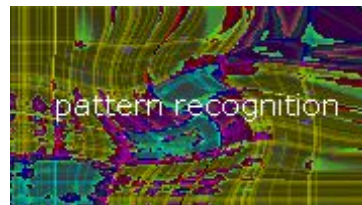




Outline



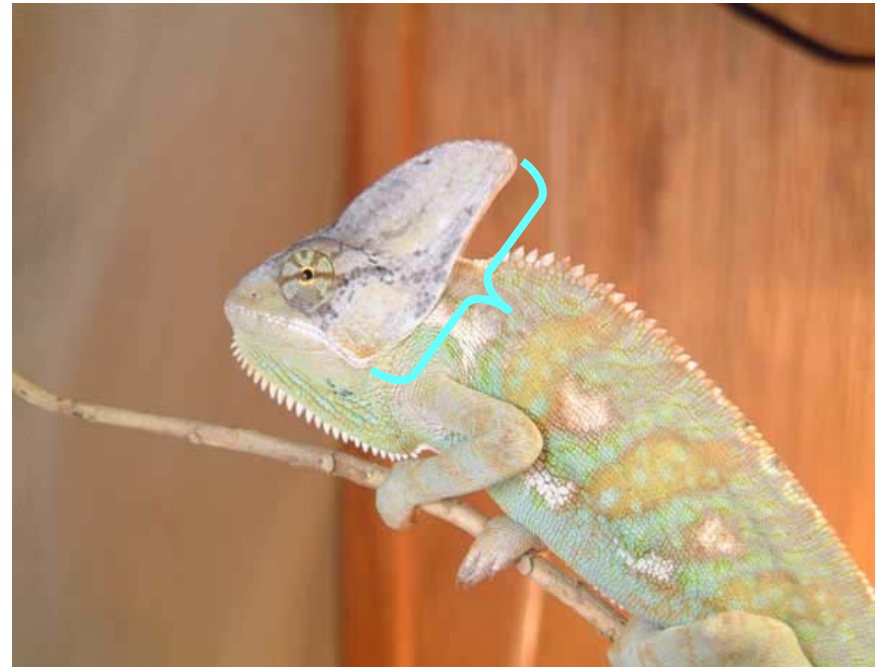
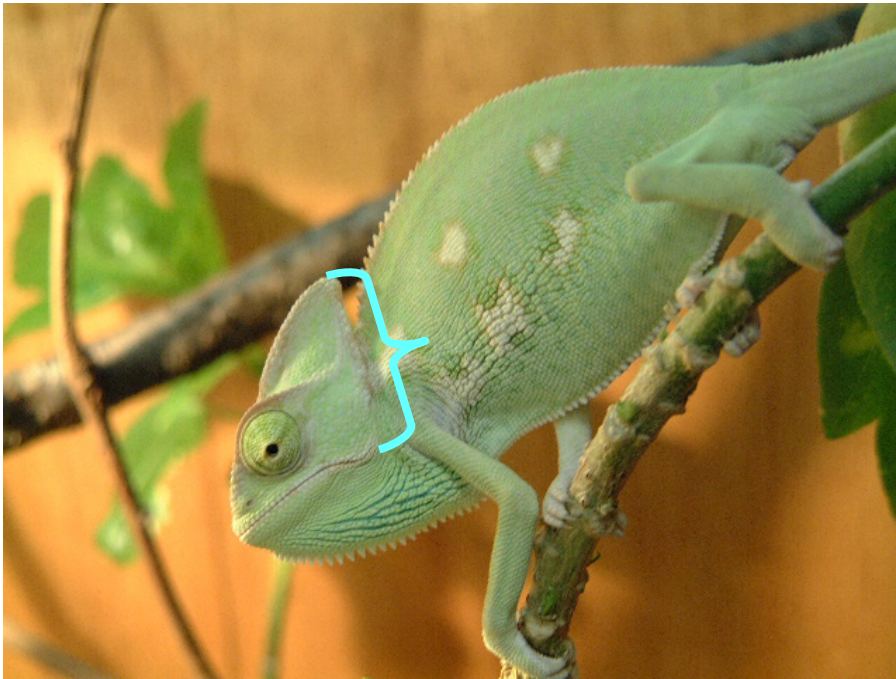
**Boosting learning algorithms
for supervised/unsupervised learning
U-loss functions**



Which chameleon wins?



Pattern recognition...



What is pattern recognition?

- There are a lot of examples for pattern recognition.**
- In principle pattern recognition is a prediction problem for class label which would classify a human interestingness and importance.**
- Originally human brain wants to label phenomena a few words , for example (good, bad), (yes, no), (dead, alive), (success, failure), (effective, no effect)....**
- Brain intrinsically predicts the class label from empirical evidence.**

Practice

- Character recognition
- voice recognition
- Image recognition
- face recognition
- finger print recognition
- speaker recognition
- ☆ Credit scoring
- ☆ Medical screening
- ☆ Default prediction
- ☆ Weather forecast
- ☆ Drug response
- ☆ Treatment effect
- ☆ Failure prediction
- ☆ Infectious disease

Binary classification

Feature vector

Class label

$$\mathbf{x} = (x_1, \dots, x_p) \longrightarrow y \in \{-1, +1\}$$

score $F : \mathbf{x} \rightarrow z$

classifier $h_F(\mathbf{x}) = \text{sgn}\{F(\mathbf{x})\}$

In a binary class $y = -1, 1$ ($G=2$)

$$F(\mathbf{x}) = F(\mathbf{x}, +1) - F(\mathbf{x}, -1)$$

$$\text{sgn}(F(\mathbf{x})) = 1 \Leftrightarrow F(\mathbf{x}, +1) > F(\mathbf{x}, -1)$$

$$\text{sgn}(F(\mathbf{x})) = -1 \Leftrightarrow F(\mathbf{x}, +1) < F(\mathbf{x}, -1)$$

$$\text{sgn}(F(\mathbf{x})) = \arg \max_{y \in \{-1, 1\}} F(\mathbf{x}, y)$$

Probability distribution

Let $p(\mathbf{x}, y)$ be a pdf of random vector (\mathbf{x}, y) .

$$p(B, C) = \int_B \left\{ \sum_{y \in C} p(\mathbf{x}, y) \right\} d\mathbf{x}$$

Marginal density $p(y) = \int_{\mathbf{x}} p(\mathbf{x}, y) d\mathbf{x}$ $p(\mathbf{x}) = \sum_{y \in \{1, \dots, G\}} p(\mathbf{x}, y)$

Conditional density $p(y | \mathbf{x}) = \frac{p(\mathbf{x}, y)}{p(\mathbf{x})}$ $p(\mathbf{x} | y) = \frac{p(\mathbf{x}, y)}{p(y)}$

$$p(\mathbf{x}, y) = p(\mathbf{x} | y)p(y) = p(y | \mathbf{x})p(\mathbf{x})$$

$$\frac{p(\mathbf{x}, y = 1)}{p(\mathbf{x}, y = -1)} = \frac{p(y = 1 | \mathbf{x})}{p(y = -1 | \mathbf{x})}$$

Error rate

Feature vector $\mathbf{x} \in \mathbf{R}^p$, class label $y \in \{1, \dots, G\}$

Classifier $h_F(\mathbf{x})$ has error rate

$$\text{Err}(h_F) = \Pr(h_F(\mathbf{x}) \neq y)$$

$$\text{Err}(h_F) = \sum_{i \neq j} \Pr(h_F(\mathbf{x}) = i, y = j) = 1 - \sum_{i=1}^G \Pr(h_F(\mathbf{x}) = i, y = i)$$

Training error

$$\text{Err}^{\text{train}}(h_F) = \frac{\#\{i : h_F(\mathbf{x}_i) \neq y_i\}}{n} \text{ for } D_{\text{train}} = \{(\mathbf{x}_i, y_i) : i=1, \dots, n\}$$

Test error

$$\text{Err}^{\text{test}}(h_F) = \frac{\#\{j : h_F(\mathbf{x}_j^{\text{test}}) \neq y_j^{\text{test}}\}}{m} \text{ for } D_{\text{test}} = \{(\mathbf{x}_j^{\text{test}}, y_j^{\text{test}}) : j=1, \dots, m\}$$

False negative/positive

	$y = +1$	$y = -1$
$h_F(\mathbf{x}) = +1$	True Positive	False Positive
$h_F(\mathbf{x}) = -1$	False Negative	True Negative

FPR $FP(h_F) = \text{pr}(h_F(\mathbf{x}) = +1 | y = -1)$

FNR $FN(h_F) = \text{Pr}(h_F(\mathbf{x}) = -1 | y = +1)$

$$\text{Err}(h_F) = FN(h_F) \text{pr}(y = +1) + FP(h_F) \text{pr}(y = -1)$$

Bayes rule

Let $p(y|\mathbf{x})$ be a conditional probability of y given \mathbf{x} .

$$\text{Define } F_0(\mathbf{x}) = \log \frac{p(y = 1 | \mathbf{x})}{p(y = -1 | \mathbf{x})}.$$

The classifier $h_{\text{Bayes}}(\mathbf{x}) = \text{sgn}(F_0(\mathbf{x}))$ leads to Bayes rule.

Theorem 1

For any classifier h

$$\text{Err}(h_{\text{Bayes}}) \leq \text{Err}(h)$$

Note: The optimal classifier is equivalent to the likelihood ratio. However, in practice $p(y|\mathbf{x})$ is unknown, so we have to learn $h_{\text{Bayes}}(\mathbf{x})$ based on the training data set.

Error rate for Bayes rule

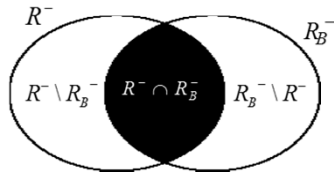
Discriminant space associated with Bayes classifier is given by

$$R_B^+ = \{\mathbf{x} \in \mathbf{R}^P : h_{\text{Bayes}}(\mathbf{x}) = +1\} = \left\{ \mathbf{x} \in \mathbf{R}^P : p(y = +1 | \mathbf{x}) \geq p(y = -1 | \mathbf{x}) \right\}$$

$$R_B^- = \{\mathbf{x} \in \mathbf{R}^P : h_{\text{Bayes}}(\mathbf{x}) = -1\} = \left\{ \mathbf{x} \in \mathbf{R}^P : p(y = +1 | \mathbf{x}) < p(y = -1 | \mathbf{x}) \right\}$$

In general, when a classifier h associates with spaces $\{R^+, R^-\}$

$$\text{Err}(h) - \text{Err}(h_{\text{Bayes}}) = \left(\int_{R^-} - \int_{R_B^-} \right) p(\mathbf{x}) p(y = +1 | \mathbf{x}) d\mathbf{x} + \left(\int_{R^+} - \int_{R_B^+} \right) p(\mathbf{x}) p(y = -1 | \mathbf{x}) d\mathbf{x}$$



$$= \left(\int_{R^- \setminus R_B^-} - \int_{R_B^- \setminus R^-} \right) p(\mathbf{x}) p(y = +1 | \mathbf{x}) d\mathbf{x} + \left(\int_{R^+ \setminus R_B^+} - \int_{R_B^+ \setminus R^+} \right) p(\mathbf{x}) p(y = -1 | \mathbf{x}) d\mathbf{x}$$

$$\geq \left(\int_{R^- \setminus R_B^-} - \int_{R_B^- \setminus R^-} \right) p(\mathbf{x}) p(y = -1 | \mathbf{x}) d\mathbf{x} + \left(\int_{R^+ \setminus R_B^+} - \int_{R_B^+ \setminus R^+} \right) p(\mathbf{x}) p(y = +1 | \mathbf{x}) d\mathbf{x}$$

$$= \left(\int_{R^-} - \int_{R_B^-} \right) p(\mathbf{x}) p(y = -1 | \mathbf{x}) d\mathbf{x} + \left(\int_{R^+} - \int_{R_B^+} \right) p(\mathbf{x}) p(y = +1 | \mathbf{x}) d\mathbf{x}$$

$$= \{1 - \text{Err}(h)\} - \{1 - \text{Err}(h_{\text{Bayes}})\}$$

$$\text{Err}(h) \geq \text{Err}(h_{\text{Bayes}})$$

Cf. McLachlan, 2002

Boost learning

Weak learners can be combined into a strong learner?

Weak learner (classifier) = error rate is slightly less than .5

**Strong learner (classifier) = error rate is slightly less than
that of Bayes classifier**

Boost by filter (Schapire, 1990)

Bagging, Arching (bootstrap)
(Breiman, Friedman, Hasite)

AdaBoost (Schapire, Freund, Batrlett, Lee)

Set of weak learners

Decision stumps

$$F_{\text{stamp}} = \left\{ f_j(\mathbf{x}, a, b) = a \operatorname{sgn}(x_j - b) : j \in \{1, \dots, p\}, a \in \{-1, +1\}, b \in \mathbb{R} \right\}$$

Linear classifiers

$$F_{\text{linear}} = \left\{ f(\mathbf{x}, \boldsymbol{\beta}) = \operatorname{sgn}(\boldsymbol{\beta}_1^T \mathbf{x} + \beta_0) : \boldsymbol{\beta} = (\boldsymbol{\beta}_1, \beta_0) \in \mathbb{R}^{p+1} \right\}$$

Neural net

SVM

k -nearest neighbor

Note: not strong but a variety of characters

$$F_{\text{stamp}} \subseteq F_{\text{linear}}$$

Exponential loss function

Let $D_{\text{train}} = \{ (\mathbf{x}_i, y_i) : i=1, \dots, n \}$ be a training (example) set.

Empirical exponential loss function for a score function $F(\mathbf{x})$

$$L_{\text{exp}}^D (F) = \frac{1}{n} \sum_{i=1}^n \exp\{ -y_i F(\mathbf{x}_i) \}$$

Expected exponential loss function for a score function $F(\mathbf{x})$

$$L_{\text{exp}}^E (F) = \int_{\mathbf{X}} \left\{ \sum_{y \in \{+1, -1\}} \exp\{ -y F(\mathbf{x}) \} q(y | \mathbf{x}) \right\} q(\mathbf{x}) d\mathbf{x}$$

where $q(y|\mathbf{x})$ is the conditional distribution given \mathbf{x} , $q(\mathbf{x})$ is the pdf of \mathbf{x} .

Adaboost

1. Initial: $w_1(i) = \frac{1}{n}$ ($i = 1 \cdots n$), $F_0(\mathbf{x}) = 0$

2. For $t = 1, \dots, T$ $\varepsilon_t(f) = \frac{\sum \mathbf{I}(y_i \neq f(\mathbf{x}_i)) w_t(i)}{\sum w_t(i)}$,

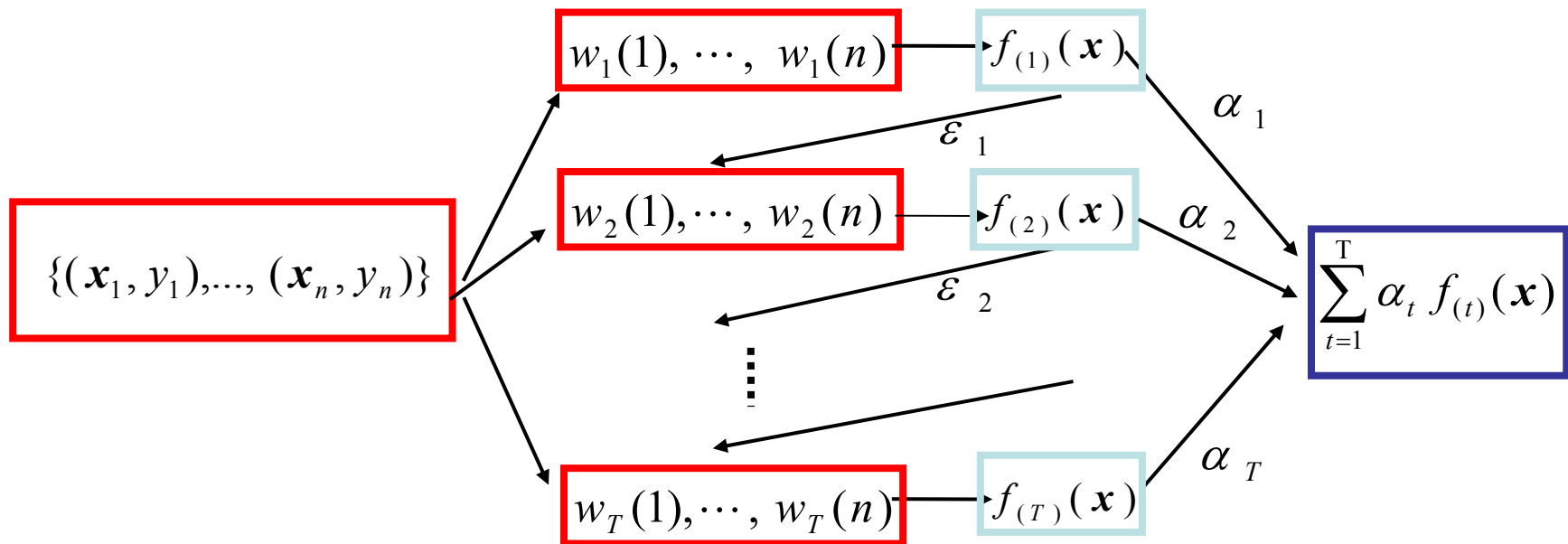
(a) $\varepsilon_t(f_{(t)}) = \min_{f \in \mathcal{F}} \varepsilon_t(f)$

(b) $\alpha_t = \frac{1}{2} \log \frac{1 - \varepsilon_t(f_{(t)})}{\varepsilon_t(f_{(t)})}$

(c) $w_{t+1}(i) = w_t(i) \exp(-\alpha_t f_{(t)}(\mathbf{x}_i) y_i)$

3. $\text{sign}(F_T(\mathbf{x}))$, where $F_T(\mathbf{x}) = \sum_{t=1}^T \alpha_t f_{(t)}(\mathbf{x})$

Learning algorithm



Final learner

$$F_T(\mathbf{x}) = \sum_{t=1}^T \alpha_t f_{(t)}(\mathbf{x})$$

Update weight

update $w_t(i) \rightarrow w_{t+1}(i)$

$f_{(t)}(\mathbf{x}_i) \neq y_i \Rightarrow$ Multiply e^{α_t}

$f_{(t)}(\mathbf{x}_i) = y_i \Rightarrow$ Multiply $e^{-\alpha_t}$

Weighted error rate $\varepsilon_t(f_{(t)}) \rightarrow \boxed{\varepsilon_{t+1}(f_{(t)})} \rightarrow \varepsilon_{t+1}(f_{(t+1)})$

$\varepsilon_{t+1}(f_{(t)}) = \frac{1}{2}$ The worst case

$$\varepsilon_{t+1}(f_{(t)}) = \frac{1}{2}$$

$$\begin{aligned} \varepsilon_{t+1}(f_t) &= \frac{\sum_{i=1}^n I(f_{(t)}(\mathbf{x}_i) \neq y_i) \frac{w_{t+1}(i)}{\sum_{i'=1}^n w_{t+1}(i')}}{\sum_{i=1}^n I(f_{(t)}(\mathbf{x}_i) \neq y_i) \exp\{-\alpha_t y_i f_{(t)}(\mathbf{x}_i)\} w_t(i)} \\ &= \frac{\sum_{i=1}^n \exp\{-\alpha_t y_i f_{(t)}(\mathbf{x}_i)\} w_t(i)}{\sum_{i=1}^n \exp\{\alpha_t\} I(f_{(t)}(\mathbf{x}_i) \neq y_i) w_t(i)} \\ &= \frac{\exp\{\alpha_t\} \sum_{i=1}^n I(f_{(t)}(\mathbf{x}_i) \neq y_i) w_t(i)}{\exp\{\alpha_t\} \sum_{i=1}^n I(f_{(t)}(\mathbf{x}_i) \neq y_i) w_t(i) + \exp\{-\alpha_t\} \sum_{i=1}^n I(f_{(t)}(\mathbf{x}_i) = y_i) w_t(i)} \\ &= \frac{\sqrt{\frac{1 - \varepsilon_t(f_{(t)})}{\varepsilon_t(f_{(t)})}} \varepsilon_t(f_{(t)})}{\sqrt{\frac{1 - \varepsilon_t(f_{(t)})}{\varepsilon_t(f_{(t)})}} \varepsilon_t(f_{(t)}) + \sqrt{\frac{\varepsilon_t(f_{(t)})}{1 - \varepsilon_t(f_{(t)})}} \{1 - \varepsilon_t(f_{(t)})\}} = \frac{1}{2} \quad \blacksquare \end{aligned}$$

Update in exponential loss

$$L_{\text{exp}}(F) = \frac{1}{n} \sum_{i=1}^n \exp\{-y_i F(\mathbf{x}_i)\}$$

Consider $F(\mathbf{x}) \rightarrow F(\mathbf{x}) + \alpha f(\mathbf{x})$

$$\begin{aligned} L_{\text{exp}}(F + \alpha f) &= \frac{1}{n} \sum_{i=1}^n \exp\{-y_i F(\mathbf{x}_i)\} \exp\{-\alpha y_i f(\mathbf{x}_i)\} \\ &= \frac{1}{n} \sum_{i=1}^n \exp\{-y_i F(\mathbf{x}_i)\} \left[e^{\alpha} \mathbf{I}(f(\mathbf{x}_i) \neq y_i) + e^{-\alpha} \mathbf{I}(f(\mathbf{x}_i) = y_i) \right] \\ &= L_{\text{exp}}(F) \{ e^{\alpha} \varepsilon(f) + e^{-\alpha} (1 - \varepsilon(f)) \} \end{aligned}$$

where
$$\varepsilon(f) = \frac{\sum_{i=1}^n \mathbf{I}(f(\mathbf{x}_i) \neq y_i) \exp\{-y_i F(\mathbf{x}_i)\}}{L_{\text{exp}}(F)}$$

Sequential optimization

$$L_{\text{exp}}(F + \alpha f) = L_{\text{exp}}(F) \{ \varepsilon(f) e^{\alpha} + (1 - \varepsilon(f)) e^{-\alpha} \}$$

$$\begin{aligned} & \varepsilon(f) e^{\alpha} + (1 - \varepsilon(f)) e^{-\alpha} \\ &= \left\{ \sqrt{\frac{1 - \varepsilon(f)}{e^{\alpha}}} - \sqrt{\varepsilon(f) e^{\alpha}} \right\}^2 + 2\sqrt{\varepsilon(f) \{1 - \varepsilon(f)\}} \end{aligned}$$

$$\geq 2\sqrt{\varepsilon(f) \{1 - \varepsilon(f)\}}$$

Equality holds if and only if $\alpha_{\text{opt}} = \frac{1}{2} \log \frac{1 - \varepsilon(f)}{\varepsilon(f)}$

AdaBoost = Seq-min exponential loss

$$\min_{\alpha \in \mathbf{R}} L_{\text{exp}}(F_{t-1} + \alpha f_{(t)}) = L_{\text{exp}}(F_{t-1}) \sqrt{\varepsilon(f_{(t)}) \{1 - \varepsilon(f_{(t)})\}}$$

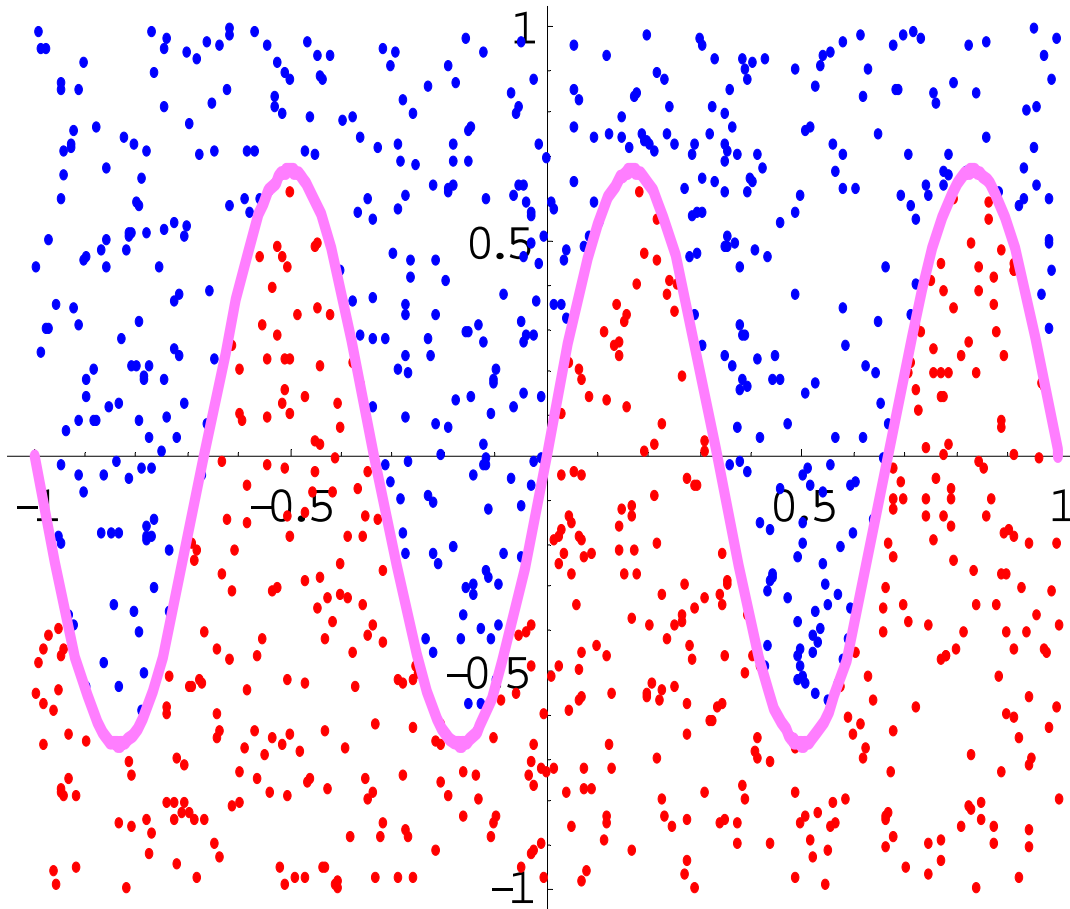
$$\alpha_{\text{opt}} = \frac{1}{2} \log \frac{1 - \varepsilon(f_{(t)})}{\varepsilon(f_{(t)})}$$

(a) $f_{(t)} = \operatorname{argmin}_{f \in F} \varepsilon_t(f)$

(b) $\alpha_t = \operatorname{argmin}_{\alpha \in \mathbf{R}} L_{\text{exp}}(F_{t-1} + \alpha f_{(t)})$

(c) $w_{t+1}(i) \propto w_t(i) \exp\{\alpha_t y_i f_t(x_i)\}$

Simulation (complete separable)



Feature space

$$[-1,1] \times [-1,1]$$

Decision boundary

$$x_2 = \sin(2\pi x_1)$$

$$\{(\mathbf{x}_i, y_i) : i=1, \dots, 1000\}$$

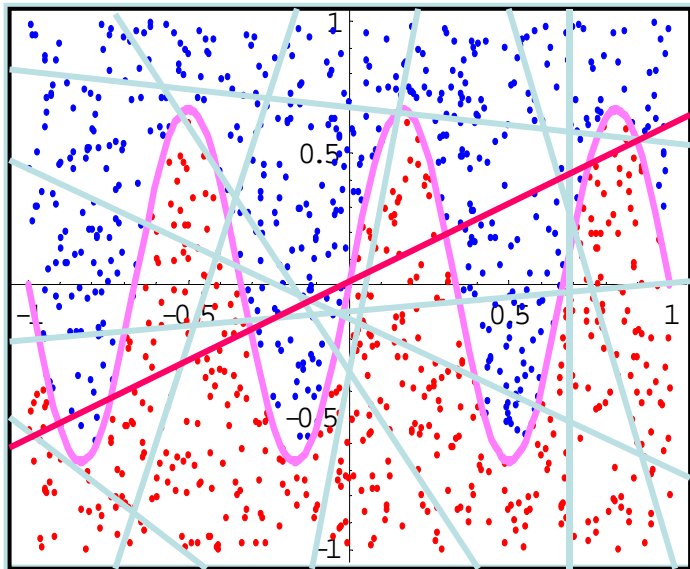
$$\mathbf{x}_i \in [-1, 1] \times [-1, 1]$$

$$y_i \in \{-1, +1\}$$

Set of linear classifiers

Linear classification machines

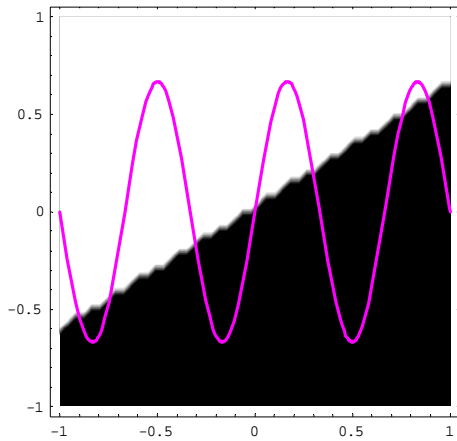
$$f(x_1, x_2) = \text{sgn}(r_1 x_1 + r_2 x_2 + r_3) = \begin{cases} +1 & \text{if } r_1 x_1 + r_2 x_2 + r_3 \geq 0 \\ -1 & \text{if } r_1 x_1 + r_2 x_2 + r_3 < 0 \end{cases}$$



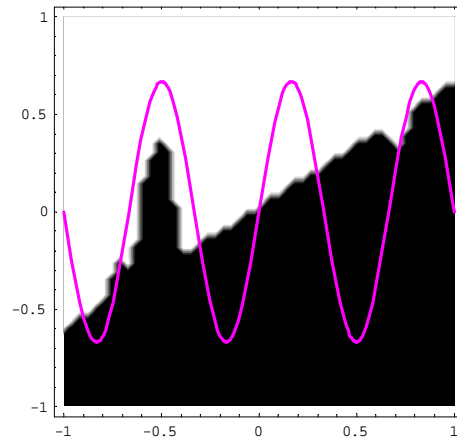
Random generation

$$\{r_1, r_2, r_3\} \sim U(-1, 1)^3$$

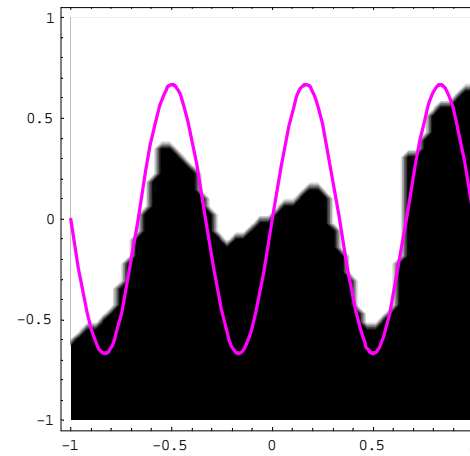
Learning process



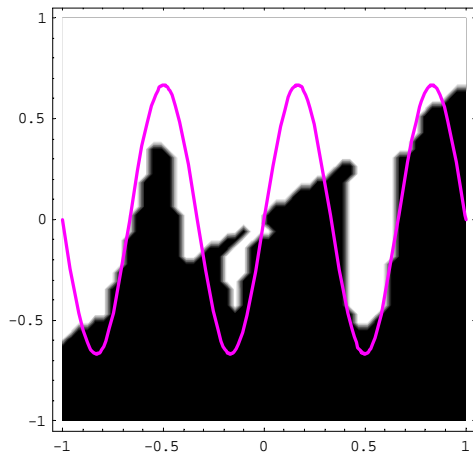
Iter = 1, train err = 0.21



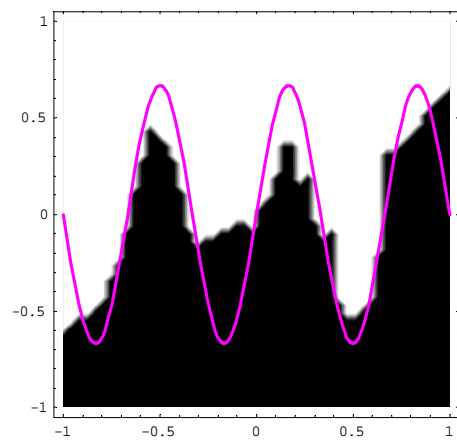
Iter = 13, train err = 0.18



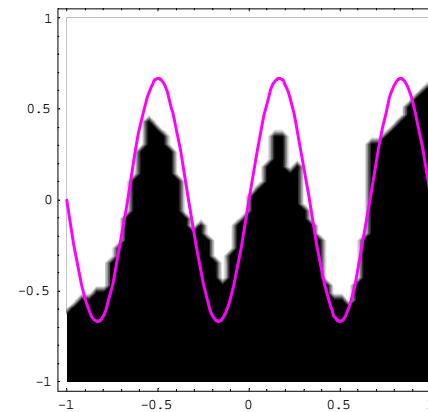
Iter = 17, train err = 0.10



Iter = 23, train err = 0.10

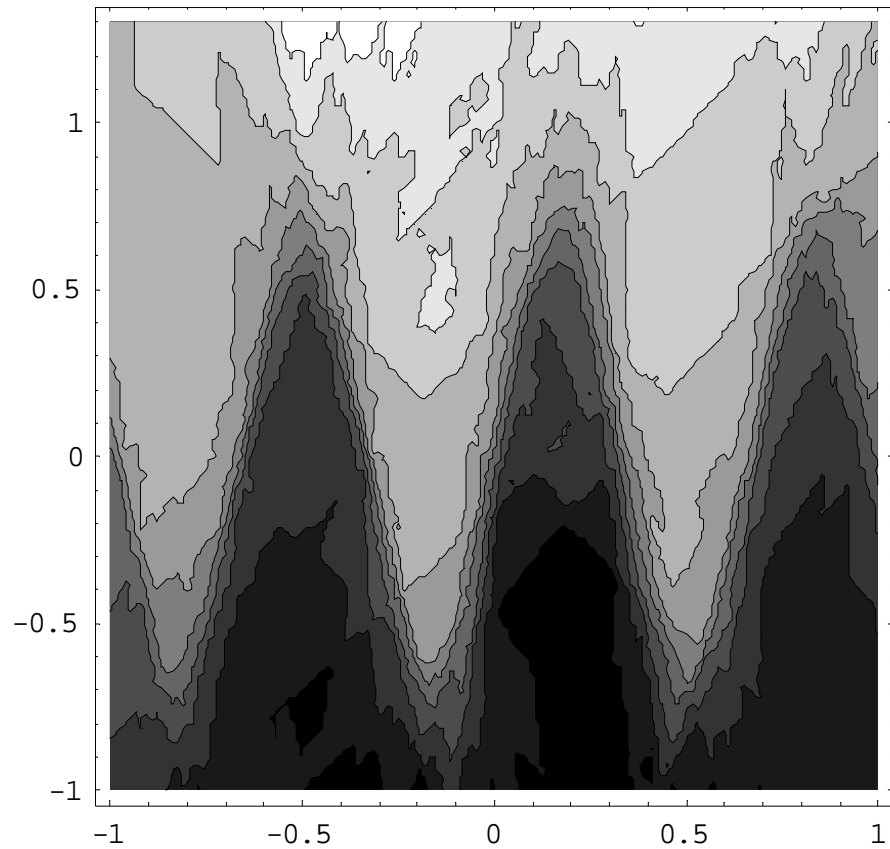


Iter = 31, train err = 0.095

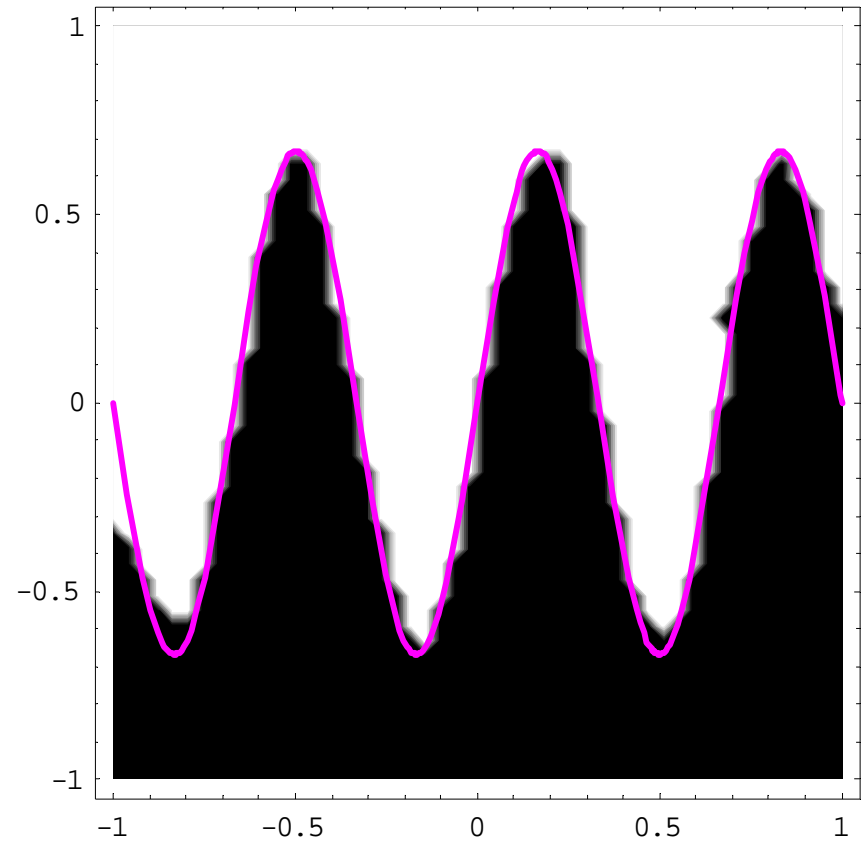


Iter = 47, train err = 0.08

Final decision boundary

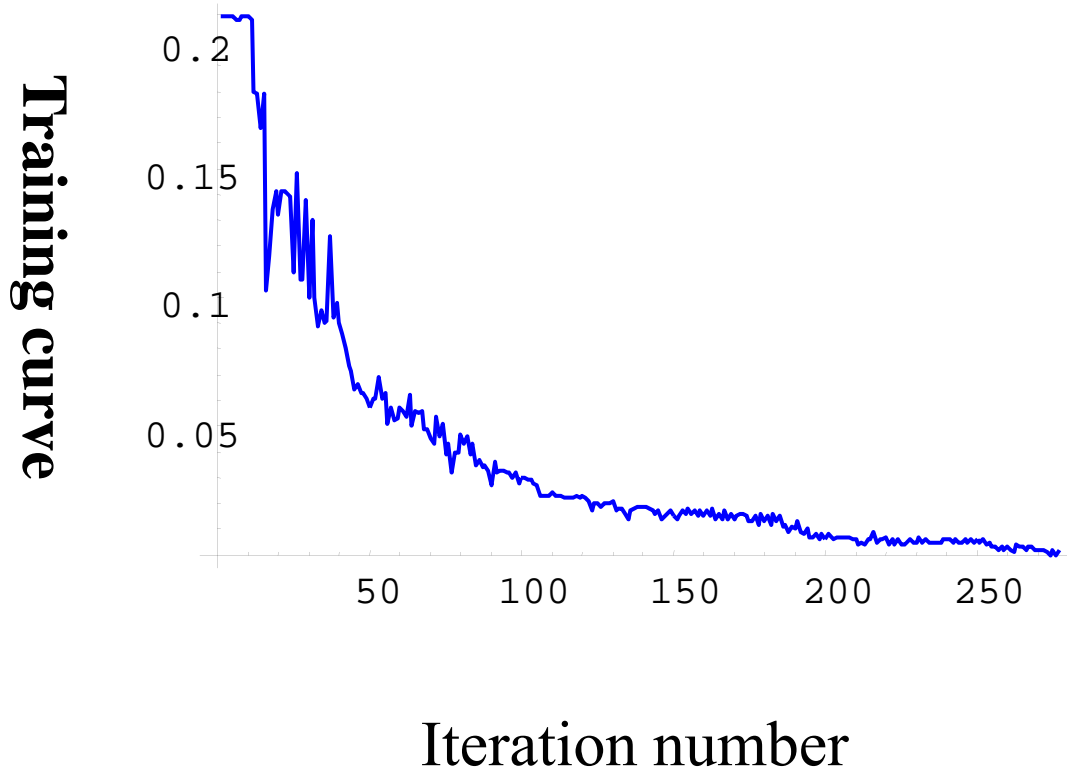


Contour of $F(x)$



Sign($F(x)$)

Learning curve



KL divergence

Label set $Y = \{1, \dots, G\}$ Feature space $X \subseteq \mathbf{R}^p$

Nonnegative function $m(\mathbf{x}, y), \mu(\mathbf{x}, y)$ ($\mathbf{x} \in X, y \in Y$)

$$D_{\text{KL}}(m, \mu) = \int_X \sum_{y=1}^G \left\{ m(\mathbf{x}, y) \log \frac{m(\mathbf{x}, y)}{\mu(\mathbf{x}, y)} - m(\mathbf{x}, y) + \mu(\mathbf{x}, y) \right\} d\mathbf{x}$$

Note

For conditional distribution $p(y|\mathbf{x}), q(y|\mathbf{x})$ given \mathbf{x} with common marginal density $p(\mathbf{x})$ we can write

$$m(\mathbf{x}, y) = p(y | \mathbf{x}) p(\mathbf{x}), \mu(\mathbf{x}, y) = q(y | \mathbf{x}) p(\mathbf{x})$$

Then,

$$D_{\text{KL}}(m, \mu) = \int_X \left\{ \sum_{y=1}^G p(y | \mathbf{x}) \log \frac{p(y | \mathbf{x})}{q(y | \mathbf{x})} \right\} p(\mathbf{x}) d\mathbf{x}$$

Twin of KL loss function

For data distribution $q(\mathbf{x}, y) = q(\mathbf{x}) q(y|\mathbf{x})$ we model as

$$m_1(y | \mathbf{x}) = \sum_{g=1}^G \exp\{F(\mathbf{x}, g) - F(\mathbf{x}, y)\}$$

$$m_2(y | \mathbf{x}) = \frac{\exp\{F(\mathbf{x}, y)\}}{\sum_{g=1}^G \exp\{F(\mathbf{x}, g)\}}$$

Then

$$D_{\text{KL}}(q, m) = \int_{\mathbf{x}} \left\{ \sum_{y \in \{+1, -1\}} \underbrace{q(y | \mathbf{x}) \log \frac{q(y | \mathbf{x})}{m(y | \mathbf{x})}}_{\text{Log loss}} - q(y | \mathbf{x}) + \underbrace{m(y | \mathbf{x})}_{\text{Exp loss}} \right\} q(\mathbf{x}) d\mathbf{x}$$

Log loss

Exp loss

Bound for exponential loss

Empirical exp loss $\bar{L}_{\text{exp}}(F) = \frac{1}{n} \sum_{i=1}^n \exp(-y_i F(\mathbf{x}_i))$

Expected exp loss $L_{\text{exp}}(F) = \mathbb{E}\{\exp(-YF(\mathbf{X}))\}$

Theorem.

Let F be a space of all discriminant functions and

$$F_{\text{opt}} = \arg \min_{F \in F} L_{\text{exp}}(F).$$

Then $F_{\text{opt}}(\mathbf{x}) = \frac{1}{2} \log \frac{p(y = +1 | \mathbf{x})}{p(y = -1 | \mathbf{x})}$.

Variational calculus

$$\begin{aligned}\frac{\partial}{\partial \delta} L_{\text{exp}}(F) &= \mathbb{E} \left[\frac{\partial}{\partial F} \exp(-YF(\mathbf{X})) \right] = \mathbb{E} [Y \exp(-YF(\mathbf{X}))] \\ &= \mathbb{E} [\exp(-F(\mathbf{X})) p(y = +1 | \mathbf{X}) - \exp(F(\mathbf{X})) p(y = -1 | \mathbf{X})] \\ &= \mathbb{E} \left[\exp(-F(\mathbf{X})) p(y = +1 | \mathbf{X}) \left\{ \exp(2F(\mathbf{X})) - \frac{p(y = +1 | \mathbf{X})}{p(y = -1 | \mathbf{X})} \right\} \right]\end{aligned}$$

Hence $F_{\text{opt}}(\mathbf{x}) = \frac{1}{2} \log \frac{p(y = +1 | \mathbf{x})}{p(y = -1 | \mathbf{x})}$.

On AdaBoost

FDA, or Logistic regression

$$F(\mathbf{x}) = \alpha_1 \mathbf{x} + \alpha_0 = \sum_{j=1}^p \alpha_{1j} x_j + \alpha_0$$

- Parametric approach to Bayes classifier
-

AdaBoost

$$F(\mathbf{x}) = \sum_{t=1}^T \alpha_{1t} f_t(\mathbf{x})$$

- Parametric approach to Bayes classifier, but dimension and basis function are flexible
- Each $f_t(\mathbf{x})$ itself is a classifier, cf. real AdaBoost.
- The stopping time T can be selected according to the state of learning.

U-Boost

U-empirical loss function $L_U^{\text{emp}}(\theta) = -\frac{1}{n} \sum_{i=1}^n \xi(q_\theta(\mathbf{x}_i)) + \langle U(\xi(q_\theta)) \rangle$

In a context of classification

$$L_U^{\text{emp}}(F) = -\frac{1}{n} \sum_{i=1}^n F(\mathbf{x}_i, y_i) + \frac{1}{n} \sum_{i=1}^n \sum_{g=1}^G U(F(\mathbf{x}_i, g))$$

Unnormalized U-loss

$$L_U^{(0)}(F) = \frac{1}{n} \sum_{i=1}^n \sum_{g=1}^G U(F(\mathbf{x}_i, g) - F(\mathbf{x}_i, y_i))$$

Normalized U-loss

$$L_U^{(1)}(F) = -\frac{1}{n} \sum_{i=1}^n F(\mathbf{x}_i, y_i) + \frac{1}{n} \sum_{i=1}^n \sum_{g=1}^G U(F(\mathbf{x}_i, g))$$

subject to $\sum_{g=1}^G u(F(\mathbf{x}, g)) = 1$

U-Boost (binary)

Unnormalized U -loss

$$L_U^{(0)}(F) = \frac{1}{n} \sum_{i=1}^n U(-y_i F(\mathbf{x}_i))$$

Note $\sum_{g=\pm 1} U(F(\mathbf{x}_i, g) - F(\mathbf{x}_i, y_i)) = U(-y_i F(\mathbf{x}_i))$ where $F(\mathbf{x}) = \frac{1}{2} \{F(\mathbf{x}, 1) - F(\mathbf{x}, -1)\}$

Bayes risk consistency

$$\frac{\partial}{\partial F(\mathbf{x})} \sum_{y=\pm 1} U(-yF(\mathbf{x})) = \frac{\partial}{\partial F(\mathbf{x})} \{U(-F(\mathbf{x}))p(y=1|\mathbf{x}) + U(F(\mathbf{x}))p(y=-1|\mathbf{x})\}$$

$$\frac{u(F^*(\mathbf{x}))}{u(-F^*(\mathbf{x}))} = \frac{p(y=1|\mathbf{x})}{p(y=-1|\mathbf{x})} \quad \text{where } F^* = \arg \min_F E\{U(-yF(\mathbf{x}))\}$$

$F^*(\mathbf{x})$ is Bayes risk consistent because $\frac{\partial}{\partial F} \frac{u(F)}{u(-F)} = \frac{\dot{u}(F)u(-F) + \dot{u}(-F)u(F)}{\{u(-F)\}^2} > 0$

AdaBoost with margin

Eta-loss function

$$L_{\eta}(F) = (1 - \eta)L_{\text{exp}}(F) + \eta L_{\text{naive}}(F),$$

with generator $U(F) = (1 - \eta) \exp(F) + \eta F$

$$L_{\text{exp}}(F) = \sum_{i=1}^N \exp(-y_i F(\mathbf{x}_i))$$

$$L_{\text{naive}}(F) = - \sum_{i=1}^N y_i F(\mathbf{x}_i)$$

regularized



EtaBoost for Mislabeled

Expected Eta-loss function

$$L_\eta(F) = (1 - \eta)E\{\exp(-yF(\mathbf{x})) - \eta yF(\mathbf{x})\}$$

Optimal score $F^* = \arg \min_F L_\eta(F)$

The variational argument leads to

$$\begin{aligned} p(y = 1 | \mathbf{x}) &= \frac{(1 - \eta)e^{\frac{1}{2}F^*(\mathbf{x})} + \eta}{(1 - \eta)(e^{\frac{1}{2}F^*(\mathbf{x})} + e^{-\frac{1}{2}F^*(\mathbf{x})}) + 2\eta} \\ &= (1 - \varepsilon(\mathbf{x})) \frac{e^{F^*(\mathbf{x})}}{1 + e^{F^*(\mathbf{x})}} + \varepsilon(\mathbf{x}) \frac{1}{1 + e^{F^*(\mathbf{x})}} \end{aligned}$$

Mislabeled modeling

$$\text{where } \varepsilon(\mathbf{x}) = \frac{\eta}{(1 - \eta)(e^{\frac{1}{2}F^*(\mathbf{x})} + e^{-\frac{1}{2}F^*(\mathbf{x})}) + 2\eta}$$

EtaBoost

1. Initial settings : $w_1(i) = \frac{1}{n}$ ($i = 1 \cdots n$), $F_0(\mathbf{x}) = 0$

2. For $m = 1, \dots, T$ $\varepsilon_m(f) \propto \sum_{i=1}^n \mathbf{I}(y_i \neq f(\mathbf{x}_i)) w_m(i)$,

(a) $\varepsilon_m(f_{(m)}) = \min_f \varepsilon_m(f)$

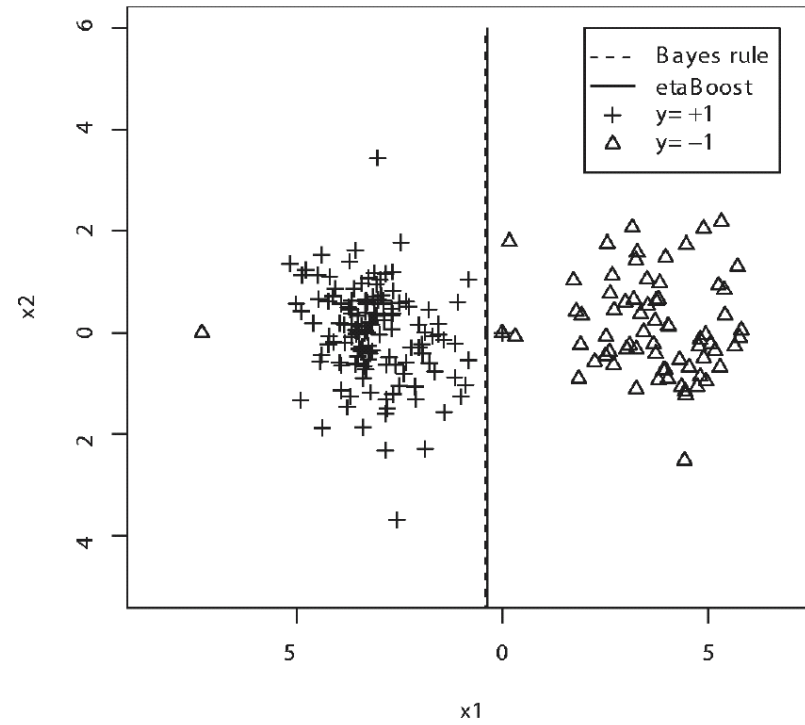
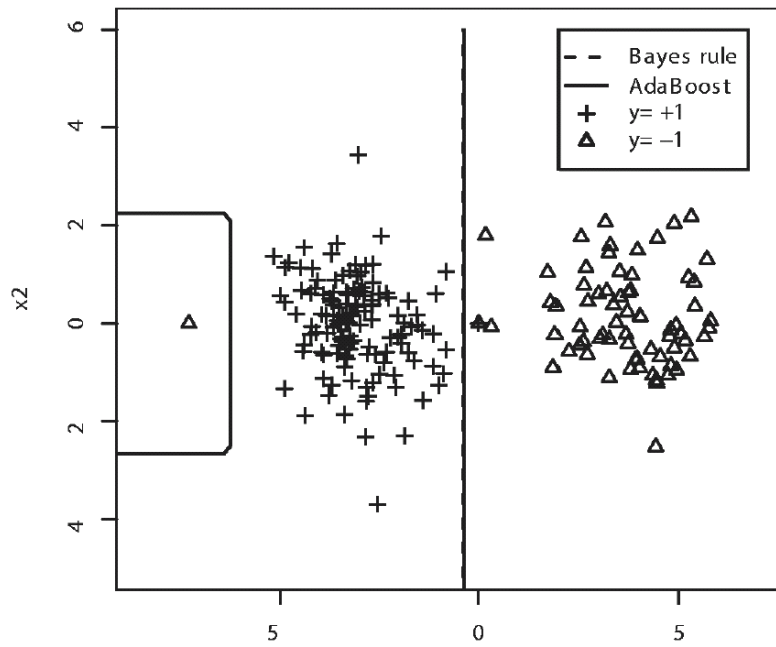
(b)
$$\alpha_m^* = \log \frac{\sqrt{1 - \varepsilon_m(f_m)} + (\eta K_m)^2 + \eta K_m}{\sqrt{\varepsilon_m(f_m)}}$$

where $\varepsilon_m(f)$ is defined in AdaBoost, $K_m = \frac{(1 - 2\varepsilon_1(f_m))}{2\sqrt{\varepsilon_m(f_m)}} \left(\frac{(1 - \eta)Z_m}{N} \right)^{-1}$

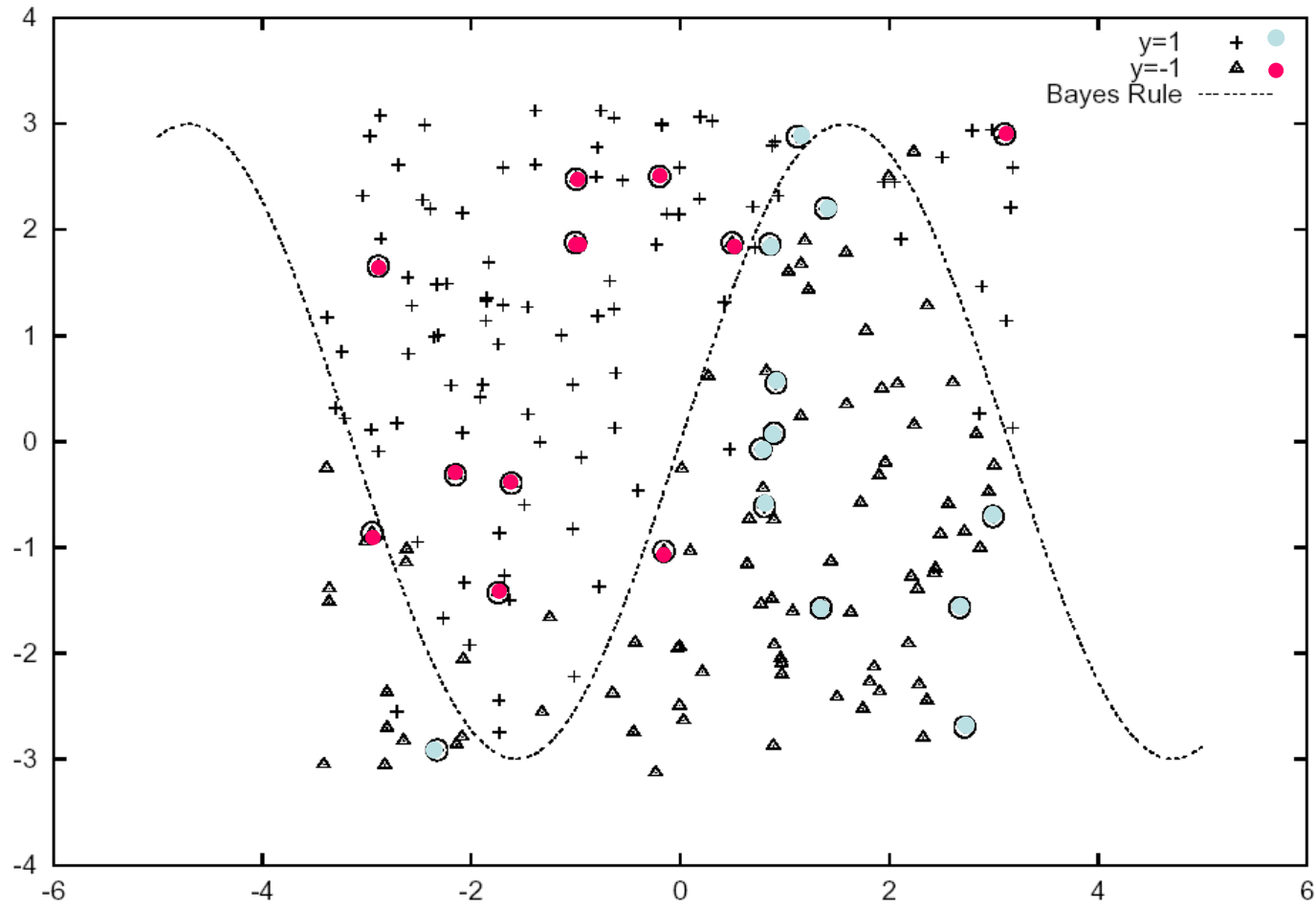
(c) $w_{m+1}(i) \propto w_m(i) \exp(-\alpha_m^* f_{(m)}(\mathbf{x}_i) y_i)$

3. $\text{sign}(F_T(\mathbf{x}))$, where $F_T(\mathbf{x}) = \sum_{t=1}^T \alpha_t f_{(t)}(\mathbf{x})$

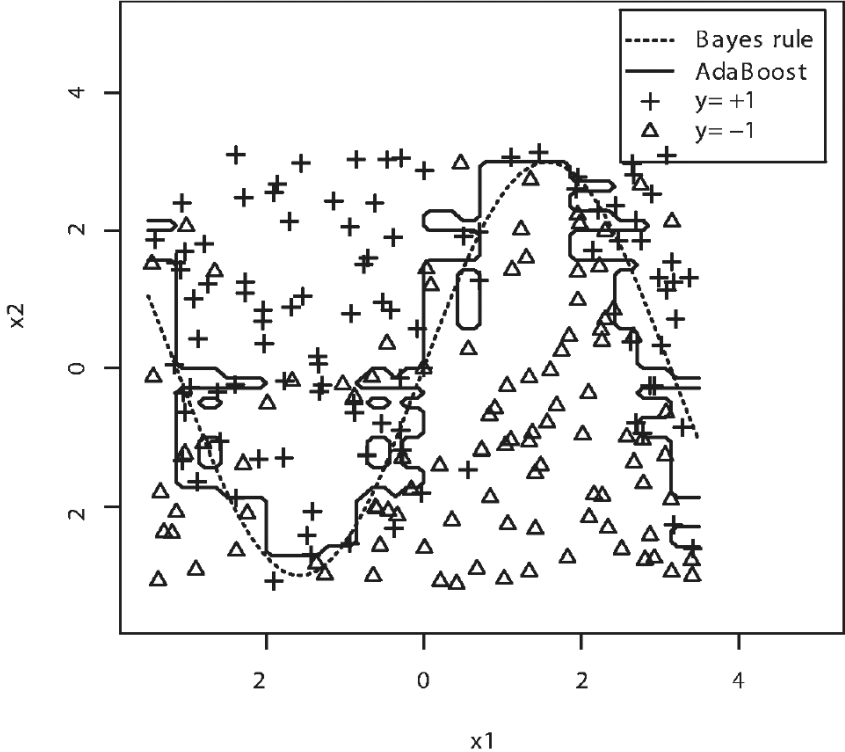
A toy example



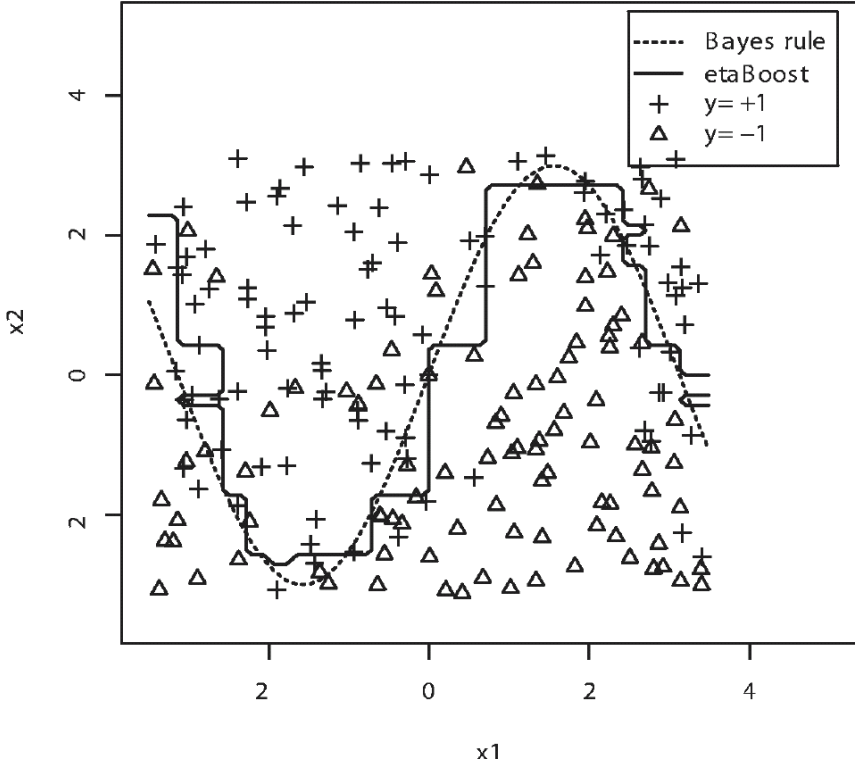
Examples partly mislabeled



AdaBoost vs. EtaBoost



AdaBoost



EtaBoost

ROC curve, AUC

$(\mathbf{X}, Y), \mathbf{X} \in \mathbb{R}^p, Y \in \{0, 1\}$

判別関数: $F(\mathbf{X})$, 閾値: c

{	$F(\mathbf{X}) > c \Rightarrow$ 陽性 ($Y = 1$)
	$F(\mathbf{X}) \leq c \Rightarrow$ 陰性 ($Y = 0$)

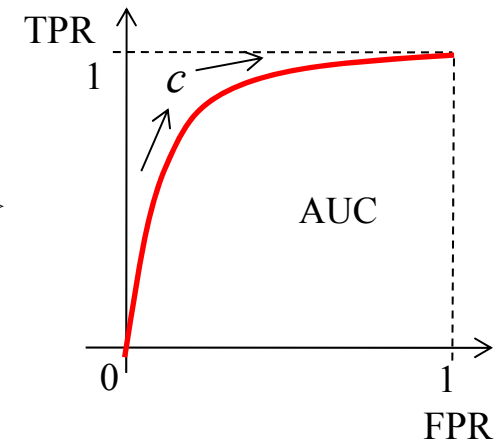
偽陽性確率 $\text{FPR}(c) = P(F(\mathbf{X}) > c | Y = 0)$

真の陽性確率 $\text{TPR}(c) = P(F(\mathbf{X}) > c | Y = 1)$

受信者動作特性 (ROC) 曲線 $\text{ROC}(F) = \{(\text{FPR}(c), \text{TPR}(c)) | c \in \mathbb{R}\}$

Cf. Bamber (1975), Pepe et al. (2003)

下側面積 (AUC) $\text{AUC}(F, \alpha) = \int_{-\infty}^{-\alpha} \text{TPR}(c') d\text{FPR}(c')$



Approximate AUC

Using probability density function g_0 and g_1 for class 0 and class 1, AUC is given as

$$\text{AUC}(F) = \int \int H(F(\mathbf{x}_1) - F(\mathbf{x}_0))g_0(\mathbf{x}_0)g_1(\mathbf{x}_1)d\mathbf{x}_0d\mathbf{x}_1,$$

where H is the Heaviside function. Similarly, the approximated AUC is given as

$$\text{AUC}_\sigma(F) = \int \int H_\sigma(F(\mathbf{x}_1) - F(\mathbf{x}_0))g_0(\mathbf{x}_0)g_1(\mathbf{x}_1)d\mathbf{x}_0d\mathbf{x}_1,$$

where $H_\sigma(x) = \Phi(\frac{x}{\sigma})$, with Φ being the standard normal distribution function.

Relationship between AUC approximate AUC

Theorem 1.

Let

$$\Psi(c) = \text{AUC}_\sigma(F + c m(\Lambda)),$$

where $\Lambda(x) = g_1(x)/g_0(x)$, and m is a strictly increasing function. Then, $\Psi(c)$ is a strictly increasing function of $c \in \mathcal{R}$, and

$$\sup_F \text{AUC}_\sigma(F) = \lim_{c \rightarrow \infty} \Psi(c) = \text{AUC}(\Lambda)..$$

Proof.

$$\frac{\partial}{\partial c} \Psi(c) = \frac{1}{2} \int \int (\zeta(x_1) - \zeta(x_0)) H'_\sigma(F(x_1) + c \zeta(x_1) - F(x_0) - c \zeta(x_0)) g_0(x_0) g_0(x_1) (\Lambda(x_1) - \Lambda(x_0)) dx_0 dx_1$$

where, $\zeta(x) = m(\Lambda(x))$. Then, we have

$$\begin{aligned} \text{AUC}_\sigma(F) &< \lim_{c \rightarrow \infty} \Psi(c) \\ &= \lim_{c \rightarrow \infty} \text{AUC}_\sigma \left[c \left\{ \frac{F}{c} + \zeta \right\} \right] \\ &= \lim_{c \rightarrow \infty} \text{AUC}_\sigma \left(\frac{F}{c} + \zeta \right) \\ &= \text{AUC}(\zeta) \\ &= \text{AUC}(\Lambda), \end{aligned}$$

Objective function of AUCBoost

At first we prepare a set \mathcal{F} , from which we choose weak classifiers to construct $F(\mathbf{x})$:

$$\mathcal{F} = \bigcup_{k=1}^p \{f(\mathbf{x}) = aH(x_k - b_k) + (1 - a)/2 \mid a \in \{-1, 1\}, b_k \in \mathcal{B}_k\},$$

Then, the objective function we propose is:

$$\overline{\text{AUC}}_{\sigma, \lambda}(F) = \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} H_{\sigma}(F(\mathbf{x}_{1j}) - F(\mathbf{x}_{0i})) - \lambda \sum_{k=1}^p \sum_{x_k \in \mathcal{B}_k} \{F_k^{(2)}(x_k)\}^2,$$

where λ is a smoothing parameter and $F_k^{(2)}$ denotes the second-order difference of F_k .

Special relation between smoothing parameter (λ) and scale parameter (σ)

Without loss of generality, the objective function is equivalent to:

Objective function of AUCBoost

$$\overline{\text{AUC}}_{\lambda}(F) = \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} \Phi(F(\mathbf{x}_{1j}) - F(\mathbf{x}_{0i})) - \lambda \sum_{k=1}^p \sum_{x_k \in \mathcal{B}_k} \left\{ F_k^{(2)}(x_k) \right\}^2$$

It is because the objective function is rewritten as:

$$\overline{\text{AUC}}_{\sigma, \lambda}(F) = \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} H_{\sigma}(F(\mathbf{x}_{1j}) - F(\mathbf{x}_{0i})) - \lambda \sigma^2 \sum_{k=1}^p \sum_{x_k \in \mathcal{B}_k} \left\{ \frac{F_k^{(2)}(x_k)}{\sigma} \right\}^2.$$

Hence we have

$$\overline{\text{AUC}}_{\sigma, \lambda}(F) = \overline{\text{AUC}}_{\sigma', \lambda'} \left(\frac{\sigma'}{\sigma} F \right),$$

where $\lambda \sigma^2 = \lambda' \sigma'^2$. This implies that the maximization of $\overline{\text{AUC}}_{\sigma, \lambda}(F)$ is equivalent to that of $\overline{\text{AUC}}_{1, \lambda \sigma^2} \left(\frac{F}{\sigma} \right)$. Therefore, we have

$$\max_{\sigma, \lambda, F} \overline{\text{AUC}}_{\sigma, \lambda}(F) = \max_{\lambda, F} \overline{\text{AUC}}_{1, \lambda}(F).$$

AUCBoost algorithm

1. Start with a score function $F_0(\mathbf{x})$.
2. For $t = 1, \dots, T$
 - a. Find the best weak classifier f_t and calculate the coefficient α_t as

$$f_t(\mathbf{x}) = \operatorname{argmax}_{f \in \mathcal{F}} \left. \frac{\partial}{\partial \alpha} \overline{\text{AUC}}_{\lambda}(F_{t-1} + \alpha f) \right|_{\alpha=0},$$

$$\alpha_t = \operatorname{argmax}_{\alpha > 0} \overline{\text{AUC}}_{\lambda}(F_{t-1} + \alpha f_t).$$

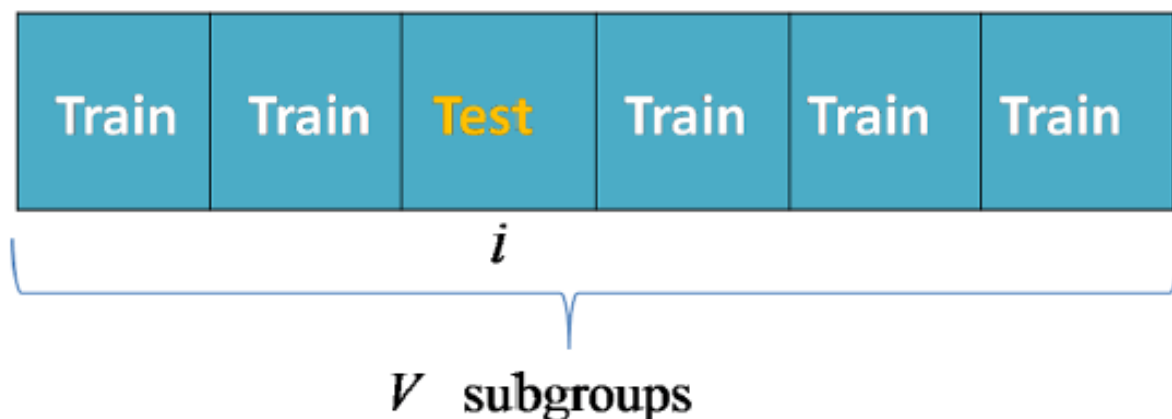
- b. Update the score function as

$$F_t(\mathbf{x}) = F_{t-1}(\mathbf{x}) + \alpha_t f_t(\mathbf{x}).$$

3. Finally, output the final score function

$$F(\mathbf{x}) = F_0(\mathbf{x}) + \sum_{t=1}^T \alpha_t f_t(\mathbf{x})$$

V -fold cross validation

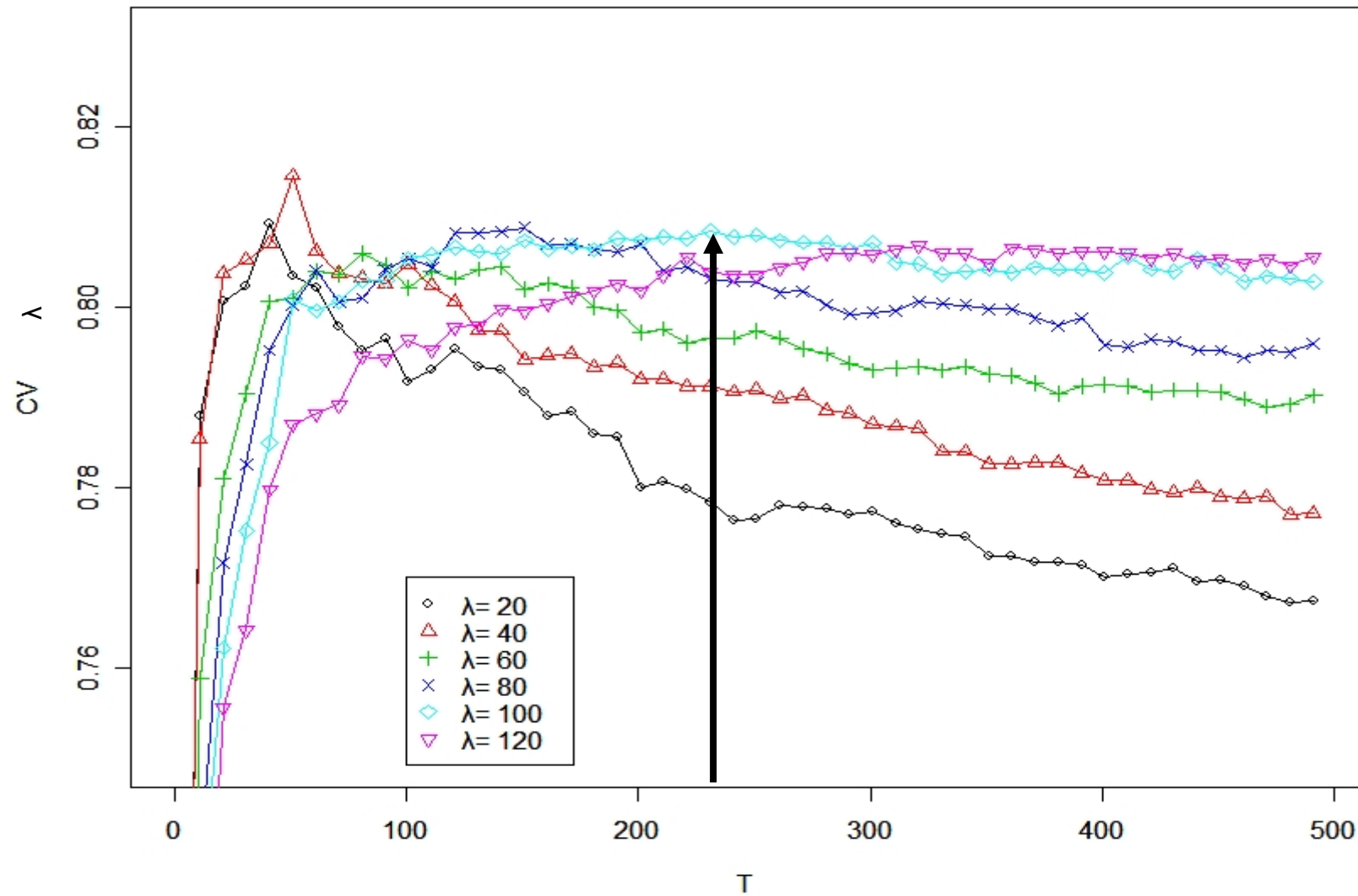


$$\text{AUC}_{\text{cv}}(\lambda, T) = \frac{1}{V} \sum_{i=1}^V \overline{\text{AUC}}_{\lambda}^{(i)}(F^{(-i)})$$

$\overline{\text{AUC}}_{\lambda}^{(i)}$: $\overline{\text{AUC}}_{\lambda}$ calculated based on only the i -th part of the data

$F^{(-i)}$: score function constructed with the i -th part of the data removed

5-fold CVの結果 (AUCBoost)



Score plot

Score plots are useful graphical tools for representing the contribution of each marker to a total score function $F(x)$. When $F(x)$ is decomposed as

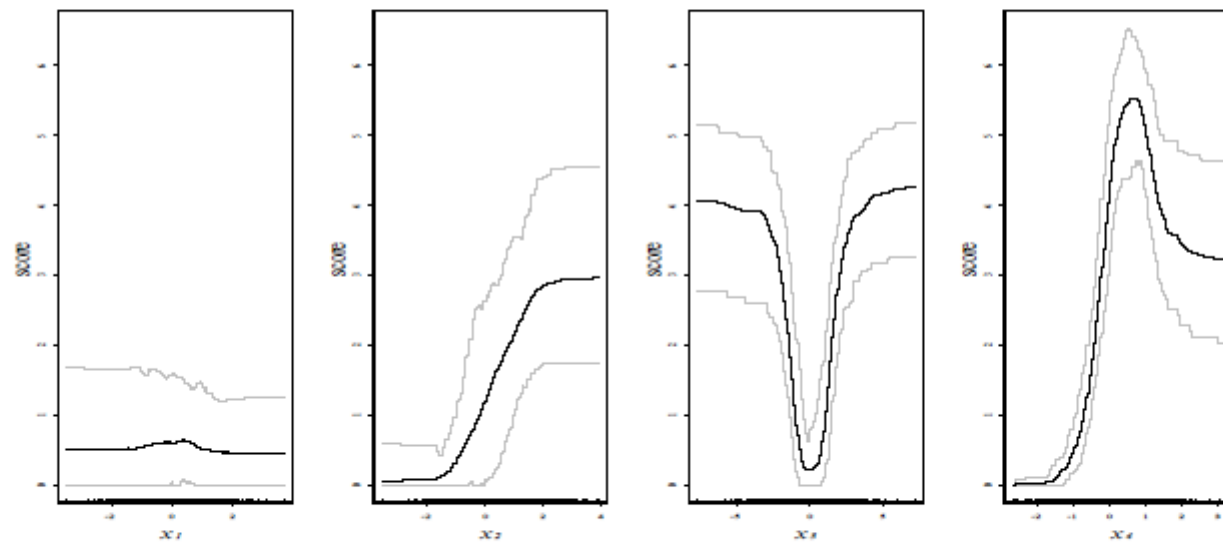
Definition

$$F(x) = \sum_{k=1}^p F_k(x_k), \quad x (= (x_1, \dots, x_p)) \in \mathbf{R}^p$$

Each plot of $F_k(x_k)$ against x_k is called a score plot. See Friedman *and others* (2000) and Kawakita *and others* (2005).

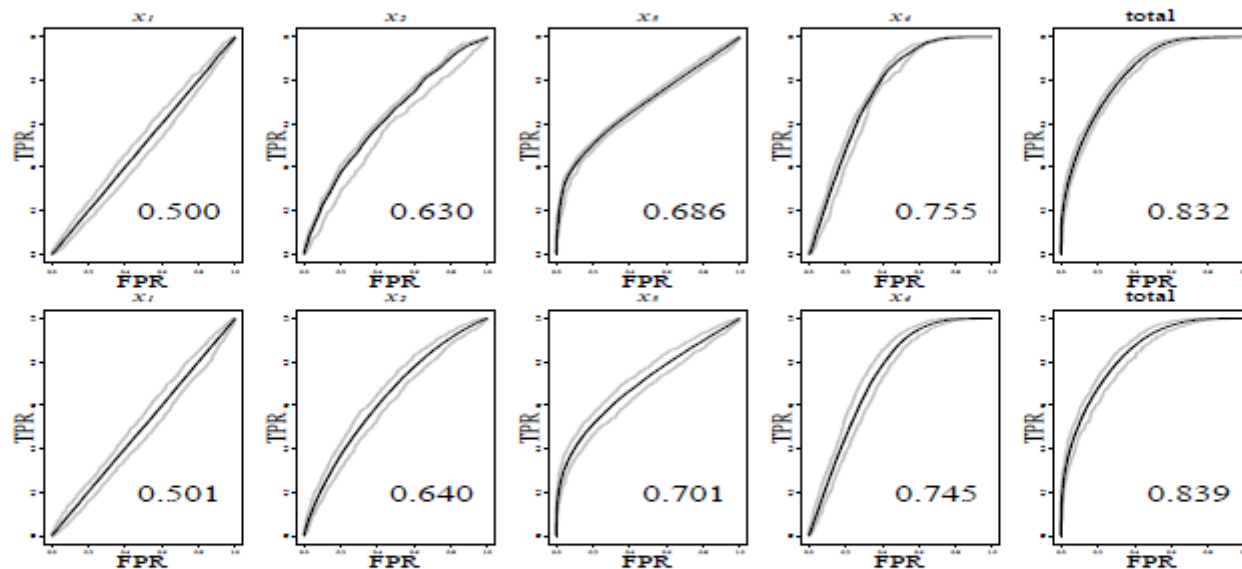
Example of score plots under normality assumption ($X_0 \sim N(\mu_0, \Sigma_0)$, $X_1 \sim N(\mu_1, \Sigma_1)$)

$\mu_0 = (0, 0, 0, 0)'$, $\mu_1 = (0, 0.5, 0, 0.5)'$, $\Sigma_0 = \text{diag}(1, 1, 1, 1)$ and
 $\Sigma_1 = \text{diag}(1, 1, 4, 0.25)$



Score plots for AUCBoost. The black lines indicate mean score plots and the gray lines indicate the 95 percent confidence bands.

Corresponding score AUCs under normality assumption ($X_0 \sim N(\mu_0, \Sigma_0)$, $X_1 \sim N(\mu_1, \Sigma_1)$)



Score ROC curves and the ROC curve constructed by AUCBoost (upper panels) and by \hat{F}_N (lower panels). The black lines indicate mean ROC curves and the gray lines indicate the 95 percent confidence bands. The corresponding mean values of AUCs are also calculated.

2標本問題とAUCブーストの関係

AUCブーストの学習アルゴリズムは逐次的に各 t -ステップで得られた判別関数 $F_t(x)$ のウィルコクソン検定統計量

$$\hat{C}(F_t) = \frac{1}{n_0 n_1} \sum_{k=1}^{n_1} \sum_{i=1}^{n_0} I(F_t(\mathbf{x}_{1k}) > F_t(\mathbf{x}_{0i}))$$

を $(t+1)$ -ステップ

$$F_{t+1}(\mathbf{x}) = F_t(\mathbf{x}) + \alpha_{t+1} f_{t+1}(\mathbf{x})$$

で最大方向に増大するように設計されている.

このように

$$\dots < \hat{C}(F_t) < \hat{C}(F_{t+1}) < \dots$$

となる中で全ての単点解析が統合されている.

この意味で全ての単点解析はブースティングに

組むことは可能である. Eguchi-Komori (2009)

The partial area under the ROC curve

Setting

$\mathbf{x} \in \mathbf{R}^p$: marker vector

$y \in \{0, 1\}$: class label

$F(\mathbf{x})$: score function, c : threshold

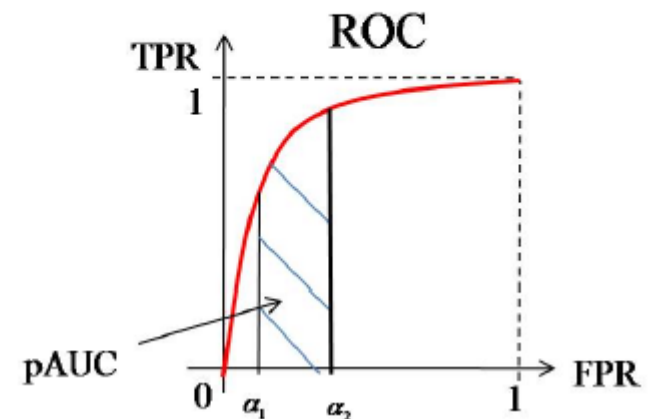
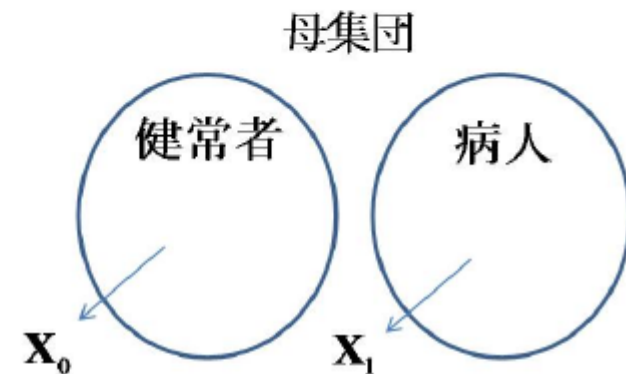
$$\begin{cases} F(\mathbf{x}) > c \Rightarrow \text{positive} \\ F(\mathbf{x}) \leq c \Rightarrow \text{negative} \end{cases}$$

$$\text{FPR}(c) = P(F(\mathbf{X}) > c | y = 0)$$

$$\text{TPR}(c) = P(F(\mathbf{X}) > c | y = 1)$$

⇓

$$\text{ROC} = \{(\text{FPR}(c), \text{TPR}(c)) | c \in \mathbf{R}\}$$

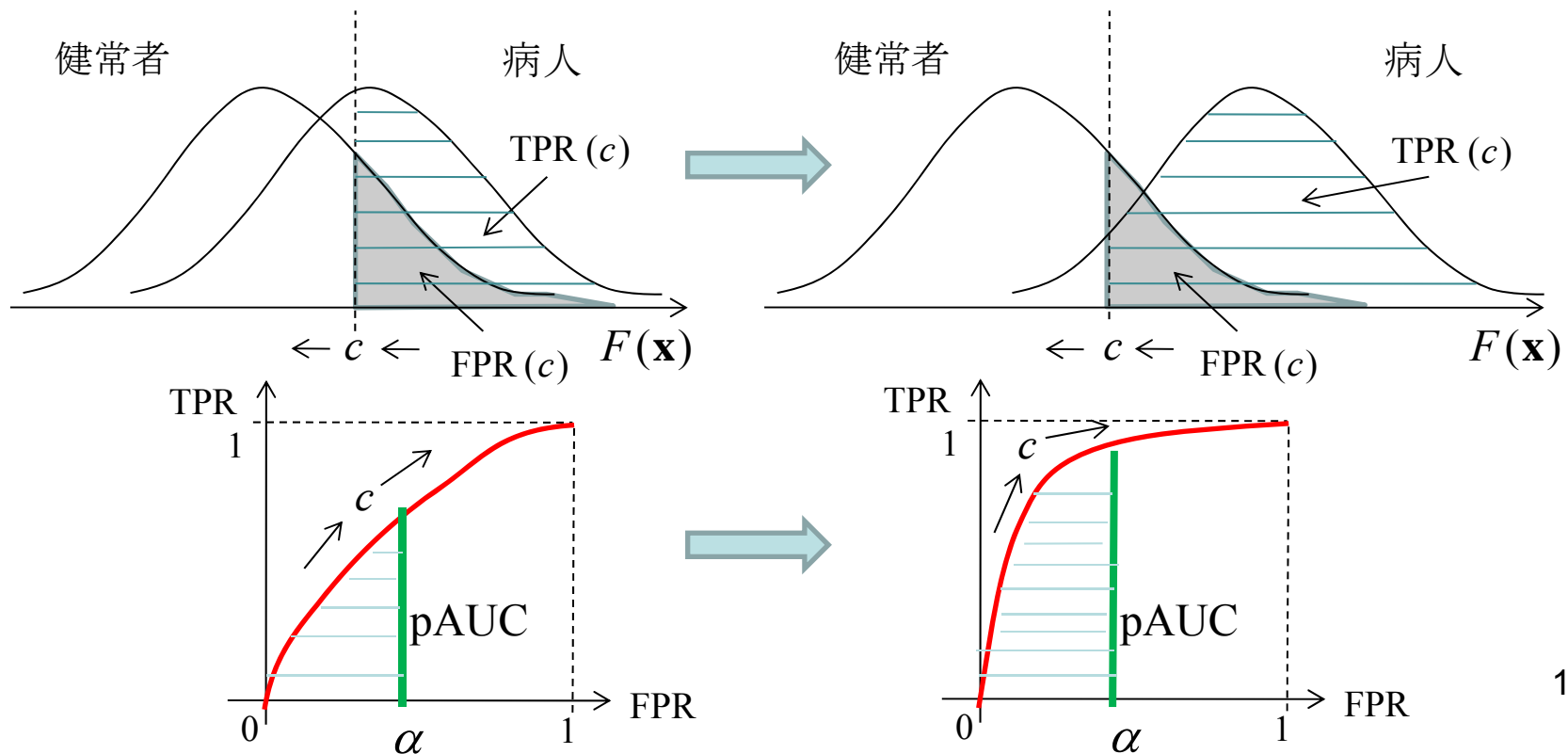


From AUC to pAUC

- The AUC has been severely criticized for inconsistencies arising between the statistical significance and the corresponding clinical significance (Cook, 2007).
- Recently, Pencina *and others* (2008) proposed a criterion termed integrated discriminant improvement, and showed its advantage over the AUC in the assessment of a new marker.

Partial AUC

$$\begin{aligned} \text{pAUC} &= \int_c^{-\infty} \text{TPR}(c) d\text{FPR}(c) \\ &= P(F(X_0) < F(X_1), c \leq F(X_0)) \end{aligned}$$



pAUCBoost algorithm

1. Set the initial value of score function as $F_0(\mathbf{x}) = 0$, and the absolute value of the initial coefficient as $|\beta_0(f)| = 1$.
2. For $t = 1, \dots, T$,
 - a. For all f 's $\in \mathcal{F}$, calculate the values of thresholds \bar{c}_1 and \bar{c}_2 of $F_{t-1} + \beta_{t-1}(f)f$
 - b. Update the coefficient with one-step Newton-Raphson iteration:

$$\beta_{t-1}(f) \rightarrow \beta_t(f)$$

- c. Find the best weak classifier f_t .

$$f_t = \operatorname{argmax}_{f \in \mathcal{F}} \overline{\text{pAUC}}_{\lambda}(F_{t-1} + \beta_t(f)f, \alpha_1, \alpha_2)$$

- d. Update the score function as.

$$F_t(\mathbf{x}) = F_{t-1}(\mathbf{x}) + \beta_t(f_t)f_t(\mathbf{x}).$$

3. Finally, output a final score function as

Robust for noninformative genes

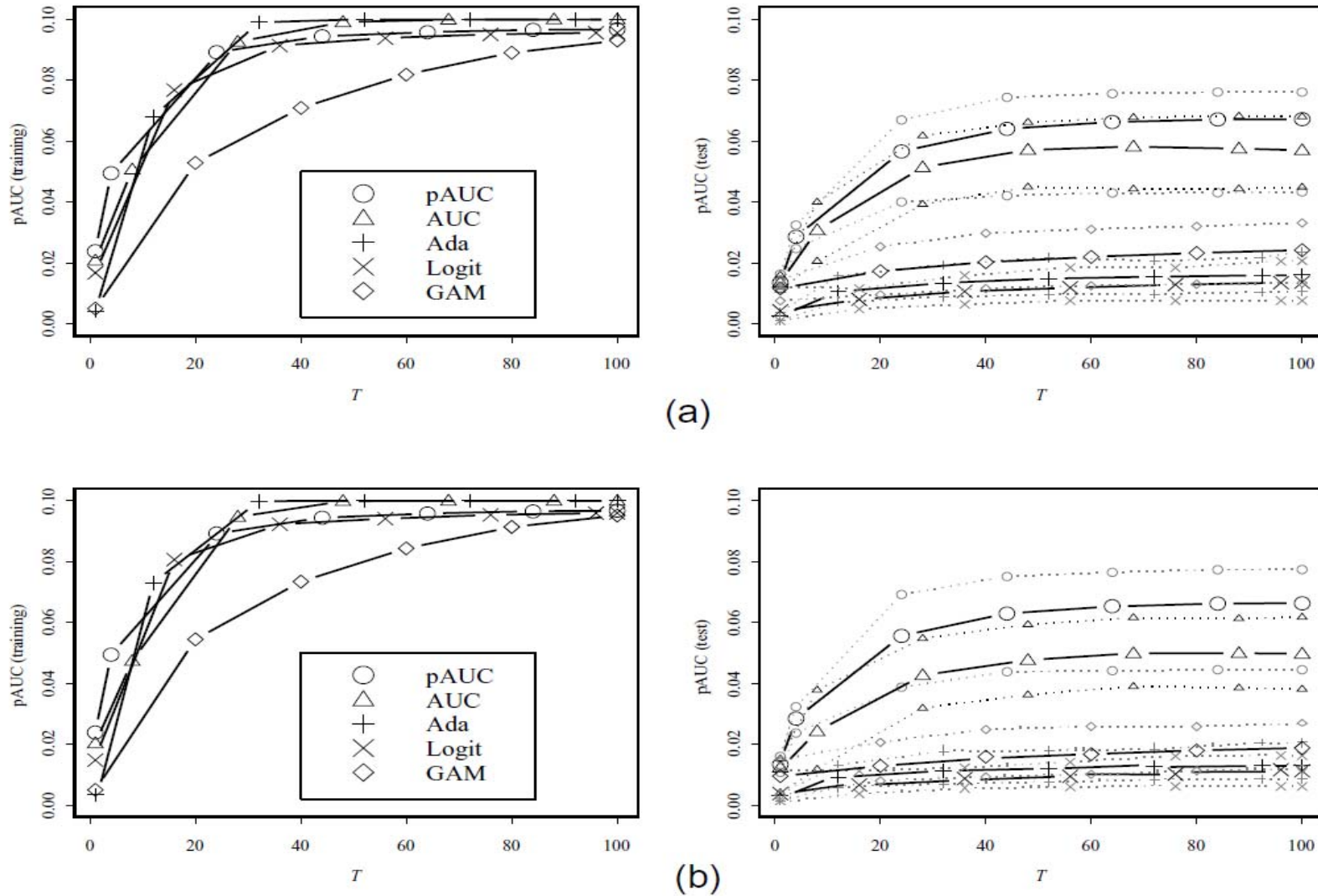


Figure 3 Results of simulation study based on the values of the pAUC. (a) The results of the pAUC with FPR between 0 and 0.1 for training data (left panel) and test data (right panel) with only informative genes. The gray dashed lines indicate the 95% confidence bands. (b) the results of the pAUC with noninformative genes added.

Robustness for noninformative genes

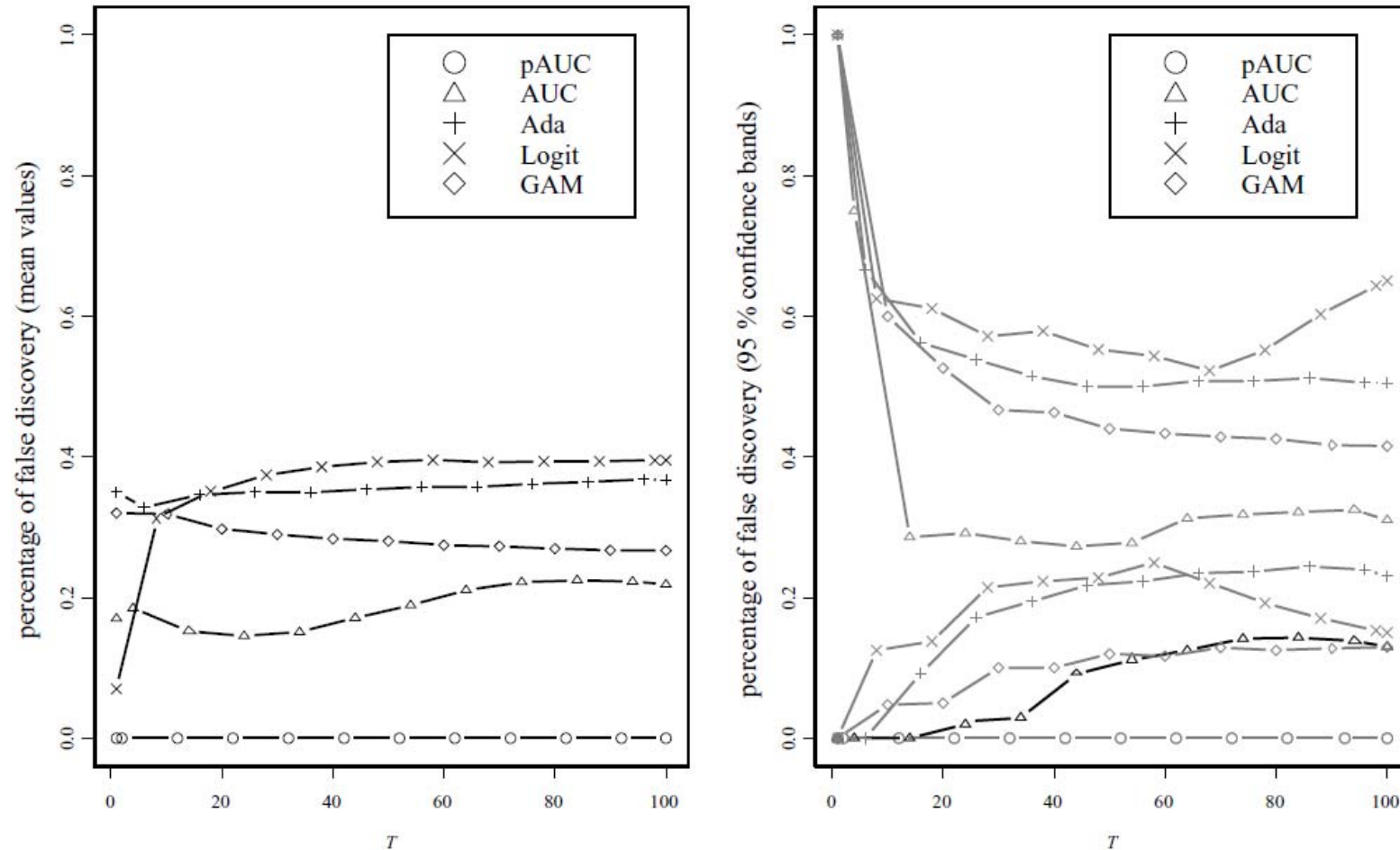


Figure 4 Results of simulation study based on the marker selection. The mean values of percentage of false discovery (left panel) and 95% confidence bands (right panel) for each boosting method. The horizontal axis denotes the iteration number of T . The lower sides of the 95% confidence bands of AUCBoost are shown by the heavy black line to emphasize the difference from those of pAUCBoost.

An Extension of the Receiver Operating Characteristic Curve and AUC-Optimal Classification

Takashi Takenouchi Osamu Komori Shinto Eguchi

Neural Computation 24, 2789–2824 (2012)

生成関数

$$U(0) = 0, \quad U(z) \leq 1, \quad \lim_{z \rightarrow \infty} U(z) = 1.$$

真の U 陽性確率

$$\text{TPU}_U(F, c) = \int_x U(F(x) - c) p_{+1}(x) dx,$$

$$\begin{aligned} A_U(F) &= \int_{c: \infty}^{-\infty} \text{TPU}_U(F, c) d\text{FPR}(F, c) \\ &= \int_x \int_{x'} U(F(x) - F(x')) p_{+1}(x) p_{-1}(x') dx dx' \end{aligned}$$

$$\tilde{A}_U(F) = \frac{1}{n_{+1} n_{-1}} \sum_{i: y_i = +1} \sum_{j: y_j = -1} U(F(x_i) - F(x_j))$$

Boost learning algorithm

1. Initialize a weight for a pair of examples as $w_1(i, j) = \frac{1}{n_{+1}n_{-1}}$
2. For $t = 1, \dots, T$,
 - a. Find a base classifier as

$$f_t(x) = \operatorname{argmin}_{f \in \mathcal{F}} (\operatorname{err}_t(f) - \operatorname{acc}_t(f)),$$

$$\operatorname{acc}_t(f) = \sum_{i:y_i=+1} \sum_{j:y_j=-1} w_t(i, j) \mathbb{I}(f(x_i) = +1) \mathbb{I}(f(x_j) = -1),$$

$$\operatorname{err}_t(f) = \sum_{i:y_i=+1} \sum_{j:y_j=-1} w_t(i, j) \mathbb{I}(f(x_i) = -1) \mathbb{I}(f(x_j) = +1).$$

- b. $\alpha_t = \operatorname{argmax}_{\alpha'} \tilde{A}_U(F_{t-1} + \alpha' f_t).$

- c. Update the weight as

$$w_{t+1}(i, j) = \frac{U'(F_t(x_i) - F_t(x_j))}{Z_t},$$

3. Output the discriminant function $F_T(x)$.

Bayes risk consistency

Let $F^* = \operatorname{argmax}_F A_U(F)$.

Then the likelihood ratio is

$$\frac{p_{+1}(x)}{p_{-1}(x)} = \exp(2\Psi_{F^*}(F^*(x))) \quad \text{where} \quad \Psi_F(z) = \frac{1}{2} \log \frac{\int U'(F(x') - z) p_{+1}(x') dx'}{\int U'(z - F(x')) p_{-1}(x') dx'}.$$

Remark. The conditional probability of $Y = y$ given \mathbf{x} is in logistic model

$$p(y|x) = \frac{1}{1 + \exp\left(-2y \left(\Psi_{F^*}(F^*(x)) + \log \sqrt{\frac{\pi_{+1}}{\pi_{-1}}}\right)\right)}$$

ベンチマークデータ

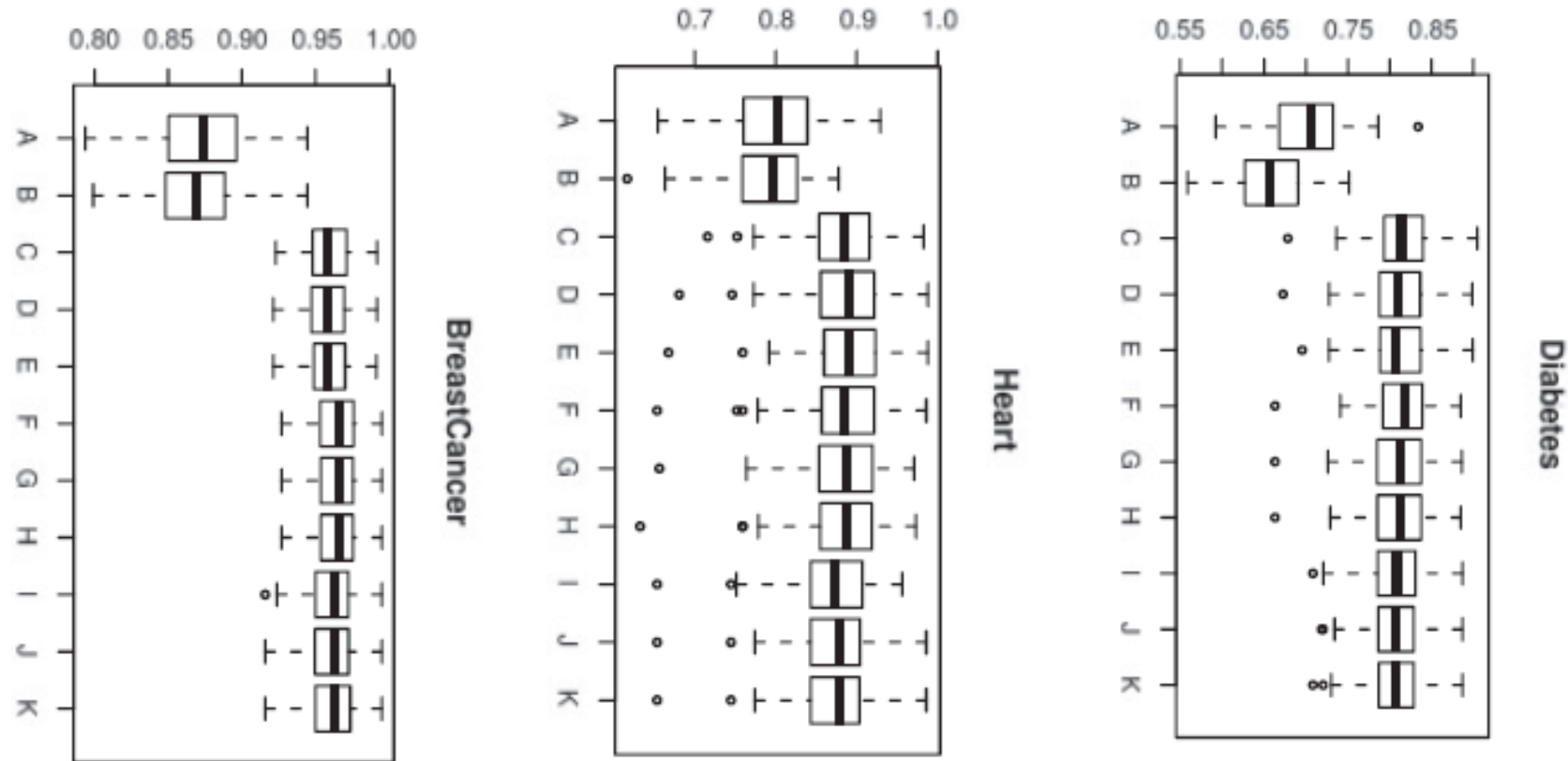
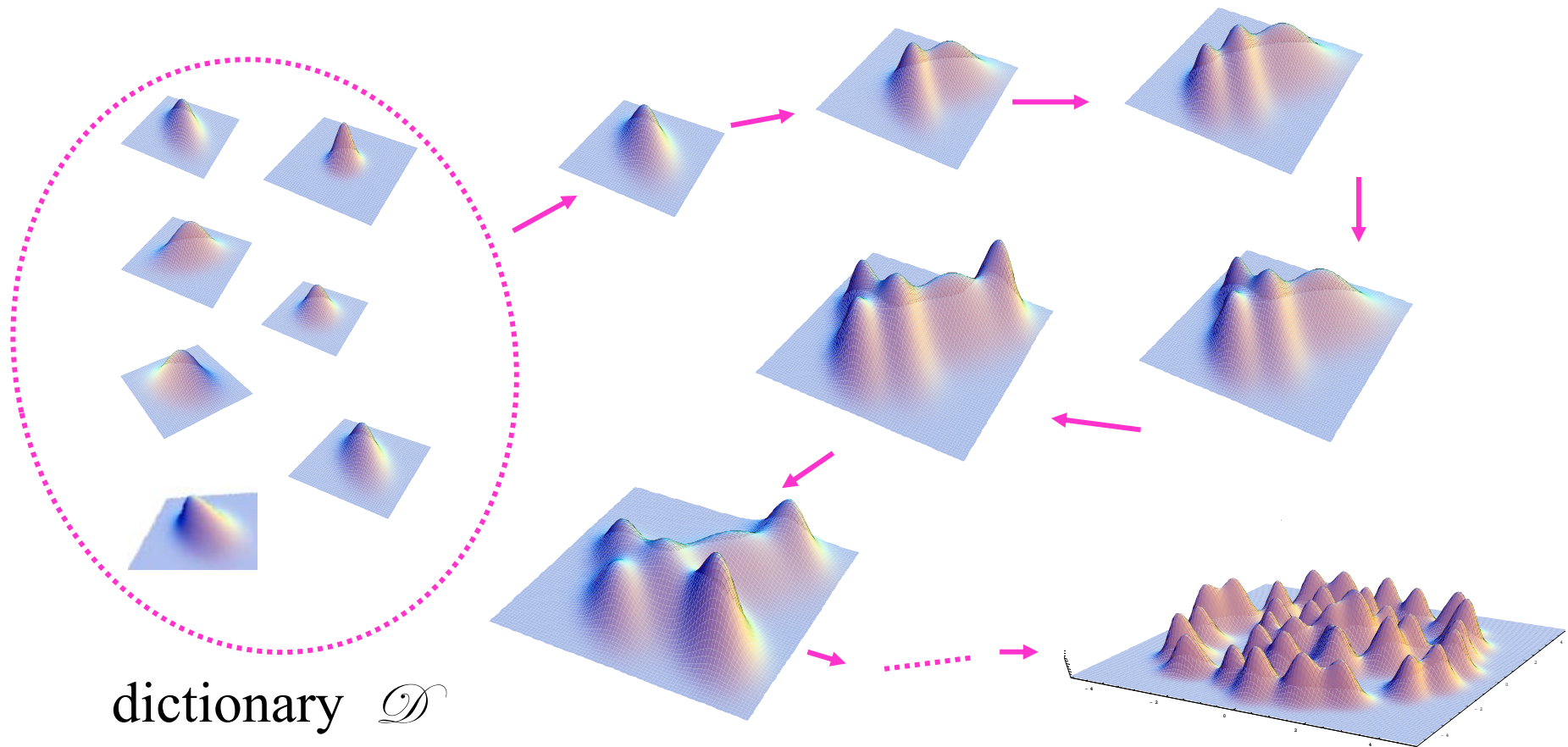


Figure 2: Box plots of AUC values for each data set. (A) SVM with RBF kernel. (B) SVM with linear kernel. (C) Boosting with U_{rank} . (D) Boosting with U_{logit} . (E) Boosting with U_{mada} . (F) U -AUCBoost with U_{rank} . (G) U -AUCBoost with U_{logit} . (H) U -AUCBoost with U_{mada} . (I) pU -AUCBoost with U_{rank} . (J) pU -AUCBoost with U_{logit} . (K) pU -AUCBoost with U_{mada} .

Density estimation with minimization of U -divergence

Kanta Naito · Shinto Eguchi



U-Boost learning

U-loss function $L_U(f) = -\frac{1}{n} \sum_{i=1}^n \xi(f(\mathbf{x}_i)) + \int_{\mathbb{R}^p} U(\xi(f(\mathbf{x}))) d\mathbf{x}$

Dictionary of density functions

$$\mathcal{D} = \{ g_\lambda(\mathbf{x}) : g_\lambda(\mathbf{x}) \geq 0, \int g_\lambda(\mathbf{x}) d\mathbf{x} = 1, \lambda \in \Lambda \}$$

Learning space = U-model

$$\mathcal{D}_U^* = \xi^{-1}(\text{co}(\xi(\mathcal{D}))) = \left\{ \xi^{-1} \left(\sum_{\lambda \in \Lambda} \pi_\lambda \xi(g_\lambda(\mathbf{x})) \right) \right\}$$

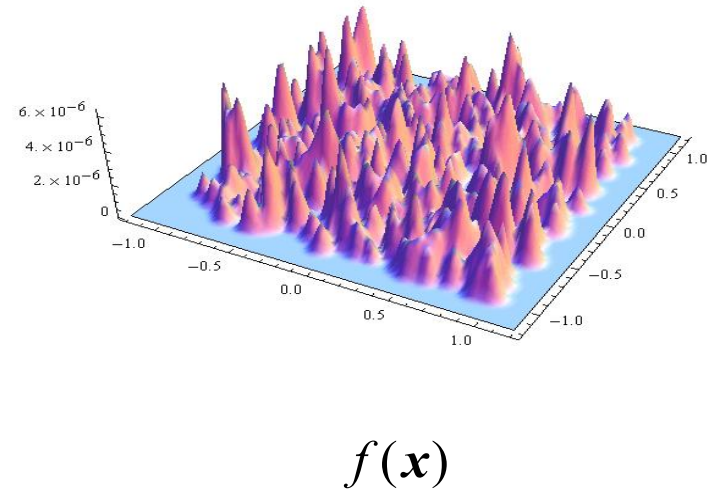
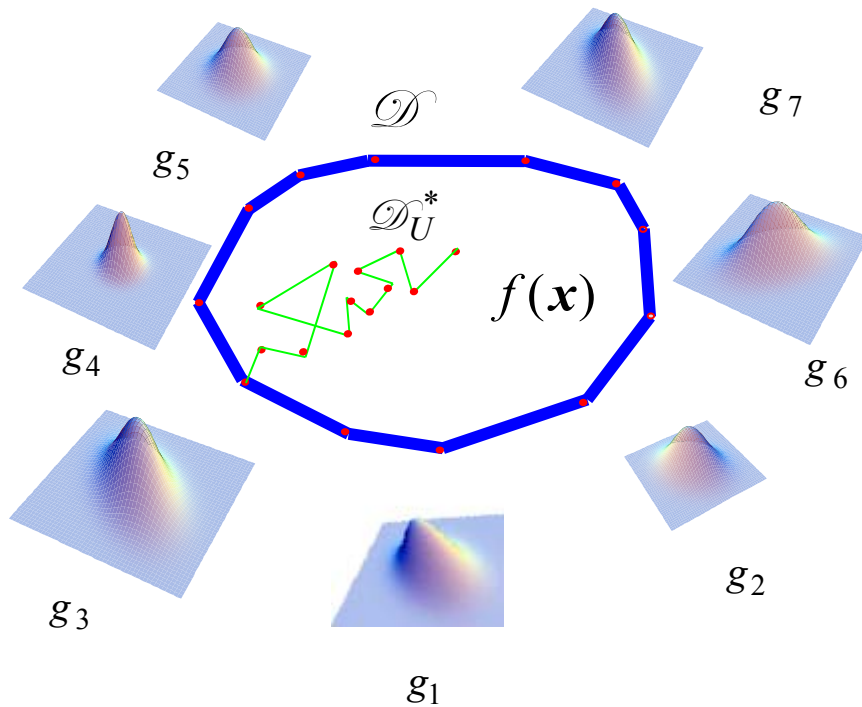
- Let $f(\mathbf{x}, \boldsymbol{\pi}) = \xi^{-1} \left(\sum_{\lambda \in \Lambda} \pi_\lambda \xi(g_\lambda(\mathbf{x})) \right)$. Then $f(\mathbf{x}, (0, \dots, 1, \dots, 0)) = g_\lambda(\mathbf{x})$
(λ)
- $\mathcal{D}_U^* \supseteq \mathcal{D}$

Goal : find $f^* = \underset{f \in \mathcal{D}_U^*}{\text{argmin}} L_U(f)$

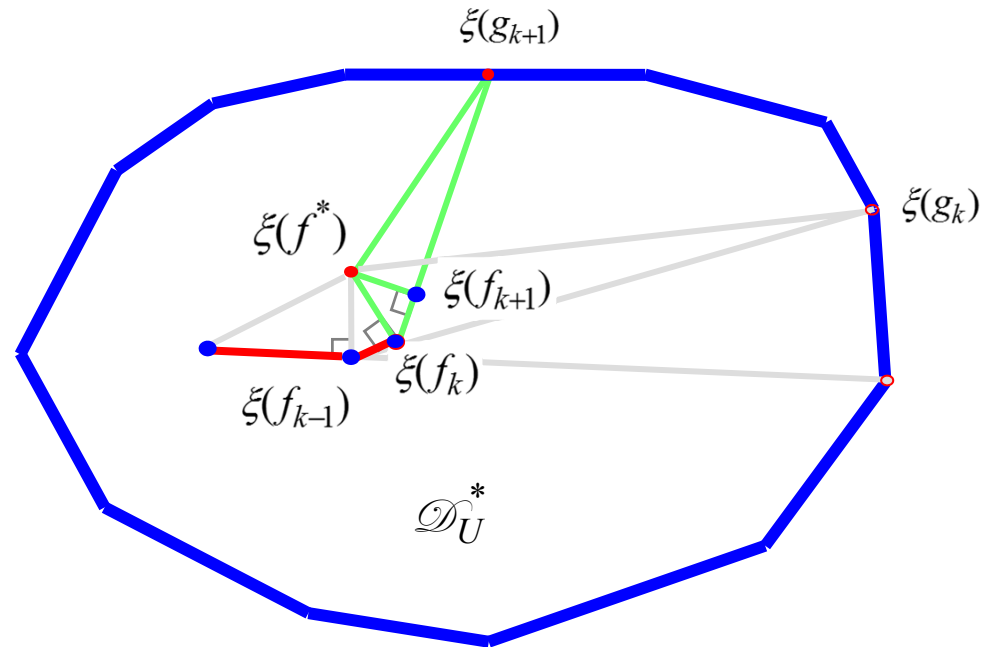
Inner step in the convex hull

$$\mathcal{D} = \{\phi(x, \lambda) : \lambda \in \Lambda\} \longrightarrow \mathcal{D}_U^* = \xi^{-1}(\text{co}\xi(\mathcal{D}))$$

Goal : $f^* = \underset{f \in \mathcal{D}_U^*}{\text{argmin}} L_U(f)$ $f^*(\mathbf{x}) = \xi^{-1}(\hat{\pi}_1 \xi(\hat{f}_1(\mathbf{x})) + \dots + \hat{\pi}_{\hat{k}} \xi(\hat{f}_{\hat{k}}(\mathbf{x})))$



Convergence of Algorithm



$$L_U(f_k) = L_U(f_{k+1}) + D_U(f_k, f_{k+1})$$

$$\xi(f_1) \rightarrow \dots \rightarrow \xi(f_{k-1}) \rightarrow \xi(f_k) \rightarrow \xi(f_{k+1}) \rightarrow \dots \rightarrow \xi(f^*)$$

$$L_U(f_1) = L_U(f_{k+1}) + \sum_{j=1}^k D_U(f_j, f_{j+1})$$

ξ -mixture

$$\xi(f_K) = (1 - \alpha_{K-1})\xi(f_{K-1}) + \alpha_{K-1}\xi(g_K)$$

$$= \pi_1 \xi(g_1) + \dots + \pi_K \xi(g_K) \rightarrow \xi(f^*)$$

Simplified boosting

$$\hat{f}_K = \xi^{-1}\left(\sum_{k=1}^K \pi_k \xi(g_k)\right)$$

For $k = 0, \dots, K-1$,

$$f_k \longrightarrow f_{k+1} = \xi^{-1}\left((1 - \alpha_{k+1})\xi(f_k) + \alpha_{k+1}\xi(g_{k+1})\right)$$

such that

$$g_{k+1} = \arg \min_{g \in \mathcal{D}} L_U(\xi^{-1}\left((1 - \alpha)\xi(f_k) + \alpha\xi(g)\right)),$$

$$\alpha_{k+1} = \frac{c}{k+c}$$

with the initial $f_1 = \arg \min_{g \in \mathcal{D}} L_U(g)$

Non-asymptotic bound

Theorem. Assume that a data distribution has a density $g(\mathbf{x})$ and that

$$(A) \quad \sup_{(\psi, \phi, \varphi) \in \text{co}(\xi(\mathcal{D})) \times \mathcal{D} \times \mathcal{D}} \int \ddot{U}(\psi) \{\xi(\phi) - \xi(\varphi)\}^2 \leq b_U$$

Then we have

$$\mathbf{E}_g D_U(g, \hat{f}_K) \leq \text{FA}(g, \mathcal{D}_U^*) + \text{EE}(g, \mathcal{D}) + \text{IE}(K, \mathcal{D}),$$

where

$$\text{FA}(g, \mathcal{D}_U^*) = \inf_{f \in \mathcal{D}_U^*} D_U(g, f) \quad (\text{Functional approximation})$$

$$\text{EE}(g, \mathcal{D}) = 2 \mathbf{E}_g \left\{ \sup_{f \in \mathcal{D}} \left| \frac{1}{n} \sum_{i=1}^n \xi(f(\mathbf{x}_i)) - \mathbf{E}_p \xi(f) \right| \right\} \quad (\text{Estimation error})$$

$$\text{IE}(K, \mathcal{D}) = \frac{c^2 b_U^2}{K + c - 1} \quad (c: \text{step-length constant}) \quad (\text{Iteration effect})$$

Remark. Trade between $\text{FA}(p, \mathcal{D}_U^*)$ and $\text{EE}(p, \mathcal{D})$

Lemmas

Lemma 1 Let \tilde{f} be a density estimator and $f^* = \xi^{-1}(\sum_{j=1}^N p_j \xi(\phi_j(\mathbf{x})))$. Then

$$\sum_{j=1}^N p_j L_U(\xi^{-1}((1-\pi)\xi(\tilde{f}(\mathbf{x})) + \pi\xi(\phi_j(\mathbf{x})))) - L_U(f^*) \leq (1-\pi)\{L_U(\tilde{f}) - L_U(f^*)\} + \Delta_U(\tilde{f}, f^*, \pi).$$

Lemma 2 Let \tilde{f} be in $\text{co}(\xi(\mathcal{D}))$. Then

$$\Delta_U(f, f^*, \pi) \leq \pi^2 b_U^2 \quad (\forall \pi \in [0,1])$$

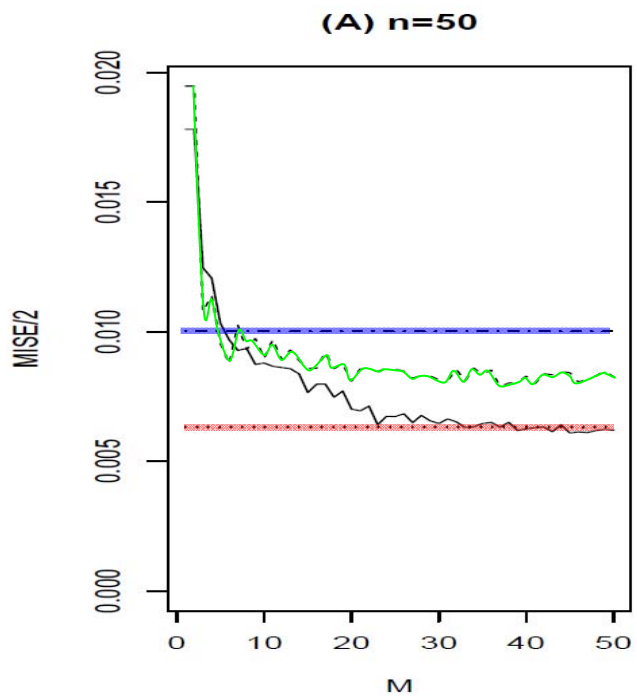
Lemma 3 Under assumption (A)

$$L_U(\hat{f}_K) \leq \inf_{\lambda} L_U(f_{\lambda}) + \frac{c^2 b_U^2}{K+c-1} \quad \text{where } f_{\lambda} = \xi^{-1}(\sum \lambda_j \xi(\phi_j))$$

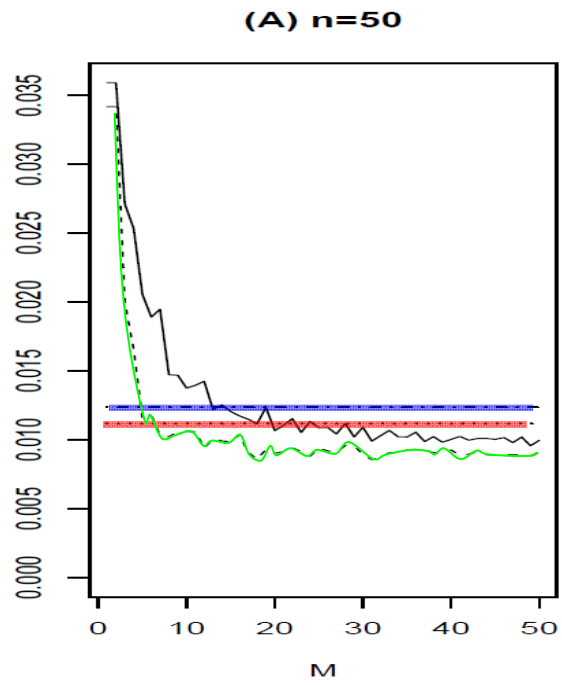
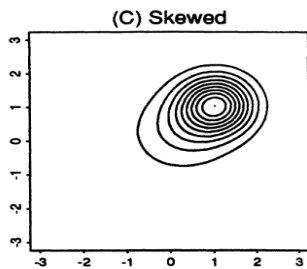
Lemma 4 If $L_U(\hat{f}_K) \leq \inf_{\lambda} L_U(f_{\lambda}) + \tau$, then

$$\mathbf{E}_g D_U(g, \hat{f}_K) \leq \inf D_U(g, f_{\lambda}) + \tau + 2 \mathbf{E}_g \sup_{f \in \mathcal{D}} \left\| \frac{1}{n} \sum \xi(f(\mathbf{x}_i)) - \mathbf{E}_g \xi(f) \right\|$$

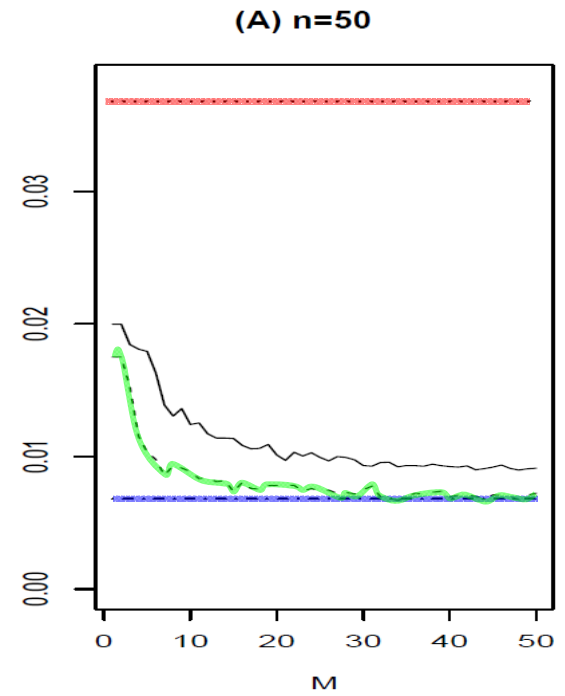
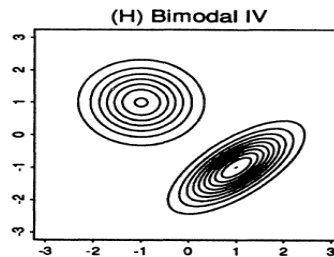
----- β -Boost ($\beta=1$)
----- β -Boost ($\beta=1/2$)
----- KDE
----- RSDE



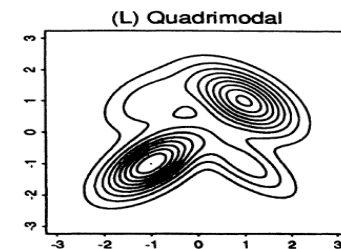
C (skewed-unimodal)



H (bimodal)

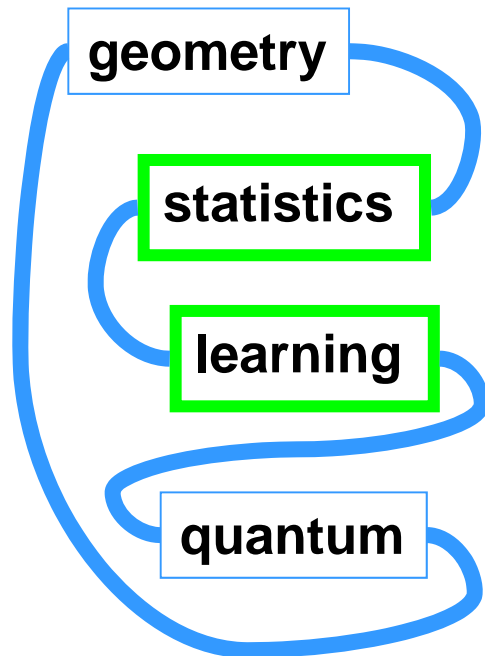


L (quadrimodal)



Future of information geometry

Deepening



Poincare conjecture

Optimal transform

Free probability

Compressed sensing

Semiparametrics

Dynamical systems

Mathematical evolutionary theory

Dempster-Shafer theory

Poincaré conjecture 1904

Every simply connected, closed 3-manifold is
homeomorphic to the 3-sphere

Ricci flow by Richard Hamilton in 1981

Perelman 2002, 2003

Optimal control theory due to Pontryagin and Bellman

Wasserstein space

$$d_W(f, g) = \min\{\sqrt{E\|X - Y\|^2} : X \sim f, Y \sim g\}$$

Optimal transport

$$\phi_{f, g} = \arg \min\{\sqrt{E\|X - \phi(X)\|^2} : X \sim f, \phi(X) \sim g\}$$

Optimal transport

Theorem (Brenier, 1991) There exists a convex function Φ such that $\nabla\Phi = \phi_{f,g}$

Assume that $\exists K > 0$ st $-\mathbf{u}^\top \left\{ \frac{\partial}{\partial \mathbf{x} \partial \mathbf{x}^\top} \log g(\mathbf{x}) \right\} \mathbf{u} \geq K \|\mathbf{u}\|^2$

Talagrand inequality $D(f, g) \geq \frac{K}{2} d_W(f, g)^2$

Log Sobolev inequality $D(f, \Phi^2 g) \leq \frac{2}{K} E_g \|\nabla\Phi\|^2$

Optimal transport theory is extending on a general manifold

C. Villani (2010) Optimal transport theory as Fields medal

Geometers consider not a space, but a distribution family on the space.

Thank you