

総研大夏期大学院 2012年9月19日 20日, 統数研

# 最尤理論、誕生から百年の統計学と 今後の方向への一考察

江口 真透

統計数理研究所

Summer-semester in Guas, 19-20 September, 2012, ISM

**Maximum likelihood theory  
from Ronald Fisher in 1912  
with discussion on future directions**

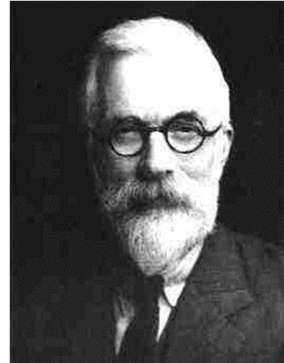
**Shinto Eguchi**

**Institute of Statistical Mathematics**

# Outline

- **MLE 100 years from Fisher 1912**
- **Information geometry**
- **Power entropy and divergence**
- **Two cultures in statistics and AdaBoost**
- **Poincaré conjecture and Optimal transport theory**
- **Likelihood for equation model**

# R. A. Fisher



Statistics

Population genetics

**Ronald A. Fisher**  
**(1890-1962)**

analysis of variance

maximum likelihood estimation

Design of experiment

Randomization test

Linear discriminant analysis

Genetical Theory of Natural Selection

Fundamental theorem of natural selection

Wright-Fisher model

Fisher's equation



## ON AN ABSOLUTE CRITERION FOR FITTING FREQUENCY CURVES.

By *R. A. Fisher*, Gonville and Caius College, Cambridge.

1. IF we set ourselves the problem, in its essence one of frequent occurrence, of finding the arbitrary elements in a function of known form, which best suit a set of actual observations, we are met at the outset by an arbitrariness which appears to invalidate any results we may obtain. In the general problem of fitting a theoretical curve, either to an observed curve, or to an observed series of ordinates, it is, indeed, possible to specify a number of different standards of conformity between the observations and the theoretical curve, which definitely lead to different though mutually approximate results. This mutual approximation, though convenient in practice in that it allows a computer to make a legitimate choice of the method which is arithmetically simplest, is harmful from the theoretical standpoint as tending to obscure the practical discrepancies, and the theoretical indefiniteness which actually exist.

# Maximum likelihood



**Fisher (1912)**

$$\int \{f(x, \theta) - y(x)\}^2 dx \quad \sum_{i=1}^n \{f(x_i, \theta) - y(x_i)\}^2$$

---

$$\log P = \int y(x) \log f(x, \theta) dx$$

$$\log P' = \sum_{i=1}^n \log f(x_i, \theta)$$

The most probable set of values for the  $\theta$ 's will make  $P$  a maximum.

**MLE 1912-2012**, Aldrich (1997)

Efficiency, sufficiency, Fisher information,...

Boltzmann-Shannon entropy, Kullback-Leibler divergence

Max entropy distribution (exponential model)

# Parametric Statistical Modeling by Minimum Integrated Square Error

David Scott Technometrics (2001)

The likelihood function plays a central role in parametric and Bayesian estimation, as well as in nonparametric function estimation via local polynomial modeling. However, integrated square error has enjoyed a long tradition as the goodness-of-fit criterion of choice in nonparametric density estimation. In this article, I investigate the use of integrated square error, or  $L_2$  distance, as a theoretical and practical estimation tool for a variety of *parametric* statistical models. I show that the asymptotic inefficiency of the parameters estimated by minimizing the integrated square error or  $L_2$  estimation ( $L_2E$ ) criterion

$$\int \{f(x, \theta) - y(x)\}^2 dx \leftarrow \frac{?}{?} \rightarrow \int y(x) \log f(x, \theta) dx$$

# Asymptotic efficiency

**Statistical model**

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

**Maximum likelihood estimator**

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \sum_{i=1}^n \log f(X_i, \theta)$$

**Asymptotic normal**

If  $\{X_i\}_{i=1}^n \stackrel{\text{iid}}{\sim} f(x, \theta)$ , then

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{D} N(0, I_{\theta}^{-1})$$

**Asymptotic efficient**

$\text{var}_A(\hat{\theta}) \leq \text{var}_A(\tilde{\theta})$  for any CAN estimator  $\tilde{\theta}$

**Estimative distribution**

$$f(x, \hat{\theta})$$

**Confidence interval**

$$\Pr\{C(r_{\alpha}) \mid f(\cdot, \theta)\} = 1 - \alpha$$

$$C(r_{\alpha}) = \{ \theta : n(\theta - \hat{\theta})^T I(\theta)(\theta - \hat{\theta}) \leq r_{\alpha}^2 \}$$

# Most cited statisticians

1 [robert tibshirani](#)  
Professor of Health Research and Policy, a  
Verified email at stanford.edu  
Cited by 116885

[Michael I. Jordan](#)  
Professor of EECS and Professor of Statistic  
Verified email at cs.berkeley.edu  
Cited by 53311

2 [Donald B Rubin](#)  
John L. Loeb Professor of Statistics, Harva  
Verified email at stat.harvard.edu  
Cited by 115695

[William G Cochran 1909-1980](#)  
Professor of Statistics, Harvard University  
Verified email at stat.harvard.edu  
Cited by 48114

3 [B Efron](#)  
Professor of statistics, Stanford University  
Verified email at stat.stanford.edu  
Cited by 69104

[Leo Breiman 1928-2005](#)  
Professor of Statistics, UC Berkeley  
Verified email at stat.berkeley.edu  
Cited by 46238

[Arthur P Dempster](#)  
Professor Emeritus of Statistics, Harvard Un  
Verified email at stat.harvard.edu  
Cited by 39427

**Cf. Physics, Medicine, Mathematics**



# DEFINING THE CURVATURE OF A STATISTICAL PROBLEM (WITH APPLICATIONS TO SECOND ORDER EFFICIENCY)

BY BRADLEY EFRON

*Stanford University*

Statisticians know that one-parameter exponential families have very nice properties for estimation, testing, and other inference problems. Fundamentally this is because they can be considered to be "straight lines" through the space of all possible probability distributions on the sample space. We consider arbitrary one-parameter families  $\mathcal{F}$  and try to quantify how nearly "exponential" they are. A quantity called "the statistical curvature of  $\mathcal{F}$ " is introduced. Statistical curvature is identically zero for exponential families, positive for nonexponential families. Our purpose is to show that families with small curvature enjoy the good properties of exponential families. Large curvature indicates a breakdown of these properties. Statistical curvature turns out to be closely related to Fisher and Rao's theory of second order efficiency.

# Statistical curvature

$$g(x) = \frac{\Gamma\left(\frac{f+1}{2}\right)}{f^{\frac{1}{2}}\Gamma\left(\frac{1}{2}\right)\Gamma\left(\frac{f}{2}\right)} \left(1 + \frac{x^2}{f}\right)^{-(f+1)/2}$$

$$\gamma_{\theta}^2 = \frac{6[3f^2 + 18f + 19]}{f(f+1)(f+5)(f+7)} :$$

$f$	1	2	5	10	20	$\rightarrow \infty$
$\gamma_{\theta}^2$	2.5	1.063	.306	.107	.0334	$\sim 18/f^2$

# Information geometry

**Statistical model**  $M = \{f(x, \theta), \theta \in \Theta\}$

**Information metric**  $g_{ij}(\theta) = -E_{f(\cdot, \theta)} \left\{ \frac{\partial}{\partial \theta_i} \log f(x, \theta) \frac{\partial}{\partial \theta_j} \log f(x, \theta) \right\}$

**Exponential connection**

$$\overset{e}{\Gamma}_{ij,k}(\theta) = E_{f(\cdot, \theta)} \left\{ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(x, \theta) \frac{\partial}{\partial \theta_k} \log f(x, \theta) \right\}$$

**Mixture connection**

$$\overset{m}{\Gamma}_{ij,k}(\theta) = \overset{e}{\Gamma}_{ij,k}(\theta) + E_{f(\cdot, \theta)} \left\{ \frac{\partial}{\partial \theta_i} \log f(x, \theta) \frac{\partial}{\partial \theta_j} \log f(x, \theta) \frac{\partial}{\partial \theta_k} \log f(x, \theta) \right\}$$

**Rao (1949) Amari (1982)**



# Natural Gradient Works Efficiently in Learning

Shun-ichi Amari

*RIKEN Frontier Research Program, Saitama 351-01, Japan*

When a parameter space has a certain underlying structure, the ordinary gradient of a function does not represent its steepest direction, but the natural gradient does. Information geometry is used for calculating the natural gradients in the parameter space of perceptrons, the space of matrices (for blind source separation), and the space of linear dynamical systems (for blind source deconvolution). The dynamical behavior of natural gradient online learning is analyzed and is proved to be Fisher efficient, implying that it has asymptotically the same performance as the optimal batch estimation of parameters. This suggests that the plateau phenomenon, which appears in the backpropagation learning algorithm of multilayer perceptrons, might disappear or might not be so serious when the natural gradient is used. An adaptive method of updating the learning rate is proposed and analyzed.

# m-geodesic and e-geodesic

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

e-geodesic  $p_t^{(e)}(\mathbf{x}) = c_t \exp\{(1-t)\log p(\mathbf{x}) + t\log q(\mathbf{x})\}$

A model  $M$  is **m-geodesic** if

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

A model  $M$  is **e-geodesic** if

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

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

**Mean match model**  $M^{(m)} = \{ p(\mathbf{x}) : E_p \{ \mathbf{t}(\mathbf{x}) \} = E_{p_0} \{ \mathbf{t}(\mathbf{x}) \} \}$

# Kullback-Leibler divergence

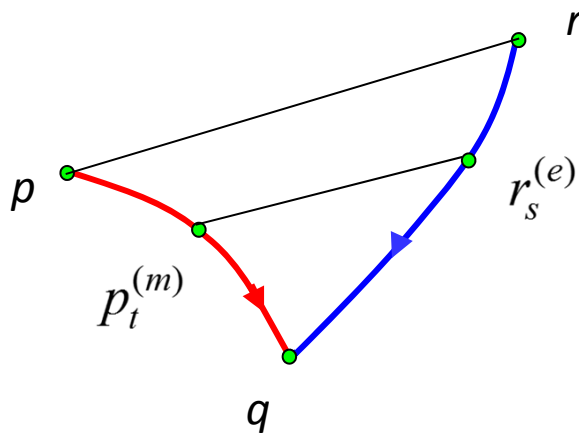
**KL divergence (1951)**  $D(g, f) = E_g \log \frac{g(X)}{f(X)}$

**AIC**  $AIC(M) = -2 \max_{\theta \in \Theta} L_n(\theta) + 2 \dim(\theta)$  (Akaike, 1973)

$$D(p, r) = D(p, q) + D(q, r)$$

$$p_t^{(m)} = (1-t)p + tq \quad r_s^{(e)}(x) = c_s \exp\{(1-s)\log r + s\log q\}$$

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



Nagaoka-Amari (1982)

# MLE on Exponential family

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

**Cumulant function**  $\kappa(\boldsymbol{\theta}) = \log \{ E_{p_0} e^{\boldsymbol{\theta}^T \mathbf{t}(\mathbf{x})} \}$

$$\boldsymbol{\eta} := \frac{\partial}{\partial \boldsymbol{\theta}} \kappa(\boldsymbol{\theta}) = E_{\boldsymbol{\theta}} \{ \mathbf{t}(\mathbf{x}) \} \quad \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \kappa(\boldsymbol{\theta}) = \text{var}_{\boldsymbol{\theta}} \{ \mathbf{t}(\mathbf{x}) \}$$

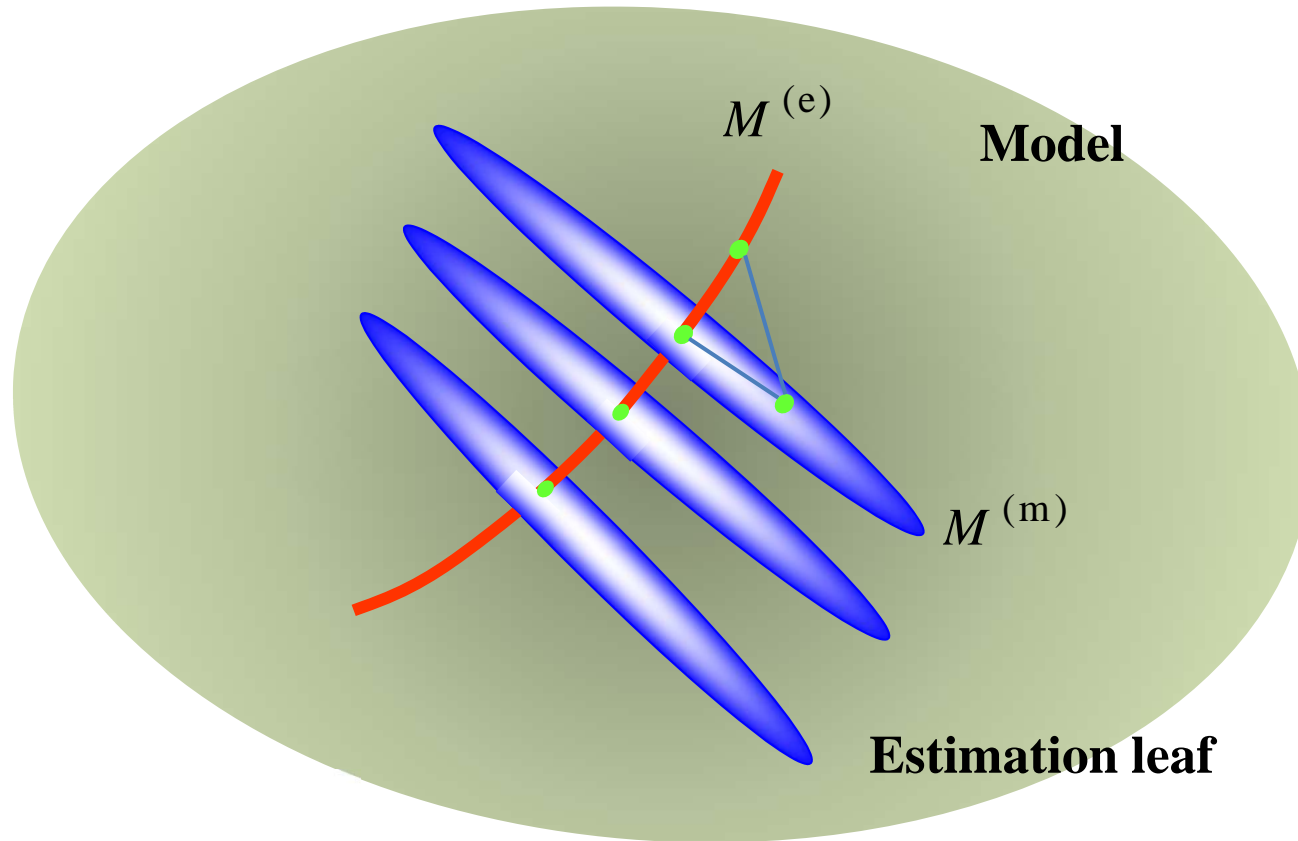
**Log likelihood function**  $L(\boldsymbol{\theta}) = \sum_{i=1}^n \log p(\mathbf{x}_i, \boldsymbol{\theta}) = n \{ \boldsymbol{\theta}^T \bar{\mathbf{t}} - \kappa(\boldsymbol{\theta}) \}$

**Maximum likelihood**  $\max_{\boldsymbol{\theta} \in \Theta} \frac{L(\boldsymbol{\theta})}{n} = \max_{\boldsymbol{\theta} \in \Theta} \{ \boldsymbol{\theta}^T \bar{\mathbf{t}} - \kappa(\boldsymbol{\theta}) \} = \kappa^*(\bar{\mathbf{t}})$

**Likelihood equation**  $\frac{\partial}{\partial \boldsymbol{\theta}} \kappa(\boldsymbol{\theta}) = \bar{\mathbf{t}} \quad \text{or} \quad \boldsymbol{\eta}_{\text{ML}} = \bar{\mathbf{t}}$

**Mean match model**  $M^{(m)} = \{ p(\mathbf{x}) : E_p \{ \mathbf{t}(\mathbf{x}) \} = \boldsymbol{\eta} \}$

# Pair of (model, estimation) Amari (1982)



$$\mathcal{F} = \bigcup_{\theta \in \Theta} \mathcal{L}_\theta$$

$$\mathcal{L}_\theta = \{g : \hat{\theta}(g) = \theta\}$$

# Non-exponential family

## t-distribution

Univariate

$$f(x, \mu, \sigma^2, \nu) = \frac{\Gamma(\frac{\nu+1}{2})}{(\pi\nu)^{\frac{1}{2}} \Gamma(\frac{\nu}{2})} \left\{ 1 + \frac{1}{\nu} \frac{(x - \mu)^2}{\sigma^2} \right\}^{-\frac{\nu+1}{2}}$$

$p$ -variate

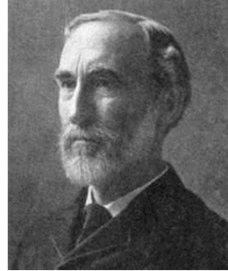
$$f_{\gamma}(x, \mu, \Sigma, \nu) = \frac{\Gamma(\frac{\nu+p}{2})}{(\pi\nu)^{\frac{p}{2}} \Gamma(\frac{\nu}{2})} \left\{ 1 + \frac{1}{\nu} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}^{-\frac{\nu+p}{2}}$$

## Generalized Pareto distribution

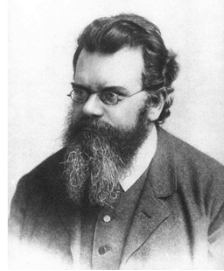
$$f(x, \theta, \xi) = \theta (1 + \xi \theta x)^{-\frac{1}{\xi} - 1}$$

## Unified view to power exponential family

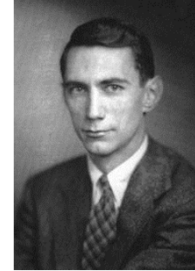
# Gibbs-Boltzmann-Shannon



**Josiah Willard Gibbs**  
(1839 - 1903)



**Ludwig E. Boltzmann**  
(1844 - 1906)



**Claude E. Shannon**  
(1916- 2001)

**GBS-entropy**

$$H_{\text{GBS}}(p) = -E_p \log p(X)$$

**Log-likelihood**

$$-\frac{1}{n} \sum_{i=1}^n \log p(\mathbf{x}_i, \theta) \approx H_{\text{GBS}}(p(\cdot, \theta))$$

**A new estimator if we extend GBS-entropy ?**

# Possible Generalization of Boltzmann–Gibbs Statistics

Constantino Tsallis<sup>1</sup>

*Journal of Statistical Physics, Vol. 52, Nos. 1/2, 1988*

---

With the use of a quantity normally scaled in multifractals, a generalized form is postulated for entropy, namely  $S_q \equiv k[1 - \sum_{i=1}^W p_i^q]/(q-1)$ , where  $q \in \mathbb{R}$  characterizes the generalization and  $\{p_i\}$  are the probabilities associated with  $W$  (microscopic) configurations ( $W \in \mathbb{N}$ ). The main properties associated with this entropy are established, particularly those corresponding to the microcanonical and canonical ensembles. The Boltzmann–Gibbs statistics is recovered as the  $q \rightarrow 1$  limit.

$$S_q(\mathbf{p}) = \frac{k}{q-1} \left( 1 - \sum_{j=1}^W p_j^q \right)$$



# Hill's diversity index in 1973

Let  $p_j$  be a relative frequency of a speice  $j$  for  $j = 1, \dots, S$

$$D_a(\mathbf{p}) = \left( \sum_{j=1}^S p_j^a \right)^{\frac{1}{1-a}}$$

$$\begin{aligned} \lim_{a \rightarrow 0} \log \{ D_a(\mathbf{p}) \} &= \lim_{a \rightarrow 0} \frac{1}{1-a} \log \left( \sum_{j=1}^S p_j^a \right) \\ &= - \sum_{j=1}^S p_j \log p_j = H(\mathbf{p}) \end{aligned}$$

# An Analysis of Transformations

By G. E. P. Box                      and                      D. R. Cox

*University of Wisconsin*

*Birkbeck College, University of London*

## SUMMARY

In the analysis of data it is often assumed that observations  $y_1, y_2, \dots, y_n$  are independently normally distributed with constant variance and with expectations specified by a model linear in a set of parameters  $\theta$ . In this paper we make the less restrictive assumption that such a normal, homoscedastic, linear model is appropriate after some suitable transformation has been applied to the  $y$ 's. Inferences about the transformation and about the parameters of the linear model are made by computing the likelihood function and the relevant posterior distribution. The contributions of normality, homoscedasticity and additivity to the transformation are separated. The relation of the present methods to earlier procedures for finding transformations is discussed. The methods are illustrated with examples.

$$y^{(\lambda)} = \begin{cases} \frac{y^\lambda - 1}{\lambda} & (\lambda \neq 0), \\ \log y & (\lambda = 0), \end{cases}$$

# Projective Entropy

$$\gamma\text{-cross entropy} \quad C_\gamma(g, f) = -\frac{\int f(x)^\gamma g(x) dx}{\left(\int f(x)^{1+\gamma} dx\right)^{\frac{\gamma}{1+\gamma}}}$$

$$\gamma\text{-diagonal entropy} \quad H_\gamma(f) = C_\gamma(f, f) = -\left(\int f^{1+\gamma}\right)^{\frac{1}{1+\gamma}}$$

$$\gamma\text{-divergence} \quad D_\gamma(g, f) = C_\gamma(g, f) - H_\gamma(g)$$

Cf. Good (1972) Fujisawa-Eguchi (2008) Eguchi-Kato (2010) Ferrari-Yang (2010)

Remark: The diagonal entropy is equivalent to Tsallis entropy

$$S_q(f) = \frac{1}{q-1} \left(1 - \int f^q\right)$$

$$S_q(f) = \frac{1}{\gamma} \{H_\gamma(f)^q + 1\} \quad \text{where } q = 1 + \gamma$$

# Three properties

$$C_\gamma(g, f) = -\int \left( \frac{f}{\|f\|_q} \right)^\gamma g, \quad (q = 1 + \gamma)$$

Linearity

$$C_\gamma(\alpha g + \beta h, f) = \alpha C_\gamma(g, f) + \beta C_\gamma(h, f)$$

scale invariance

$$C_\gamma(g, \lambda f) = C_\gamma(g, f) \quad (\forall \lambda > 0)$$

Lower Bound

$$C_\gamma(g, f) \geq C_\gamma(g, g)$$

These properties lead to uniqueness for  $C_\gamma$  ?

# Theorem

Let  $\Gamma(g, f) = \Phi\left(\int \rho(f)\right) \int \psi(f)g$

If  $\Gamma(g, f)$  satisfies

$$(i) \quad \Gamma(g, \lambda f) = \Gamma(g, f) \quad (\forall \lambda > 0)$$

$$(ii) \quad \Gamma(g, f) \geq \Gamma(g, g)$$

then there is a constant  $\gamma > 0$  such that  $\Gamma(g, f) = C_\gamma(g, f)$

**Remark:** If  $(\Phi(Y), \rho(f), \psi(f)) = (Y^{-\frac{\gamma}{1+\gamma}}, f^{1+\gamma}, f^\gamma)$ ,  
then  $\Gamma(g, f) = \left(\int f(x)^{1+\gamma} dx\right)^{-\frac{\gamma}{1+\gamma}} \int f^\gamma g = C_\gamma(g, f)$

# Proof

$$(i) \Rightarrow \frac{\partial}{\partial \lambda} \Gamma(g, \lambda f) = 0$$

$$\frac{\partial}{\partial \lambda} \Gamma(g, \lambda f) = \int \left[ \left\{ \frac{\partial}{\partial \lambda} \Phi \left( \int \rho(\lambda f) \right) \right\} \psi(\lambda f(x)) + \Phi \left( \int \rho(\lambda f) \right) \frac{\partial \psi(\lambda f(x))}{\partial \lambda} f(x) \right] g(x) dx$$

$$\frac{\partial \log \psi(X)}{\partial X} = -c X^{-1} \quad \therefore \psi(X) = X^\gamma \quad (\gamma = c),$$

$$(ii) \Rightarrow \frac{\partial}{\partial f_i} \Gamma(g, f) \Big|_{f=g} = 0 \quad \text{if } \Gamma(g, f) = \Phi \left( \sum_{j=1}^m \rho(f_j) \right) \sum_{j=1}^m g_j (f_j)^\gamma$$

$$(*) \quad \frac{\partial}{\partial f_i} \Gamma(g, f) \Big|_{f=g} = \left\{ \dot{\Phi} \left( \sum_{j=1}^m \rho(g_j) \right) \sum_{j=1}^m (g_j)^{\gamma+1} \right\} \dot{\rho}(g_i) + \left\{ \Phi \left( \sum_{j=1}^m \rho(g_j) \right) \right\} \gamma (g_i)^\gamma = 0 \quad [\because (ii)]$$

$$\frac{\partial \rho(g)}{\partial g} = -c \gamma g^\gamma \quad \therefore \rho(g) = -c \frac{\gamma}{\gamma+1} g^{\gamma+1}$$

Taking a sum of (\*) over  $i$ ,

$$\frac{\dot{\Phi}(Y)}{\Phi(Y)} = -\frac{\gamma+1}{\gamma} Y, \quad \therefore \Phi(Y) = Y^{-\frac{\gamma+1}{\gamma}}$$

# Shannon-Khinchin axiom

$$S_g[p] = \sum_{i=1}^W g(p_i)$$

**K1:** The requirement that  $S$  depends continuously on  $p$  implies that  $g$  is a continuous function.

**K2:** The requirement that the entropy is maximal for the equi-distribution  $p_i = 1/W$  (for all  $i$ ) implies that  $g$  is a concave function.

**K3:** The requirement that adding a zero-probability state to a system,  $W + 1$  with  $p_{W+1} = 0$ , does not change the entropy implies  $g(0) = 0$ .

$$S_{c,d}[p] = \frac{e \sum_i^W \Gamma(1+d, 1-c \ln p_i)}{1-c+cd} - \frac{c}{1-c+cd}$$

**Cf. Hanel and Thurner (2011)**

# Max $\gamma$ -entropy

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

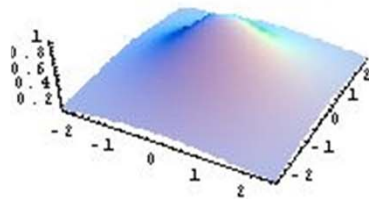
$\gamma$ -entropy  $H_\gamma(f) = -\|f\|_q \quad (q = 1 + \gamma)$

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

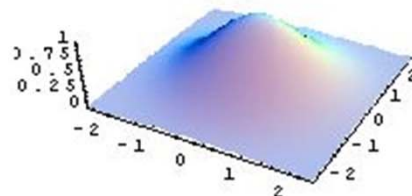
Eguchi-Komori-Kato (2011)

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

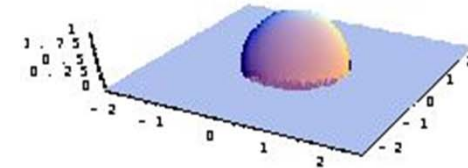
$$\mathcal{F}_\gamma(\mu, \Sigma) = \{f \in \mathcal{F}(\mu, \Sigma) : \text{Support}(f) \subseteq \{x : (x - \mu)^\top \Sigma^{-1}(x - \mu) \leq \frac{1}{2}\kappa_\gamma\}\}$$



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



$\gamma = 0$  (Gaussian)



$\gamma = 2$  (Wigner)



# $\gamma$ -Estimation

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

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

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

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

Gaussian

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

$\gamma$ -estimator

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

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

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

remark If  $\gamma > 0$ ,  $(\hat{\mu}_\gamma, \hat{\Sigma}_\gamma)$  is robust gainst much oulying

# Escort normal family

$$f_\gamma(\cdot, \mu, \Sigma) \in \mathcal{F}(\mu, \Sigma)$$

Escort distribution  $e_q(f(x)) = \frac{f(x)^{1+\gamma}}{\int f^{1+\gamma}} \quad (q=1+\gamma)$

$$\begin{aligned} e_q(f(x, \mu, \Sigma)) &= c_\gamma^* \det(\Sigma)^{-\frac{1}{2}} \left\{ 1 - \frac{\kappa_\gamma}{2} (x - \mu) \Sigma^{-1} (x - \mu) \right\}^{\frac{1+\gamma}{\gamma}} \\ &= c_\gamma^* \left\{ \det(2\pi\Sigma)^{-\frac{1}{2}} \frac{\gamma}{2^{1+\gamma}} - \frac{\kappa_\gamma}{2} (x - \mu) \left( \det(\Sigma)^{-\frac{1}{2}} \frac{\gamma}{2^{1+\gamma}} \Sigma^{-1} \right) (x - \mu) \right\}^{\frac{1+\gamma}{\gamma}} \\ &= c_\gamma^* \left\{ \det(\Xi)^{\frac{2}{p\gamma+2\gamma+2}} - \frac{\gamma}{2} (x - \mu) \Xi (x - \mu) \right\}^{\frac{1+\gamma}{\gamma}} \end{aligned}$$

$$\frac{\partial}{\partial \mu} \int e_q(f_\gamma(x, \mu, \Sigma)) = 0 \Rightarrow E_{f_\gamma(\cdot, \mu, \Sigma)}(X) = \mu$$

$$\frac{\partial}{\partial \Xi} \int e_q(f_\gamma(x, \mu, \Xi)) = 0 \Rightarrow E_{f_\gamma(\cdot, \mu, \Sigma)}\{(X - \mu)(X - \mu)^T\} = \Sigma$$

# $\gamma$ -estimation on $\gamma$ -normal model

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

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

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

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

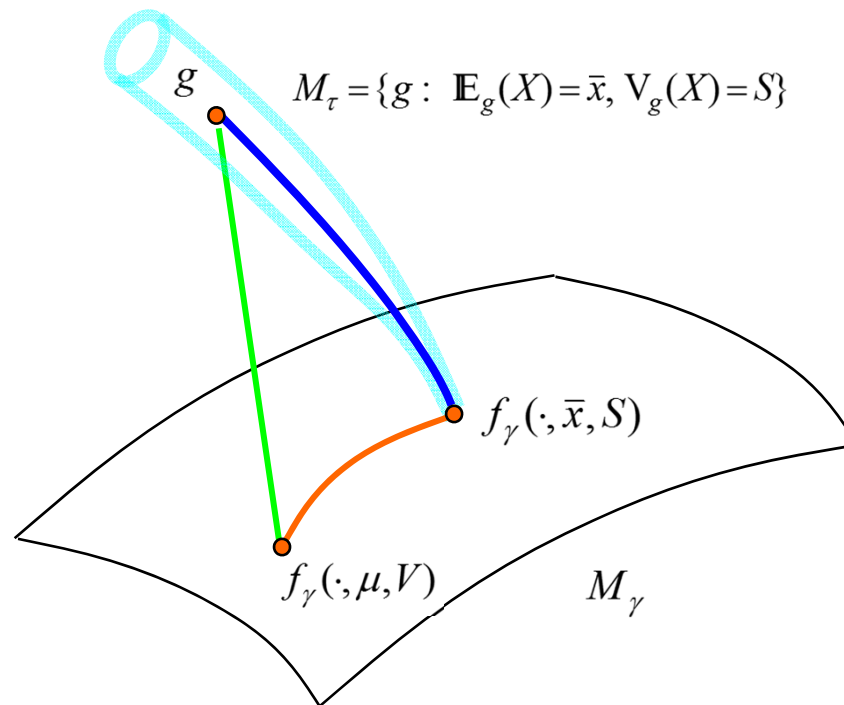
$$L_\gamma(\mu, \Sigma) - L_\gamma(\bar{x}, S) = \det(\Xi)^{\frac{\gamma}{p\gamma+2\gamma+2}} - \frac{1}{2} \kappa_\gamma \{ (\mu - \bar{x})^\top \Xi (\mu - \bar{x}) + \text{tr}(S \Xi) \}$$

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

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

# Pythagoras triangle

Loss decomposition  $L_\gamma(\mu, \Sigma) - L_\gamma(\bar{x}, S) = D_\gamma(f_\gamma(\cdot, \bar{x}, S), f_\gamma(\cdot, \mu, \Sigma))$



Pythagorean  $D_\gamma(g, f_\gamma(\cdot, \mu, V)) - D_\gamma(g, f_\gamma(\cdot, \bar{x}, S)) = D_\gamma(f_\gamma(\cdot, \bar{x}, S), f_\gamma(\cdot, \mu, V))$

# MLE for normal model

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

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

Max 0-entropy  $H_0(f_{\mu, \Sigma}) = \max_{f \in \mathcal{F}(\mu, \Sigma)} H_0(f)$  where  $H_0(f) = -\int f \log f$

Likelihood  $L_0(\mu, \Sigma) = \sum_{i=1}^n \log f_{\mu, \Sigma}(x_i)$

## Teicher's theorem (1961)

Let  $g_{\mu, \Sigma}$  be a location-scale family.

MLE  $(\hat{\mu}, \hat{\Sigma})$  is the sample mean and variance  $(\bar{x}, S)$   
if and only if  $g_{\mu, \Sigma}$  is Gaussian.

$$\text{where } (\bar{x}, S) = \left( \frac{1}{n} \sum x_i, \frac{1}{n} \sum (x_i - \bar{x})(x_i - \bar{x})^T \right)$$

## $\gamma$ -estimator leads to $\gamma$ -normal model

Let  $(\hat{\mu}_\gamma, \hat{\Sigma}_\gamma) = \arg \min_{(\mu, V)} L_\gamma(g(\cdot, \mu, \Sigma))$

where  $g(x, \mu, \Sigma) = c_h \det(2\pi \Sigma)^{-\frac{1}{2}} h\left(-\frac{1}{2}(x - \mu)\Sigma^{-1}(x - \mu)\right)$

Then

$$g(x, \mu, \Sigma) = f_\gamma(x, \mu, \Sigma) \iff (\hat{\mu}_\gamma, \hat{\Sigma}_\gamma) = (\bar{x}, S)$$

$\gamma = 0$  (Gaussian) Cf. Teicher (1961)

$(\gamma$ -estimator ,  $\gamma'$ -model)

model loss function	0-model	$\gamma$ -model
0-loss (- log-likelihood)	$(\bar{x}, S)$ MLE	M-estimator
$\gamma$ -loss	emergence $\gamma$ -estimator	$(\bar{x}, S)$ $\gamma$ -estimator

# $U$ -entropy and $U$ -divergence

$U$  is a convex and increasing function Cf. Eguchi-Kano (2001)

**$U$ -cross entropy**  $C_U(g, f) = -E_g \xi(f(X)) + \int U(\xi(f))$

**$U$ -entropy**  $H_U(f) = -E_f \xi(f(X)) + \int U(\xi(f))$

where  $\xi = (U)^{-1}$

**$U$ -divergence**  $D_U(g, f) = C_U(g, f) - H_U(g)$

Example 1. **Boltzmann-Shannon entropy.**

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

$$C_0(g, f) = -\int g \log f + 1 \quad H_0(f) = -\int f \log f + 1$$

$$D_0(g, f) = \int g(\log g - \log f)$$



Example 2. **Power entropy.** Tsallis (1988) Basu et al (1998) Cf. Box-Cox (1964)

$$(U(t), \xi(s)) = \left( \frac{(1+\beta t)^{\frac{1+\beta}{\beta}}}{1+\beta}, \frac{s^\beta - 1}{\beta} \right) \quad (0 \leq \beta \leq 1)$$

$$C_\beta(g, f) = \int -g \frac{f^\beta - 1}{\beta} + \frac{g^{\beta+1}}{\beta+1} \quad H_\beta(f) = \int \left( -\frac{f^{\beta+1}}{\beta(\beta+1)} + \frac{f}{\beta} \right)$$

$$D_\beta(g, f) = \int g \frac{g^\beta - f^\beta}{\beta} - \frac{g^{\beta+1} - f^{\beta+1}}{\beta+1} \quad \text{Tsallis entropy } (q=\beta+1)$$

Example 3. **Sigmoid entropy.** Takenouchi-Eguchi (2004)

$$(U(t), \xi(s)) = (e^t - \eta t, \log(s + \eta)) \quad (0 < \eta < \frac{1}{2})$$

$$C_\eta(g, f) = -\int (g + \eta) \log(f + \eta) \quad H_\eta(f) = -\int (f + \eta) \log(f + \eta)$$

$$D_\eta(g, f) = \int (g + \eta) \{ \log(g + \eta) - \log(f + \eta) \}$$

# max $U$ -entropy model

Let  $p_0(\mathbf{x})$  be a pdf and let  $\mathbf{t}(\mathbf{x})$  be a statistic of dimension  $d$ .

Consider a problem  $\max_{p \in \mathcal{Q}(p_0)} H_U(p)$

where  $\mathcal{Q}(p_0) = \{p : E_p \mathbf{t}(X) = E_{p_0} \mathbf{t}(X)\}$   
Mean-equal space

The Lagrangian  $L(p, \boldsymbol{\theta}, \lambda) = \int \{-p \xi(p) + U(\xi(p)) + \boldsymbol{\theta}^T (\mathbf{t} - \boldsymbol{\tau}_0) p\} d\mathbf{x}$

leads to  $\xi(p(\mathbf{x})) = \boldsymbol{\theta}^T \mathbf{t}(\mathbf{x}) + \text{const.}$

So  $p(\mathbf{x}) = \dot{U}(\boldsymbol{\theta}^T \mathbf{t}(\mathbf{x}) - \kappa_U(\boldsymbol{\theta}))$ , we call  **$U$ -model**

# U-estimation

Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be from  $p(\mathbf{x})$  and let  $p_\theta(\mathbf{x})$  be a model function.

**U-loss function**  $L_U(p_\theta) = -\frac{1}{n} \sum_{i=1}^n \xi(p_\theta(\mathbf{x}_i)) + \int U(\xi(p_\theta))$

**U-estimator**  $\hat{\theta}_U = \arg \min_{\theta \in \Theta} L_U(p_\theta)$

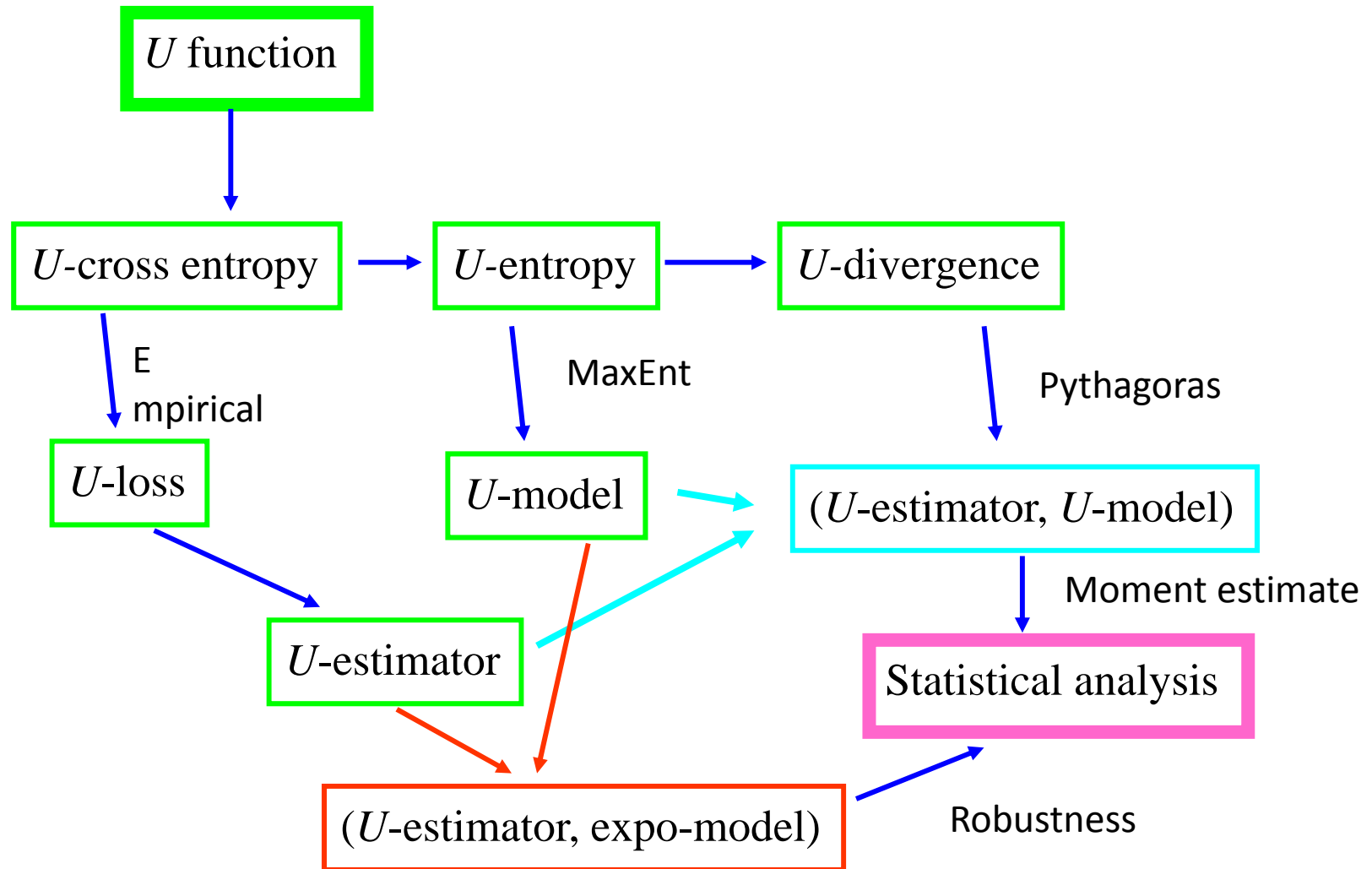
Note:  $E_p L_U(p_\theta) = C_U(p, p_\theta)$  ;  $L_U(p_\theta) \xrightarrow[\text{a.s.}]{} C_U(p, p_\theta)$  as  $n \rightarrow \infty$

**Consistency**  $\hat{\theta}_U \xrightarrow[\text{a.s.}]{} \theta(p)$  as  $n \rightarrow \infty$

where  $\theta(p) = \arg \min_{\theta \in \Theta} C_U(p, p_\theta)$

**Asymptotic normality**  $\sqrt{n}(\hat{\theta}_U - \theta(p)) \xrightarrow[D]{} N(0, V_U(\theta, p))$  as  $n \rightarrow \infty$

# $U$ -method



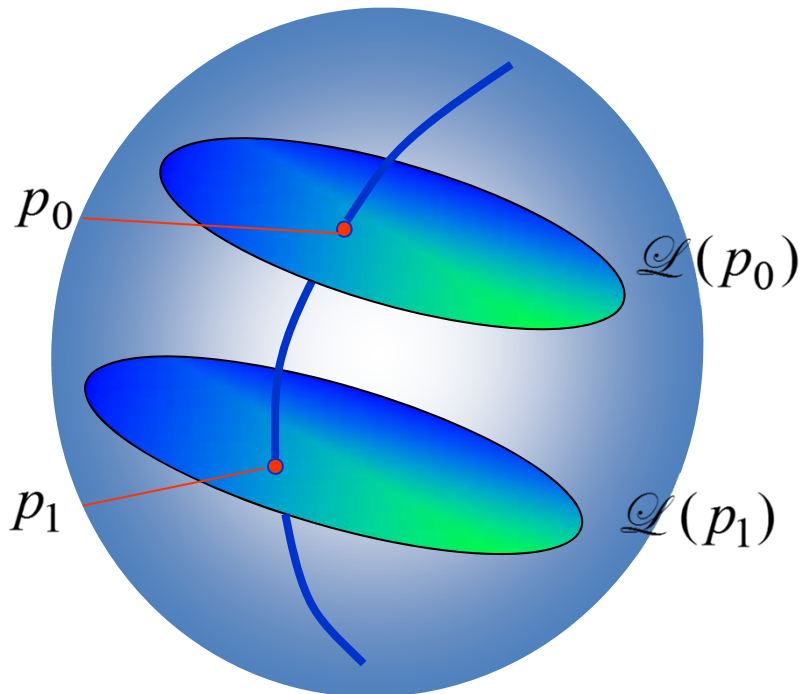
## $U_1$ -estimator under exp-model

<b>model</b>	Exponential	$U_1$ -model
<b>loss function</b>		
$U_{\text{exp}} = \text{likelihood}$	<b>MLE</b>	<b>Robust M-estimate</b>
$U_1$ -loss	<b>Robust <math>U_1</math>-estimator</b>	Moment estimator

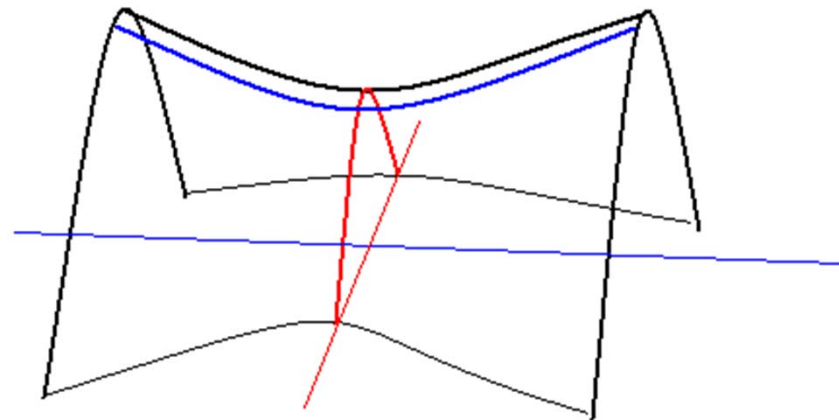
# Minimax game

$$\max_{p \in \mathcal{Q}(p_0)} \min_{q \in \mathcal{P}} C_U(p, q) = C_U(p^*, p^*) = \min_{q \in \mathcal{P}} \max_{p \in \mathcal{Q}(p_0)} C_U(p, q)$$

$$\text{where } p^* = \operatorname{argmax}_{p \in \mathcal{Q}(p_0)} H_U(p)$$



Cf. Grunwald and Dawid (2004)



# Statistical Modeling: The Two Cultures

Leo Breiman

*Abstract.* There are two cultures in the use of statistical modeling to reach conclusions from data. One assumes that the data are generated by a given stochastic data model. The other uses algorithmic models and treats the data mechanism as unknown. The statistical community has been committed to the almost exclusive use of data models. This commitment has led to irrelevant theory, questionable conclusions, and has kept statisticians from working on a large range of interesting current problems. Algorithmic modeling, both in theory and practice, has developed rapidly in fields outside statistics. It can be used both on large complex data sets and as a more accurate and informative alternative to data modeling on smaller data sets. If our goal as a field is to use data to solve problems, then we need to move away from exclusive dependence on data models and adopt a more diverse set of tools.

*Statistical Science*

2001, Vol. 16, No. 3, 199–231

# Pattern recognition

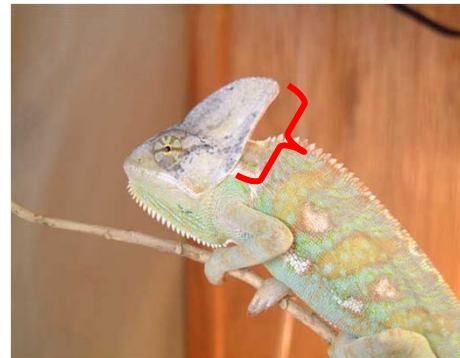
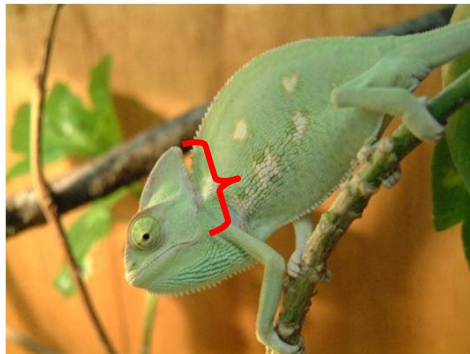
Feature vector  $\mathbf{x} = (x_1, \dots, x_p)$

Class label  $y \in \{-1, +1\}$

Classifier  $\mathbf{x} \mapsto y = f(\mathbf{x})$

Training data  $D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$

Statistical classifier  $f(\mathbf{x}) = \text{sgin}(F(\mathbf{x}))$





# Set of weak learners

## Decision stamp

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

## Linear classifiers

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

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

**No stronger, but exhaustive!**

ANNs SVMs *k*-NNs

# Adaboost

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

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

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

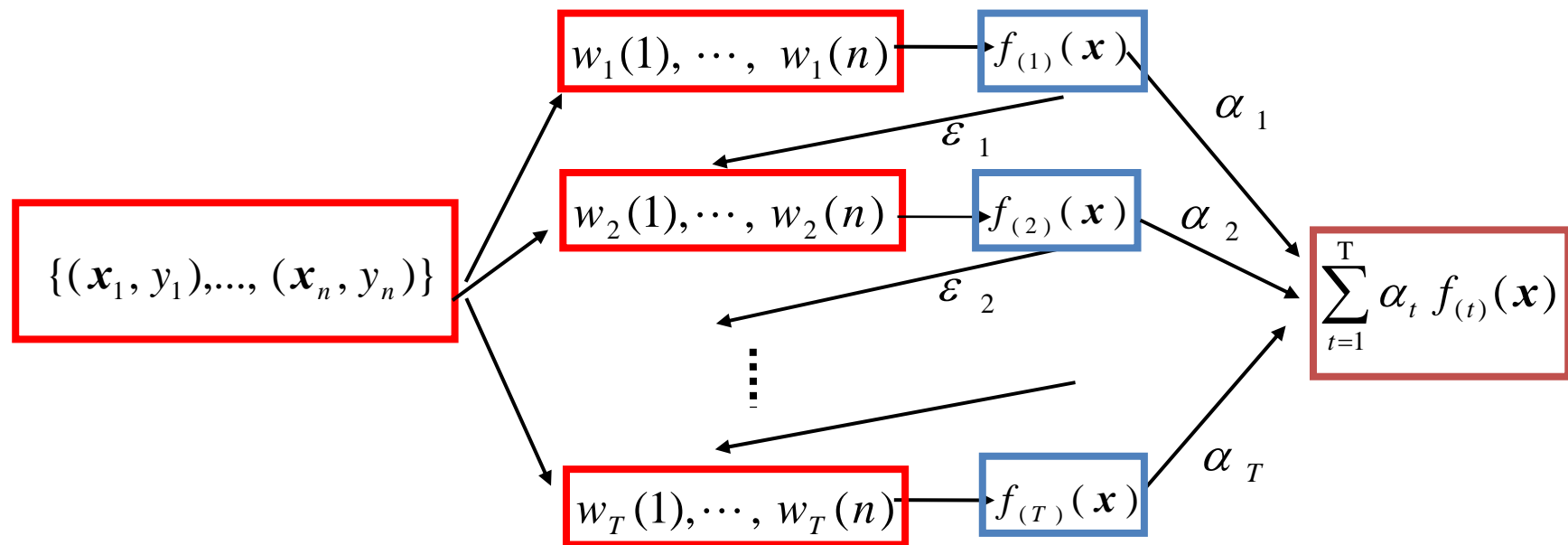
(b)       $\alpha_t = \frac{1}{2} \log \frac{1 - \varepsilon_t(f_{(t)})}{\varepsilon_t(f_{(t)})}$

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

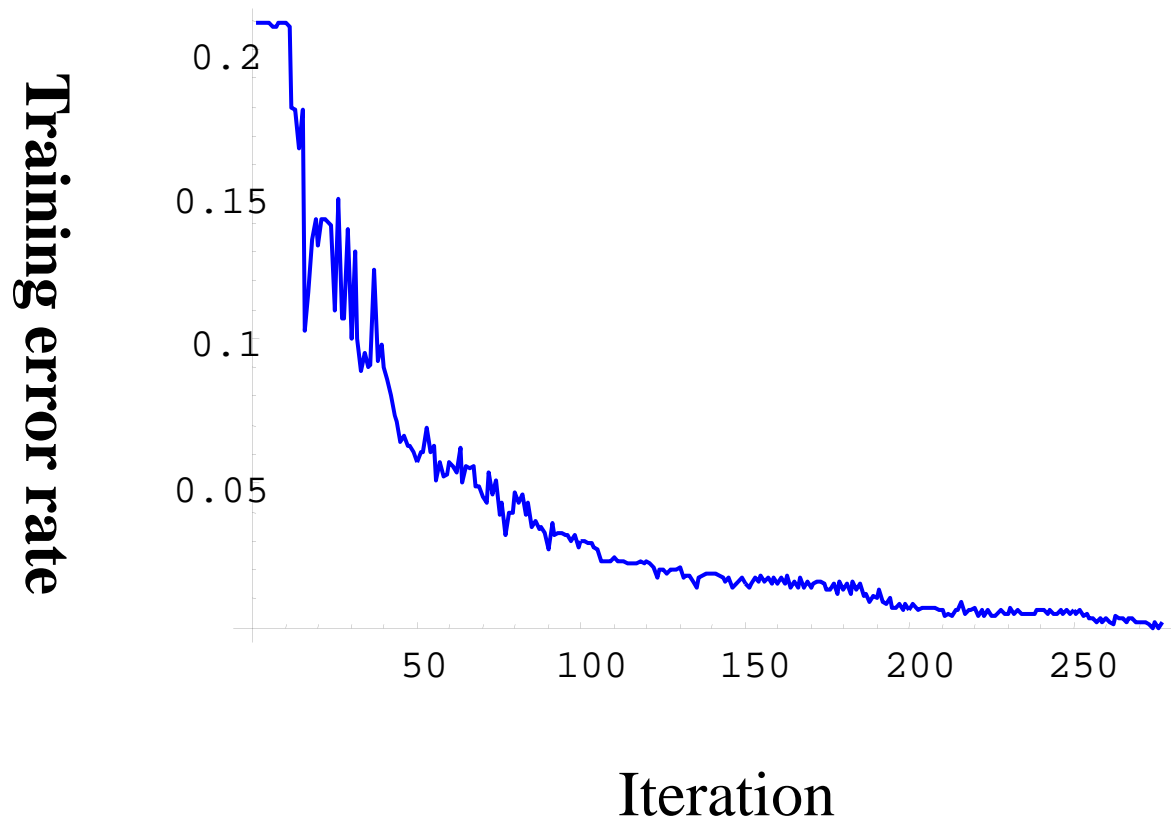
Freund-Schapire (1997).

# Learning process



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

# Learning curve



# Stopping rule

**Training data**  $D_{\text{train}} = \{ (\mathbf{x}_i, y_i) : i = 1, \dots, n \}$

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

How to find  $T$  ?

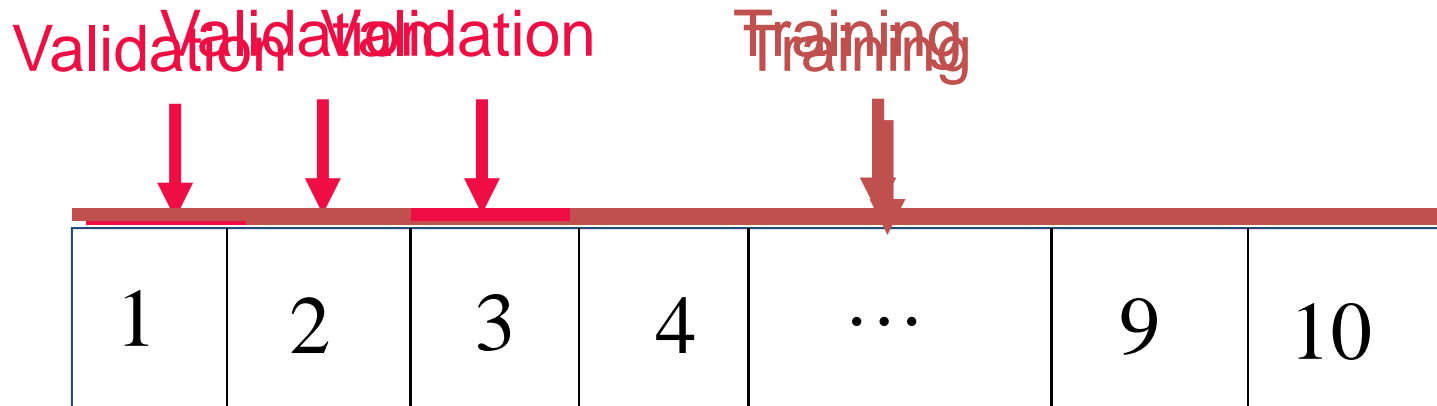
✗  $T_{\text{opt}} = \arg \min_{T > 0} \text{Err}^{\text{train}}(h_{F_T})$

✗  $T_{\text{opt}} = \arg \min_{T > 0} \text{Err}^{\text{test}}(h_{F_T})$

○  $T_{\text{opt}} = \arg \min_{T > 0} \text{CVErr}^{\text{train}}(h_{F_T})$

# 10-fold CV Error rate

$$D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n, F_T(\mathbf{x}) = \sum_{t=1}^T \alpha_t f_{(t)}(\mathbf{x})$$



$$\varepsilon(1) \varepsilon(2) \varepsilon(3)$$

$$\Rightarrow \text{Averaging by } \frac{1}{10} \sum_{k=1}^{10} \varepsilon(k)$$

# Update of exp-loss

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

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

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

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 minimization

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

$$\varepsilon(f) e^{\alpha} + (1 - \varepsilon(f)) e^{-\alpha}$$

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

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

$$\text{Equality iff } \alpha_{\text{opt}} = \frac{1}{2} \log \frac{1 - \varepsilon(f)}{\varepsilon(f)}$$



# Adaboost = sequential min expo-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)  $\alpha_t = \operatorname{arg min}_{\alpha \in \mathbf{R}} L_{\text{exp}}(F_{t-1} + \alpha f_{(t)})$

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

# Bayes risk consistency

## Expected loss functional

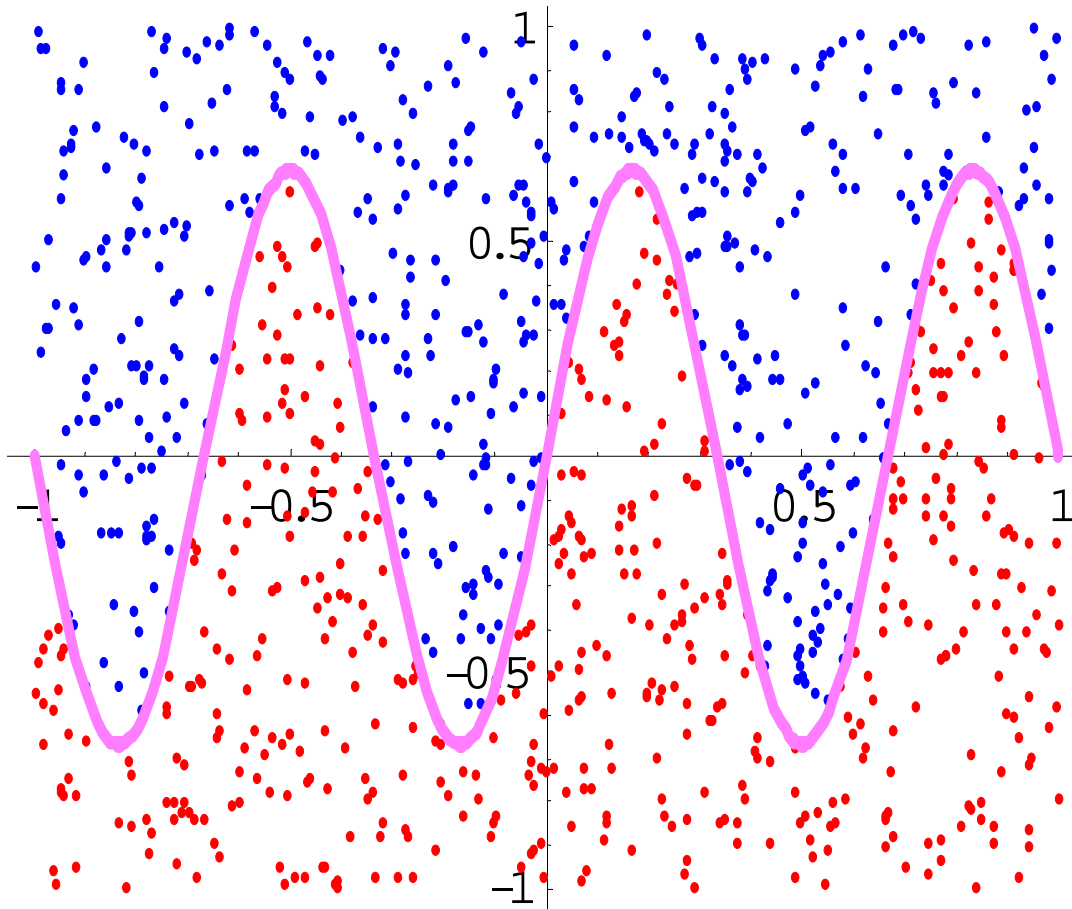
$$\begin{aligned}\mathbb{L}_{\text{exp}}(F) &= \sum_{y=\pm 1} \int \exp\{-y F(\mathbf{x})\} dP(y, \mathbf{x}) \\ &= \int \{e^{-F(\mathbf{x})} p(+1 | \mathbf{x}) + e^{F(\mathbf{x})} p(-1 | \mathbf{x})\} f(\mathbf{x}) d\mathbf{x}\end{aligned}$$

## Optimal discriminant function

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

$$\begin{aligned}&\mathbb{L}_{\text{exp}}(F) - \mathbb{L}_{\text{exp}}(F_{\text{opt}}) \\ &= \int \{e^{-F(\mathbf{x})} p(+1 | \mathbf{x}) + e^{F(\mathbf{x})} p(-1 | \mathbf{x}) - 2\sqrt{p(+1 | \mathbf{x}) p(-1 | \mathbf{x})}\} f(\mathbf{x}) d\mathbf{x} \\ &= \int \{\sqrt{e^{-F(\mathbf{x})} p(+1 | \mathbf{x})} - \sqrt{e^{F(\mathbf{x})} p(-1 | \mathbf{x})}\}^2 f(\mathbf{x}) d\mathbf{x} \geq 0\end{aligned}$$

# A simulation



**Feature space**

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

**Decision boundary**

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

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

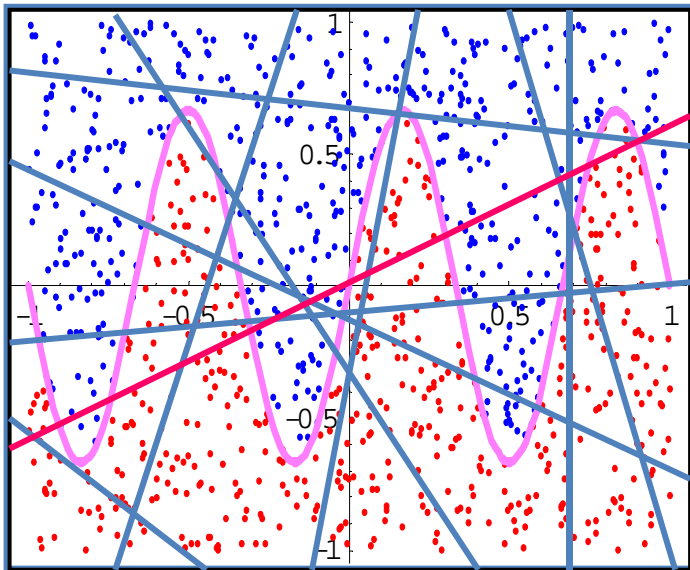
$$\mathbf{x}_i \in [-1, 1] \times [-1, 1]$$

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

# Set of linear classifiers

## Linear classifier

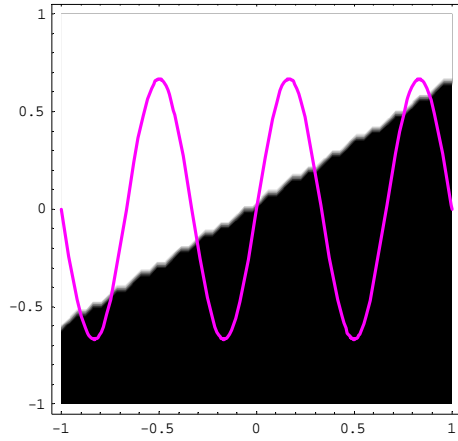
$$f(x_1, x_2) = \text{sgn}(r_1 x_1 + r_2 x_2 + r_3) = \begin{cases} +1 & \text{if } r_1 x_1 + r_2 x_2 + r_3 \geq 0 \\ -1 & \text{if } r_1 x_1 + r_2 x_2 + r_3 < 0 \end{cases}$$



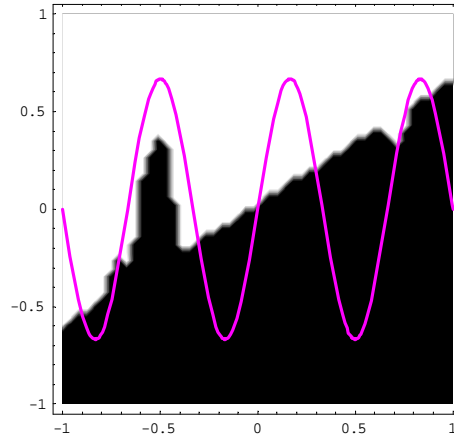
## Random generation

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

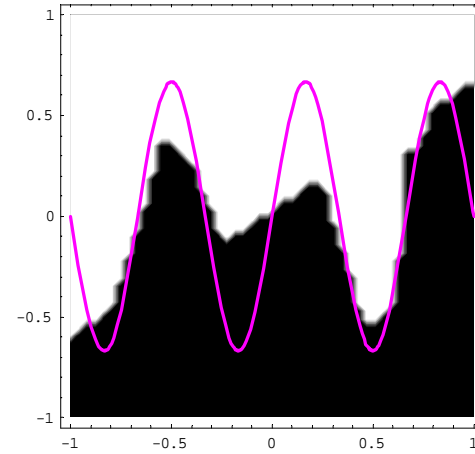
# Learning process



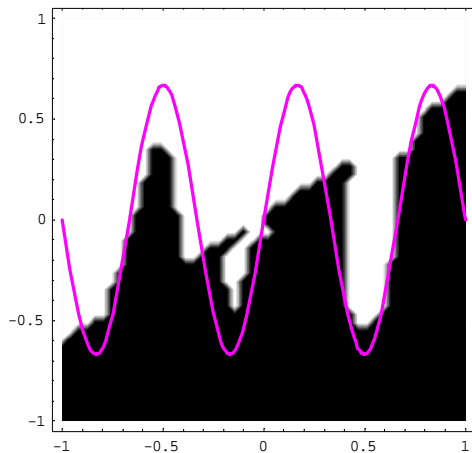
Iter = 1, train err = 0.21



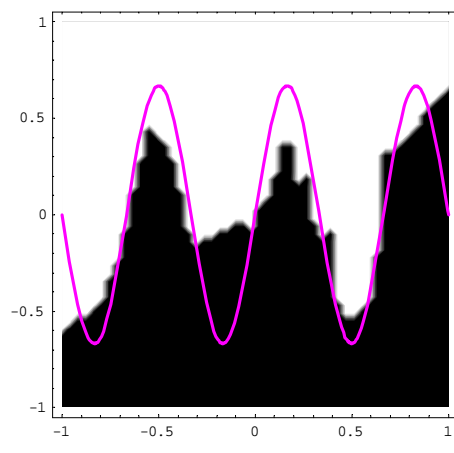
Iter = 13, train err = 0.18



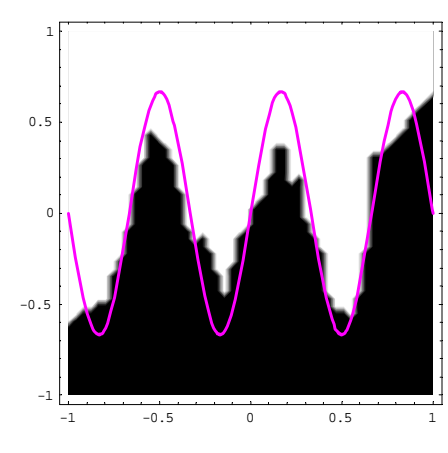
Iter = 17, train err = 0.10



Iter = 23, train err = 0.10

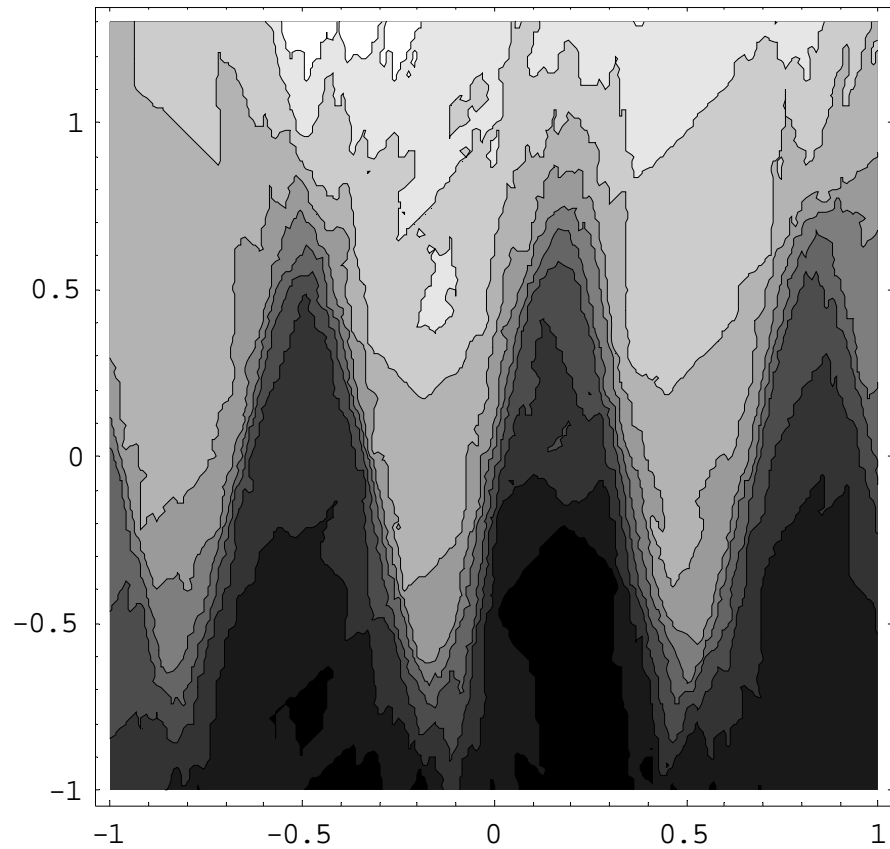


Iter = 31, train err = 0.095

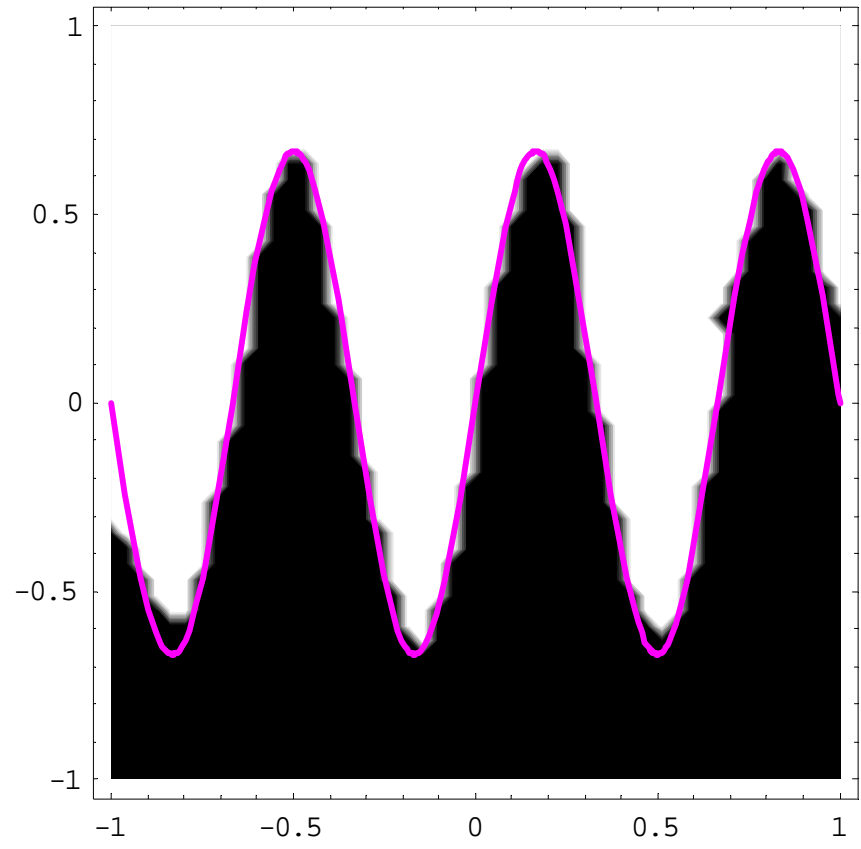


Iter = 47, train err = 0.08

# Final decision



Contour of  $F(x)$



Sign( $F(x)$ )

# $U$ -boost for classification

Let  $(\mathbf{x}, y)$  be a variable with feature vector  $\mathbf{x}$  and label  $y$

Decision rule  $h(\mathbf{x}) = \operatorname{argmax}_{y \in \{1, \dots, g\}} \hat{q}(y | \mathbf{x})$

$U$ -boost:  $q_t(y | \mathbf{x}) = u(\alpha_t f_t(\mathbf{x}, y) + \xi(q_{t-1}(y | \mathbf{x})))$

Step 1.  $f_t = \operatorname{argmin}_f \varepsilon_t(f | q_{t-1})$

Step 2.  $\alpha_t = \operatorname{argmin}_\alpha L_U^{\text{emp}}(\alpha f_t + \xi(q_{t-1}))$

$$L_U^{\text{emp}}(\xi(q_t)) - L_U^{\text{emp}}(\xi(q_{t+1})) = D_U(q_{t+1}, q_t)$$

Murata *et al* (2004)

# U-Boost for density estimation

**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  $f(\mathbf{x}, (0, \dots, 1, \dots, 0)) = g_\lambda(\mathbf{x})$   
( $\lambda$ )

•  $\mathcal{D}_U^* \supseteq \mathcal{D}$

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



# U-Boost algorithm

---

(A) Find  $f_1 = \arg \min_{g \in \mathcal{D}} L_U(g)$

(B) Update  $f_k \rightarrow f_{k+1} = \xi^{-1}((1 - \alpha_{k+1})\xi(f_k) + \alpha_{k+1}\xi(g_{k+1}))$  st  
 $(\alpha_{k+1}, g_{k+1}) = \arg \min_{(\alpha, g) \in (0,1) \times \mathcal{D}} L_U(\xi^{-1}((1 - \alpha)\xi(f_k) + \alpha\xi(g)))$

(C) Select  $K$ , and  $\hat{f} = \xi^{-1}((1 - \alpha_K)\xi(f_{K-1}) + \alpha_K\xi(g_K))$

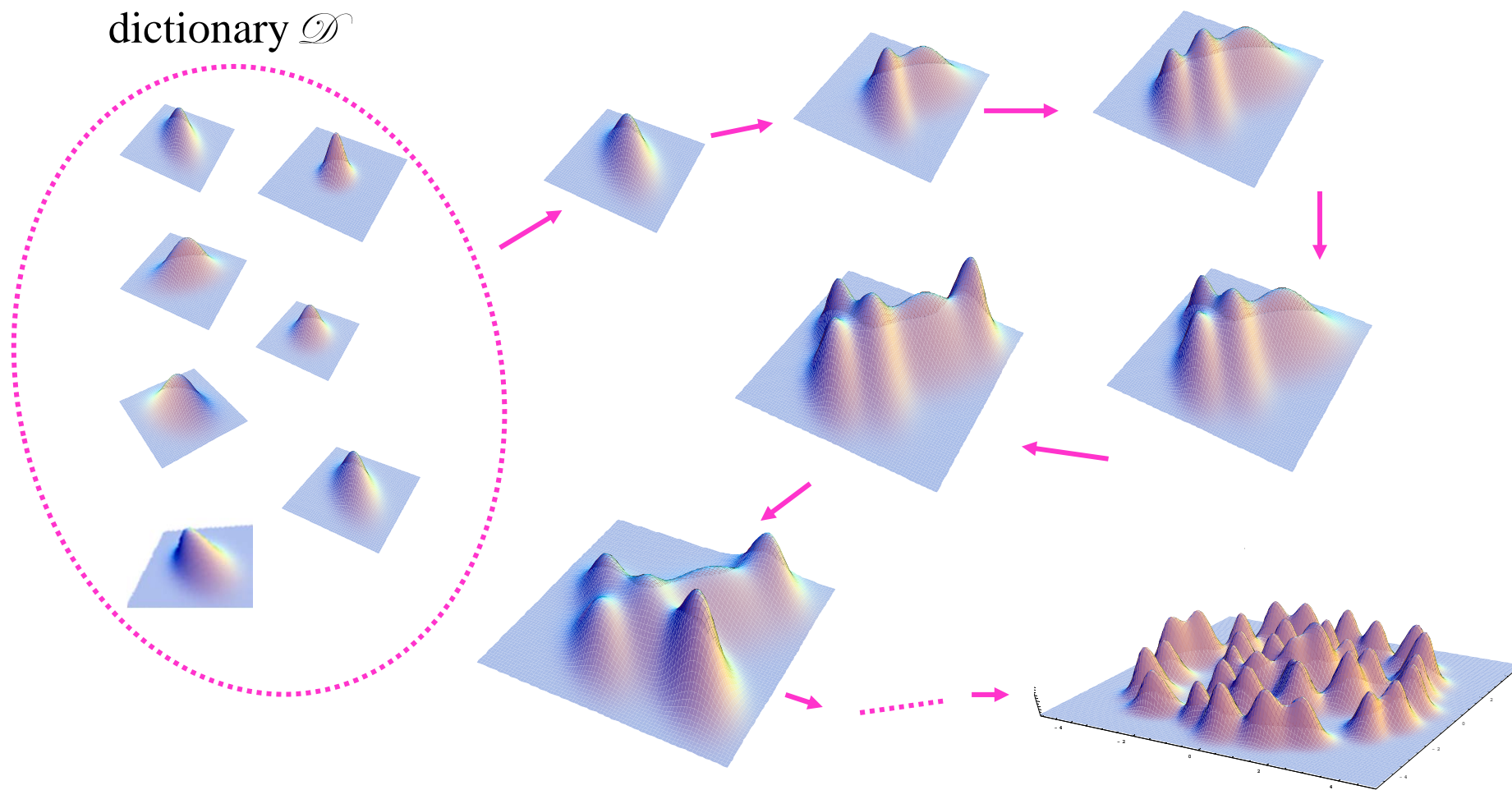
---

Example 2. **Power entropy**  $f^*(x) = \left( \sum_k \pi_k g_k(x)^\beta \right)^{\frac{1}{\beta}}$

If  $\beta = 0$ ,  $f^*(\mathbf{x}) = \exp\left(\sum_k \pi_k \log g_k(\mathbf{x})\right) = \prod_k g_k(\mathbf{x})^{\pi_k}$  [Friedman et al \(1984\)](#)

If  $\beta = 1$ ,  $f^*(\mathbf{x}) = \sum_k \pi_k g_k(\mathbf{x})$  [Klemela \(2007\)](#)

# Learning?

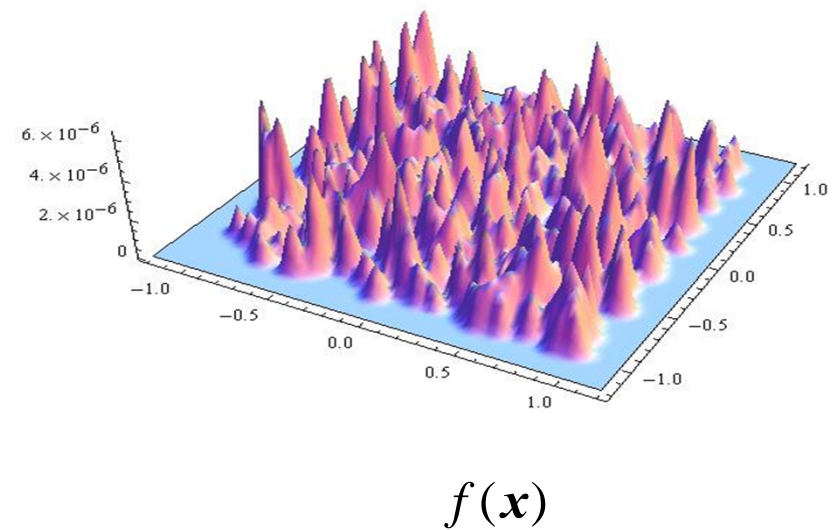
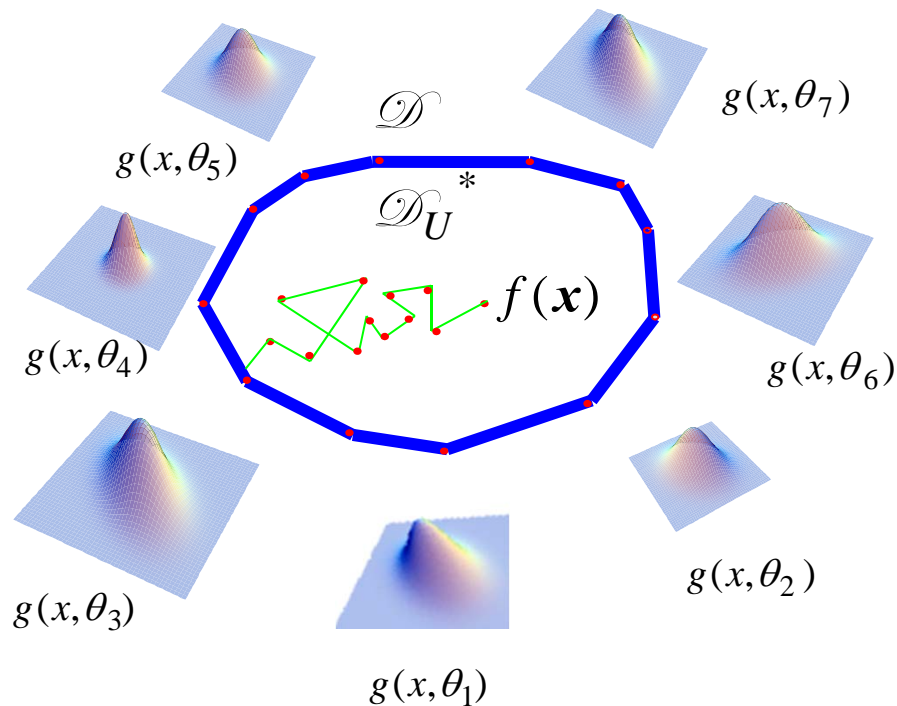


**Boosting is a forward stagewise algorithm**

# Inner step in the convex hull

$$\mathcal{D} = \{g(x, \theta) : \theta \in \Theta\} \longrightarrow \mathcal{D}^* = \text{convex}(\mathcal{D})$$

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

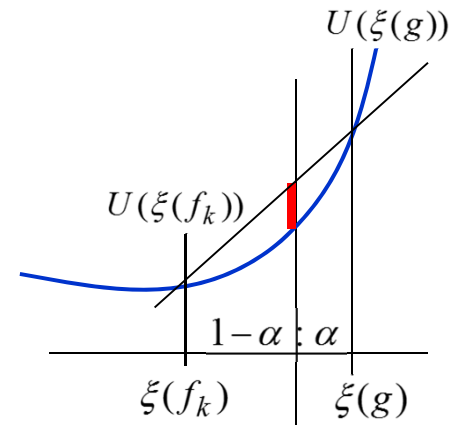


# Loss decomposition

$$L_U(\xi^{-1}((1-\alpha)\xi(f_k) + \alpha\xi(g))) = (1-\alpha)L_U(f_k) + \alpha L_U(g) - \Delta_U(f_k, g, \alpha)$$

where 
$$\Delta_U(f_k, g, \alpha) = \int \underbrace{\{(1-\alpha)U(\xi(f_k)) + \alpha U(\xi(g)) - U((1-\alpha)\xi(f_k) + \alpha\xi(g))\}}$$

- Note 1. (1)  $\Delta_U(f_k, g, \alpha) \geq 0,$   
 (2)  $\Delta_U(f_k, g, \alpha) = 0 \Leftrightarrow f_k = g \text{ (a.e.)}$



Min problem: 
$$\min_{g \in \mathcal{D}} \{ \alpha L_U(g) - \Delta_U(f_k, g, \alpha) \}$$

Make  $L_U(g)$  smaller, and  $\Delta_U(f_k, g, \alpha)$  larger in  $g$  simultaneously.

# Non-asymptotic bound

Naito and Eguchi (2012)

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

Then we have

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

where

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

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

$$\text{IE}(K) = \frac{b_U}{K + c_U} \quad (b_U, c_U \text{ are constants}) \quad (\text{Iteration effect})$$

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

# Statistical applications

(kernel) *U*-PCA   *U*-ICA (mixture)

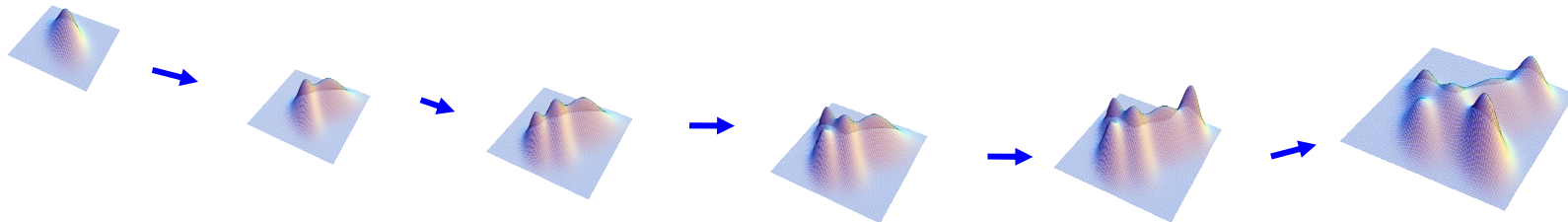
*U*-Boost   *U*-SVM   *U*-AUC   *U*-cluster

.....

Robustness (outliers) (Exp-model, *U*-loss)

Redescendency (hidden structure) (*U*-model, *U*-estimate)

***U*-Boosting density**   Evolve (*U*-model, *U*-loss)



# Drawback of Boosting

$(\hat{\delta}, \hat{f}_*)$  is suffered from over - learning by employing the same dataset  $D$

because 
$$(\hat{\delta}, \hat{f}_*) = \arg \min_{(\delta, f_*) \in F} L_1((1-\delta)\hat{f}_t + \delta f_*, D)$$

- Bagging:  $(\hat{\delta}, \hat{f}_*)$  is learned from bootstrapped samples  $D^*$ 's and averaging  
Breiman (1996)
- Predetermined mixing ratio as  $\delta = (\delta_1, \dots, \delta_t, \dots)$   
Early stopping (Zhang-Yu, 2005), Clemela (2008)
- Boosted Lasso (Zhao - Yu, 2004), Bühlmann (2006)  
Meinshausen and Bühlmann (2011), Stadler-Bühlmann-van de Geer (2010)

# 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  $\Phi$  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.**

# Fisher's equation in 1937

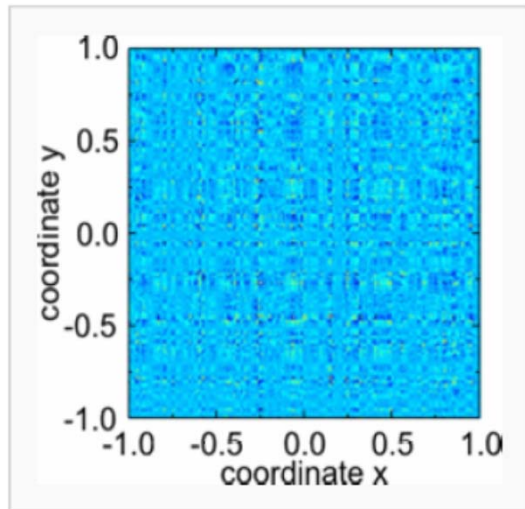
Fisher-Kolmogorov equation

$$\frac{\partial u}{\partial t} = \Delta u + u(1-u)$$

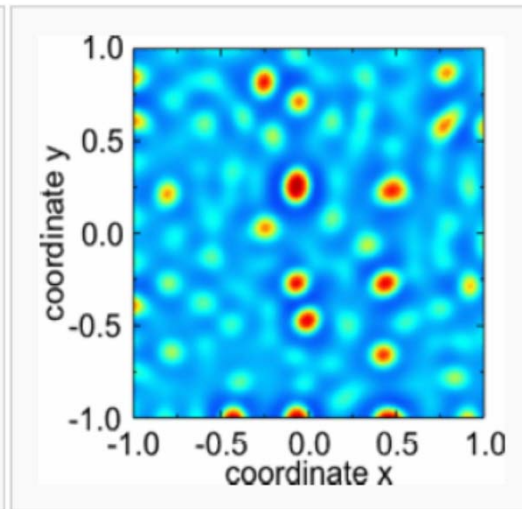
Reaction–diffusion systems

$$\frac{\partial u}{\partial t} = \Delta u + f(u)$$

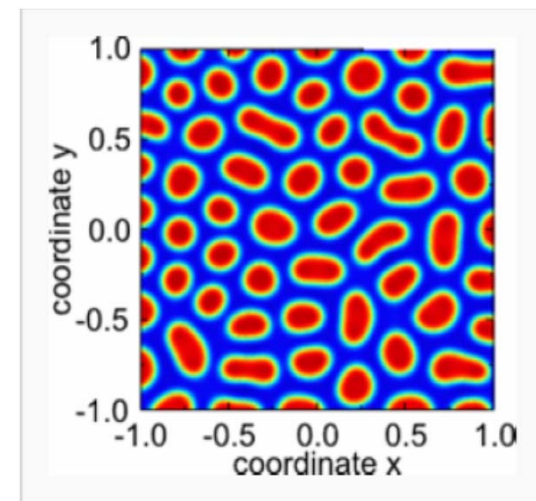
Ecology, physiology, combustion, crystallization, plasma physics, phase transition problems cf. Wiki (reaction diffusion process)



Noisy initial conditions at  $t = 0$ .



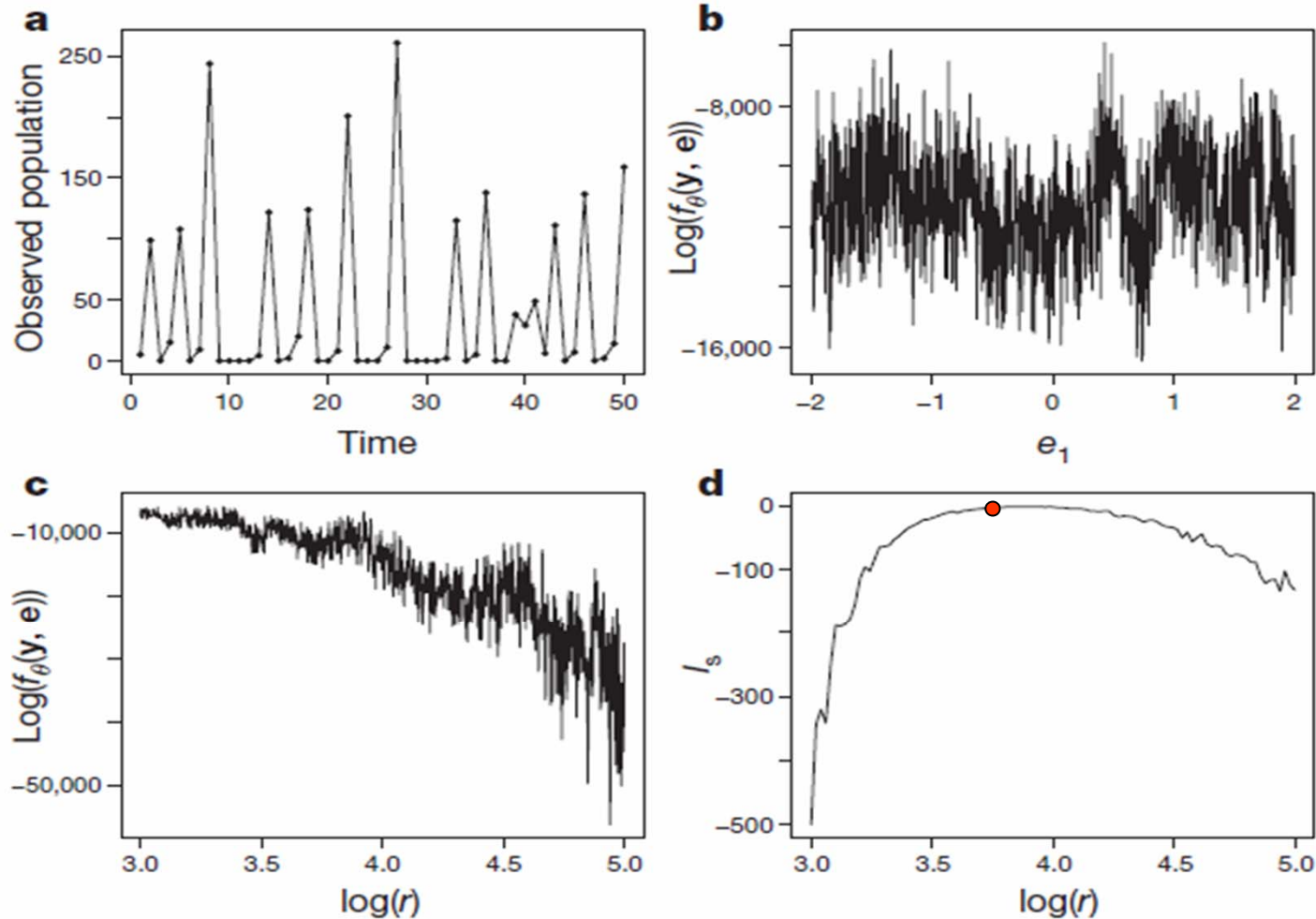
State of the system at  $t = 10$ .



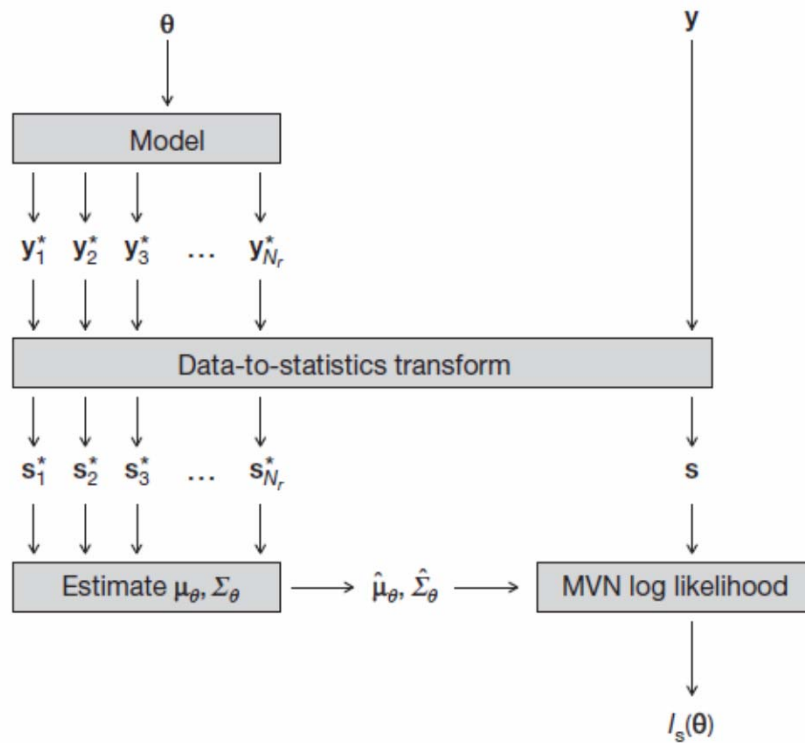
Almost converged state at  $t = 100$ .

# Statistical inference for noisy nonlinear ecological dynamic systems

Simon N. Wood



# Statistical algorithm



Sufficient statistic

$$\mathbf{s} = (s_1, s_2, s_3)^T$$

$$y_t^\alpha = s_1 y_{t-1}^\alpha + s_2 y_t^{2\alpha} + s_3 y_t^{3\alpha} + \varepsilon_t$$

# Simulated curved normal model

Ricker map

$$N_{t+1} = rN_t \exp(-N_t + e_t)$$

where  $\{e_t\}$  are independent  $N(0, \sigma_e^2)$

Sampling

$$y_t \sim \text{Poisson}(\phi N_t)$$

Statistical parameter

$\theta = (r, \sigma_e^2, \phi)$  defines  $f(\mathbf{y}, \mathbf{e}, \theta)$

$$\theta \xrightarrow{\text{simulate}} (\mathbf{y}_1^*, \dots, \mathbf{y}_M^*) \longrightarrow (\mathbf{s}_1^*, \dots, \mathbf{s}_M^*)$$

$$\hat{\boldsymbol{\mu}}_{\theta} = \sum_{i=1}^M \mathbf{s}_i^*, \quad \mathbf{V}_{\theta} = \sum_{i=1}^M (\mathbf{s}_i^* - \hat{\boldsymbol{\mu}}_{\theta})(\mathbf{s}_i^* - \hat{\boldsymbol{\mu}}_{\theta})^{\top}$$

$$L(\theta) = -\frac{1}{2}(\mathbf{s} - \hat{\boldsymbol{\mu}}_{\theta})^{\top} \mathbf{V}_{\theta}^{-1} (\mathbf{s} - \hat{\boldsymbol{\mu}}_{\theta}) - \frac{1}{2} \log \det(2\pi \mathbf{V}_{\theta})$$

# Conclusions and Future directions

Release Fisher's mind control

