

Optimization inspired by Information Theory

Flemming Topsøe

University of Copenhagen

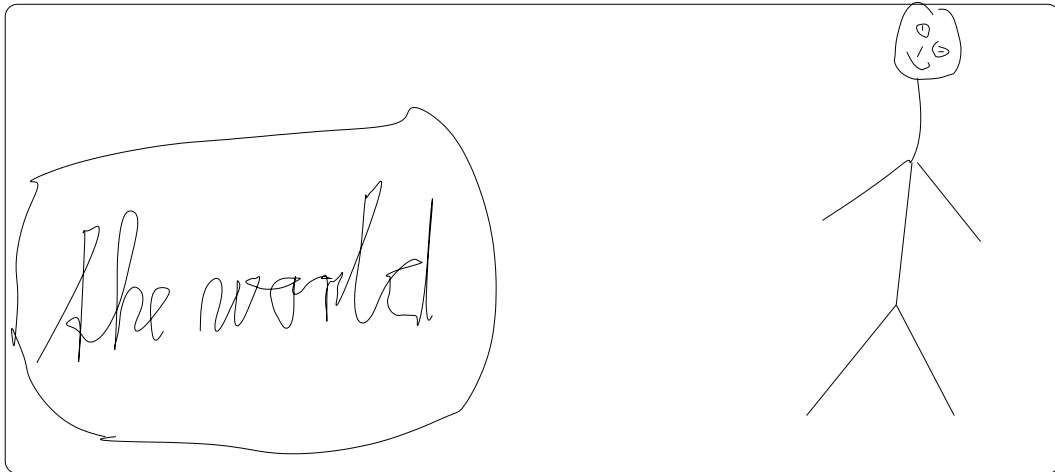
Department of Mathematical Sciences

Seminar, Statistics and Probability Theory Group,

April 9th, 2008

Inspired by previous work on information theoretical optimization, an axiomatic approach to certain special two-person zero-sum games is developed. One of the players (the statistician, physicist, investor, planner, or ...), is imagined to have a "mind", the other (the data, the physical system, the market, the costumers, ...) not. Somewhat provocatorially, it is the ambition of the speaker that after the talk, the audience will (finally!) understand what lies behind the notion of an exponential family.

A source of inspiration



Natures side

Player I

Observers side (you!)

Player II

A situation of **conflict** between two “persons”:
Player I who does not have a mind,
Player II – you – who does.

... and another

Motto:

When you seek an extremum, *don't* differentiate!

Example: $x^2 - 6x + 1 = ?$

1.st method: Differentiate

– and you are done.

2.nd method: Write expression as $(x - 3)^2 - 8$

– and you are done

Best method? 2.nd! And more clear, e.g. you do indeed find a *minimum* this way.

Motto turned positive:

If your problem is *natural*,
there is an *intrinsic method* !

Games!

(X, Y, Φ) defines a *two-person zero-sum game* with Φ as *objective function* (or *complexity*) if $\Phi : X \times Y \rightarrow \mathbb{R}$ (or $\overline{\mathbb{R}}$) and we associate notions of *optimal strategies* and *equilibrium* with this function.

Players and strategies:

Player I is a **maximizer**, chooses x

Player II is a **minimizer**, chooses y .

Specific and global values:

$$\text{val}_I(x) = \inf_{y \in Y} \Phi(x, y) = \inf \Phi_x \left(\text{entropy! } H(x) \right)$$

$$\text{val}_I = \sup_{x \in X} \text{val}_I(x),$$

$$\text{val}_{II}(y) = \sup_{x \in X} \Phi(x, y) = \sup \Phi^y \left(\text{risk! } R(y) \right)$$

$$\text{val}_{II} = \inf_{y \in Y} \text{val}_{II}(y).$$

$y \in Y$ is an *optimal response* to $x \in X$ if

$\Phi(x, y) = \text{val}_I(x)$. We put $\hat{x} = \text{resp}_{II}(x) = \text{argmin} \Phi_x$.

Terminology: y is *adapted to* x or x is a *matching strategy* to y (matches Pl.-II aim to respond optimally).

$$y \in \text{resp}_{\text{II}}(x) \Leftrightarrow x \in \text{match}_{\text{I}}(y) \Leftrightarrow y \in \text{argmin} \Phi_x,$$

$$x \in \text{resp}_{\text{I}}(y) \Leftrightarrow y \in \text{match}_{\text{II}}(x) \Leftrightarrow x \in \text{argmax} \Phi^y.$$

Visualization: e.g. as a matrix

	y_1	y_2	y_3	y_4	y_5	$\text{val}_{\text{I}}(\cdot)$
x_1	0	0	5	1	3	0
x_2	3	5	1	1	1	1
x_3	2	0	1	5	2	0
$\text{val}_{\text{II}}(\cdot)$	3	5	5	5	3	

Redundancy: Compare the *potentially possible* with the *actually achieved* to obtain *Player-I redundancy* and *Player-II redundancy*:

$$\delta_{\text{I}}(x, y) = \text{val}_{\text{II}}(y) - \Phi(x, y),$$

$$\delta_{\text{II}}(x, y) = \Phi(x, y) - \text{val}_{\text{I}}(x) \quad (\text{divergence! } D(x, y)).$$

Define $\text{span}(x, y) = \text{val}_{\text{II}}(y) - \text{val}_{\text{I}}(x)$, then:
 $\text{span}(x, y) = \delta_{\text{I}}(x, y) + \delta_{\text{II}}(x, y)$, hence:
minimax inequality holds: $\text{val}_{\text{I}} \leq \text{val}_{\text{II}}$.

Game Theoretical Equilibrium: If $\text{val}_I = \text{val}_{II} \in \mathbb{R}$. Ideally: γ_Φ is in **game theoretical equilibrium** and **has optimal strategies**, say (x_0, y_0) . Notation: $\gamma_\Phi \in \text{GTE}(x_0, y_0)$. If (x_0, y_0) not specified, write $\gamma_\Phi \in \text{GTE}^*$.

If $\Phi(x_0, y_0) \in \mathbb{R}$, then (**Player-I satisfaction**):

$$\Phi(x, y_0) \leq \Phi(x_0, y_0) \text{ for all } x \in X,$$

$$\text{val}_{II}(y_0) = \Phi(x_0, y_0),$$

$$\delta_I(x_0, y_0) = 0,$$

$$x_0 \in \text{resp}_I(y_0),$$

$$y_0 \in \text{match}_{II}(x_0)$$

are equivalent and so are (**Player-II satisfaction**):

$$\Phi(x_0, y_0) \leq \Phi(x_0, y) \text{ for all } y \in Y,$$

$$\text{val}_I(x_0) = \Phi(x_0, y_0),$$

$$\delta_{II}(x_0, y_0) = 0,$$

$$y_0 \in \text{resp}_{II}(x_0),$$

$$x_0 \in \text{match}_I(y_0).$$

saddle-value theorem: Let $(x_0, y_0) \in X \times Y$ and assume that $\Phi(x_0, y_0) \in \mathbb{R}$. Then the following conditions are equivalent:

$$\begin{aligned} & \gamma_\Phi \in \text{GTE}(x_0, y_0), \\ \rightarrow & \forall (x, y) : \Phi(x, y_0) \leq \Phi(x_0, y_0) \leq \Phi(x_0, y), \\ & \text{val}_I(x_0) = \Phi(x_0, y_0) = \text{val}_{II}(y_0), \\ & \delta_I(x_0, y_0) = \delta_{II}(x_0, y_0) = 0, \\ & x_0 \in \text{resp}_I(y_0) \text{ and } y_0 \in \text{resp}_{II}(x_0), \\ & y_0 \in \text{match}_{II}(x_0) \text{ and } x_0 \in \text{match}_I(y_0). \end{aligned}$$

If so, $\text{val}(\gamma_\Phi) = \Phi(x_0, y_0)$.

In this case, we talk about: *Nash equilibrium, Nash equilibrium pair, saddle-value inequalities, saddle point.*

Thus $\gamma_\Phi \in \text{GTE}^* \Leftrightarrow \exists$ saddle point.

...

(figures very helpful - plan to have some on the black-board)

Subgames, preparations

Introduce **subgames** by restricting strategy set for Player I: $\gamma_\Phi(X_0)$ corresponding to a **preparation**, X_0 . Let Γ be class of all subgames. Expand notation: $\text{val}_I(X_0)$, $\text{val}_{II}(y|X_0)$ and $\text{val}_{II}(X_0)$ if necessary.

Level- and **sub-level sets** become important: $L^y(h) = \{\Phi^y = h\}$, $SL^y(h) = \{\Phi^y \leq h\}$, in full:

$$SL^y(h) = \{x \in X | \Phi^y(x) \leq h\} = \{x \in X | \Phi(x, y) \leq h\}$$

From the saddle-value theorem:

(x_0, y_0) are optimal strategies for a subgame in equilibrium iff $\Phi(x_0, y_0) \in \mathbb{R}$ and y_0 is adapted to x_0 . If so, the possible preparations are all X_0 with

$$\{x_0\} \subseteq X_0 \subseteq SL^{y_0}(h) \text{ with } h = \Phi(x_0, y_0).$$

The sets $(SL^y(h))_{y,h}$ with the **level** h some finite value of Φ^y for an argument which matches y are thus the **maximal preparations**.

Typically the maximal preparations are not “practically feasible” – but level sets are. Given a preparation X_0 , define the **exponential family (Player-II domain)**, by

$$\mathcal{E}_{\text{II}}(X_0) = \{y \in Y \mid \exists h \in \mathbb{R} : X_0 \subseteq L^y(h)\}.$$

An $y \in \mathcal{E}_{\text{II}}(X_0)$ is a **robust PI.-II strategy** and h is the **level of robustness**.

If y_0 is robust and adapted to $x_0 \in X_0$, then $\gamma_\Phi(X_0) \in \text{GTE}(x_0, y_0)$ and $\text{val}(\gamma_\Phi(X_0)) = \Phi(x_0, y_0)$.

Exponential family for a **preparation family** \mathcal{X} :

$$\mathcal{E}_{\text{II}}(\mathcal{X}) = \bigcap_{X_0 \in \mathcal{X}} \mathcal{E}_{\text{II}}(X_0).$$

If \mathcal{E}_{II} is an exponential family (PI.-II domain), the corresponding **exponential family (PI.-I domain)** is

$$\mathcal{E}_{\text{I}} = \{x \mid \text{resp}_{\text{II}}(x) \cap \mathcal{E}_{\text{II}} \neq \emptyset\} = \text{match}_{\text{I}}(\mathcal{E}_{\text{II}}).$$

From \mathcal{E}_{II} define family \mathcal{X} of **associated preparations**:

$$\mathcal{E}_{\text{II}}^\perp = \bigcup_{y \in \mathcal{E}_{\text{II}}} \{L^y(h) \mid h \in \mathbb{R}\}.$$

More another time, to be worked out (geometry etc.)

Axioms for Complexity, entropy, divergence.

Strategy sets are X, Y , a map $x \curvearrowright \hat{x}$ of X into Y gives the *response*. $\text{MOL}(X)$ denotes set of *molecular measures*: $\{\alpha \in M_+^1(X) \mid \text{supp}(\alpha) \text{ finite}\}$.

Axiom 1 *Linking identity* $\Phi(x, y) = H(x) + D(x, y)$ holds, $D \geq 0$ and $D(x, y) = 0 \Leftrightarrow y = \hat{x}$.

Axiom 2 X is convex and Φ affine in first variable: For $y \in Y, \alpha \in \text{MOL}(X)$,

$$\Phi\left(\sum_{x \in X} \alpha_x x, y\right) = \sum_{x \in X} \alpha_x \Phi(x, y).$$

Axiom 3 X is topological, algebraic operations continuous and, for each $(x_0, y_0) \in X \times Y, x \curvearrowright D(x, y_0)$ and $x \curvearrowright D(x_0, \hat{x})$ are lower semi-continuous.

Axiom 4 Every *D-Cauchy sequence* has a convergent subsequence.

$(x_n)_{n \geq 1}$ *D-Cauchy* means:

$$\lim_{n, m \rightarrow \infty} D\left(x_n, \left(\frac{1}{2}x_n + \frac{1}{2}x_m\right)\right) = 0.$$

First consequences

Introduce **barycentre** $b(\alpha) = \sum_{x \in X} \alpha_x x$, and associated **information rate**

$$I(\alpha) = \sum_{x \in X} \alpha_x D(x, \widehat{b(\alpha)}).$$

Concavity and convexity properties:

Let $\alpha \in \text{MOL}(X)$. Then

$$H\left(\sum_{x \in X} \alpha_x x\right) = \sum_{x \in X} \alpha_x H(x) + I(\alpha)$$

and, if $H(b(\alpha)) < \infty$, then, for every $y \in Y$,

$$\sum_{x \in X} \alpha_x D(x, y) = D\left(\sum_{x \in X} \alpha_x x, y\right) + I(\alpha).$$

Last identity is the **compensation identity**.

Main theorem

Assume Axioms 1-4 are satisfied and let X_0 be a convex preparation.

A sequence $(x_n) \subseteq X_0$ is *asymptotically optimal* if $\lim_{n \rightarrow \infty} H(x_n) = H_{\max}$. A strategy $x \in X$ (not necessarily in X_0) is an *H_{\max} -attractor* if $D(x_n, \hat{x}) \rightarrow 0$ for every asymptotically optimal sequence (x_n) .

If X_0 is convex and $H_{\max}(X_0) < \infty$, then PI. II has a unique optimal strategy y^* , and a H_{\max} -attractor x^* exists and $y^* = \widehat{x^*}$. The game is in equilibrium and for each $x \in X_0$ and each $y \in Y$:

$$H(x) + D(x, y^*) \leq H_{\max}(X_0).$$

Creation of Information Triples

Atomic Triples, Integration

(ϕ, h, d) with $X = Y = \text{real interval}$, and response the identity leads to *atomic information triples*.

Example 1 y_0 a **prior**,

$$\phi(x, y) = (x - y)^2 - (x - y_0)^2,$$

$$h(x) = -(x - y_0)^2,$$

$$d(x, y) = (x - y)^2.$$

Example 2

$$\phi(x, y) = x \ln \frac{1}{y},$$

$$h(x) = x \ln \frac{1}{x},$$

$$d(x, y) = x \ln \frac{x}{y}.$$

Examples are of **Bregman type**: for “smooth” strictly concave h , (ϕ, h, d) with ϕ and d defined by

$$\begin{aligned}\phi(x, y) &= h(y) + (x - y) h'(y), \\ d(x, y) &= h(y) - h(x) + (x - y) h'(y),\end{aligned}$$

is an atomic information triple.

A natural process of **integration** leads to more general triples. Given measure μ on set T and then some function space $X \subseteq I^T$, take identity as response and define (Φ, H, D) by integration, i.e.

$$\Phi(x, y) = \int_T \phi(x(t), y(t)) d\mu(t)$$

and similarly for H and D

By integration, Example 1 extends to a triple over Hilbert space:

$$\begin{aligned}\Phi(x, y) &= \|x - y\|^2 - \|x - y_0\|^2, \\ H(x) &= -\|x - y_0\|^2, \\ D(x, y) &= \|x - y\|^2.\end{aligned}$$

And similarly, Example 2 leads to standard discrete information theory by integration w.r.t. counting measure over an “alphabet”.

Equivalence, Relativization

Equivalence results from adding to both Φ and to H an affine function defined on X

If (Φ, H, D) is given and you add $x \mapsto -\Phi(x, y_0)$, you obtain the **relativized triple with y_0 as prior**:

$$\tilde{\Phi}(x, y) = D(x, y) - D(x, y_0)$$

$$\tilde{H}(x) = -D(x, y_0)$$

$$\tilde{D}(x, y) = D(x, y).$$

(for this, it suffices that D satisfies the compensation identity). Leads to Kullback’s **minimum information discrimination principle**.

Randomization

Start with (Φ, H, D) . Allow *randomized strategies* $\alpha \in MOL(X)$ for Player I. Put $b(\alpha) = \sum_{x \in X} \alpha_x x$. Randomization then gives:

$$\begin{aligned}\hat{\alpha} &= \widehat{b(\alpha)}, \\ \tilde{\Phi}(\alpha, y) &= \sum_{x \in X} \alpha_x \Phi(x, y), \\ \tilde{H}(\alpha) &= \sum_{x \in X} \alpha_x \Phi(x, \widehat{b(\alpha)}), \\ \tilde{D}(\alpha, y) &= D(b(\alpha), y).\end{aligned}$$

By equivalence you obtain:

$$\begin{aligned}\tilde{\Phi}_0(\alpha, y) &= \sum_{x \in X} \alpha_x D(x, y), \\ \tilde{H}_0(\alpha) &= \sum_{x \in X} \alpha_x D(x, \widehat{b(\alpha)}), \\ \tilde{D}_0(\alpha, y) &= D(b(\alpha), y).\end{aligned}$$

Classical Information Theory

Let \mathbb{A} , the *alphabet*, be discrete, put $X = M_{+}^1(\mathbb{A})$, and $Y = K(\mathbb{A})$, the set of *code length functions* over \mathbb{A} , i.e. the set of $\kappa : \mathbb{A} \rightarrow [0, \infty]$ such that *Krafts equality*

$$\sum_{i \in \mathbb{A}} e^{-\kappa_i} = 1$$

holds. The response $P \curvearrowright \kappa$ is defined by $\kappa_i = \ln \frac{1}{p_i}$; $i \in \mathbb{A}$ and for Φ we take *average code length*, i.e.

$$\Phi(P, \kappa) = \langle \kappa, P \rangle = \sum_{i \in \mathbb{A}} p_i \kappa_i .$$

Then

$$H(P) = \sum_{i=1}^n p_i \ln \frac{1}{p_i},$$
$$D(P, Q) = \sum_{i=1}^n p_i \ln \frac{p_i}{q_i} .$$

Consider preparation of the form

$$\mathcal{P} = \{P | \langle f_k, P \rangle = a_k \text{ for } k = 1, \dots, m\}$$

and seek MaxEnt-distribution. Lagrange multipliers one possibility. Better: Seek $\kappa \in \mathcal{E}_{II}$. Clearly, κ of the form

$$\kappa = \alpha + \beta_1 f_1 + \dots + \beta_m f_m$$

works. By Krafts equality this requires that

$$\alpha = \ln Z(\beta_1, \dots, \beta_m)$$

with the **partition function** given by

$$Z(\beta_1, \dots, \beta_m) = \sum_{i \in \mathbb{A}} \exp(-\beta_1 f_1(i) - \dots - \beta_m f_m(i)).$$

Adjust coefficients so that the matching distribution is consistent ($\in \mathcal{P}$), and you are done.

.....

Further examples include separation, location theory, universal coding, determination of capacity, duality and more ... Next time!