# Information and Games

Flemming Topsøe (topsoe@math.ku.dk)
University of Copenhagen
Department of Mathematical Sciences

"Facets of Entropy" Copenhagen, October 24-26, 2007

based on
http://www.math.ku.dk/~topsoe/jogo2.pdf and
http://www.math.ku.dk/~topsoe/CTnexttopsoe.pdf
Overview: existence and identification

\*\*\*\*\*\*\*\*\*\*

So said Kolmogorov (1983):
*"Information theory
must
precede probability theory
and not be based on it"*

Let us follow him!

Nature
holds the truth!
does not have a mind!
an absolute
passive

Observer
seeks the truth
has a mind
YOU!
inventive

In so doing, I reveal nature of entropy and divergence:

- Entropy is minimal complexity
- Divergence measures discrepancy between actual and minimal complexity

– and what then is complexity ?

- Complexity is whatever complexity should be in any concrete situation.

– hmmm. So axiomatize based on game theoretical thinking! — but why games?

- because they provide a perfect setting for the modeling of conflict situations

Is that what information is about? yes(!?)
Which tools?

- the simplest are the most important: two-person zero-sum games!

Modeling of un-symmetric conflict situations between players, one having a mind, the other not - with focus: complexity ($\Phi$), entropy ($H$) and divergence ($D$)

or: pay-off ($\Psi$), maximal pay-off ($\Pi$), divergence ($D$)

- strategy sets $X$ and $Y$ (e.g. X=Y),
- connection $X \to Y$ written $x \curvearrowright \hat{x}$ (e.g. $\hat{x} = x$),
- complexity $\Phi : X \times Y \to ]-\infty, \infty]$ (or $[0, \infty]$),
- entropy H $: X \to ]-\infty, \infty]$ (or $[0, \infty]$),
- divergence D $: X \times Y \to [0, \infty]$
- and possibly a preparation $X_0 \subseteq X$ (or more!)

**Axiom 1** (linking): For $(x, y) \in X \times Y$,
$\Phi(x, y) = H(x) + D(x, y)$ (linking identity) and
$D(x, y) = 0 \Leftrightarrow y = \hat{x}$
**Axiom 2** (affinity): $X$ is convex and $\Phi$ affine in
first variable: $\Phi(\sum a_\nu x_\nu, y) = \sum a_\nu \Phi(x_\nu, y)$
**Axiom 3** (semi-continuity): reference topology $(X, \tau)$ is
Hausdorff, algebraic operations continuous and, for any
$(x_0, y_0)$, $x \curvearrowright D(x, y_0)$ and $x \curvearrowright D(x_0, \hat{x})$
are $\tau$-lower semi-continuous.
**Axiom 4** (weak completeness): For a sequence
$(x_n)$ in $X$, put $x_{n,m} = \frac{1}{2}x_n + \frac{1}{2}x_m$ and
$y_{n,m} = \widehat{x_{n,m}}$. If "Cauchy property" $D(x_n, y_{n,m}) \to 0$
as $n, m \to \infty$ then $\exists x, (x_{n_k})_{k \geq 1}$ such that $x_{n_k} \to x$

strong convergence: $x_n \twoheadrightarrow x$ if $D(x_n, \hat{x}) \to 0$

Sometimes pay-off Ψ, maximal pay-off Π and Divergence D as before are more natural to work with.
By "duality" you pass from the one system to the other:
$(\Phi, \mathsf{H}, \mathsf{D}) \leftrightarrow (\Psi, \Pi, \mathsf{D})$ with $\Phi = -\Psi$, $\mathsf{H} = -\Pi$.
We talk about information triples based on complexity and information triples based on pay-off.

**Example 1** (classical information theory) $\mathbb{A}$ a discrete alphabet, $X = M_+^1(\mathbb{A})$, set of probability distributions over $\mathbb{A}$, $Y = K(\mathbb{A})$, set of code length functions over $\mathbb{A}$, i.e. set of $\kappa : \mathbb{A} \to [0, \infty]$ such that Krafts equality holds: $\sum_{i \in \mathbb{A}} \exp(-\kappa_i) = 1$. Let $P \curvearrowright \hat{P}$ be the bijection $P \leftrightarrow \kappa$ with $\kappa_i = \ln \frac{1}{p_i}$, $p_i = \exp(-\kappa_i)$. With $\Phi$ as average code length and with

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

$$\mathsf{H}(P) = \sum_{i=1}^{n} p_i \ln \frac{1}{p_i}$$

$$\mathsf{D}(P, \kappa) = \sum_{i=1}^{n} p_i \ln \frac{p_i}{q_i} \text{ (here } Q \leftrightarrow \kappa) \qquad \square$$

**Example** (geometric version of updating)

$X = Y$, a Hilbert space, with identity as connection. Let $y_0$ be a point in $Y$, the prior. With

$$\Psi(x, y) = \|x - y_0\|^2 - \|x - y\|^2$$
$$\Pi(x) = \|x - y_0\|^2$$
$$\mathsf{D}(x, y) = \|x - y\|^2,$$

Axioms 1-4 hold. □

For standard updating you will replace $\|x - y\|^2$ by a standard divergence, say Kullback-Leibler divergence. The reason why seemingly different objects can be used for the same task lies in the common identities fulfilled by these quantities. They involve entropy, divergence and information rate $\mathrm{I}(\cdot)$ (really the same as mutual information ) and hold whenever Axioms 1 and 2 hold. Let $\mathrm{MOL}(X)$ be the set of probability distributions over $X$ with finite support. Then, for $\alpha \in \mathrm{MOL}(X)$, define

$$\mathrm{I}(\alpha) = \sum_{x \in X} \alpha_x \, \mathsf{D}(x, \widehat{\bar{x}})$$

**Theorem (identities re concavity- and convexity)**

(i) Let $\overline{x} = \sum_{x \in X} \alpha_x x$ be a convex combination of elements in $X$ corresponding to $\alpha \in \mathsf{MOL}(X)$. Then

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

(ii) With notation as in (i), assume that $\mathsf{H}(\overline{x}) < \infty$ and let $y \in Y$. Then

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

(iii) For elements $\alpha_1, \cdots, \alpha_m$ in $\mathsf{MOL}(X)$ with barycentres $\overline{x_1} \ldots, \overline{x_m}$, and for any mixture $\alpha = \sum_1^m w_k \alpha_k$ with a barycentre $\overline{x}$ of finite entropy, the following identity holds:

$$\mathsf{I}\left(\sum_{k=1}^m w_k \alpha_k\right) = \sum_{k=1}^m w_k \mathsf{I}(\alpha_k) + \sum_{k=1}^m w_k \mathsf{D}(\overline{x_k}, \overline{x}).$$

**Proof:** (i) is trivial by linking and (i) $\Rightarrow$ (ii) $\Rightarrow$. (iii) $\qquad \square$

Consider $(\Phi, H, D)$ satisfying Axioms 1-4 and a preparation $X_0$. This defines two-person zero-sum game $\gamma(X_0)$ with $\Phi$ as objective function, Observer as minimizer and Nature as optimizer, but with restriction to strategies in $X_0$. Some important notions:

- $x \in X_0$: consistent strategy
- $\sup_{x \in X_0} \inf_{y \in Y} \Phi(x, y) = \sup_{x \in X_0} H(x)$:
  maximum entropy value, $H_{max} = H_{max}(X_0)$
- $\inf_{y \in Y} \sup_{x \in X_0} \Phi(x, y) = \inf_{y \in Y} R(y)$:
  minimal risk, $R_{min} = R_{min}(X_0)$.
- minimax inequality: $H_{max} \leq R_{min}$.
- *equilibrium* means that $H_{max} = R_{min} < \infty$.
- $x$ MaxEnt-strategy: consistent $x$ with $H(x) = H_{max}$.
- A sequence $(x_n)$ of consistent strategies is
  asymptotically optimal if $\lim_{n \to \infty} H(x_n) = H_{max}$.
- $x \in X$ is a $H_{max}$-attractor if $x_n \twoheadrightarrow x$ for
  every asymptotically optimal sequence $(x_n)$.
- $y \in Y$ $R_{min}$-strategy if $R(y) = R_{min}(X_0)$.
- $(x^*, y^*)$ optimal pair: $x^*$ MaxEnt-strat., $y^*$ $R_{min}$-strat.

**Main Theorem** $X_0$ convex, $\mathsf{H}_{\max}(X_0) < \infty$. Then:

- Observer has unique optimal strategy $y^*$
- an $\mathsf{H}_{\max}$-attractor $x^*$ exists and $y^* = \widehat{x^*}$.
- $\mathsf{H}_{\max}$-attractors are equivalent, hence unique if connection is injective
- the game is in equilibrium: $\mathsf{H}_{\max}(X_0) = \mathsf{R}_{\min}(X_0)$
- for $x \in X_0$, $y \in Y$ strong inequalities hold:

$$\mathsf{H}(x) + \mathsf{D}(x, y^*) \leq \mathsf{H}_{\max}(X_0)$$
$$\mathsf{R}_{\min}(X_0) + \mathsf{D}(x^*, y) \leq \mathsf{R}(y).$$

Application to updating: Consider $(\Phi, \mathsf{H}, \mathsf{D})$ satisfying Axioms 1-4 with $x \curvearrowright \hat{x}$ injective and, to simplify, $\Phi$ finite. Let $X_0 \subseteq X$ be convex and let $y_0 \in Y$ be a "prior". Define:

$$\Psi(x, y) = \mathsf{D}(x, y_0) - \mathsf{D}(x, y)$$
$$\Pi(x) = \mathsf{D}(x, y_0).$$

Consider corresponding game. For $y \in Y$, let

$$\Gamma(y) = \inf_{x \in X_0} \Psi(x, y)$$

guaranteed updating gain associated with $y$ and let

$$\Gamma_{\mathsf{max}} = \sup_{y \in Y} \Gamma(y), \; \mathsf{D}_{\mathsf{min}} = \inf_{x \in X_0} \mathsf{D}(x, y_0)$$

A strategy $x^* \in X$ is the generalized I-projection of $y_0$ on $X_0$ if $x_n \twoheadrightarrow x^*$ for every sequence $(x_n)$ in $X_0$ which is asymptotically optimal in the sense that $\lim_{n \to \infty} \mathsf{D}(x, y_0) = \mathsf{D}_{\mathsf{min}}$.

---

**Theorem** The game is in equilibrium: $\Gamma_{\mathsf{max}} = \mathsf{D}_{\mathsf{min}}$. There is a unique generalized I-projection $x^*$ of $y_0$ on $X_0$ and $y^* = \widehat{x^*}$ is the unique optimal strategy for Observer: $\Gamma(y^*) = \Gamma_{\mathsf{max}}$. Furthermore, for $(x, y) \in X_0 \times Y$,

$$\mathsf{D}(x, y_0) \geq \mathsf{D}_{\mathsf{min}} + \mathsf{D}(x, y^*),$$
$$\Gamma(y) + \mathsf{D}(x^*, y) \leq \Gamma_{\mathsf{max}}.$$

---

This follows from applying main theorem to

$$(x, y) \curvearrowright -\Psi(x, y) = \mathsf{D}(x, y) - \mathsf{D}(x, y_0)$$
$$x \curvearrowright \mathsf{H}(x) - \Phi(x, y_0) = -\mathsf{D}(x, y_0)$$
$$(x, y) \curvearrowright \Phi(x, y) - \mathsf{H}(x) = \mathsf{D}(x, y).$$

**Theorem** Let $X_0$ be a closed convex subset of the Hilbert space $Y$ and let $y_0 \notin X_0$. Then there exists a hyperplane which separates $y_0$ from $X_0$.

**Proof:** Take $(\Psi, \Pi, \mathsf{D})$ as before:

$$\Psi(x, y) = \|x - y_0\|^2 - \|x - y\|^2$$
$$\Pi(x) = \|x - y_0\|^2$$
$$\mathsf{D}(x, y) = \|x - y\|^2,$$

Conclude from main Theorem that the associated game is in equilibrium. The value for Nature is $\inf_{x \in X_0} \|x - y_0\|^2$ which is positive by assumption. Therefore, the value of the game for Observer must also be positive, i.e.

$$\sup_{y \in Y} \inf_{x \in X_0} \left( \|x - y_0\|^2 - \|x - y\|^2 \right) > 0.$$

We conclude that there exists $y \in Y$ such that

$$\|x - y_0\| > \|x - y\| \text{ for all } x \in X_0.$$

This shows that the hyperplane of all $x$ with the same distance to $y$ as to $y_0$ separates $y_0$ and $X_0$. $\square$

Change focus from existence to identification. Key to this: Nash equilibrium! Requires pair $(x^*, y^*) \in (X_0, Y)$ such that saddle value inequalities $\Phi(x, y^*) \leq \Phi(x^*, y^*) \leq \Phi(x^*, y)$ hold for all $(x, y) \in X_0 \times Y$. If also $\Phi(x^*, y^*) < \infty$, equilibrium follows with $(x^*, y^*)$ a (MaxEnt,$R_{min}$)-pair. With our special assumptions we find:

**Theorem** Let $x^* \in X_0$, put $y^* = \widehat{x^*}$. Then the game is in equilibrium with $(x^*, y^*)$ as (MaxEnt,$R_{min}$)-pair *iff* $H(x^*) < \infty$ and Nash's inequality $\Phi(x, y^*) \leq \Phi(x^*, y^*)$ holds for all $x \in X_0$.

Nash inequality gives strong inequalities (Pythagorean ineq.).

Useful corollaries: Kuhn-Tucker type theorems (not discussed here) and robustness lemma below. Given preparation $X_0$, $y \in Y$ is robust if $\exists c < \infty \forall x \in X_0 : \Phi(x, y) = c$. Define exponential family:
$\mathcal{E} = \mathcal{E}(X_0) = \{y | y \text{ robust } \}$.

**Lemma** If $x^*$ is consistent and $y^* = \widehat{x^*}$ robust then $(x^*, y^*)$ is a (MaxEnt,$R_{min}$)-pair.

Only Axiom 1 required for this! Proof is easy.

*Return to standard probabilistic setting*:

Discrete alphabet $\mathbb{A}$, both strategy sets $= M_+^1(\mathbb{A})$, connection=identity. Now use notation: $P$, $Q$ rather than $x$, $y$ and $\mathcal{P}$ for typical preparation. We aim at discussing preparations given by linear constraints. Given set $f = (f_\nu)_{1 \leq \nu \leq k}$ of real functions on $\mathbb{A}$. The natural preparations are

$$\mathcal{P}_a = \{P | \langle f_\nu, P \rangle = a_\nu \text{ for } 1 \leq \nu \leq k\}$$

with $a = (a_\nu)_{1 \leq \nu \leq k} \in \mathbb{R}^k$.
Natural exponential family:

$$\mathcal{E} = \{Q | \ Q \text{ robust for all the } \mathcal{P}_a\}.$$

Question: which complexity measures ?
try to simplify search for distributions in $\mathcal{E}$!

key idea: Take $\Phi$ of the form

$$\Phi(P, Q) = \xi_Q\Big( \langle \overline{\kappa}(Q), P \rangle \Big) \text{ where}$$

$$\Big( \overline{\kappa}(Q) \Big)(i) = \kappa(q_i) \text{ for } i \in \mathbb{A} \text{ thus}$$

$$\langle \overline{\kappa}(Q), P \rangle = \sum_{i \in \mathbb{A}} p_i \kappa(q_i) \,.$$

with assumptions:
- The $\xi_Q$'s increasing and concave (e.g. linear)
- the coder $\kappa$ is smooth, decreasing, convex, $\kappa(1) = 0$
- $\Phi$ satisfies Axiom 1

classical: $\xi_Q$'s the identity map, $\kappa(q) = \ln \frac{1}{q}$
Then $\kappa^{-1}$ is restriction of $x \curvearrowright \exp(-x)$ to $[0, \infty]$. Entropy generated by this measure of complexity is standard BGS-entropy.

From previous page:

$$\Phi(P,Q) = \xi_Q\Big(\langle\overline{\kappa}(Q),P\rangle\Big) \text{ where}$$

$$\langle\overline{\kappa}(Q),P\rangle = \sum_{i\in\mathbb{A}} p_i\kappa(q_i)\,.$$

Trivial but key observation: any $Q$ for which $\overline{\kappa}(Q)$ is a linear combination of the constant function 1 and the given functions $f_1,\cdots,f_k$, i.e. of the form

$$\overline{\kappa}(Q) = \lambda_0 + \lambda_1 f_1 + \cdots + \lambda_k f_k = \lambda_0 + \lambda\cdot f$$

for certain constants $\lambda_0$ and $\lambda = (\lambda_1,\cdots,\lambda_k)$, is a member of $\mathcal{E}$.

Motivated by this, fix constants $\lambda = (\lambda_1,\cdots,\lambda_k)$ and ask: $\exists\lambda_0,Q : \overline{\kappa}(Q) = \lambda_0 + \lambda\cdot f$? This amounts to $q_i = \kappa^{-1}(\lambda_0+\lambda\cdot f(i))$. Summarizing what this leads to we find:

**Theorem (MaxEnt calculus)** Let $\lambda = (\lambda_1, \cdots, \lambda_k)$ be given constants. Then "normally", the equation

$$\sum_{i \in \mathbb{A}} \kappa^{-1}\Big(\lambda_0 + \lambda \cdot f(i)\Big) = 1$$

has a solution, necessarily unique, and $Q$ given by

$$q_i = \kappa^{-1}\Big(\lambda_0 + \lambda \cdot f(i)\Big) \text{ for } i \in \mathbb{A}$$

has the stipulated form, hence belongs to the exponential family $\mathcal{E}$. This distribution is the MaxEnt-distribution for $\mathcal{P}_a$ with $a = (a_1, \cdots, a_k)$ given by

$$a_\nu = \sum_{i \in \mathbb{A}} q_i f_\nu(i) \text{ for } \nu = 1, \cdots, k$$

and, for this value of $a$,

$$\mathsf{H}_{\max}(\Phi, \mathcal{P}_a) = \xi_Q(\lambda_0 + \lambda \cdot a) \,.$$

Theorem replaces and expands the standard recipe for MaxEnt-calculations. Focus on $\lambda_0$ rather than on the classical partition function.

Which complexity measures?

**Example** The complexity measures

$$\Phi^B(P,Q) = \frac{1}{q-1} + \sum \left( q_i^q - \frac{q}{q-1} p_i q_i^{q-1} \right)$$

$$\Phi^C(P,Q) = \frac{1}{1-q} \sum p_i^q (1 - q_i^{1-q})$$

$$\Phi^R(P,Q) = \frac{1}{1-q} \left( \frac{\sum p_i^q}{\sum p_i^q q_i^{1-q}} - 1 \right).$$

all give Tsallis entropy. Only $\Phi^B$ is good! □

First restrict form of $\Phi$: From

$$\Phi(P,Q) = \xi_Q \Big( \langle \overline{\kappa}(Q), P \rangle \Big) \text{ to}$$

$$\Phi(P,Q) = \langle \overline{\kappa}(Q), P \rangle + \overline{\xi}(Q) \text{ with}$$

$$\overline{\xi}(Q) = \sum_{i \in \mathbb{A}} \xi(q_i).$$

$\kappa$ is the coder, $\xi$ the corrector.

Generation of information triples à la Bregman:

Bregman generator: a strictly concave and smooth real function h defined on $[0, 1]$ with $h(0) = h(1) = 0$ and $h'(1) = -1$.

From h we generate two more functions, $\phi = \phi(p, q)$, and $d = d(p, q)$:

$$\phi(p, q) = h(q) + (p - q) \, h'(q) \,,$$
$$d(p, q) = h(q) - h(p) + (p - q) \, h'(q) \,.$$

Consider the internal functions, $\Phi = \Phi_h$, $H = H_h$ and $D = D_h$ generated by $\phi$, h and d:

$$\Phi(P, Q) = \sum_{i \in \mathbb{A}} \phi(p_i, q_i) \,,$$
$$H(P) = \sum_{i \in \mathbb{A}} h(p_i) \,,$$
$$D(P, Q) = \sum_{i \in \mathbb{A}} d(p_i, q_i) \,.$$

[ the slides ended here – orally I ended by stressing the importance of not considering entropy alone, and the essentials of using the method of generation à la Bregman ]