# BEYOND SHANNON Generalized entropies and rational inquiry

Ancona

Bayes by the Sea

September 2018

Vincenzo CRUPI

Center for Logic, Language, and Cognition Department of Philosophy and Education University of Turin <u>vincenzo.crupi@unito.it</u>

www.vincenzocrupi.com





Cognitive Science 42 (2018) 1410–1456 © 2018 Cognitive Science Society, Inc. All rights reserved. ISSN: 1551-6709 online DOI: 10.1111/cogs.12613

#### Generalized Information Theory Meets Human Cognition: Introducing a Unified Framework to Model Uncertainty and Information Search

Vincenzo Crupi,<sup>a</sup> Jonathan D. Nelson,<sup>b,c</sup> Björn Meder,<sup>c</sup> Gustavo Cevolani,<sup>d</sup> Katya Tentori<sup>e</sup>

\*Center for Logic, Language, and Cognition, Department of Philosophy and Education, University of Turin <sup>b</sup>School of Psychology, University of Surrey <sup>c</sup>Center for Adaptive Behavior and Cognition, Max Planck Institute for Human Development <sup>d</sup>IMT School for Advanced Studies, Lucca <sup>c</sup>Center for Mind/Brain Sciences, University of Trento

Received 27 January 2017; received in revised form 5 March 2018; accepted 6 March 2018

#### Abstract

Searching for information is critical in many situations. In medicine, for instance, careful choice of a diagnostic test can help narrow down the range of plausible diseases that the patient might have. In a probabilistic framework, test selection is often modeled by assuming that people's goal is to reduce uncertainty about possible states of the world. In cognitive science, psychology, and medical decision making, Shannon entropy is the most prominent and most widely used model to formalize probabilistic uncertainty and the reduction thereof. However, a variety of alternative entropy metrics (Hartley, Quadratic, Tsallis, Rényi, and more) are popular in the social and the natural sciences, computer science, and philosophy of science. Particular entropy measures have been predominant in particular research areas, and it is often an open issue whether these divergences emerge from different theoretical and practical goals or are merely due to historical accident. Cutting across disciplinary boundaries, we show that several entropy and entropy reduction measures arise as special cases in a unified formalism, the Sharma-Mittal framework. Using mathematical results, computer simulations, and analyses of published behavioral data, we discuss four key questions: How do various entropy models relate to each other? What insights can be obtained by considering diverse entropy models within a unified framework? What is the psychological plausibility of different entropy models? What new questions and insights for research on human information acquisition follow? Our work provides several new pathways for theoretical and empirical research, reconciling apparently conflicting approaches and empirical findings within a comprehensive and unified information-theoretic formalism.

Keywords: Entropy; Uncertainty; Value of information; Information search; Probabilistic models

# key rival models from the literature

	entropy measure	expected entropy reduction	key references
quadratic entropy	$I - \sum_{h \in H} P(h)^2$	$\sum_{e \in E} P(e) \left[ \sum_{h \in H} P(h \mid e)^2 - \sum_{h \in H} P(h)^2 \right]$	Gini (1912) Niiniluoto and Toumela (1973) Horwich (1982) Crupi and Tentori (2014)
Hartley entropy	$ln\sum_{h\in H} P(h)^0$	$\sum_{e \in E} P(e) \left[ ln \sum_{h \in H} P(h)^{\circ} - ln \sum_{h \in H} P(h \mid e)^{\circ} \right]$	Hartley (1928) Aczél, Forte, and Ng (1974)
Shannon entropy	$\sum_{h \in H} P(h) \ln\left(\frac{1}{P(h)}\right)$	$\sum_{e \in E} P(e) \left[ \sum_{h \in H} P(h \mid e) \ln \left( P(h \mid e) \right) - \sum_{h \in H} P(h) \ln \left( P(h) \right) \right]$	Shannon (1948) Carnap and Bar-Hillel (1953) Oaksford and Chater (1994, 2003) Nelson (2005, 2008) Nelson <i>et al.</i> (2010)
error entropy	$I - \max_{h \in H} \left[ P(h) \right]$	$\sum_{e \in E} P(e) \left[ \max_{h \in H} \left[ P(h \mid e) \right] - \max_{h \in H} \left[ P(h) \right] \right]$	Fano (1961) Baron (1985) Baron, Beattie, and Hershey (1988) Crupi, Tentori, and Lombardi(2009) Meder and Nelson (2012) Rusconi <i>et al.</i> (2014)

#### key rival models from the literature



A graphical illustration of quadratic ( ° ), Hartley (•), Shannon (+), and error entropy (solid line) as distinct measures of uncertainty over a binary hypothesis set  $H = \{h, \neg h\}$ as a function of the probability of h. (Note that Hartley entropy "jumps" to 0 for extreme values of *P*(*h*).)

Rényi (1961): entropy as a (parametric) generalized mean of Shannon's original atomic information values [i.e., ln(1/P(x))]

**Rényi (1961)**: entropy as a (parametric) <u>generalized mean</u> of Shannon's original atomic information values [i.e., *ln*(1/*P*(*x*))]



Rényi (1961): entropy as a (parametric) <u>generalized mean</u> of Shannon's original atomic information values [i.e., ln(1/P(x))] <u>Rényi entropies</u> (r = order):

$$\Rightarrow ent_{P}^{(r)}(H) = \frac{1}{1-r} ln\left[\sum_{h \in H} P(h)^{r}\right]$$

Hartley and Shannon entropy are special cases (for r = 0, I, respectively), but <u>not</u> quadratic and error

Rényi (1961): entropy as a (parametric) generalized mean of Shannon's original atomic information values [i.e., ln(1/P(x))] <u>**Rényi entropies**</u> (r = order):

$$\Rightarrow ent_{P}^{(r)}(H) = \frac{1}{1-r} ln\left[\sum_{h \in H} P(h)^{r}\right]$$

Hartley and Shannon entropy are special cases (for r = 0, I, respectively), but <u>not</u> quadratic and error

Tsallis (1988): standard weighted average of a <u>parametric generalization</u> of Shannon's original atomic information values [→Tsallis logarithm]



A graphical illustration of the generalized atomic information function  $ln_t[1/P(h)]$ based on Tsallis's generalized logarithm for four different values of the parameter t (0, 1, 2, and 5, respectively, for the curves from top to bottom).



**Rényi entropies** 
$$(r = order)$$
:

$$= \operatorname{ent}_{P}^{(r)}(H) = \frac{1}{1-r} \ln \left[ \sum_{h \in H} P(h)^{r} \right]$$

Hartley and Shannon entropy are special cases (for r = 0, I, respectively), but <u>not</u> quadratic and error

Tsallis (1988): standard weighted average of a <u>parametric generalization</u> of Shannon's original atomic information values [→Tsallis logarithm]

Tsallis entropies (
$$t = degree$$
):  

$$ent_{P}^{(t)}(H) = \frac{1}{t-1} \left[ 1 - \sum_{h \in H} P(h)^{t} \right]$$

Shannon and quadratic entropy are special cases (for t = 1, 2, respectively), but <u>not</u> Hartley and error

Rényi (1961): entropy as a (parametric) <u>generalized mean</u> of Shannon's original atomic information values [i.e., ln(1/P(x))]

Tsallis (1988): standard weighted average of a <u>parametric generalization</u> of Shannon's original atomic information values [→Tsallis logarithm] Sharma-Mittal entropies: biparametric family of measures of <u>order</u> *r* and <u>degree</u> *t* 

$$- ent_{P}^{(r,t)}(H) = \frac{1}{t-1} \left[ 1 - \left( \sum_{h \in H} P(h)^{r} \right)^{\frac{t-1}{r-1}} \right]$$
$$ent_{P}^{(r,t)}(H) = ln_{t}e_{r} \left[ \sum_{h \in H} P(h) ln_{r} \left( \frac{1}{P(h)} \right) \right]$$













A map of how different entropy measures are derived within the Sharma-Mittal framework for different settings of the order (r) and degree (t) parameters.

#### (i) irrelevant combination

<u>ex</u>.: A card was drawn from a well-shuffled stardard deck and kept hidden.  $H = \{2, 3, ..., king, ace\}$  is the set of (thirteen) hypotheses concerning the value of the card drawn, while  $K = \{\text{hearts}, \text{diamonds}, \text{clubs}, \text{spades}\}$ is the set of (four) hypotheses concerning its suit. So knowing the true value of the combined variable  $H \times K$  (value *plus* suit) amounts to knowing exactly which card has been drawn out of the whole deck of fifty-two. (We assume H and K probabilistically independent.) You can ask about H.

#### (i) irrelevant combination

<u>ex</u>.: A card was drawn from a well-shuffled stardard deck and kept hidden.  $H = \{2, 3, ..., \text{king, ace}\}$  is the set of (thirteen) hypotheses concerning the value of the card drawn, while  $K = \{\text{hearts, diamonds, clubs, spades}\}$ is the set of (four) hypotheses concerning its suit. So knowing the true value of the combined variable  $H \times K$  (value *plus* suit) amounts to knowing exactly which card has been drawn out of the whole deck of fifty-two. (We assume H and K probabilistically independent.) You can ask about H.

<u>Shannon</u>:  $R_p^{S}(H,H) = ln(I3) = R_p^{S}(H \times K,H)$ 



#### (i) irrelevant combination

<u>ex</u>.: A card was drawn from a well-shuffled stardard deck and kept hidden.  $H = \{2, 3, ..., \text{king, ace}\}$  is the set of (thirteen) hypotheses concerning the value of the card drawn, while  $K = \{\text{hearts, diamonds, clubs, spades}\}$ is the set of (four) hypotheses concerning its suit. So knowing the true value of the combined variable  $H \times K$  (value *plus* suit) amounts to knowing exactly which card has been drawn out of the whole deck of fifty-two. (We assume H and K probabilistically independent.) You can ask about H.

- <u>Shannon</u>:  $R_p^s(H,H) = ln(I3) = R_p^s(H \times K,H)$
- <u>quadratic</u>:  $R_{P}^{Q}(H,H) = 0.92 > 0.23 = R_{P}^{Q}(H \times K,H)$

(i) irrelevant combination

## (ii) <u>commutativity</u>

<u>ex</u>.: A card was drawn from a well-shuffled stardard deck and kept hidden.  $H = \{\text{red}, \text{black}\}, \text{while } K = \{\text{hearts}, \text{diamonds}, \text{clubs}, \text{spades}\}$ 

(i) irrelevant combination

## (ii) <u>commutativity</u>

<u>ex</u>.: A card was drawn from a well-shuffled stardard deck and kept hidden.  $H = \{\text{red}, \text{black}\}, \text{while } K = \{\text{hearts}, \text{diamonds}, \text{clubs}, \text{spades}\}$ 

<u>Shannon</u>:  $R_{P}^{S}(H,K) = ln(2) = R_{P}^{S}(K,H)$ 



(i) irrelevant combination

## (ii) <u>commutativity</u>

<u>ex</u>.: A card was drawn from a well-shuffled stardard deck and kept hidden.  $H = \{\text{red}, \text{black}\}, \text{while } K = \{\text{hearts}, \text{diamonds}, \text{clubs}, \text{spades}\}$ 

<u>Shannon</u>:  $R_{p}^{S}(H,K) = ln(2) = R_{p}^{S}(K,H)$ <u>quadratic</u>:  $R_{p}^{Q}(H,K) = 0.75 > 0.19 = R_{p}^{Q}(K,H)$ 

#### (i) <u>additivity of what</u>?

for <u>all and only</u> Rényi entropies (i.e., SM entropies of degree t = 1):  $ent_{p}(H \times K) = ent_{p}(H) + ent_{p}(K)$  if  $H \perp_{P} K$ 

while for degree t = 2 (including quadratic):

 $ent_{P}(H \times K) = ent_{P}(H) + ent_{P}(K) - ent_{P}(H)ent_{P}(K)$  if  $H \perp_{P} K$ 

but for <u>ALL</u> SM entropies (of any degree t and order r):  $R_{p}(H, E \times F) = R_{p}(H, E) + R_{p}(H, F)$ if  $E \perp_{p} F \mid H$ 

#### (i) <u>additivity of what</u>?

## (ii) order vs. concavity?

order (r) is an index of the *imbalance* of the entropy function: it indicates how much the entropy measure discounts minor hypotheses order-0 measures  $\rightarrow$  the actual probability distribution is neglected: live hypothesis are just counted, as if they were all equally important order- $\infty$  measures  $\rightarrow$  only the most likely hypothesis matters,

and all minor hypotheses are disregarded altogether

intermediate values of  $r \rightarrow$  more likely hypotheses count more, but less likely hypotheses do retain some weight: the higher (lower) r is, the more (less) it so happens that more likely hypotheses are regarded and less likely hypotheses are discounted



#### (i) <u>additivity of what</u>?

#### (ii) order vs. concavity?

but most Rényi entropies (r > 1) are <u>non-concave</u>, thus allowing for *informationally detrimental* experiments!

#### <u>ex</u>:

 $prior P(H) = \{0.63, 0.185, 0.185\}$   $P(e) = 0.44 \qquad P(H|e) = \{1, 0, 0\}$   $P(not-e) = 0.56 \qquad P(H|not-e) = \{0.33..., 0.33..., 0.33...\}$   $P(f) = 0.9 \qquad P(H|f) = \{0.63, 0.185, 0.185\}$   $P(not-f) = 0.1 \qquad P(H|not-f) = \{0.63, 0.185, 0.185\}$ which test is most useful,  $E = \{e, not-e\}$  or  $F = \{f, not-f\}$ ?

#### (i) <u>additivity of what</u>?

## (ii) order vs. concavity?

but most Rényi entropies (r > 1) are <u>non-concave</u>, thus allowing for *informationally detrimental* experiments!

order WITH concavity  $\rightarrow$  power entropies (quadratic for order r = 2)

$$ent_{P}^{(r)}(H) = I - \left(\sum_{h \in H} P(h)^{r}\right)^{\frac{1}{r-1}}$$

#### entropy, surprise, and evidential support

Shannon and quadratic entropy of H amount to the expected surprise of finding out the true element of H, with surprise =  $ln_t[1/P(h)]$ , for t = 1 and t = 2, respectively

suprise  $\rightarrow$  probababilistic evidential support as reduction of surprise:

$$C_{P}(h,e) = ln_{t}\left(\frac{I}{P(h)}\right) - ln_{t}\left(\frac{I}{P(h \mid e)}\right)$$

- $t = I \rightarrow \log probability ratio measure$
- $t = 2 \rightarrow$  probability *difference* measure

#### entropy, inaccuracy, and divergence

Shannon and quadratic entropy of *H* also amount to the *expected inaccuracy* of *P*, with

$$lnacc(P, w_{i}) = t \cdot ln_{t} \left(\frac{1}{P(h_{i})}\right) - ln_{t} \left(\frac{1}{\sum_{h \in H} P(h)^{t}}\right)$$

for 
$$t = 1$$
 and  $t = 2$ , respectively

$$t = I \rightarrow logarithmic score$$

 $t = 2 \rightarrow Brier score$ 

if divergence = reduction (difference) of expected inaccuracy from a posterior perspective (e.g., Roche & Shogenji), then  $t = 1 \rightarrow$  KL-divergence, while  $t = 2 \rightarrow$  squared Euclidean distance and their expectations are identical to the expected reduction of Shannon and quadratic entropy, respectively

## conclusions: with quadratic instead of Shannon...

- irrelevant combination and communitativity problems are solved...
- additive behaviour of (independent) experiments with regards to a target hypothesis space is retained...
- a better behaved (because uniformly concave) continuum of entropy measures of varying order is obtained...
- a connection with a much better behaved probabilistic measure of evidential support is established...
- and the advantages of a close link with sound (in)accuracy and divergence measures are also preserved