

# IMPLICATION ENTROPIQUE ET CAUSALITE

Régis GRAS<sup>1</sup>, Raphaël COUTURIER<sup>2</sup>

## TITLE

### Entropic implication and causality

## RESUME

La méthode d'analyse de données, l'Analyse Statistique Implicative (A.S.I.), a pour objectif, à partir du croisement d'une population et de variables, l'extraction et l'étude de certaines relations d'association, appelées règles, de type implicatif donc non symétrique. Un indice statistique, l'intensité d'implication, permet d'attribuer un indicateur numérique de qualité à une quasi-règle de type  $a \Rightarrow b$  (si  $a$  alors en général  $b$ ) où  $a$  et  $b$  sont des variables de natures diverses. A travers différentes applications, est apparue la nécessité d'adapter le concept d'intensité à des situations où les populations en jeu deviennent très importantes. Nous fournissons ici une réponse à cette demande par une autre intensité d'implication formalisée sur des bases utilisant le concept d'entropie de Shannon et sensible aux variations de cardinaux.

*Mots-clés : contre-exemple, implication classique, implication entropique, indice d'inclusion, règle*

## ABSTRACT

The Statistical Implicative Analysis (SIA) aims at extracting and studying association relationship, called rules, from the crossing of a population and some variables. Those rules are non symmetrical. A statistical index, the implicative intensity, provides a quality measure of a quasi-rule of the form  $a \Rightarrow b$  (if  $a$  then in general  $b$ ) where  $a$  and  $b$  are variables of different natures. With different applications, it appears that the requirement to adapt the concept of intensity to situations where the populations at stake become large. In this paper we meet this requirement with a new implicative intensity based on the Shannon's entropy concept.

*Keywords : counter-example, classic implication, entropic implication, rule*

## 1 Introduction

L'analyse implicative classique établit une mesure de qualité - l'intensité d'implication  $\varphi(a,b)$  - à la règle  $a \Rightarrow b$

- d'une part, en comparant le nombre de contre-exemples à cette règle observés dans la contingence à celui que l'on observerait en toute hypothèse si les deux variables  $a$  et  $b$  étaient indépendantes ;
- d'autre part, en lui associant une mesure qui valorise la probabilité de l'écart entre les contre-exemples contingents à la règle et l'observable selon un modèle de distribution du nombre de contre-exemples aléatoires sous cette hypothèse.

---

<sup>1</sup> Ecole Polytechnique de l'Université de Nantes, Equipe Connaissance et Décision, Laboratoire d'Informatique de Nantes-Atlantique (LINA), UMR 6241, [regisgra@club-internet.fr](mailto:regisgra@club-internet.fr)

<sup>2</sup> FEMTO-ST, département DISC, Université de Franche-Comté, [raphael.couturier@univ-fcomte.fr](mailto:raphael.couturier@univ-fcomte.fr)

L'objectif de l'ASI est d'extraire d'un ensemble de données issu du croisement sujets x variables, des lois, des règles qui ne peuvent être ni des croyances, ni des connaissances justifiées (fonction épistémique des lois) mais plutôt des états de choses récurrents, permanents responsables de la vérité ou la validité des croyances et des connaissances (Encyclopedia Universalis, 1990, corpus 13). L'ASI a donc essentiellement une fonction descriptive, voire prédictive dans certains cas. Pour ce faire, on lui confère la qualité d'explicabilité de relations causales. La causalité, dont il est question ici, est une façon d'organiser systématiquement des faits d'expérience et, à travers cette organisation, de leur donner du sens. Ainsi, la causalité n'est pas, dans notre pouvoir et notre volonté, un moyen rationnel d'accéder à une vérité ni à une raison d'être des choses (Encyclopedia Universalis, 1990, corpus 5). La mise en évidence d'une stabilité relationnelle entre un fait a (par exemple un attribut, simple ou complexe) et un autre fait b par une règle  $a \Rightarrow b$  conduit à exprimer la possibilité de désigner le premier comme une cause et le second comme un effet. Soulignons qu'au même effet peuvent être associées conséquemment plusieurs causes disjointes ou une cause conjointe de plusieurs faits. De la même façon, un phénomène peut être cause de plusieurs conséquences. Mais la **signification causale** ne peut pas s'extraire ni de la succession temporelle, ni de la concomitance dont rend compte la corrélation. S'impose, pour cette dernière, une dissymétrie de la relation entre les faits. C'est le facteur qui distingue l'analyse implicative de la ressemblance.

Deux réserves ont été émises par les utilisateurs de l'approche classique exprimée ci-dessus, à travers la variété des applications qu'ils ont rencontrées :

- lorsque les tailles des ensembles d'individus traités augmentent, atteignant des effectifs de l'ordre de plusieurs centaines ou plus, l'intensité d'implication  $\varphi(a,b)$  a tendance à ne plus être suffisamment discriminante car ses valeurs peuvent être très voisines de 1, alors que l'inclusion dont elle cherche à modéliser la qualité, est loin d'être satisfaite. Ce phénomène a été déjà signalé par A. Bodin (1997) dont les travaux traitent des ensembles de grande taille d'élèves impliqués dans des enquêtes internationales;
- le modèle classique de la quasi-implication retient essentiellement la mesure, extraite d'une échelle de probabilité, de l'intensité de la quasi-règle  $a \Rightarrow b$ . Or comme nous l'avons abordé en introduction à propos de la causalité, la prise en compte complémentaire de la contraposée de l'implication de  $\bar{b} \Rightarrow \bar{a}$  (si pas d'« effet » alors pas de « cause ») est indispensable pour renforcer l'évaluation de la qualité suffisamment bonne de la relation de quasi-implication, voire quasi-causale, de a sur b<sup>3</sup>. En même temps, elle pourrait permettre de corriger la difficulté évoquée en relation à la taille des ensembles en jeu. En effet si A et B sont des ensembles de petite taille par rapport à E, leurs complémentaires seront importants et réciproquement. Cependant, si avec l'approche classique nous évaluons de la même façon la rareté des contre-exemples à la règle  $\bar{b} \Rightarrow \bar{a}$ , nous constatons que l'intensité d'implication est la même que celle tirée de la rareté des contre-exemples à la règle  $a \Rightarrow b$ . Ainsi, la prise en compte de la contraposée n'apporte pas de nouvelle information quant à la qualité causale de la règle.

---

<sup>3</sup> Ce phénomène est signalé par Y. Kodratoff dans son article publié dans les Actes du Colloque « Fouille dans les données par la méthode implicative », IUFM de Caen, juin 2000. Nous avons aussi abordé cette question du paradoxe de Hempel en Introduction de l'ouvrage

Deux critères épistémologiques nous sont alors imposés pour pallier ces réserves observées dans les applications et améliorer le modèle formalisé par l'intensité d'implication dans l'approche classique :

- formaliser l'implication de telle façon que la règle directe  $a \Rightarrow b$  et sa contraposée  $\bar{b} \Rightarrow \bar{a}$  soient mesurées de façon spécifique par un indice qui ainsi les distinguera ;
- conserver la sensibilité de la qualité de la règle aux effectifs de toutes les populations en jeu.

Une première solution a été apportée (Gras et Kuntz, 2001) sous le nom d'implication entropique. Elle associait un indice entropique ou d'inclusion, que nous allons rappeler dans le § 2.1, et l'intensité d'implication classique. Or, il est apparu que cette association présentait un caractère quelque peu artificiel, ce qui nous conduit à une révision de cette formalisation. Elle conduira alors à une nouvelle intensité entropique associant cette fois l'indice entropique et un coefficient statistique (§ 2.2). Quelques simulations permettront de comparer dans le § 2.3 ses variations à celles de l'intensité classique en fonction des paramètres intervenant dans les observations. Une application illustrera dans le § 3 cette nouvelle intensité entropique. Dans l'annexe, nous établirons une relation probabiliste entre les deux formes de l'intensité d'implication.

## 2 Formalisation de la quasi-règle implicative dans l'approche entropique.

La solution<sup>4</sup> que nous apportons utilise un indice qui rend compte de la dissymétrie entre les situations  $S_1 = (a \text{ et } b)$  et  $S'_1 = (a \text{ et non } b)$  qui concerne la quasi-règle  $a \Rightarrow b$  ainsi que celle entre les situations  $S_2 = (\text{non } a \text{ et non } b)$  et  $S'_2 = (a \text{ et non } b)$  qui concerne la quasi-règle contraposée. Notons que ce sont les mêmes instances qui contredisent la quasi-implication et sa contraposée. Les valeurs relatives de ces instances sont fondamentales dans notre approche.

### 2.1 Construction d'indice d'inclusion

Pour rendre compte de l'incertitude liée à un éventuel pari de l'appartenance à une des deux situations  $S_1$  ou  $S'_1$ , (resp.  $S_2$  ou  $S'_2$ ), nous avons choisi le concept d'entropie de Shannon (1949). Ainsi nous déterminons l'entropie conditionnelle relative à  $S_1$  et  $S'_1$  lorsque  $a$  est réalisée, incertitude sur  $b$  lorsque  $a$  est connu

$$H(b/a) = -\frac{n_{a \wedge b}}{n_a} \log_2 \frac{n_{a \wedge b}}{n_a} - \frac{n_{a \wedge \bar{b}}}{n_a} \log_2 \frac{n_{a \wedge \bar{b}}}{n_a}$$

puis l'entropie conditionnelle relative à  $S_2$  et  $S'_2$ , incertitude sur non  $a$  lorsque non  $b$  est réalisé :

<sup>4</sup> J. Blanchard apporte dans (Blanchard *et al.*, 2005) une réponse à ce problème par une mesure de « l'écart à l'équilibre ».

$$H(a/b) = -\frac{n_{a\bar{b}}}{n_{\bar{b}}} \log_2 \frac{n_{a\bar{b}}}{n_{\bar{b}}} - \frac{n_{\bar{a}\bar{b}}}{n_{\bar{b}}} \log_2 \frac{n_{\bar{a}\bar{b}}}{n_{\bar{b}}}$$

Ces entropies conditionnelles, contenant des informations mutuelles, sont à valeurs dans  $[0,1]$  et devraient être simultanément faibles. En conséquence, les dissymétries entre les situations  $S_1$  et  $S'_1$  (resp.  $S_2$  et  $S'_2$ ) devraient être simultanément fortes si l'on souhaite disposer d'un bon critère d'inclusion de A dans B. En effet les entropies conditionnelles représentent l'*incertitude* moyenne des expériences qui consistent à observer si b (resp. non a) est réalisé lorsque l'on a observé a (resp. non b). Le complément à 1 de cette incertitude représente donc l'*information* moyenne recueillie par la réalisation de ces expériences. Plus cette information est importante, plus forte est la garantie de la qualité simultanée de l'implication et de sa contraposée. Nous devons maintenant adapter ce critère numérique entropique au modèle attendu dans les différentes situations cardinales.

Pour que le modèle ait la signification attendue, il doit satisfaire, selon nous, les contraintes épistémologiques complémentaires suivantes :

1° il devra intégrer les valeurs de l'entropie et même, pour amplifier les contrastes, prendre le carré de ces valeurs ;

2° ce carré varie aussi de 0 à 1, pour rendre compte du déséquilibre, c'est-à-dire d'une tendance à l'inclusion en s'opposant à l'entropie, c'est-à-dire à l'incertitude. La valeur retenue sera le complément à 1 de son carré ;

3° afin de prendre en compte les deux informations propres à la quasi-implication et à sa contraposée, c'est le produit des valeurs que nous retiendrons. Le produit a la propriété de s'annuler dès que l'un de ses termes s'annule, i.e. dès que cette qualité s'efface ;

4° enfin, le produit ayant une dimension 2 par rapport à l'entropie, nous prendrons sa racine carrée pour revenir à la même dimension et représenter une moyenne géométrique.

5° pour ne pas précipiter le rejet dès les premiers contre-exemples et renforcer le changement de comportement de l'intensité d'implication, dès que leur effectif dépasse la moitié du nombre d'observations potentielles ( $n_a/2$ ), on a jugé préférable d'avoir d'une part, aux bornes une tangente horizontale et d'autre part, au milieu une tangente verticale (point d'inflexion afin de ralentir la décroissance vers 0). D'où une petite transformation analytique des entropies conditionnelles  $H(b/a)$  et  $H(\bar{a}/\bar{b})$  et disposant de caractères analytiques adéquats aux attentes dont cette dernière exigence.

Posons  $\alpha = \frac{n_a}{n}$  la fréquence de a,  $\bar{\beta} = \frac{n_{\bar{b}}}{n}$  la fréquence de non b et  $t = \frac{n_{a\bar{b}}}{n}$  la fréquence des contre-exemples. Nous construisons alors, en utilisant la transformation analytique évoquée, deux fonctions  $h_1$ , qui traduit l'entropie conditionnelle  $H(b/a)$ , et  $h_2$ , qui traduit l'entropie conditionnelle  $H(\bar{a}/\bar{b})$ . Elles sont définies sur  $[0,1]$ , où  $\mathbf{1}(\cdot)$  est une fonction indicatrice et le logarithme log est en base 2.

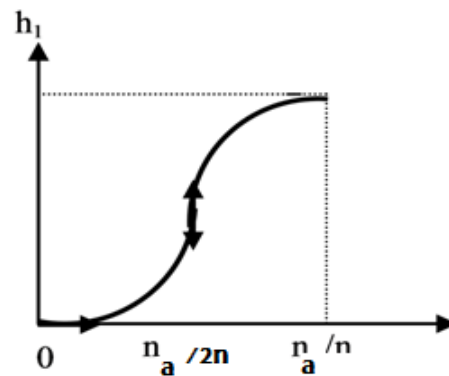


FIGURE 1 – Graphe de la fonction  $h_1$  qui admet une tangente horizontale en  $t=0$  et en  $t=\alpha$ , et une tangente verticale en  $t= \alpha/2$  comme souhaité

De là, nous proposons la définition suivante permettant de déterminer le critère entropique.

**Définition 1:** L'indice d'inclusion de A, support de a, dans B, support de b, est le nombre :

$$i(a, b) = \left( \left[ 1 - h_1(t) \right] \left[ 1 - h_2(t) \right] \right)^{\frac{1}{2}}$$

qui intègre l'information délivrée par la réalisation du nombre de contre-exemples, d'une part à la quasi-règle  $a \Rightarrow b$  et, d'autre part, à la quasi-règle  $\bar{b} \Rightarrow \bar{a}$ .

## 2.2 Formalisation de l'intensité d'implication entropique

L'indice entropique que nous venons de définir est malheureusement indépendant de toute homothétie des effectifs des populations puisqu'il est basé sur des fréquences relatives. Il rejoint ainsi la plupart des indices d'implication de la littérature (Gras & Couturier, 2010) Afin de le rendre sensible aux effectifs, il nous faut introduire un coefficient en  $n$  qui devra valoriser, sur le plan statistique, les observations sur de grandes populations. Il devra factoriser l'indice  $i(a,b)$ , croître avec  $n$ , donc minorer cet indice entropique lorsque  $n$  sera petit, et tendre vers 1 lorsque  $n$  tendra vers l'infini. De plus, il devra permettre de conserver une tangente horizontale à la nouvelle intensité entropique, caractère qui assure une décroissance lente au voisinage d'un nombre de contre-exemples faible.

Considérons la suite des expériences qui consistent à vérifier sur  $n$  tirages indépendants du couple  $(a,b)$  s'il y a ou non contre-exemple à  $a \Rightarrow b$  de probabilité  $p$  d'absence de contre-exemple à chaque tirage. La variable aléatoire, centrée et réduite, du nombre total de contre-exemples suit asymptotiquement la loi normale  $N(0,1)$ . Examinant la probabilité pour que le nombre de contre-exemples aléatoire soit inférieur au nombre observé, le test associé est unilatéral. L'intervalle de confiance de  $p$  au seuil de 0,95 a pour amplitude unilatérale la valeur  $\frac{1}{2\sqrt{n}}$  après majoration de l'écart-type par  $\frac{1}{2}$ . L'amplitude de l'intervalle où se produisent peu ou pas de contre-exemples au-delà

du seuil de 0,05 est donc  $1 - \frac{1}{2\sqrt{n}}$ . C'est un indicateur croissant avec  $n$ , indicateur de confiance qui ne peut que renforcer le critère de qualité de la règle  $a \Rightarrow b$  dont rend compte l'indice d'inclusion. Plus il est grand, plus fiable est le rôle de l'indice entropique pour signifier la qualité de la règle. Aussi, nous adopterons comme indice d'implication entropique le nombre :

$$\Psi(a,b) = \left(1 - \frac{1}{2\sqrt{n}}\right) i(a,b)$$

Ainsi, lorsque le nombre  $n$  d'observations est faible, l'affectation à l'indice entropique du coefficient modérateur relativise la valeur de cet indice pour représenter une relation dont on attendrait la capacité prédictive et la fonction causale.

TABLEAU 1 – **Exemple** : On dispose des données d'observations des variables binaires  $a$  et  $b$  (attributs par exemple) sur 100 individus

	b	Non b	
a	28	2	30
Non a	40	30	70
	68	32	100

On obtient alors :  $h_1 = 0,00739$  ;  $h_2 = 0,00565$  ;  $i(a,b) = 0,99348$  ;  $\Psi(a,b) = 0,9438$

Si l'on divise toutes les valeurs par 10, on obtient  $\Psi(a,b) = 0,8364$

Si l'on multiplie toutes les valeurs par 10, on obtient  $\Psi(a,b) = 0,97777$

Si l'on multiplie toutes les valeurs par 100, on obtient  $\Psi(a,b) = 0,98885$

## 2.3 Quelques simulations

Afin d'observer les comportements respectifs comparés entre l'implication classique  $\phi(a,b)$ , l'indice d'implication entropique  $i(a,b)$  et l'intensité d'implication entropique  $\Psi(a,b)$ , nous simulons plusieurs situations et comparons les courbes obtenues afin de valider les critères exigés.

**Figure 2** :  $n$ ,  $n_a$  et  $n_b$  faibles ; le nombre de contre-exemples varie de 0 à  $n_a$ . Nous observons bien une tangente horizontale au voisinage de 0, un point d'inflexion en  $n_a/2$  (point d'équilibre)<sup>5</sup> pour  $\Psi(a,b) = 0,5$  aux approximations des tracés près, et une valeur nulle en  $n_{a,b} = n_a$ . Nous passons ainsi quant à notre degré de confiance en la règle de l'état « adhérer », aux états successifs : « accepter », « douter », « négliger » et « réfuter ».

Notons de plus que si le nombre de contre-exemples est nul la valeur de  $\Psi(a,b)$  n'est pas égale à 1, ce qui peut choquer le bon sens. Or nous avons dit par ailleurs que l'intensité d'implication mesurait un certain étonnement statistique d'observer un

<sup>5</sup> Nous retrouvons ainsi une propriété dite d'écart à l'équilibre étudiée dans (Blanchard et al, 2005)

certain nombre de contre-exemples à la règle stricte. Il n'est donc pas étonnant que ce nombre soit nul, même si la valeur de l'étonnement, corrigé de la faiblesse de la population ( $n=10$ ) est faible (0,7).

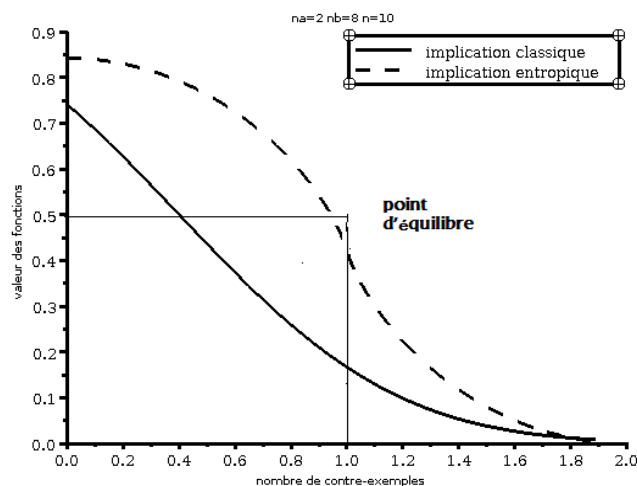


FIGURE 2 –

**Figure 3 :** Tous les paramètres précédents subissent une homothétie x10. Par suite, les valeurs de l'intensité entropique sont améliorées, en particulier pour  $n_{a,b} = 0$  où  $\Psi(a,b) = 0,85$ . On note les mêmes propriétés que celles signalées pour la figure 1 en  $n_a/2$ .

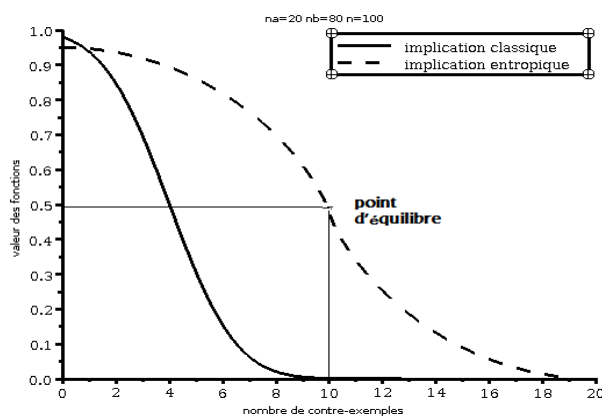


FIGURE 3 –

**Figure 4 :** Cette fois, le coefficient homothétique est 1000 par rapport à la Figure 2. On constate encore une amélioration de l'intensité entropique et surtout une différence entre  $\Psi(a,b)$  et  $\phi(a,b)$  : l'intensité d'implication classique « tarde » à décroître lorsque les contre-exemples se multiplient. Sa valeur, même pour  $n_a/5$  est égale à 1 puis décroît trop rapidement sans nuance.

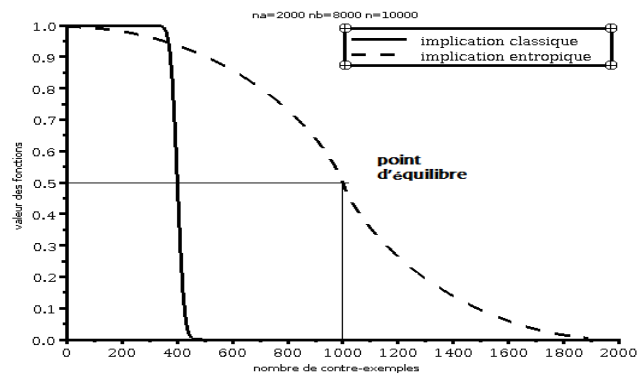


FIGURE 4 –

**Figure 5 :** On retrouve les mêmes phénomènes qu'avec la **Figure 2** avec cependant une meilleure résistance à la décroissance du fait d'un paramètre  $n_a$  acceptant plus de contre-exemples.

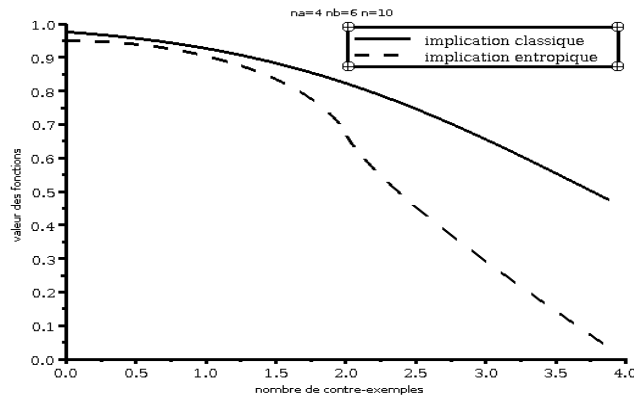


FIGURE 5 –

**Figures 6 et 7 :** Les observations conduisent aux mêmes conclusions que précédemment prenant en compte l'effet homothétie et dilatation de  $n_a$ .

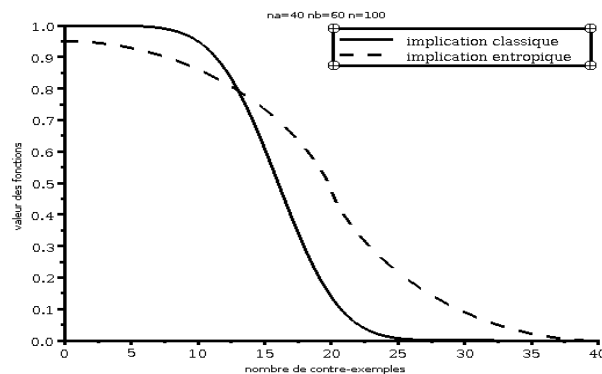


FIGURE 6 –



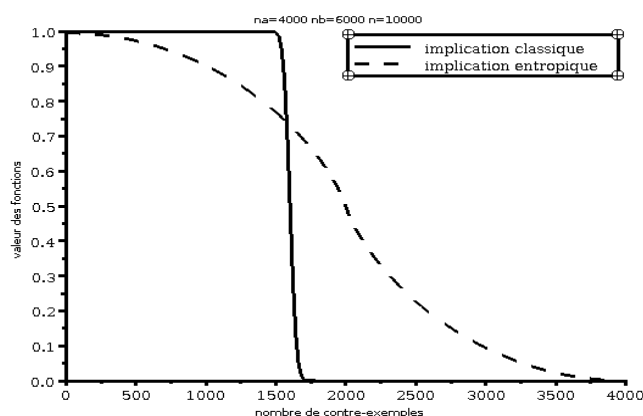


FIGURE 7 –

### 3 Application

Nous utilisons un fichier de données relatives à un questionnaire présenté à plus de 300 enseignants (Bodin et al, 1999) et portant sur les objectifs qu'ils assignent à l'enseignement des mathématiques à des élèves de 16 à 18 ans. L'enseignant doit choisir parmi 15 objectifs les 6 qu'il considère majeurs dans son rôle didactique, il doit les hiérarchiser (de 1 à 6) puis choisir un degré d'accord avec 10 opinions communément émises par ses collègues. Une analyse implicative des réponses conduit au graphe partiel suivant obtenu avec l'option dite « cône ».



E- développement de l'imagination et la créativité

I - développement de l'esprit critique

(OP8) - .pouvoir donner un exemple ou un contre-exemple personnels à l'affirmation : "si deux applications  $f$  et  $g$  sont strictement croissantes sur un intervalle, l'application produit  $fxg$  y est également croissante".

(OP7) - pouvoir reconnaître si un nombre entier écrit dans la base 10 est divisible par 4.

FIGURE 8 – *graphe implicatif*

Analysons ce chemin implicatif : les deux objectifs E et I, dénotant une volonté de l'enseignant d'associer l'élève à des processus où il dispose de la liberté de créer et critiquer impliquent la capacité à fabriquer un exemple (créer) et à critiquer une affirmation (contre-exemple). L'analyse en terme de causalité nous permet d'avancer l'hypothèse très vraisemblable suivante : les choix didactiques E et I de l'enseignant auraient une influence (causale) sur la compétence à résoudre des exercices ouverts où

créativité et esprit critique peuvent s'exprimer. Si ces objectifs sont atteints alors les items OP8 et OP7 ne doivent pas poser de problème à l'élève. On peut aussi exprimer cette dépendance, peut-être cette subordination, de la façon suivante : des exercices ouverts se résolvent d'autant plus aisément par un élève que l'on a développé chez lui des qualités de créativité et d'esprit critique.

## 4 Conclusion

Nous avons défini un nouvel indice d'implication. Il semble répondre aux attentes des chercheurs qui, utilisant l'intensité d'implication classique, se trouvent embarrassés par sa trop rapide convergence vers 1. Mais ces deux intensités d'implication ne sont pas, bien entendu, étrangères l'une à l'autre comme le montre l'étude présentée en annexe : *l'implication classique révèle la qualité de l'étonnement statistique du faible nombre de contre-exemples à la règle ; l'implication entropique décèle le déséquilibre cardinal autour du nombre de contre-exemples*. Cependant et en plus, cette dernière prend en compte la contraposée de la règle et, par conséquent, est un très bon indicateur d'une hypothèse de causalité. L'examen de quelques simulations où l'on a fait varier certains paramètres montre une bonne adéquation du résultat obtenu avec les exigences épistémologiques énoncées. Des applications attendues devraient confirmer cette tendance.

## Références

- [1] Blanchard J., Kuntz P., Guillet F. & Gras R. (2002) : Améliorer la mesure de l'étonnement statistique des règles à l'aide de la contraposée. *Rapport au C.N.R.S. du STIC GafoQualité*
- [2] Blanchard J., Guillet F., Briand H. & Gras R. (2005) Ipee : Indice probabiliste d'écart à l'équilibre pour l'évaluation de la qualité des règles, *Extraction et Gestion des Connaissances : état et perspectives, RNTI-E-5, Cépaduès*.391-395.
- [3] Bodin, A. (1997), Modèles sous-jacents à l'analyse implicative et outils complémentaires. *Prépublication IRMAR. n°97-32, 1-24*
- [4] Bodin, A. & Gras, R.(1999) Analyse du préquestionnaire enseignants avant EVAPM-Terminales, *Bulletin n°425 de l'Association des Professeurs de Mathématiques de l'Enseignement Public, 772-786, Paris*
- [5] Couturier, R. & Ag Almouloud, S. (2009) CHIC : historique et fonctionnalités. In R. Gras, J.C. Régnier & F. Guillet (Eds). *Analyse Statistique Implicative. Une méthode d'analyse de données pour la recherche de causalités* Toulouse : Ed. Cépaduès. p. 279-291.
- [6] Gras, R., Kuntz, P., Couturier, R. & Guillet F [2001]: Une version entropique de l'intensité d'implication pour les corpus volumineux, *Proceedings des Journées Extraction et Gestion des Connaissances EGC.'2001 de Nantes (18-19 janvier 2001)*, Hermès, p 69-80, ISBN 2-7462-0216-6

- [7] Gras, R., Kuntz, P. & Briand H. (2002) : Les fondements de l'analyse statistique implicative et quelques prolongements pour la fouille des données, *Mathématiques et Sciences Humaines*, n° 154-155, p 9-30, ISSN 0987 6936
- [8] Gras, R., & Couturier, R. (2010). Spécificités de l'Analyse Statistique Implicative (A.S.I.) par rapport à d'autres mesures de qualité de règles d'association. In J.C. Régnier, R. Gras, F. Spagnolo & B. Di Paola (Eds) *Quaderni di Ricerca in Didattica - GRIM* (ISSN on-line 1592-4424), Palerme : Université de Palerme, p.19-57
- [9] Shannon, C.E. & Weaver W. (1949) *The mathematical theory of communication*. Univ. of Illinois Press..

## Ouvrages

- [1] *Statistical Implicative Analysis*, R. Gras, E. Suzuki, F. Guillet and F. Spagnolo (Eds) Berlin Heidelberg : Springer-Verlag, ISBN 978-3-540-78982-6, 2008.
- [2] *Analyse Statistique Implicative ; une méthode d'analyse de données pour la recherche de causalités*, sous la direction de R. Gras, réd. invités : R. Gras, J.-C. Régnier & F. Guillet, (2009), RNTI E-16, Cépaduès Editions Toulouse. ISSN : 1764.1667, ISBN : 978.2.85428.897.1.
- [3] Régnier, J.C., Gras, R., Spagnolo, F. & Di Paola B. (Eds) (2011) *Analyse Statistique Implicative, Objet de recherche et de formation en analyse des données, outil pour la recherche multidisciplinaire, Prolongement des débats* [http://math.unipa.it/~grim/QRDM\\_20\\_Suppl\\_1.htm](http://math.unipa.it/~grim/QRDM_20_Suppl_1.htm) *QRDM - Quaderni di Ricerca in Didattica - GRIM* (ISSN on-line 1592-4424) Université de Palerme,

## Annexe

### Relation entre l'intensité d'implication classique et l'entropie conditionnelle

Nous montrons ci-dessous qu'il existe une relation étroite entre les lois des entropies conditionnelles  $H(b/a)$ ,  $H(\text{non } a/\text{non } b)$  et  $\varphi(a,b)$  ce qui prouve que ces notions ne sont pas étrangères tout en mesurant de façon différente des caractères de qualité des règles directe et contraposée.

Soit  $X$  la v.a. représentant le nombre aléatoire de contre-exemples  $N_{a \wedge \bar{b}}$ . Alors la probabilité pour que l'entropie conditionnelle de  $b$  sachant  $a$  soit inférieure à  $\alpha$  s'écrit :

$$\forall \alpha \in [0,1], \Pr[H(b/a) < \alpha] = \Pr \left\{ \frac{1}{\ln 2} \left[ -\frac{n_a - X}{n_a} \ln \frac{n_a - X}{n_a} - \frac{X}{n_a} \ln \frac{X}{n_a} \right] < \alpha \right\} =$$

$$\Pr \left\{ -[(n_a - X) \ln \frac{n_a - X}{n_a} + X \ln \frac{X}{n_a}] < \alpha \cdot n_a \ln 2 \right\}$$

Posons  $\frac{X}{n_a} = X' \in [0, \frac{1}{2}]$  Alors :

$$\Pr[H(b/a) < \alpha] = \Pr \{ -[(1-X') \ln(1-X') + X' \ln X'] < \alpha \ln 2 \} =$$

$$\Pr [H(b/a) < \alpha] = \Pr \{ -\ln X'^{X'} (1-X')^{1-X'} < \alpha \ln 2 \} = \Pr [ \ln X'^{X'} (1-X')^{1-X'} > -\ln 2^\alpha ] =$$

$$\Pr [ X'^{X'} (1-X')^{1-X'} > \frac{1}{2^\alpha} ]$$

Or la fonction  $f: x$  associe  $f(x) = x^x(1-x)^{1-x}$  est inversible car bijective sur l'intervalle considéré pour  $X'$  (propriété de l'entropie). D'où :

$$\Pr[H(b/a) < \alpha] = \Pr [X' < f^{-1}(\frac{1}{2^\alpha})] = \Pr[X < n_a f^{-1}(\frac{1}{2^\alpha})]$$

Choisissant  $\alpha$  tels que  $r = n_{a \wedge \bar{b}} = n_a f^{-1}(\frac{1}{2^\alpha})$ , on a donc, par exemple lorsque  $X$  suit

la loi de Poisson de paramètre  $\lambda = \frac{n_a n_b}{n}$  :

$$\Pr[H(b/a) < \alpha] = \sum_{X=0}^r \frac{e^{-\lambda} \lambda^X}{X!}$$

$\Pr[H(b/a) < \alpha] = \Pr [X < n_{a \wedge \bar{b}}] = 1 - \varphi_\alpha(a,b)$  où  $\varphi_\alpha(a,b)$  est l'intensité d'implication classique pour une observation de contre-exemples égale à  $n_{a \wedge \bar{b}}$

Autrement dit, il existe une relation étroite entre entropie conditionnelle et intensité d'implication.

### Remarque

Si nous nous intéressons à l'entropie conditionnelle  $H(\bar{a}/\bar{b})$  au lieu de  $H(b/a)$ , nous montrerions, par analogie des calculs, que :

$$\Pr [H(\bar{a}/\bar{b}) < \beta] = \Pr [X < n_{\bar{a}} f^{-1}(\frac{1}{2^\beta})]$$