

La ressemblance mathématisée

JEAN-PAUL DELAHAYE

Être intelligent c'est voir des ressemblances, mais qu'est-ce qu'une ressemblance ?

Que la tête de Louis Philippe ressemble à une poire, c'est ce que les caricatures ont montré. Que la carte de l'Italie ressemble à une botte, cela saute aux yeux. Que le bruit de la grosse caisse ressemble au tonnerre, le compositeur le sait et l'utilise. Que l'histoire de *La bicyclette bleue* ressemble à celle d'*Autant en emporte le vent*, c'est ce que les juges ont dû évaluer. La reconnaissance de la ressemblance est essentielle à notre survie : sans une notion instinctive de « situation équivalente », sans une maîtrise des « visages semblables », des « animaux de la même espèce », etc., nous ne pourrions nous orienter, nous reconnaître, identifier les aliments comestibles, les bêtes dangereuses.

La difficulté est que la même eau ne coule jamais deux fois sous le pont : deux situations ne sont jamais identiques en tout point, jamais deux animaux ne sont exactement pareils, jamais le visage de nos connaissances ne présente le même aspect à deux instants différents, etc.

Nous savons pourtant juger que deux situations, deux objets, deux images sont analogues, et cette capacité est une composante essentielle de l'intelligence. L'un des problèmes majeurs de l'Intelligence

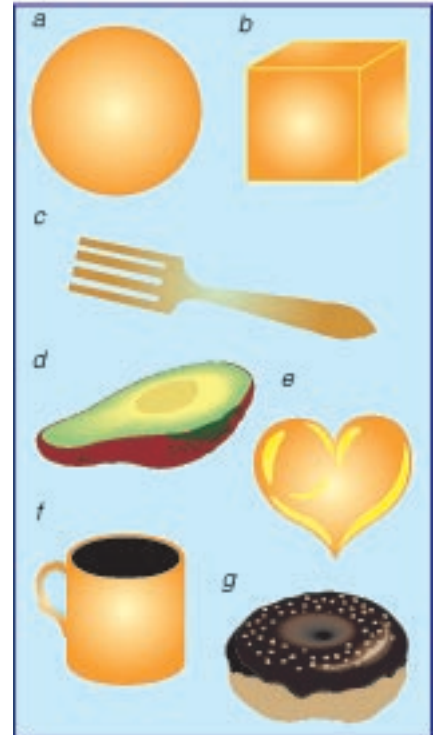
artificielle est de confectionner des programmes qui, au moins partiellement, reconnaissent des objets approchants. Le problème est difficile (et tout progrès en Intelligence artificielle est difficile), mais faut-il pour cela renoncer à le résoudre et penser qu'il est définitivement hors de portée ?

Les mathématiques qui prétendent fournir des outils pour comprendre le monde, le modéliser, le maîtriser proposent-elles de bonnes théories de la ressemblance ?

Nous verrons que oui en examinant plusieurs de ces théories. La dernière d'entre elles a été créée récemment par un groupe de physiciens, de mathématiciens et d'informaticiens : elle est fascinante par les liens qu'elle établit avec la thermodynamique et parce qu'on peut la considérer comme la théorie ultime de la ressemblance.

LES SEMBLABLES GÉOMÉTRIQUE ET TOPOLOGIQUE

L'idée mathématique la plus ancienne est celle de semblable géométrique : deux objets sont semblables si, en leur faisant effectuer des translations, des rotations, des symétries et des homo-



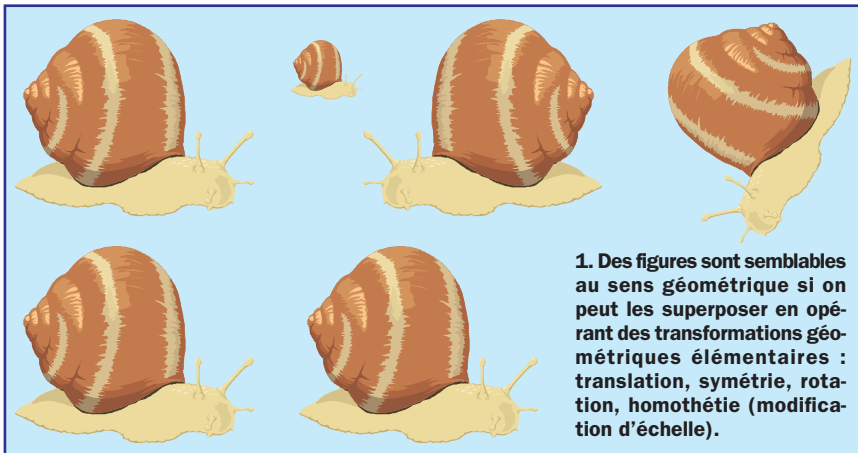
2. La surface d'un ballon peut se déformer continûment et prendre la forme d'un cube, d'un demi-avocat, d'un cœur, ou même d'une fourchette. Ces objets sont topologiquement semblables et pourtant, il est inacceptable de les considérer « semblables ». La tasse et le gâteau torique sont topologiquement équivalents entre eux mais pas à la sphère.

thétiques (agrandissements et rapetissements bien proportionnés), on peut les superposer.

Cette notion de figures géométriques semblables est utile pour faire des plans de voitures ou de maisons, mais elle est trop limitée, car trop rigide. Ce n'est pas elle qui pourra, par exemple, nous indiquer si deux photographies différentes représentent le même visage. Les mathématiciens, conscients de cette limite, ont proposé une notion bien plus souple de ressemblance : l'homéomorphie.

Deux surfaces sont homéomorphes si, en imaginant qu'elles sont fabriquées avec un caoutchouc parfait, on peut déformer l'une en l'autre sans faire de déchirure. Cette fois, la notion de ressemblance obtenue est trop molle ! La surface d'un cube est topologiquement semblable à celle d'un ballon ou même d'une fourchette. La *ressemblance topologique* conduit d'ailleurs à des situations paradoxales : un système d'anneaux enlacés est transformable progressivement en le même système d'anneaux libérés. Bien d'autres notions de topologie généralisent et étendent celle d'homéomorphisme, mais chacune ne saisit qu'une part bien mince de l'idée d'analogie.

Une raison en est sans doute que,



1. Des figures sont semblables au sens géométrique si on peut les superposer en opérant des transformations géométriques élémentaires : translation, symétrie, rotation, homothétie (modification d'échelle).

pour traiter de figures semblables, il faut renoncer à donner une réponse OUI ou NON et confectionner une *mesure de ressemblance*. Cette évaluation numérique va être donnée par ce qu'on appelle des distances : deux objets identiques seront à distance nulle ; deux objets très semblables seront à petite distance l'un de l'autre ; la distance entre deux objets très différents sera grande.

LA DISTANCE DE HAUSDORFF

Une première distance satisfaisante dans un assez grand nombre de cas est la distance de Hausdorff, nommée ainsi pour honorer Félix Hausdorff, mathématicien allemand mort en 1942. En 1914, Hausdorff a formulé les axiomes généraux de la topologie.

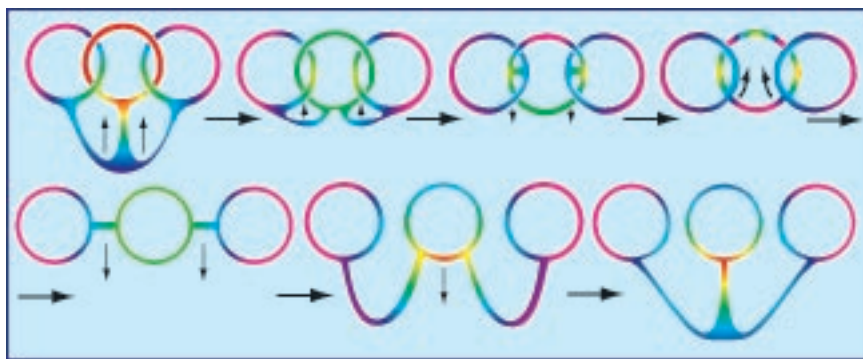
La distance de Hausdorff mesure la distance entre objets dessinés en noir et blanc sur une feuille de papier (des généralisations aux espaces à trois dimensions ou plus, ou prenant en compte les couleurs sont possibles). Pour définir cette distance entre figures, on utilise la notion usuelle de distance entre deux points.

Par définition, deux objets A et B dessinés sur une feuille de papier (c'est-à-dire deux ensembles de points A et B) sont à une distance de Hausdorff l'un de l'autre de moins de r unités, si chaque point de A est à moins de r unités d'au moins un point de B , et si, réciproquement, chaque point de B est distant de moins de r unités d'au moins un point de A .

Un procédé simple permet de visualiser ce qu'est la distance de Hausdorff entre deux ensembles A et B : on remplace chaque point de l'ensemble A par une petite tache ronde dont on fait grossir le rayon r jusqu'à ce que B soit recouvert par ce « A gonflé». On imagine aussi le processus réciproque (le recouvrement de A par un « B gonflé»). La distance de Hausdorff entre A et B est le plus petit r permettant simultanément les deux recouvrements.

Cette distance saisit correctement de nombreux aspects de l'idée intuitive de ressemblance. On considérera que deux dessins ne différant que par l'épaisseur des traits, ou par le fait que l'un est dessiné en pointillé, sont très proches l'un de l'autre : dès que l'on fait un peu gonfler l'un, il recouvre l'autre.

De même, la version brouillée d'un dessin A à travers un miroir dépoli sera peu distante du dessin A : de même encore, si, dans un dessin, vous remplacez les zones noires par des zones finement hachurées ou par des zones emplies de points aléatoires, alors la distance de la figure initiale à sa transformée sera petite.



3. Les transformations continues des surfaces réalisées en caoutchouc idéal permettent de passer d'une forme enlacée à une autre déenlacée, et «apparemment» sans rapport. Le semblable topologique apparaît trop tolérant.

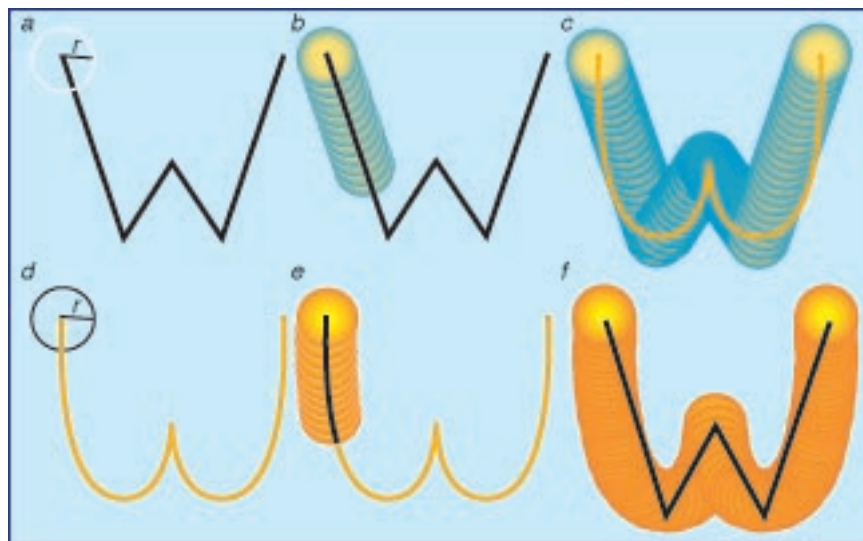
Un disque et une ellipse peu aplatie de rayons comparables seront proches pour la distance de Hausdorff. Un carré et un cercle, en revanche, seront plus éloignés. De façon générale, changer la texture ou déplacer légèrement les points d'un dessin en en préservant l'allure générale produit un dessin proche pour la distance de Hausdorff, alors que changer les formes crée des dessins éloignés.

Un cas intéressant est celui de l'approximation finie d'ensembles infinis. Un dessin comportant un nombre infini de points – par exemple, un disque plein – peut être approché de plus en plus finement par une série de figures ayant chacune un nombre fini de points. Ici l'infini, même non dénombrable, se laisse approcher facilement par du fini.

Cette distance est d'ailleurs implicitement utilisée par tous ceux qui dessinent des fractales : jamais personne ne dessine vraiment une fractale,

celle-ci étant, par définition, infiniment découpée, ce qu'aucun écran, aucun crayon, aucune imprimante ne réussit jamais à reproduire (notons aussi que personne ne dessine non plus de vrais cercles ou de vraies lignes). En revanche, les représentations qu'on donne des fractales sont bien des objets proches au sens de la distance de Hausdorff : comme Monsieur Jourdain faisait de la prose sans le savoir, nous utilisons la notion de ressemblance à la manière de Hausdorff sans jamais en avoir entendu parler.

Puisque deux dessins proches au sens intuitif sont mesurés proches au sens de la distance de Hausdorff, on comprend pourquoi cette distance est un outil essentiel de la morphologie mathématique, la discipline dont le but est la compréhension et l'analyse des images et des formes, et dont les applications concernent la vision artificielle, le traitement informatisé des images, la biologie et la géologie.



4. La distance de Hausdorff est le rayon r le plus petit tel qu'en remplaçant chaque point de A par un disque de rayon r on recouvre B , et inversement. C'est une vraie distance : $d(A,A) = 0$, $d(A,B) = d(B,A)$, $d(A,C) \leq d(A,B) + d(B,C)$ lorsque A , B et C sont des ensembles non vides et fermés (X est fermé si toute suite convergente de points de X a sa limite dans X).

INSUFFISANCE DE LA DISTANCE DE HAUSDORFF

Pourtant la distance de Hausdorff ne reconnaît pas certaines proximités entre formes que l'esprit humain identifie instinctivement. Ainsi la distance de Hausdorff entre un dessin et sa version en négatif (ou sa version où seul le contour des formes est gardé) est grande, alors que bien sûr nous, nous voyons très rapidement que le même objet est représenté. Pire, ajouter un unique point situé à l'extérieur d'un dessin donne un dessin très éloigné pour la distance de Hausdorff, alors que les deux objets nous apparaissent presque identiques. De même encore, la distance entre un objet et le même objet rapetissé d'un facteur constant est grande alors que nous percevons presque instantanément leur structure commune.

De même encore, la distance de Hausdorff ne permet pas de reconnaître la similitude entre des images d'un objet tridimensionnel (un visage, par exemple) vu sous divers angles ou des éclairages différents, alors que notre cerveau perçoit très vite l'analogie des formes.

Un défaut mineur, mais à signaler aussi de la distance de Hausdorff est qu'elle ne donne rien de bon avec des objets infinis en taille. Deux droites passant par un même point, ou bien sont confondues, et alors distantes de 0 (ce qui est satisfaisant), ou bien sont distantes de $+\infty$ même si leur angle est très faible (ce qui n'est pas acceptable).

Les déformations «caoutchouteuses» ne sont pas bien mesurées non plus : les montres molles de Dali ne sont pas proches des vraies montres pour la distance de Hausdorff, alors que, pour nous, elles le sont.

Remarquons qu'une combinaison de la notion de distance de Hausdorff avec la notion de similitude géométrique arrange certains de ces défauts : pour mesurer la distance entre A et B avec cette distance (appelée Hausdorff-bis), on envisage toutes les figures obtenues par application d'une similitude géométrique à B et l'on regarde celle qui, au sens de la distance de Hausdorff, est la plus proche de A ; c'est elle alors qui détermine la distance entre A et B pour Hausdorff-bis. Cette variante corrige les défauts concernant les déformations d'échelle et les déplacements, mais ne change rien aux insuffisances de la distance de Hausdorff concernant le passage au négatif, l'ajout d'un point à l'extérieur, les photos d'un même visage ou les déformations caoutchouteuses.

La dernière notion de ressemblance correspond à un progrès important sur le plan théorique, et en même temps

ouvre la porte à certaines applications pratiques. Cette notion, que nous appellerons distance informationnelle, est le résultat du travail commun de cinq chercheurs provenant de disciplines différentes : les physiciens Wojciech Zurek, du Santa Fe Institute, et Charles Bennett, du Centre de recherches IBM à New York ; le mathématicien Peter Gacs, de l'Université de Boston, aux États-Unis, et les informaticiens Ming Li, de l'Université de Waterloo, au Canada, et Paul Vitanyi, de l'Université d'Amsterdam.

LA DISTANCE INFORMATIONNELLE

La distance informationnelle présuppose que tous les objets que l'on considère sont des ensembles finis de points pris dans un ensemble discret. Il faut par exemple imaginer qu'il s'agit des pixels (noirs ou blancs) d'une image. Pour la distance de Hausdorff, une telle hypothèse de finitude n'était pas nécessaire : pour affiner un concept, on est parfois obligé de renoncer à l'infini, ce qui n'est peut-être pas étonnant : selon certains philosophes des mathématiques, l'infini n'est qu'une illusion. Des généralisations de la distance informationnelle à plus de deux dimensions et prenant en compte les couleurs sont possibles, mais nous ne les examinons pas ici.

L'idée fondamentale est que, si deux objets sont semblables, on passe facilement de la description de l'un à la description de l'autre, et réciproquement. En revanche, plus le passage de l'une à l'autre est long à détailler, plus les objets doivent être considérés comme différents, c'est-à-dire éloignés.

La difficulté du passage d'une description de l'objet A à la description de l'objet B est mesurée par la longueur du plus court programme qui transforme la donnée des points de A (supposés fournis pas énumération) à la donnée des points de B . Et donc la distance informationnelle est la somme des longueurs du plus court programme permettant de transformer A en B , et du plus court programme permettant de transformer B en A . On montre que cette distance varie peu quand on change le langage de programmation utilisé pour mesurer la taille des programmes.

L'idée d'utiliser une notion de plus court programme provient de la théorie algorithmique de l'information (appelée aussi théorie de la complexité de Kolmogorov), où l'on mesure la complexité d'un objet A par la taille du plus court programme qui produit A . La théorie de la distance informationnelle utilise d'ailleurs

de nombreux résultats de la théorie algorithmique de l'information.

Pour illustrer la notion de distance informationnelle, considérons le dessin de la tour Eiffel et le même dessin en négatif (nous avons dit que la distance entre ces deux objets pour la distance de Hausdorff est grande). Ces deux objets sont-ils proches pour la distance informationnelle ?

Oui, car on écrit facilement un programme de «passage au négatif» et que ce programme est court : le plus court programme transformant le premier dessin en le second (qui est le même ici que celui transformant le second en le premier) sera encore plus court (notons qu'il n'est pas indispensable de connaître précisément le plus court programme), et donc les deux dessins seront proches au sens de la distance informationnelle.

Ajouter un point à l'objet A donne un objet B que notre cerveau considère proche. En accord avec cette perception (et contrairement encore à la distance de Hausdorff, qui se trompait dans un tel cas), la distance informationnelle trouve que A et B sont très proches, car le programme le plus court ajoutant un point, et le programme le plus court enlevant un point (pour passer de B à A) sont tous les deux petits.

Plus remarquable et plus intéressant est l'exemple suivant : si l'on considère une image de la tour Eiffel prise selon un autre angle de vue, alors là encore on pourra par l'utilisation d'un programme assez court (qui contiendra des informations sur la position tridimensionnelle de la tour Eiffel) passer d'une image à l'autre. La parenté entre deux images provenant d'un même objet pris sous différents angles, ou même sous différents éclairages, est correctement détectée par la distance informationnelle.

Si l'on considère deux textes, même très longs, le second étant le résultat d'une remise en page du premier (changement des polices de caractères, des sauts des lignes, etc.), ils seront proches pour la distance informationnelle, car le passage de l'un à l'autre est décrit – et donc programmé –, de manière concise.

De nombreux autres exemples prouvent que cette mesure de distance fondée sur la taille des programmes satisfait l'essentiel de nos attentes, et que la grande majorité des défauts de la distance de Hausdorff disparaissent avec la distance informationnelle. L'utilisation pratique de cette distance semble cependant poser de graves problèmes.

Peut-on mesurer facilement la distance informationnelle entre objets ? Malheureusement, non. On majore sans mal

la distance informationnelle entre objets : à chaque fois qu'on trouve des programmes qui transforment A en B et B en A , on en déduit une telle majoration. En revanche, trouver des programmes dont on puisse prouver qu'ils sont les plus courts possibles est exceptionnel, et donc mesurer la distance informationnelle exactement est très rare.

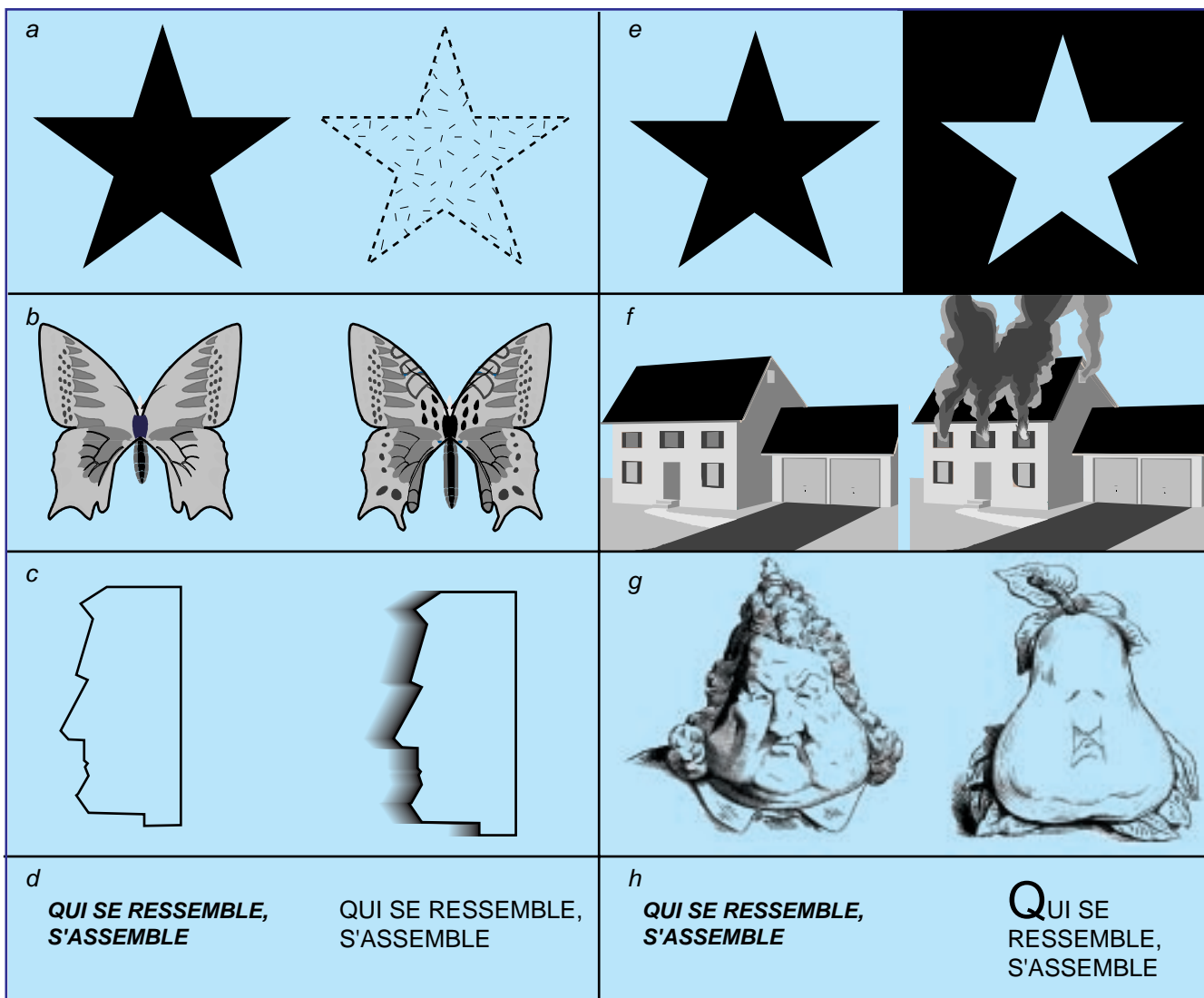
DISTANCE INFORMATIONNELLE ET INDÉCIDABILITÉ

Cette difficulté pratique est liée à l'indécidabilité logique dont le problème de l'arrêt des programmes proposé en 1936 par Alan Turing a été le premier

exemple. Le résultat de Turing nous dit que jamais on ne pourra faire un programme qui puisse, pour chaque programme qu'on lui soumettrait, calculer si oui ou non il s'arrête au bout d'un temps fini. Le résultat d'indécidabilité concernant la distance informationnelle nous dit que jamais nous ne pourrions concevoir un programme pour calculer, à chaque fois que nécessaire, la distance informationnelle de deux objets. Certains programmes pourront peut-être traiter quelques cas, mais aucun ne pourra les traiter tous, et donc aucun mécanisme général de raisonnement ou de calcul ne pourra servir d'instrument de mesure général pour la distance informationnelle.

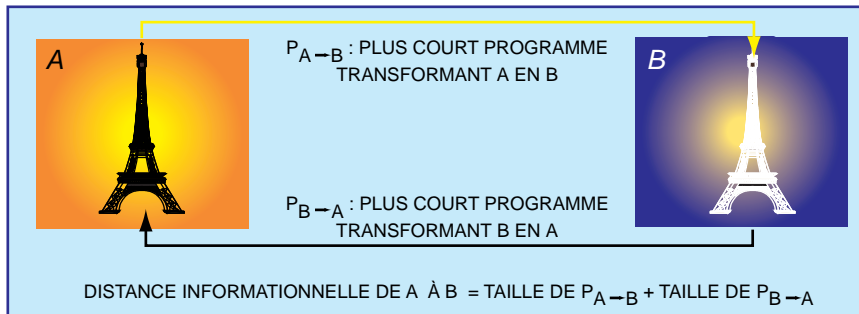
Cette impraticabilité reflète sans doute l'analyse intelligente nécessaire à la formulation des jugements de ressemblance : certains nous sont immédiats, d'autres ne sont perçus qu'après de longues analyses, d'autres encore ne le sont que par certains rares esprits. Est-il alors étonnant que la mesure du degré d'analogie entre des objets quelconques soit difficilement mécanisable et ne puisse pas, en définitive, l'être complètement.

L'intelligence est la cristallisation dans nos gènes, et ensuite dans notre cerveau, des très longs calculs réalisés par l'évolution et par le développement des cultures humaines. L'évolution biologique et le développement culturel sont



5. Pour nous, dans chaque paire les deux dessins se ressemblent clairement. La distance de Hausdorff saisit une partie de ces ressemblances (celles de la partie gauche (a), (b), (c), (d), mais passe à côté de celles de la partie droite (e), (f), (g), (h). La distance de Hausdorff est un moyen efficace mais pas infallible d'appréciation des ressemblances. En revanche, les cas où la distance de Hausdorff était impuissante à voir les ressemblances sont cette fois correctement traités par la distance informationnelle. Le programme pour transformer un dessin en son négatif (e) est très court, les

deux dessins sont donc proches pour la distance informationnelle. Le programme pour ajouter un peu de fumée et celui pour l'effacer (f) sont courts donc les deux dessins sont proches pour la distance informationnelle. Un programme permet de passer d'un dessin à l'autre (g), et ils sont relativement proches (et en tout cas bien plus proches l'un de l'autre que deux dessins quelconques). Quelques indications de remise en forme permettent de passer d'un texte imprimé à l'autre et réciproquement, qui sont donc proches l'un de l'autre pour la distance informationnelle (h).



6. La distance informationnelle entre deux objets A et B s'obtient en additionnant la taille des plus courts programmes permettant de passer de A à B et de B à A . Elle possède des propriétés de minimalité qui lui confère une importance particulière en théorie algorithmique de l'information. De plus les rapports qu'elle a avec la thermodynamique du calcul lui donnent un sens physique que les autres distances ne possèdent pas.

en effet de très longues séquences d'événements assimilables à des opérations élémentaires de calculs dont les produits sont le langage, les sciences et plus généralement nos facultés d'analyse. Ces très longs calculs, sans que nous sachions précisément comment, ont produit dans nos cerveaux une notion d'analogie très fine et très efficace, essentielle à notre survie. Cette cristallisation est incomplète (il n'y a pas de raisons sérieuses de croire que notre intelligence échappe aux limitations formulées par la logique), mais elle est bien meilleure que celle que nous réussissons à insérer dans nos programmes : aujourd'hui nous ne savons pas écrire des programmes qui égalent notre capacité à percevoir des analogies.

LA DISTANCE ENTRE SÉQUENCES GÉNÉTIQUES

L'indécidabilité rend-elle la distance informationnelle stérile? Heureusement non, car les majorations que l'on peut faire la rendent utilisable dans des cas sans doute assez nombreux.

Un exemple d'application de la distance informationnelle est la distance définie dans l'équipe de bio-informatique de Lille, composée de Max Dauchet, Éric Rivals, Jean-Stéphane Varré et moi, pour évaluer la distance entre séquences génétiques. Depuis longtemps les généticiens, et tout particulièrement ceux qui souhaitent faire de la reconstitution d'arbres phylogénétiques (ce sont les arbres qui indiquent les parentés entre espèces animales et végétales), utilisent des mesures de ressemblance entre séquences génétiques. L'idée à la base des distances utilisées jusqu'à présent était le décompte des mutations, délétions ou insertions de nucléotides (les lettres de l'alphabet génétique) : plus le nombre des changements ponctuels nécessaires pour passer d'une séquence à l'autre

est grand, plus les séquences sont considérées comme éloignées. Récemment nous avons proposé une nouvelle distance qui mesure l'éloignement entre séquences en considérant d'autres événements possibles comme les déplacements de morceaux de séquences ou leur duplication. Cette distance, que nous avons conçue comme une approximation calculable de la distance informationnelle, est plus précise que les distances utilisées classiquement (qu'elle généralise). On peut donc espérer qu'elle fournira des arbres phylogénétiques plus exacts.

Notons encore que les algorithmes modernes de compression de données pour les images de films sont fondés sur l'idée que deux images successives A et B d'une séquence sont proches et que, pour gagner de l'espace, il suffit de ne coder que leur différence : ce qui change dans A pour obtenir B . Cette quantité d'informations permettant de passer de A à B est en fait la distance informationnelle entre A et B (ou presque, puisque l'on ne s'occupe pas du passage de B à A).

INTERPRÉTATION THERMODYNAMIQUE

L'interprétation thermodynamique de la distance informationnelle est une autre confirmation de son intérêt. Les cinq chercheurs physiciens, mathématiciens et informaticiens ont en effet démontré un résultat remarquable concernant la distance informationnelle : elle est liée aux coûts thermodynamiques minimaux de la transformation de A en B et de B en A .

On sait depuis quelques années, à la suite des travaux de Rolf Landauer et de Charles Bennett, qu'il est possible, en théorie, de réaliser des calculs sans dépense d'énergie et n'entraînant donc aucun accroissement de l'entropie physique. Lors d'un calcul, les seules dépenses d'énergie inévitables proviennent de l'utilisation d'opérations irréver-

sibles comme l'effacement d'une mémoire dont on peut se passer si l'on accepte de garder des informations inutiles en fin de calcul.

Il est alors intéressant de considérer le flux minimum d'informations (ajout d'informations à A au départ, effacement d'informations à la fin du calcul une fois B obtenu) pour transformer A en B de manière réversible. Cette distance déduite de considérations thermodynamiques est équivalente à la distance informationnelle : elle en diffère d'un facteur deux au plus. D'autres propriétés de minimalité confèrent, parmi toutes les distances envisageables, un statut privilégié à la distance informationnelle.

La distance informationnelle entre A et B est donc une quantité ayant un sens en thermodynamique, ce qu'aucune autre distance ou notion mathématique de similitude ne possédait jusqu'à présent. L'adéquation entre ce que donne la nouvelle distance et l'intuition, le lien avec le monde physique, ainsi que d'autres résultats sur la distance informationnelle trouvés par le groupe des cinq chercheurs autour de Charles Bennett montrent clairement qu'une notion profonde a été identifiée.

Jean-Paul DELAHAYE est directeur adjoint du Laboratoire d'informatique fondamentale de Lille du CNRS.

e-mail : delahaye@lilfl.fr

C. BENNETT, et al. *Thermodynamics of Computation and Information Distance*. Proc. 25th ACM Symp. Theory of Computation, pp. 21-30, 1993.

C. BERGE, *Espaces topologiques et fonctions multivoques*, Éditions Dunod, Paris, 1966 (voir le chapitre 6).

J.-P. DELAHAYE, *Information, complexité et hasard*, éditions Hermès, Paris, 1994.

M. LI et P. M. B. VITANYI, *An Introduction to Kolmogorov Complexity and Its Applications*, Springer-Verlag, 1997 (voir le chapitre 8).

J. SERRA, *Image Analysis and Mathematical Morphology*, Academic Press, New York, 1982.

J.-S. VARRÉ, *Phylogénie et compression de données*, Publication du Laboratoire d'informatique fondamentale de Lille, URA CNRS 369, juillet 1996.

J.S. VARRÉ, E. RIVALS, M. DAUCHET et J.-P. DELAHAYE, *Les distances transformationnelles et applications à la phylogénie*. Journées Analyse des séquences génomiques, École Polytechnique (Palaiseau), 20 et 21 juin 1996.

Articles de P. Vitanyi sur Internet :
<http://www.cwi.nl/~paulv/publications.html>