

L'intelligence artificielle au service des données de santé

T. Perennec¹ & R. Bourgade², PA. Gourraud³, S. Supiot¹

1. *Institut de Cancérologie de l'Ouest, Service de Radiothérapie, Saint-Herblain*
2. *CHU de Nantes, Service d'Anatomie Pathologique, CHU de Nantes, Nantes*
3. *CHU de Nantes, INSERM, CIC 1413, Pôle Hospitalo-Universitaire 11 : Santé Publique, Clinique des données, Nantes, France*

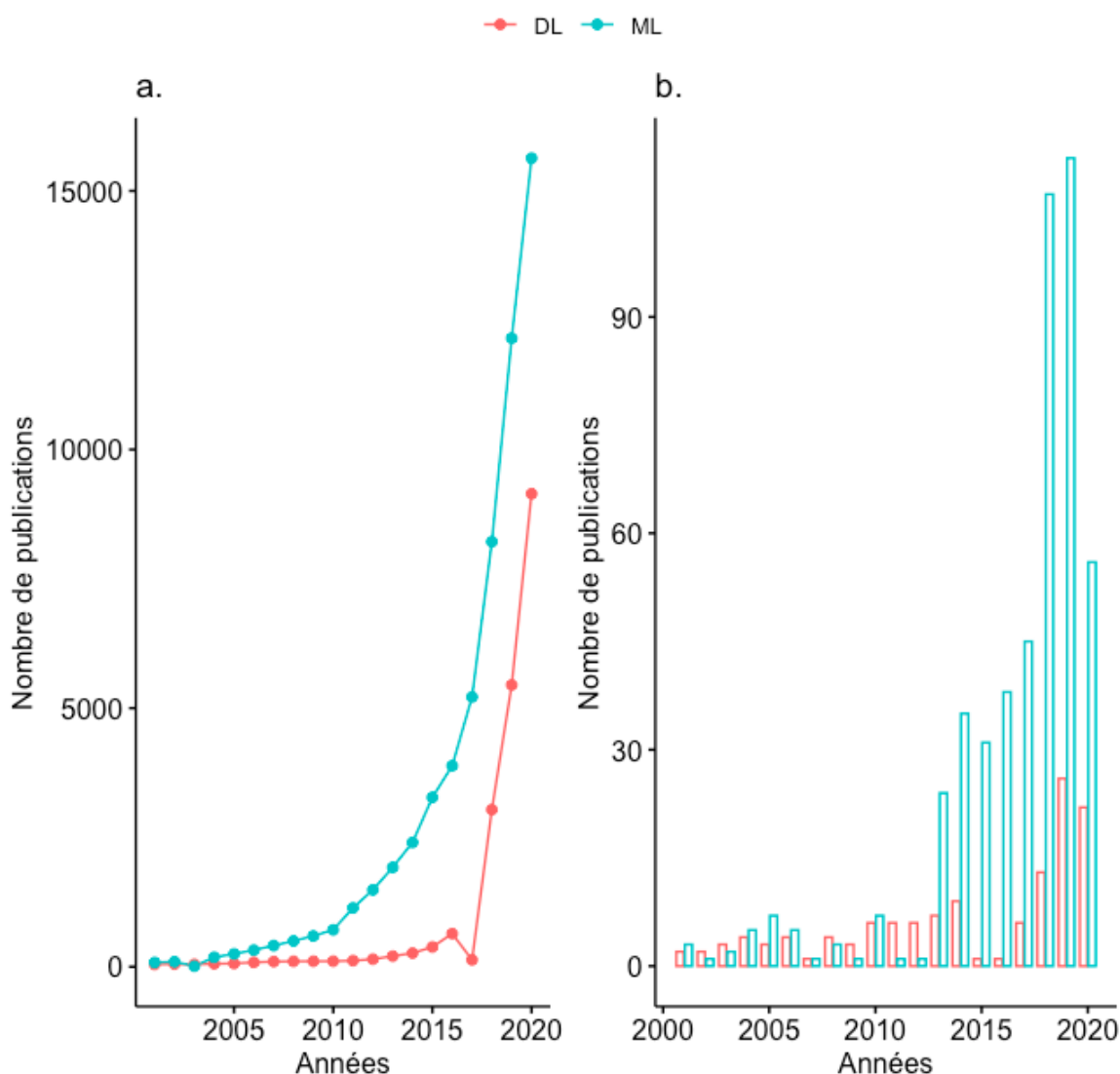
Introduction

Historiquement, le concept d'intelligence artificielle (IA) fait son apparition dès la fin des années 1950, avec un algorithme rudimentaire de classification binaire appelé "perceptron" (1). Désormais d'une très haute complexité, l'IA, discipline en constante évolution, façonne la société de demain. Elle se définit traditionnellement comme l'ensemble des théories et des techniques mises en œuvre dans le but de réaliser des machines capables de simuler ou dépasser l'intelligence humaine. Parmi ces techniques, on retrouve le *Machine Learning*, également appelé "Apprentissage automatique", développé à partir des années 1970. Ce dernier consiste à doter les ordinateurs de la capacité d'apprendre à réaliser une tâche à partir d'exemples, sans y avoir été explicitement programmé. Parmi les méthodes existantes, le *Deep Learning* repose sur l'utilisation de réseaux de neurones profonds. Il a été démocratisé au début des années 2010, et s'impose désormais comme la méthode de référence pour la résolution de problèmes complexes (2). De la conduite de véhicule autonome à la prédiction météorologique, en passant par la reconnaissance vocale, le deep learning fait dorénavant partie intégrante de notre quotidien. Cette omniprésence est due tant à l'augmentation exponentielle des données qu'aux récents progrès informatiques avec l'avènement des calculs parallèles, rendus possibles grâce aux processeurs multi-coeurs, puis plus récemment aux processeurs graphiques. La santé, secteur novateur en constante évolution, fait pourtant exception à la règle. Le sujet passionne et de très nombreux articles se sont attachés à démontrer l'intérêt du deep learning et du machine learning dans la médecine (Figure 1a). Assistance diagnostique, prédiction de réponse thérapeutique ou mise au point de nouveaux marqueurs théranostiques, les applications potentielles sont nombreuses et convergent vers une médecine personnalisée de précision, plus efficace et plus sûre (3). Pourtant, les applications concrètes en pratique quotidienne demeurent très rares et les essais cliniques publiés sur le sujet sont pratiquement inexistantes, représentant moins de 1% des articles publiés (Figure 1b). Ce paradoxe repose notamment sur le caractère hautement confidentiel des données de santé, pourtant indispensables à la mise au point de ces algorithmes. Ces données sensibles sont régies par un ensemble de textes tels que le *Règlement Général sur la Protection des Données (RGPD)* adopté par l'Union Européenne en 2016 et entré en vigueur en 2018. La

France s'est particulièrement illustrée sur la protection des données avec la création de la CNIL (Commission Nationale de l'Informatique et des Libertés) et de la loi "Informatique et Libertés" du 6 janvier 1978.

Cette réglementation, indispensable à la protection du patient, ne doit cependant pas se faire au détriment du progrès scientifique et technologique. Nous allons ainsi passer en revue certaines solutions visant à faciliter et encourager les recherches axées sur l'intelligence artificielle, dans le plus strict respect des patients et de leurs données.

Figure 1 : Nombre annuel de publications sur Pubmed concernant le Machine Learning (ML) et le Deep learning (DL), quel que soit le type d'article (a) ou ne concernant que les essais cliniques (b)



Etat des lieux des données de santé en France

Les données de santé sont définies par la CNIL comme des données à caractère personnel, relatives à la santé physique ou mentale, passée, présente ou future, d'une personne physique, qui révèlent des informations sur son état de santé. Cette définition élargie est le fruit du règlement

européen sur la protection des données personnelles (RGPD) entré en application le 25 mai 2018. Cette réglementation s'applique à toute organisation publique ou privée qui traite des données personnelles, dès lors qu'elle est établie sur le territoire de l'Union européenne ou que son activité cible directement des résidents européens. Le RGPD s'inscrit ainsi dans la continuité de la Loi française Informatique et Libertés de 1978, s'adaptant à l'évolution numérique et technologique de notre santé. En vue de garantir une utilisation optimale de ces données à des fins de recherche scientifique, de nombreux acteurs ont été mis en place. Parmi eux, le Health Data Hub, traduction de la Plateforme des Données de Santé (PDS) succède à l'Institut National des Données de Santé (INDS) le 30 novembre 2019. Ses missions consistent notamment en l'organisation et la mise à disposition des données, issues entre autres du Système Nationale des Données de Santé (SNDS), tout en garantissant le respect des droits exercés par les patients.

Le crowdsourcing, la solution collaborative

Le crowdsourcing repose sur la mise à disposition de données, afin qu'un nombre plus ou moins important de personnes puissent proposer une réponse à un problème. Ce mode de fonctionnement collaboratif n'est pas rare dans le milieu de l'intelligence artificielle (4) et plusieurs plateformes se sont spécialisées dans la mise en place de compétitions de data science telles que DrivenData, CrowdANALYTIX ou encore Kaggle. Cette dernière a été fondée en 2010 par Anthony Goldbloom et rachetée 4 ans plus tard par Alphabet, maison mère de Google. Elle met ses ressources à disposition de toute société ou groupe universitaire souhaitant résoudre un problème par le biais d'une compétition de machine learning. Les organisateurs doivent proposer une base de données labellisées et définir une méthode servant à l'évaluation des modèles, appelée "métrique". Les prix de ces compétitions varient de plusieurs milliers à 1,5 million de dollars. Cette somme avait été offerte par le gouvernement américain pour mettre au point un algorithme de reconnaissance d'objets menaçants dans les aéroports, détectés par caméra. Si les sujets traités sur ces plateformes sont variés, la santé et notamment l'oncologie y occupent une place importante. Par exemple, l'université de Radboud en partenariat avec l'institut Karolinska a organisé une compétition durant l'été 2020, visant à déterminer le grade ISUP de biopsies de prostates, à partir d'une dizaine de milliers de lames anatomopathologiques numérisées et anonymisées. La Société Française de Pathologie, quant à elle, a récemment organisé une compétition hébergée par DrivenData dont le but était la classification de lésions tumorales du col de l'utérus. Ces compétitions permettent l'émergence d'algorithmes performants visant à terme à optimiser les performances des praticiens.

Ces plateformes sont également des lieux d'échange et de partage. Ainsi, comme l'a démontré une étude récente, elles réunissent les meilleurs data scientists mondiaux, dont les travaux constituent l'état de l'art en machine learning (4). Il est à noter que certaines compétitions aux données sensibles sont réservées aux meilleurs participants de la plateforme. Enfin, elles

permettent et encouragent le partage de bases de données. Si cette opportunité paraît adaptée à la prédiction du cours d'actifs financiers ou à la reconnaissance de maladies affectant les plantes, elle le semble nettement moins concernant des données de patients potentiellement réidentifiables. Dès lors, il convient de trouver des solutions de protection à l'utilisation de ces données sensibles, afin d'exploiter l'intelligence artificielle à son plein potentiel, comme cela est possible dans d'autres disciplines.

Le leurre de la pseudo-anonymisation

La majorité des analyses statistiques réalisées dans les études cliniques porte sur des bases de données dites "pseudo-anonymisées" ou "pseudonymisées". Compromis entre données brutes et entièrement anonymisées, elles doivent permettre de répondre à une question scientifique en excluant toute donnée directement identifiante (nom, prénom etc...). La pseudonymisation est cependant un processus réversible, dont les données demeurent soumises au RGPD. Ainsi, contrairement aux données anonymisées, la réidentification de données pseudonymisées est toujours possible. C'est notamment ce qui a poussé le Conseil d'Etat à mandater la CNIL pour expertiser les moyens de pseudonymisation mis en œuvre lors de la crise sanitaire du COVID-19 à travers l'ordonnance du 19 juin 2020. En effet, certaines données indirectement identifiantes telles que les dates de naissance ou de décès, couplées à des caractéristiques rares peuvent trahir l'identité d'un patient. *L. Rocher et al* ont ainsi proposé un modèle capable d'identifier n'importe quel patient quelque soit la cohorte avec une aire sous la courbe ROC de 0,84 à 0,97. Ils ont également prouvé qu'ils pouvaient réidentifier jusqu'à 99,98% des patients américains dans n'importe quelle base de données, si au moins 15 de leurs attributs démographiques étaient connus (5).

En imagerie, le visage d'un patient peut être reconstitué à partir de coupes scannographiques, grâce aux algorithmes de reconstruction tridimensionnelle classiques.

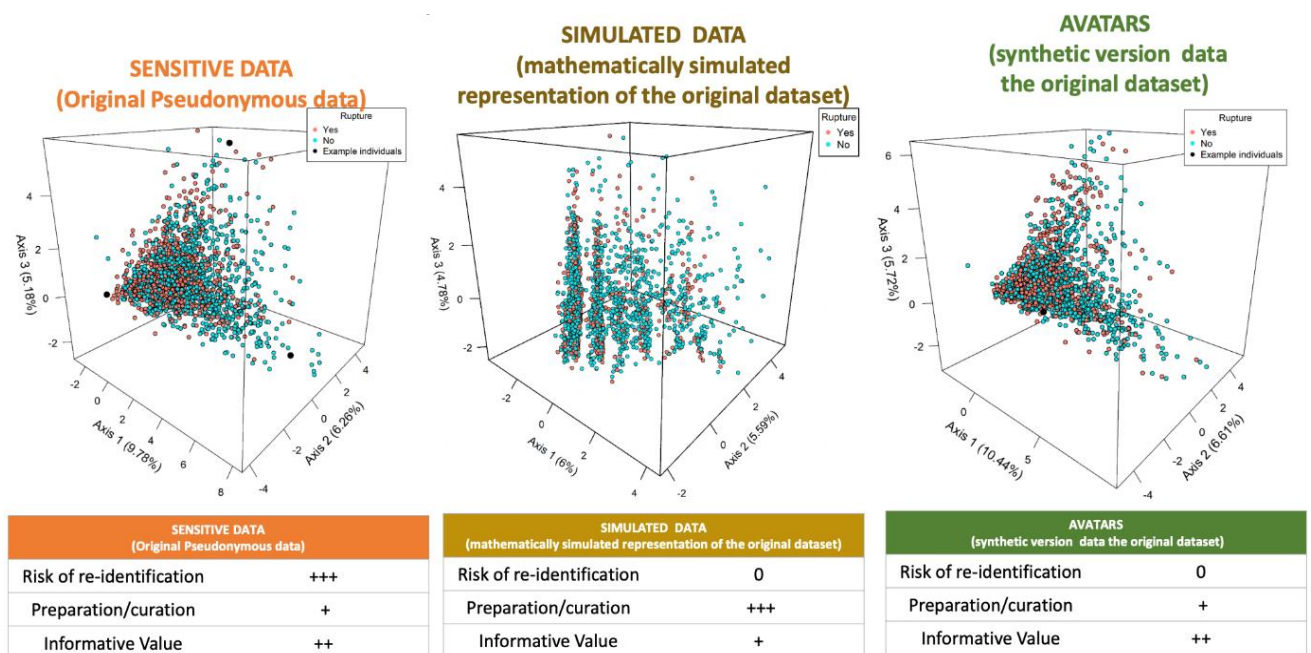
Face à ces éventualités, différentes techniques existent pour préserver la confidentialité des patients et deviendront vraisemblablement un prérequis obligatoire à l'utilisation des données de santé dans la décennie à venir.

"Avatarisation" des données

L'avatarisation est une méthode permettant d'anonymiser les données de façon plus sûre. Elle consiste à remplacer les données d'un patient par les données d'un avatar qui lui "ressemble", créé artificiellement au moyen d'un algorithme. En fonction du degré de confidentialité que l'on souhaite appliquer aux données, le logiciel renverra une base de données plus ou moins proche de la réalité, dont les analyses statistiques obtenues demeureront plus ou moins superposables. Il devient alors possible d'exploiter publiquement cette base de données, et ensuite de reproduire

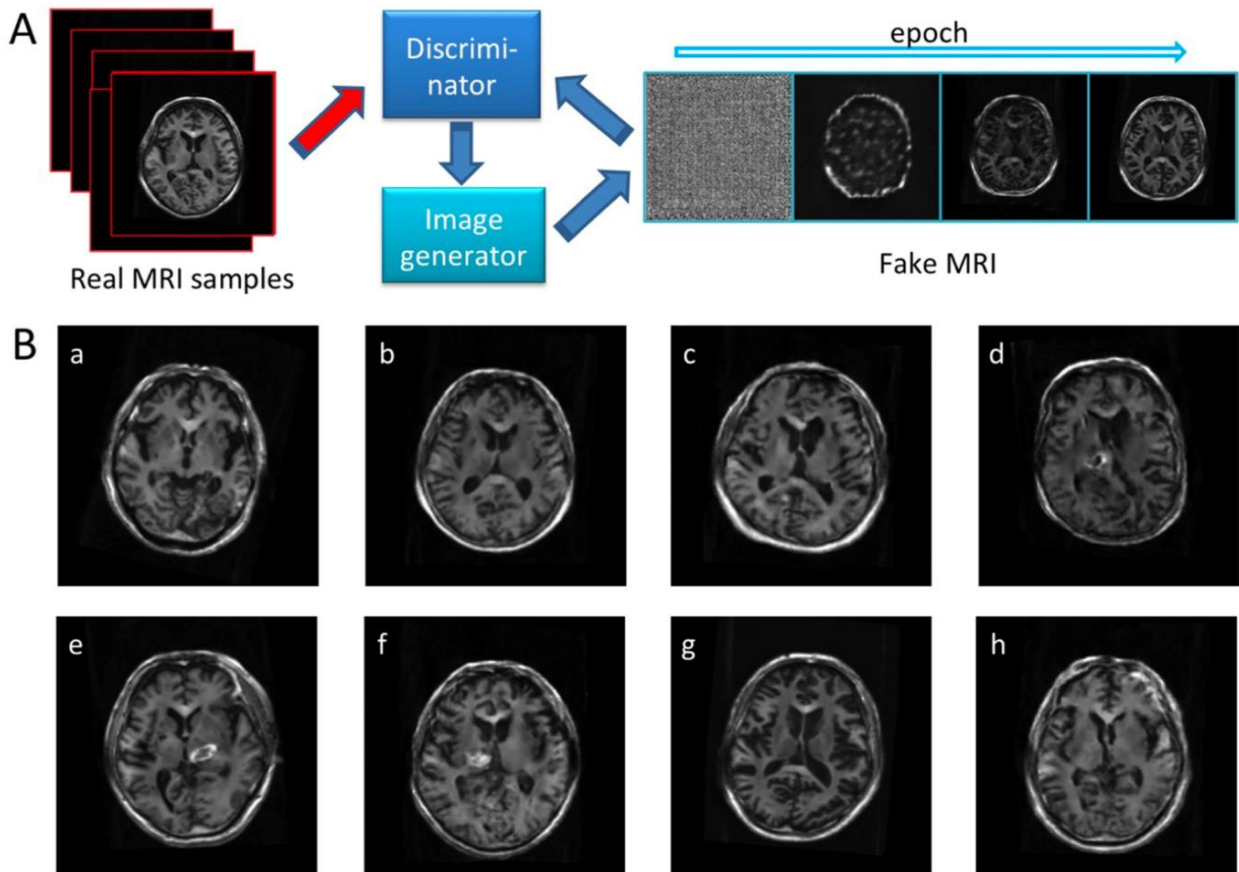
l'ensemble des analyses sur la base de données originale, en milieu interne sécurisé. Ce type de transformation synthétique a été utilisé dans l'étude ICAN afin de déterminer quels étaient les facteurs de risque de rupture d'anévrisme intracrânien. La représentation graphique des différentes sources de données (originale, simulée et avatarisée) est présentée dans la Figure 2 (6).

Figure 2 : Représentations graphiques de l'analyse factorielle de données mixtes (AFDM) de la population globale du projet ICAN comparant le jeu de données original pseudonymisé, une simulation mathématique de sa représentation et pour finir son avatarisation. En rouge les patients présentant un anévrisme intracrânien rompu et en bleu, non rompu. Chaque axe correspond à une composante principale de l'AFDM et la position de chaque patient en fonction de ces axes dépend de ses caractéristiques. (Figure adaptée de la figure 1 de Rousseau et al. (6))



Concernant les données plus complexes, comme l'imagerie radiologique ou les lames numérisées d'anatomie pathologique, il faut avoir recours à d'autres solutions. Les GAN (*Generative Adversarial Networks*), algorithmes d'apprentissage profond non supervisé, sont entraînés à générer des images de synthèse à partir d'images bien réelles, avec un fort degré de réalisme. Ce procédé a notamment été médiatisé par le deepfake (publication de discours fallacieux de personnalités politiques), dont la publication avait généré une inquiétude sur leur possible utilisation à des fins de désinformation, ou encore par la génération d'œuvres artistiques imitant le style littéraire de Shakespeare ou artistique de Van Gogh (7). Désormais utilisés en médecine, ces algorithmes permettent actuellement de générer des images supplémentaires pour pallier au manque de données disponibles (Figure 3) (8).

Figure 3 : A : Les réseaux antagonistes génératifs (GANs) sont des modèles génératifs mettant en compétition deux algorithmes. Le premier est chargé de la génération d'images artificielles, tandis que le second apprend continuellement à différencier celles qui sont réelles, de celles qui ne le sont pas. B : Des images d'IRM ont été générées grâce à des réseaux antagonistes génératifs convolutionnels profonds (DCGAN). Les images synthétiques sont les images b, c, d, f et h. (Figure adaptée de la figure 1 de Kazuhiro et al. (8))



Blockchain et propriété des données

Issue d'un processus multifactoriel impliquant de nombreux acteurs, la propriété de la donnée de santé ne peut être clairement définie. Profitant de l'engouement pour le dossier médical partagé, certaines entreprises telle la start-up française Embleema, considèrent qu'elles sont la propriété du patient. Dès lors, le partage des données de santé à un tiers est accepté, tant qu'il est consenti et rémunéré. Ces entreprises se présentent alors comme des tiers de confiance, chargés de mettre en relation le patient avec les investigateurs d'une étude. Ces derniers s'engagent à rémunérer le patient, à hauteur de l'intérêt et de la rareté de ses données partagées. Le risque de cette pratique est l'introduction d'un biais de sélection, similaire aux études rémunérées, par la probable plus forte participation de populations défavorisées. En plus de la rétribution des patients, Embleema propose la sécurisation du dossier médical personnel, par l'utilisation d'une blockchain.

Mise en lumière par le protocole Bitcoin, une crypto-monnaie issue des travaux d'un certain Satoshi Nakamoto en 2008, cette dernière pourrait jouer un rôle majeur dans la santé de demain. Il s'agit d'une technologie de stockage et de transmission d'informations entièrement décentralisée et sécurisée par cryptographie. La blockchain peut être représentée comme une base de données dans laquelle chaque information est inscrite de façon immuable, sans pouvoir être effacée ni même éditée. Cette propriété peut aussi rendre le partage de données périlleux : il est en effet ensuite impossible de les supprimer, tout partage étant irrévocable. Néanmoins, des informations telles que le droit d'accès au dossier médical ou le consentement d'inclusion dans une étude peuvent y être inscrites et garantir ainsi une parfaite traçabilité sécurisée du dossier médical, dont seul le patient détient la clef privée. Comme illustré dans la Figure 4, de nombreuses entreprises tendent à adopter cette nouvelle technologie.

	Blockchain company		
	Name	Country	Website
EMR data management	PokitDoc	USA	http://pokitdoc.com
	Gem	USA	http://enterprise.gem.co/health
	YouBase	USA	http://www.youbase.io
EHR data management	Medicalchain	USA	http://www.medicalchain.com
	HealthWizz	USA	http://www.healthwizz.com
	Curisium	USA	http://www.curisium.com
	Hearthly	Spain	http://hearthy.co
	Iryo	Slovenia	http://iryio.io
	Robomed	Russia	http://www.robomed.io
PHR data management	Medcredits	USA	https://medcredits.io
	MyClinic	UK	https://myclinic.com
Point-of-care genomics	Nebula Genomics	USA	http://www.nebula.org
	Genomes.io	USA	http://www.genomes.io
	TimiCoin	USA	http://www.timicoin.io
	Shivom	Switzerland	http://shivom.io
Oncology patients network	OncoPower	USA	http://oncopower.org
Pharma & drug development	Embleema	France	http://www.embleema.com
	BlockPharma	France	http://www.blockpharma.com
	Chronicled MediLedger	USA	http://www.mediledger.com

EMR: Electronic Medical Record, EHR: Electronic Health Record, PHR: Personal Health Record.

Table 1 : Listes d'entreprises utilisant une blockchain selon le secteur d'activité. (Adaptée de la table 1 de Dimitrov (9))

Conclusions

Prérequis indispensable à l'intelligence artificielle, l'utilisation des données de santé balance entre perspective de progrès scientifiques d'une part et respect accru de la confidentialité d'autre part. Les données constituent désormais une véritable réserve de valeur, attirant de nombreuses convoitises, parfois au détriment du patient. Il convient ainsi de définir un juste équilibre qui facilite une mise à disposition des données, tout en garantissant une sécurité optimale, sous peine de manquer le "virage IA" au profit de nations moins protectionnistes, dont la conception de la propriété des données diffère de la nôtre.

Références

1. Rosenblatt F. The perceptron: a perceiving and recognizing automaton. Project Para Report. No. 85-460-1. 1958;
2. LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature*. 2015; 521:436–44.
3. Ching T, Himmelstein DS, Beaulieu-Jones BK, Kalinin AA, Do BT, Way GP, et al. Opportunities and obstacles for deep learning in biology and medicine. *J Royal Soc Interface* 2018;15:20170387.doi:10.1098/rsif.2017.0387.
4. Tauchert C, Darmstadt T. Crowdsourcing Data Science: A qualitative analysis of organizations' Usage of Kaggle Competitions.
5. Rocher L, Hendrickx JM, de Montjoye Y-A. Estimating the success of re-identifications in incomplete datasets using generative models. *Nat Commun*. 2019 Dec;10(1):3069.
6. Rousseau O, Karakachoff M, Gaignard A, Bellanger L, Bijlenga P, Constant Dit Beaufils P, et al. Location of intracranial aneurysms is the main factor associated with rupture in the ICAN population. *J Neurol Neurosurg Psychiatry*. 2021;92:122–8.
7. Zhu J-Y, Park T, Isola P, Efros AA. Unpaired Image-to-Image Translation Using Cycle-Consistent Adversarial Networks. In: 2017 IEEE International Conference on Computer Vision (ICCV). Venice: IEEE; 2017. p. 2242–51.
8. Kazuhiro K, Werner RA, Toriumi F, Javadi MS, Pomper MG, Solnes LB, et al. Generative Adversarial Networks for the creation of realistic artificial brain magnetic resonance images. 2018;4:159–63.
9. Dimitrov DV. Blockchain Applications for Healthcare Data Management. *Healthc Inform Res*. 2019;25:51.