



GOVERNEMENT

*Liberté
Égalité
Fraternité*



Innovation - Écologie - Territoires
Commissariat Général au Développement Durable



**SANTÉ
ENVIRONNEMENT**

JUILLET 2024

ÉTAT DES CONNAISSANCES
L'IA : SYNTHÈSE
DES CONNAISSANCES ET PERSPECTIVES
POUR LA SANTÉ ENVIRONNEMENTALE
À SANTÉ PUBLIQUE FRANCE

Résumé

L'IA : synthèse des connaissances et perspectives pour la santé environnementale à Santé publique France

Si l'intérêt grandissant pour l'intelligence artificielle (IA) est manifeste, il semble parfois difficile d'identifier précisément le périmètre, les outils et les méthodes sous-tendues par la notion d'« intelligence artificielle ».

Ce travail de synthèse des connaissances centré sur les utilisations des méthodes d'IA pour la santé environnementale (SE), basé sur la littérature et l'interrogation d'experts du sujet, a permis de faire la lumière sur un nouvel univers de mots et expressions qu'il est nécessaire de définir tant les interprétations sont nombreuses.

L'IA est mobilisée dans le domaine de la santé environnementale à différentes étapes des études : de la collecte de données à l'analyse en passant par la prédiction. Les capacités de traitement des données massives par l'IA offrent des possibilités intéressantes pour faciliter l'exploitation des données d'observation de la terre, simuler des données d'exposition, faire ressortir des profils de territoires, croiser les données environnementales et sanitaires, faciliter la recherche documentaire et communiquer des résultats de façon dynamique et adaptée, etc.

Les opportunités offertes par l'IA dans le domaine de la santé environnementale sont immenses. Cependant, il est essentiel de considérer plusieurs enjeux avant de mobiliser cette technologie. Ces enjeux incluent des aspects juridiques, éthiques, écologiques et techniques. Chacun de ces aspects représente une source d'incertitude qu'il faut prendre en compte.

MOTS CLÉS : INTELLIGENCE ARTIFICIELLE, SANTÉ ENVIRONNEMENTALE, INNOVATION, OUTILS, MÉTHODE

Citation suggérée : Chaperon L, Deplanche M, Grignon P, Haroutunian L, Stempfelet M. L'IA : synthèse des connaissances et perspectives pour la santé environnementale à Santé publique France. Saint-Maurice : Santé publique France, 2024. 38 p. Disponible à partir de l'URL : www.santepubliquefrance.fr

ISSN : 2609-2174 - ISBN-NET : 979-10-289-0920-8 - RÉALISÉ PAR LA DIRECTION DE LA COMMUNICATION, SANTÉ PUBLIQUE FRANCE - DÉPÔT LÉGAL : JUILLET 2024

Abstract

AI: Synthesis of Knowledge and Perspectives for Environmental Health at Santé publique France

While the growing interest in Artificial Intelligence (AI) is clear, it sometimes seems difficult to identify precisely the scope, tools and methods underlying the notion of 'artificial intelligence'.

This review of knowledge on the uses of AI methods for environmental health (EH), based on the literature and interviews with experts on the subject, has shed light on a new universe of words and expressions that need to be defined, given the wide range of interpretations.

AI is used in the field of environmental health at various stages of studies : from data collection to analysis and prediction. AI's capacity to process massive data offers interesting possibilities for facilitating the use of earth observation data, simulating exposure data, highlighting territorial profiles, cross-referencing environmental and health data, facilitating documentary research and communicating results in a dynamic and appropriate way, etc.

The opportunities are huge, but legal, ethical, ecological and technical issues need to be taken into account before AI can be used to improve environmental health.

KEY WORDS: ARTIFICIAL INTELLIGENCE, ENVIRONMENTAL HEALTH, INNOVATION, METHOD, TOOL

Auteurs

Laura Chaperon

Direction santé environnement travail, Santé publique France

Margaux Deplanche

Direction santé environnement travail, Santé publique France

Paul Grignon

Ecolab, ministère de la transition écologique et de la cohésion des territoires (MTECT)/
Commissariat général au développement durable (CGDD)/Service de recherche et de
l'innovation (SRI)

Laetitia Haroutunian

Direction scientifique et international, Santé publique France)

Morgane Stempfelet

Direction santé environnement travail, Santé publique France)

Relecteurs

Johnny Platon

Direction appui, traitement et analyse des données, Santé publique France

Clémence Fillol

Direction santé environnement, Santé publique France

Remerciements

Nous remercions chaleureusement les personnes interrogées lors d'entretiens et qui ont accepté de partager leurs expériences :

Juliette Froppier (Ecolab), Marie Ramon-Daré (Ecolab), Emmanuel Bacri (*Health Data Hub* - HDH), Maria-René Palomo Arbizu (HDH), Matthieu Porte et Marie Gombert (Institut national de l'information géographique et forestière - IGN), Grégoire Etot (Office français de la biodiversité - OFB), Guillaume Boulanger (Santé publique France), Édouard Chatignoux (Santé publique France), Johnny Platon (Santé publique France), Matthieu Brient (OpenDataFrance), Julie Letrertre (*Copernicus European Centre for Medium-Range Weather Forecasts* - ECMWF), Claire Monteleoni (Institut national de recherche en sciences et technologies du numérique - Inria), Régis Chatellier et Alexis Léautier (Commission nationale de l'informatique et des libertés - Cnil), Adeline Martin (Institut de médecine et physiologie spatiales du Centre national d'études spatiales - Cnes Medes), Laure Malherbe et Jean-Yves Chatelier (Institut national de l'environnement industriel et des risques - Ineris), Marc Malenfer (Institut national de recherche et de sécurité - Inrs)

TABLE DES MATIÈRES

Résumé	2
Abstract	3
Auteurs/Relecteurs/Remerciements	4
Table des figures	6
Lexique – Acronymes	6
1. INTRODUCTION	8
2. DÉFINITIONS ET PANORAMA DES MÉTHODES D'IA	10
2.1 Définitions de l'IA	10
2.1.1 Générale	10
2.2.2 La GéolA	12
2.2 Cartographie des méthodes qui constituent l'IA	12
2.3 Les approches sous-jacentes à l'IA	14
2.4 Les facteurs facilitant la mise en œuvre de projets d'IA	14
2.5 Focus : <i>Machine Learning</i>	15
2.6 Focus : l'apprentissage profond	16
3. IA POUR LA SANTÉ ENVIRONNEMENT : LES MÉTHODES ADAPTÉES	17
3.1 Le potentiel de l'IA	17
3.2 Types de méthodes et cas d'usage pour mobiliser l'IA et les données en santé- environnement	17
3.2.1 Collecte et extraction de données	18
3.2.2 Prétraitement et consolidation des données	19
3.2.3 Analyse de données	21
3.2.4 Valorisation et communication des données et des analyses	23
4. LES MOYENS DONT A BESOIN LA DATA SCIENCE ?	24
5. DISCUSSION	26
5.1 Critique et perspectives	26
5.2 Données	26
5.3 Considérations éthiques et juridiques associées à l'IA	27
Protection des données personnelles et des libertés	27
Transparence et explicabilité des algorithmes et des systèmes d'IA	27
Les biais et discriminations algorithmiques	28
Bouleversement du travail, santé mentale et exploitation des personnes pour la production de données	28
5.4 Sobriété numérique : l'IA frugale	29
Protéger les données sensibles, protéger l'environnement ?	29
6. PERSPECTIVES	31
7. CONCLUSION	33
8. RÉFÉRENCES	34
ANNEXES	36
Annexe 1 : recherche bibliographique	36
Annexe 2 : liste des organismes interrogés	37
Annexe 3 : liste non exhaustive de conférences en ligne	38

Table des figures

Figure 1. Chronologie de l'IA (alimentée par des sources diverses)	11
Figure 2. Cartographie schématique des méthodes qui constituent l'IA.....	13
Figure 3. Structure nécessaire à la mise en place d'une démarche scientifique innovante ...	24

Lexique - Acronymes

3IA	Instituts interdisciplinaires d'intelligence artificielle
ACP	Analyse en composante principale
Afnor	Association française de normalisation
AIS	Agence de l'innovation en Santé
Amama	Alphabet, Meta, Amazon, Microsoft, Apple
Basias	Base de données des anciens sites industriels et activités de services. (Base de données française diffusée publiquement depuis 1999)
BATX	Baidu, Alibaba, Tencent, Xiaomi
Chatbot	Agent conversationnel
Cnil	Commission nationale de l'informatique et des libertés
CNN	<i>Convolutional Neural Network</i>
CNNum	Conseil national du numérique (France)
Copernicus	Copernicus est le programme d'observation de la Terre de la Commission européenne qui gère six services thématiques liés à l'atmosphère, au milieu marin, aux terres, au changement climatique, à la sécurité et aux urgences.
Data mining	Analyse et fouille de données (permet d'analyser un grand volume de données et d'en faire ressortir des modèles, des corrélations, des tendances)
Dinum	Direction interministérielle du numérique
Ecolab	Pôle Intelligence artificielle du Commissariat général au développement durable du ministère de la transition écologique
ECMWF	<i>European Centre for Medium-Range Weather Forecasts</i> . Centre européen pour les prévisions météorologiques à moyen terme
ESIL	<i>Environmental Data Science Innovation & Inclusion Lab</i> (Laboratoire d'innovation et d'inclusion en science des données environnementales)
Gafam	Google, Apple, Facebook, Amazon et Microsoft (acronyme des géants du Web américains)
GeoIA	L'intelligence artificielle géospatiale (GeoAI) est l'application de l'intelligence artificielle aux données, sciences et technologies géospatiales (dont les SIG)
Giec	Groupe d'experts intergouvernemental sur l'évolution du climat (de l'Organisation des Nations Unies – ONU)
GPT	<i>Generative Pretrained Transformer</i> (à la base de l'intelligence artificielle générative)
GPU	<i>Graphics Processing Units</i> (mémoire graphique)
HDH	Health Data Hub : plateforme de données de santé mise en place par le gouvernement pour combiner les bases de données de santé existantes et faciliter leur utilisation à des fins de recherche et

développement. « Nous garantissons l'accès aisé et unifié, transparent et sécurisé, aux données de santé pour améliorer la qualité des soins et l'accompagnement des patients ».

IA	Intelligence artificielle
IA Act	Règlement européen sur l'intelligence artificielle
IGN	Institut national de l'information géographique et forestière
Inria	Institut national de recherche en sciences et technologies du numérique
Inserm	Institut national de la santé et de la recherche médicale
LiDAR	<i>Light Detection And Ranging</i>
LLM	<i>Large Language Models</i> (en français grand modèle de langage) : type de programme d'intelligence artificielle capable de reconnaître et de générer du texte
Machine learning	L'apprentissage automatique est un champ d'étude de l'intelligence artificielle qui vise à donner aux machines la capacité d'« apprendre » à partir de données, via des modèles mathématiques. Il s'agit du procédé par lequel les informations pertinentes sont tirées d'un ensemble de données d'entraînement (Glossaire de l'intelligence artificielle (IA) Cnil)
Multimodal learning	Évolution du <i>machine learning</i> qui combine plusieurs sources de données simultanément de type texte, image ou son pour résoudre des tâches plus complexes. C'est l'IA générative multimodale.
Naïades	Base de données sur la qualité des eaux de surface
NLP	<i>Natural Language processing</i>
OpenDataFrance	Association créée en 2013 qui accompagne et fédère les acteurs publics territoriaux engagés dans l'ouverture et l'usage de la donnée
Random forest	Les forêts aléatoires sont une méthode d'apprentissage automatique ensembliste, se basant sur de multiples arbres de décision entraînés sur des sous-ensembles de données légèrement différents.
Réseau de neurones	<i>Artificial neural network</i> en anglais. Dans le domaine de l'intelligence artificielle, un réseau de neurones artificiels est un ensemble organisé de neurones interconnectés permettant la résolution de problèmes complexes tels que la vision par ordinateur ou le traitement du langage naturel.
RGPD	Règlement général sur la protection des données. Mis en place en 2016, il encadre le traitement des données personnelles sur le territoire de l'Union européenne.
RNNs	<i>Recurrent Neural Networks</i> , réseaux de neurones récurrents en français sont des modèles d'apprentissage automatique qui permettent d'analyser des séquences de données, telles que du texte, de la parole ou des séries temporelles.
SIG	Système d'information géographique
SNDS	Système national des données de santé
SQL	<i>Structured Query Language</i> (langage informatique normalisé servant à exploiter des bases de données relationnelles)
TAL	Traitement automatique du langage (traduction de NLP)
Web Scraping	Collecte de données de sites Web par l'utilisation d'un script

1. INTRODUCTION

L'intelligence artificielle (IA) est souvent synonyme d'innovation et porte des enjeux qui dépassent les méthodes de calcul et les concepts génériques. Intégrer l'IA dans un programme de travail, l'afficher dans une feuille de route projette une structure et des équipes dans la prospective et l'innovation. Cela les rend potentiellement attractives pour le financement, le recrutement, et d'autres opportunités. Si l'IA est déjà très présente dans la recherche depuis une cinquantaine d'années, son déploiement à grande échelle, dans l'industrie notamment, a démarré en 2020 [1].

En France, la stratégie nationale pour l'intelligence artificielle a été lancée en 2017 par le président François Hollande [2].

Le **Plan France IA**, qui fait suite au [rapport Stratégie France IA](#) [2] présenté en mars 2017, vise à « *promouvoir la recherche, l'innovation et l'adoption de l'IA en France* ». Ce plan comprend des investissements dans la recherche fondamentale, la création de chaires universitaires, le soutien aux *start-ups* et la promotion d'une éthique de l'IA.

Dans le même élan, le **Rapport Villani** de 2018 « *Donner un sens à l'intelligence artificielle* », incite au développement d'une IA dont la coordination des recherches est assurée par l'État. Dès lors, ont été créés des instituts interdisciplinaires d'intelligence artificielle (3IA) répartis sur l'ensemble du territoire national. La **direction interministérielle du numérique** (Dinum) a quant à elle été mise en place en 2019 pour élaborer la stratégie numérique de l'État.

Le **Programme d'investissement d'avenir** (PIA) est quant à lui lancé en 2021. Il est doté de 20 milliards d'euros sur cinq ans dont une enveloppe de 12,5 milliards d'euros pour financer des investissements dans les filières technologiques et des projets collaboratifs.

La stratégie nationale pour l'IA s'inscrit finalement dans la suite du Rapport Villani et dans le cadre du **Plan France 2030** [4] qui vise à accélérer la transformation des secteurs clés de l'économie française par l'innovation, pour mettre l'IA au service de l'économie et de la société. Un coordinateur national pour l'IA est nommé depuis janvier 2023 pour articuler et piloter la stratégie nationale pour l'IA au niveau interministériel.

Aujourd'hui, il existe beaucoup de documents stratégiques pour structurer le développement de l'IA dans les organisations publiques et privées. On peut citer la feuille de route de l'intelligence artificielle et de la transition écologique du pôle ministériel 2021-2024 [5] ou encore la feuille de route IA de l'IGN 2022-2024 [6].

En 2022, l'**Agence de l'innovation en Santé** (AIS)¹ est créée, représentant ainsi une des mesures phares du plan « innovation Santé 2030 » qui a vocation à piloter, en lien avec les ministères et les opérateurs concernés, la mise en œuvre du volet santé de France 2030. La France est aujourd'hui clairement identifiée comme lieu attractif pour les entreprises de la Tech, en particulier en matière d'IA. Le 15 février 2024, a été inauguré à Paris le hub Google dédié à l'IA. Celui-ci est ouvert à tout l'écosystème français du domaine.

Encore très lointaine il y a quelques années, l'IA est aujourd'hui citée dans beaucoup de domaines tels que l'industrie ou encore la santé, et constitue un sujet de préoccupation central dans les secteurs où sont manipulées et analysées des données. Le domaine de l'épidémiologie environnementale s'y intéresse de plus en plus notamment pour augmenter

¹ [Agence de l'innovation en santé : la feuille de route et les douze travaux prioritaires présentés | enseignementsup-recherche.gouv.fr](#)

les capacités de montée en qualité et d'enrichissement des données environnementales avec des données de mesures, des données d'observation de la terre (occupation du sol, températures, etc.), des données de modélisation (notamment pour la prédiction en termes d'expositions environnementales – par exemple concernant les vagues de chaleur, les inondations, les feux de forêts, etc.) ou des données sur des événements de santé.

Si l'intérêt grandissant pour l'IA est manifeste, il semble parfois difficile d'identifier précisément le périmètre, les outils et les méthodes sous-tendues par la notion d'« intelligence artificielle ».

Les concepts, les méthodes, les champs d'application, les limites et les perspectives de l'IA font l'objet d'un nombre important de conférences, webinaires, formations scientifiques qui permettent de se familiariser avec le sujet, d'acquérir des connaissances et ainsi de préciser les contours de ce champ (Annexe 3).

Ce travail de synthèse des connaissances centré sur les utilisations des méthodes d'IA pour la santé environnementale (SE), basé sur la littérature (Annexe 1) et l'interrogation d'experts du sujet (Annexe 2), a permis de faire la lumière sur un nouvel univers de mots et expressions qu'il est nécessaire de définir tant les interprétations sont nombreuses. En effet, l'IA est un sujet aux acceptions multiples et parfois floues. Le terme peut faire référence à diverses méthodes telles que l'automatisation, l'apprentissage statistique ou *machine learning*, l'apprentissage profond ou *deep learning*...

Compte tenu de la nécessité de préciser ce que le terme d'IA recouvre ici, le présent document propose d'abord de revenir, en introduction, sur la définition des grands concepts et termes structurants. La première partie dresse ensuite un panorama des types de méthodes caractérisant l'intelligence artificielle, des types d'approches sous-jacentes, et des facteurs facilitant et limitant la mise en œuvre de projets mobilisant l'IA. Une focalisation sur certaines méthodes, pour en comprendre les grands principes, est ensuite proposée.

Une seconde partie développe la façon dont l'IA peut être mobilisée dans le cadre d'études dans le domaine de la santé-environnement. Des exemples de projets menés sont dans ce cadre mentionnés à titre illustratif, et des suggestions de champs prometteurs à l'intersection de l'IA et de la santé-environnement sont formulées en complément.

Enfin, le document conclut sur les opportunités et différents types de contraintes (notamment juridiques, éthiques et techniques) pour la mobilisation de l'IA au service de la santé-environnement.

2. DÉFINITIONS ET PANORAMA DES MÉTHODES D'IA

2.1 Définitions de l'IA

2.1.1 Générale

Il n'existe pas de définition unique de l'intelligence artificielle. Celle-ci a sans doute évolué au cours du temps, au gré des avancées technologiques et des découvertes scientifiques dans le domaine (Figure 1).

Dans le dictionnaire *Le Robert* (version en ligne), elle est définie comme « l'ensemble des théories et des techniques développant des programmes informatiques complexes capables de simuler certains traits de l'intelligence humaine (raisonnement, apprentissage...). Cette définition laisse un flou autour des deux composantes clés de l'IA que sont « des programmes informatiques complexes » et « une simulation de certains traits de l'intelligence ». Il est donc indispensable d'y associer d'autres visions pour éclaircir le sens de l'association des deux notions « intelligence » et « artificielle ».

Cédric Vasseur, spécialiste en IA et contributeur du Rapport Villani de 2018, la définit en découpant l'expression. Derrière le terme « artificiel », il s'agit d'une intelligence « créée par l'Homme ». L'intelligence est plus difficile à définir mais il est communément admis que le concept de résolution de problème en fait partie. On pourra ainsi reprendre la définition de Dartmouth en 1956 : l'IA est une machine qui résout des problèmes habituellement résolus par des êtres humains ou des animaux. De façon assez similaire, le Conseil de l'Europe, dans son glossaire de l'IA², définit l'IA comme l'« ensemble des sciences, théories et techniques dont le but est de reproduire par une machine des capacités cognitives d'un être humain. Les développements actuels visent à pouvoir confier à une machine des tâches complexes auparavant déléguées à un humain. »

L'intelligence artificielle est définie par le *United Nations Environment Programme* comme l'ensemble des processus qui donnent aux machines la possibilité d'apprendre de l'expérience à mesure qu'elles recueillent davantage de données pour effectuer des tâches comme des humains. Ces processus incluent l'apprentissage (acquisition d'informations et les règles d'utilisation de ces informations), le raisonnement (utiliser des règles pour aboutir à des conclusions approximatives ou définitives) et l'autocorrection [7].

Des contours qui ont évolué au cours du temps...

Malgré ce manque de consensus de définition, l'IA a une histoire qu'il est important de rappeler pour comprendre les bases de la construction du domaine : la conférence de Dartmouth est le lieu de naissance de l'intelligence artificielle. Organisée par John Mc Carthy et Marvin Minsky en 1956, elle a mis en lumière les premiers éléments des réseaux de neurones ou du *machine learning*.

S'en suivent des périodes de faste mais aussi de régression de l'intelligence artificielle. On parle alors d'hivers de l'IA au gré des avancées technologiques, des nouvelles données et du développement des puissances de calcul.

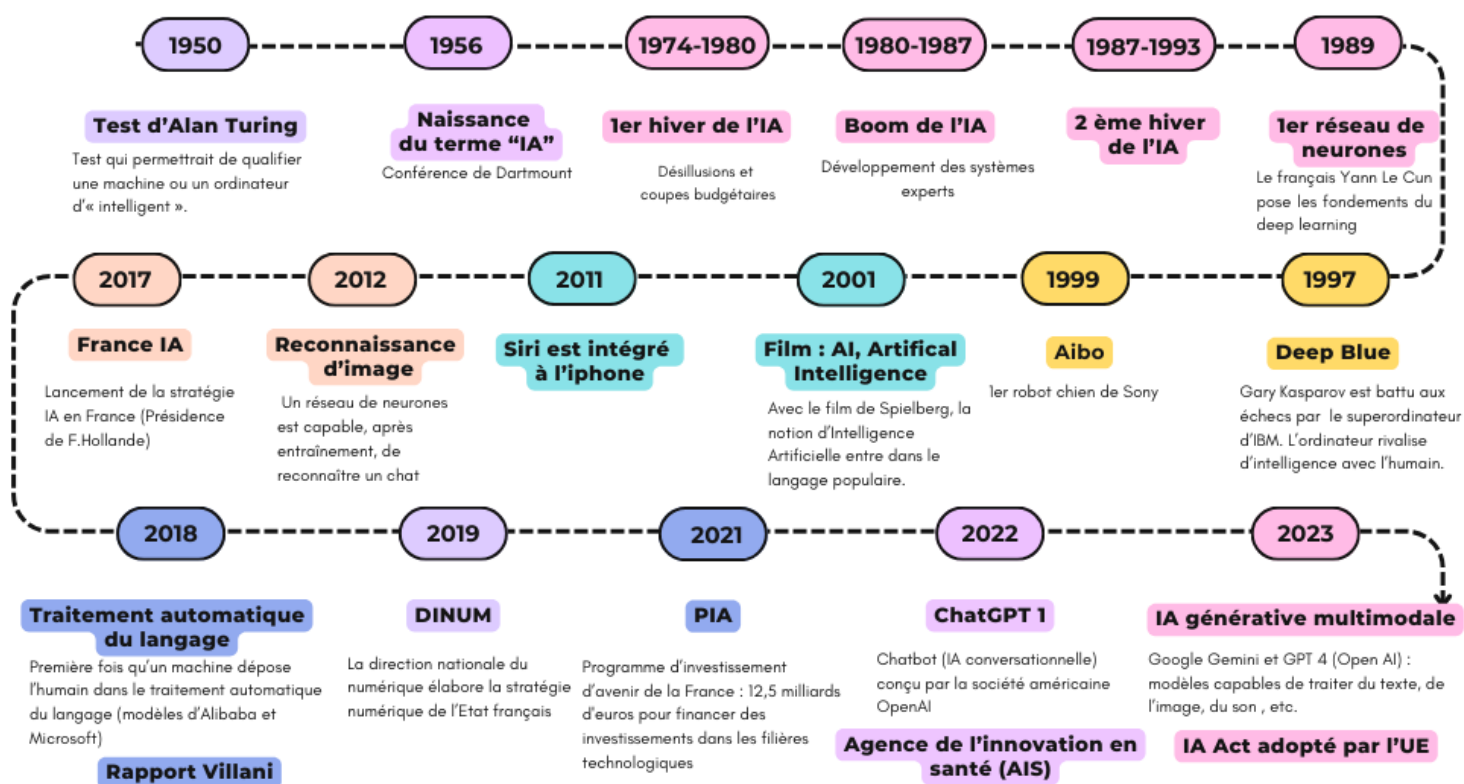
² <https://www.coe.int/fr/web/artificial-intelligence/glossary>

Un véritable tournant pour l'IA s'engage en 2012 avec les progrès liés au *deep learning* amenés par l'équipe de Geoffrey Hinton qui surpasse alors les technologies de l'époque dans la reconnaissance d'objets et les IA génératives³ (ChatGPT, Midjourney, etc.).

Le *deep learning* ou « apprentissage profond » est une invention française. En effet c'est Yann Le Cun qui pose les fondements du *deep learning* en 1989 puis qui mettra au point en 2010 un système de réseaux de neurones basé sur le fonctionnement du cerveau humain pour la vision.

Figure 1 : Chronologie de l'IA (alimentée par des sources diverses⁴)

Chronologie de l'IA dans le contexte français



Dans le nouveau paysage que dessine l'IA comme ensemble d'outils et de méthodes, il y a également l'idée d'une communauté de l'innovation. Faire partie de la communauté de l'IA est un aspect extrêmement valorisant. À l'inverse, se sentir exclu de ce « monde innovant » apparaît comme décevant voire dégradant. L'IGN (Institut national de l'information géographique et forestière) vise par exemple à démocratiser l'IA en son sein, mais aussi pour l'ensemble de la société, en accompagnant pédagogiquement la capacité d'agir qu'elle offre sans laisser personne sur le bord de la route [6].

³ 2012: A Breakthrough Year for Deep Learning <https://medium.com/neuralmagic/2012-a-breakthrough-year-for-deep-learning-2a31a6796e73>

⁴ https://fr.wikipedia.org/wiki/Histoire_de_l%27intelligence_artificielle

2.2.2 La GéolA

Depuis quelques années, on entend parler d'IA appliquée aux données spatiales : **la GeolA ou Geo Intelligence**. Il s'agit de l'IA appliquée aux données géographiques, aux méthodes d'analyse spatiale et aux outils d'intelligence géographique (systèmes d'information géographique, SIG) dans le but d'accélérer notre compréhension du monde (dont les expositions environnementales)⁵. Les spécificités des données spatiales (en termes de taille, de format, etc.), qui limitent assez vite les traitements avec des outils et des méthodes SIG classiques, en font de bonnes candidates pour l'utilisation de méthodes d'IA [8]. On parle de géo intelligence et la notion intègre trois domaines : l'intelligence décisionnelle, l'intelligence artificielle et la technologie des systèmes d'information géographique (SIG) pour l'aide à la décision à l'échelle des territoires.

Les objectifs du couple SIG/IA sont généralement : la simulation de données géographiques manquantes, la simulation de scénarios (simulation de l'évolution de l'occupation des sols, simulation d'accidents industriels en 3D, etc.) ou encore la prédiction pour l'aide à la décision et la gestion des risques (déforestation, sécheresses, inondations) [8].

La spécificité de la GeolA est d'autant plus notable qu'aujourd'hui, l'information géolocalisée est omniprésente (navigation, capteurs embarqués, suivi des colis, etc.) et constitue un enjeu pour la surveillance des environnements et par conséquent pour l'estimation des expositions des populations aux pollutions environnementales qui représentent un déterminant de santé majeur [9]. C'est aussi une source précieuse d'information sur les comportements individuels spatialisés (consommations alimentaires, pratiques du sport, type de mobilité, etc.).

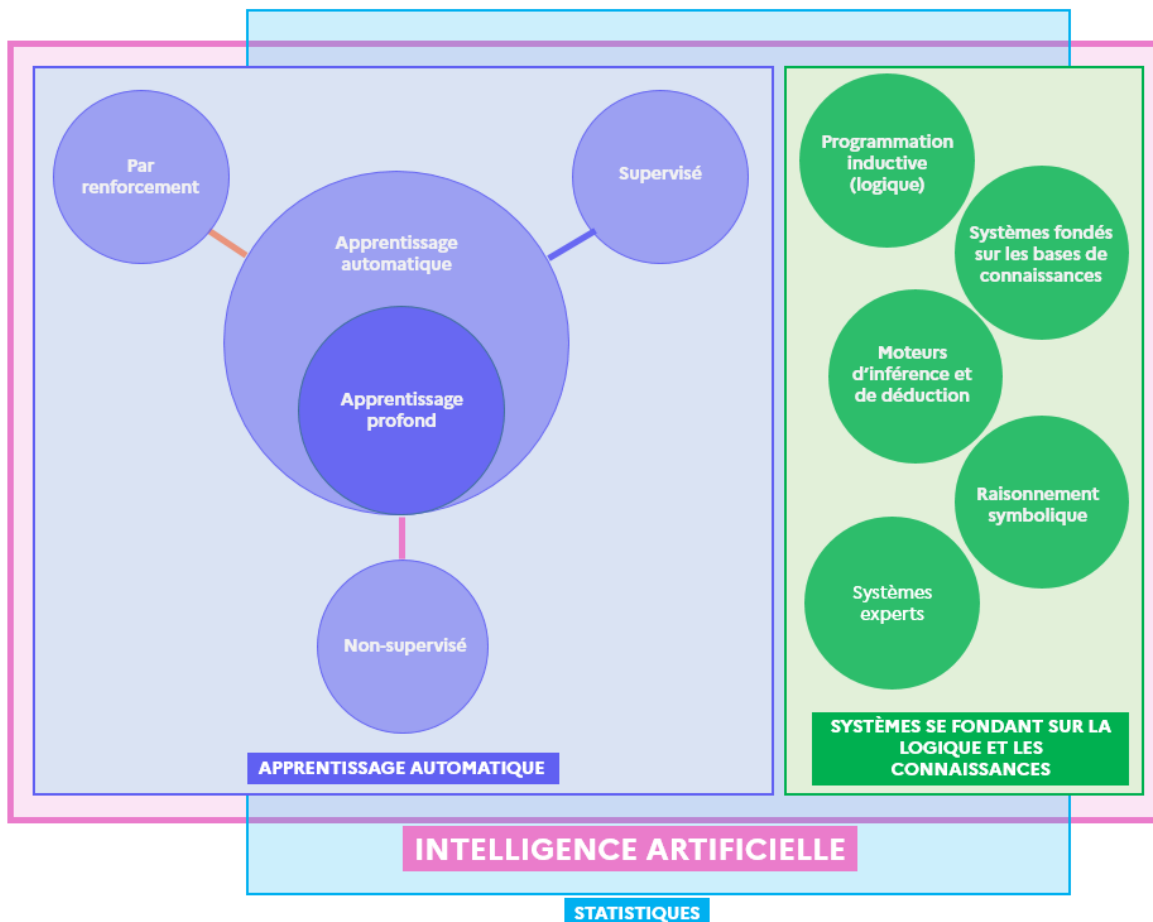
La GeolA s'appuie totalement sur les différentes méthodes d'IA connues qui sont balayées ci-après : *machine learning*, *deep learning* en premier lieu.

2.2 Cartographie des méthodes qui constituent l'IA

Au-delà de l'apprentissage statistique (*Machine learning*) et de l'apprentissage profond (*Deep learning*) que l'on connaît parfois mieux, l'IA est un domaine très vaste qui regroupe beaucoup de méthodes : la reconnaissance textuelle, les systèmes experts, la robotique, le traitement automatique du langage, l'analyse de données massives, etc. Les méthodes d'intelligence artificielle se nourrissent notamment des approches statistiques plus « classiques » (Figure 2).

⁵ [https://www.esri.com/fr-fr/capabilities/geoai/overview#:~:text=L'intelligence%20artificielle%20g%C3%A9ospatiale%20\(GeoAI,environnement%20et%20les%20risques%20op%C3%A9rationnels.](https://www.esri.com/fr-fr/capabilities/geoai/overview#:~:text=L'intelligence%20artificielle%20g%C3%A9ospatiale%20(GeoAI,environnement%20et%20les%20risques%20op%C3%A9rationnels.)

Figure 2 : cartographie schématique des méthodes qui constituent l'IA



Une distinction est souvent opérée entre systèmes d'apprentissage stochastiques ou statistiques d'une part, et systèmes déterministes d'autre part.

Systèmes stochastiques ou statistiques : différents **types d'apprentissage automatique** peuvent être distingués⁶ :

- l'apprentissage supervisé (à partir de données annotées) ;
- l'apprentissage non supervisé (à partir de données non labellisées) ;
- l'apprentissage par renforcement⁷.

Systèmes déterministes : **systèmes se fondant sur la logique et les connaissances**. La Cnil⁸ distingue ainsi :

- la programmation inductive (logique) ;
- les systèmes fondés sur les bases de connaissance ;
- les moteurs d'inférence et de déduction ;
- le raisonnement symbolique ;
- Les systèmes experts.

⁶ Voir notamment : <https://www.cnil.fr/fr/quel-est-le-perimetre-des-fiches-pratiques-sur-lia>

⁷ Par apprentissage par renforcement, on entend un algorithme pour lequel l'apprentissage est basé sur les différentes expériences rencontrées préalablement. L'optimisation de la prédiction s'opère ainsi au cours du temps, sur la base des expériences passées (<https://www.cnil.fr/fr/definition/apprentissage-par-renforcement>).

⁸ <https://www.cnil.fr/fr/quel-est-le-perimetre-des-fiches-pratiques-sur-lia>

2.3 Les approches sous-jacentes à l'IA

L'IA fait appel à trois types de techniques et d'approches :

- Une approche statistique. En effet, l'apprentissage statistique (*Machine learning*) se base intrinsèquement sur une approche statistique (régression linéaire, régression logistique, les règles de Bayes et les tests d'hypothèses pour découvrir des modèles et des relations dans les données), comme le laisse transparaître le terme lui-même.
- Une approche basée sur l'apprentissage automatique (à l'aide de données d'entraînement). Par exemple, une tâche d'apprentissage statistique (*machine learning*) de type apprentissage supervisé permet d'apprendre une fonction de prédiction à partir d'exemples annotés : des données d'entraînement sont ainsi utilisées.
- Une approche fondée sur la logique et la connaissance. Ainsi, les méthodes d'intelligence artificielle comprennent des modèles plus ou moins complexes, parfois basés sur un ensemble de connaissances et règles logiques sous-jacentes.

2.4 Les facteurs facilitant la mise en œuvre de projets d'IA

Les techniques d'IA peuvent être appliquées grâce à la conjonction de plusieurs facteurs ⁹:

- L'existence de plus en plus de données numériques ;
- Des technologies permettant les flux d'échanges de données ;
- Des algorithmes de plus en plus complexes inspirés des neurosciences ;
- La progression exponentielle de la puissance des processeurs (loi de Moore).

Ainsi, les sciences des données englobent à la fois l'IA telle que définie précédemment, soit la dimension de modélisation (apprentissage automatique, statistique ou basé sur la logique et la connaissance), et cette seconde dimension relative à la gestion et à la manipulation des données massives (à laquelle est notamment associée la capacité computationnelle) qui rend possible la réalisation de certains projets d'IA.

FOCUS : IA et Big Data : quelles différences ?

Il existe de nombreux domaines d'application de l'intelligence artificielle qui s'appuient désormais sur les données massives et les processus nécessaires à leur exploitation : le Big Data (parfois désigné par les termes « données massives » ou mégadonnées).

Quatre phases peuvent être distinguées dans le cycle de vie des données massives [10] :

- L'acquisition et l'enrichissement des données ;
- Leur intégration et leur homogénéisation (qui renvoient étroitement aux questions d'interopérabilité entre bases de données) ;
- L'analyse des données, la visualisation et l'interprétation en résultant.

Le concept de Big Data implique ainsi à la fois le stockage, le traitement et l'analyse des grands volumes de données. La nature et les sources de provenance des données massives sont nombreuses : données de navigation ou de mobilité géolocalisées, données extraites d'internet par *web scraping*¹⁰, données satellites issues de la télédétection, mesures issues de capteurs individuels. Il peut ainsi notamment s'agir de textes, d'images ou encore de vidéos.

Généralement, les mégadonnées sont générées en temps réel et à grande échelle, ce qui explique leur caractère volumineux.

L'intelligence artificielle n'est pas synonyme de Big Data, mais les méthodes d'intelligence artificielle se basent sur les données massives, notamment pour l'entraînement des algorithmes.

⁹ [Le développement de l'intelligence artificielle : risque ou opportunité | vie-publique.fr](http://vie-publique.fr)

¹⁰ *Web scraping* : collecte de données de sites Web par l'utilisation d'un script.

2.5 Focus : *Machine Learning*

Le paragraphe suivant détaille la démarche générale associée à la réalisation d'un projet de *Machine Learning*. L'objectif d'une tâche de *Machine Learning* est de trouver un modèle/une fonction qui prédit la variable de sortie à partir des variables d'entrée avec un risque espéré le plus faible. La démarche est explicitée ci-après :

- **Identification des variables explicatives (espace d'entrée) et des variables expliquées (espace de sortie).** Chaque observation est définie par des valeurs associées aux variables explicatives et par une valeur prise pour la variable expliquée.
- **Définition d'une fonction de perte, d'une famille de modèles, et d'une méthode d'estimation.**
 - La fonction de perte permet d'évaluer la mesure dans laquelle le modèle correspond à la valeur attendue. Le choix de la fonction de perte dépend de la nature de la tâche de modélisation (par exemple classification, régression), de la famille de modèles parmi laquelle la recherche est effectuée, et de la procédure d'optimisation retenue pour la sélection du modèle au sein de cette famille.
 - On distingue plusieurs familles des modèles : les modèles linéaires, pour lesquels la prédiction est obtenue selon une combinaison linéaire de variables explicatives, les modèles polynomiaux, ou encore les modèles non-linéaires.
 - Il convient également de retenir une méthode d'estimation du modèle. On distingue ainsi :
 - Les méthodes d'apprentissage non supervisé : l'apprentissage non supervisé recouvre des méthodes d'apprentissage automatique où l'apprentissage se fait sur des données qui ne présentent pas de variable à prédire (ou label). L'objectif est l'étude de la structure sous-jacente des données. Il peut notamment s'agir de mettre en évidence des classes d'observations présentant des caractéristiques communes. L'algorithme des K-moyennes constitue un exemple de tâche d'apprentissage non supervisé.

FOCUS : l'algorithme des K-moyennes

Objectif : mettre en évidence des ensembles homogènes à partir d'un ensemble de données.

Si l'on cherche à obtenir 3 groupes les plus homogènes possibles, on cherche 3 points centraux (centroïdes) qui vont être utilisés pour définir les 3 sous-groupes de données.

1. Au départ, les 3 centroïdes sont placés de manière aléatoire au sein de l'ensemble des points.
2. L'algorithme parcourt ensuite les points de données et la distance entre chaque point de l'ensemble des points et les 3 centroïdes est mesurée. Chaque point de données est alors regroupé avec le centroïde le plus proche. Chaque point se trouve alors affecté à l'un des 3 groupes selon le centroïde auquel il a été rattaché. Sur la base des nouveaux groupes créés, un nouveau centroïde est calculé pour chaque groupe en prenant pour périmètre l'ensemble des points le constituant.
3. La démarche explicitée en 2. est répétée avec les nouveaux centroïdes, et des itérations du processus ont lieu jusqu'à ce que les centroïdes soient stabilisés.

- Les méthodes d'apprentissage supervisé : l'apprentissage supervisé recouvre des méthodes d'apprentissage automatique où l'apprentissage se fait sur des données qui présentent une variable à prédire (ou « label »).

FOCUS : arbres de décision et forêts aléatoires

Les arbres de décision

Ils reposent sur une suite de décisions pour prédire un résultat (appartenance à une classe/catégorie). L'objectif est ainsi de déterminer de manière optimale des classes/catégories de sorte à attribuer chaque observation à l'une des classes de manière la plus précise possible. Les arbres de décision visent ainsi à décomposer un problème de classification en une série de tests structurés par un ensemble de classes. Ces classes sont alors considérées comme des sous-régions homogènes. Les arbres de décision n'apparaissent pas toujours robustes : de petites variations de données peuvent induire des arbres de décision totalement différents. Ils peuvent notamment mener à un problème de surapprentissage. Cela peut mener à générer des arbres trop complexes, pouvant ainsi compromettre la propension à les réappliquer pour prédire de nouvelles données en dehors de l'ensemble d'apprentissage.

Les forêts aléatoires (ou forêts d'arbres de décision)

Ce type de méthode permet notamment de remédier à la forte variabilité pouvant caractériser les arbres de décision. La méthode des forêts aléatoires permet ainsi d'obtenir une plus grande stabilité des prédictions. Les forêts aléatoires consistent à effectuer un vote des arbres (en classification) ou une moyenne des arbres (cas de régression) obtenus sur des échantillons d'apprentissage. On retient ainsi comme prédiction la valeur la plus fréquente (cas de classification) ou la moyenne des prédictions de tous les arbres (cas de régression). C'est une méthode utilisée, par exemple, pour faire de l'imputation de données manquantes, avec des algorithmes comme MissForest¹¹ [11].

2.6 Focus : l'apprentissage profond

Le *deep learning* recouvre des méthodes d'apprentissage automatique statistique utilisant des réseaux de neurones qui sous-tendent plusieurs couches de neurones cachées. Ils nécessitent des capacités computationnelles importantes et un nombre important de données d'entraînement car de très nombreux paramètres doivent être estimés¹².

Autrement dit, le *deep learning* recouvre un ensemble de méthodes basées sur des architectures complexes combinant différentes transformations non-linéaires. Les méthodes d'apprentissage profond sont basées sur les réseaux de neurones « classiques », qui sont combinés et enrichis pour former des architectures complexes telles que les réseaux de neurones convolutifs ou les réseaux de neurones récurrents. Les réseaux de neurones convolutifs tendent à être particulièrement utilisés pour le traitement d'images, et les réseaux de neurones récurrents ont permis des avancées dans le domaine du langage naturel notamment.

¹¹ [Analyse de la performance de la méthode d'imputation de données manquantes missForest et application à des données environnementales - Espace ETS \(etsmtl.ca\)](#)

¹² <https://www.cnil.fr/fr/definition/apprentissage-profond-deep-learning>

3. IA POUR LA SANTÉ ENVIRONNEMENT : LES MÉTHODES ADAPTÉES

3.1 Le potentiel de l'IA

Intelligence épidémiologique : un potentiel sous-exploité ?

Malgré l'importance forte accordée à l'IA comme voie d'innovation pour la santé publique, la littérature compte assez peu de travaux mobilisant l'IA pour l'épidémiologie comme le montre Rodriguez-Gonzales en 2019 [12].

Dans la littérature, les méthodes d'IA sont assez peu mises en avant dans les études épidémiologiques, *a fortiori* pour des sujets santé environnement. Pourtant les méthodes d'IA balayées en amont sont autant de nouvelles possibilités de traitement complexe des données collectées (spécifiquement le *machine learning* et le *data mining*) pour évaluer la vulnérabilité des territoires, l'exposition des populations dans le temps et dans l'espace grâce à la modélisation ou la reconstitution des données historiques par exemple [9].

Dans un objectif de croisement des données environnementales et de santé pour la santé publique, les épidémiologistes ont recours à des données de surveillance environnementale ou utilisées comme telles. Il s'agit notamment de données géo référencées, de type imagerie satellite ou aérienne, données de mesures, grilles de modélisation de polluants dans l'air, etc.

Ces données spécifiques (en termes de format, de taille) sont spécialement visées par les méthodes d'IA qui facilitent leurs traitements par l'augmentation des capacités de calcul notamment. Ainsi la GeolA, définie plus haut, est une sorte de clé de voûte de l'épidémiologie environnementale.

L'IGN a d'ailleurs fait de l'IA son alliée et plus précisément « *les techniques d'apprentissage machine (...) permettent de généraliser, d'extrapoler, de systématiser et d'accélérer la production de descriptions des évolutions du territoire. C'est grâce à notre capacité à extraire de l'information sur les images aériennes, les photos de constellations de satellites ou d'autres types de capteurs et nos avancées en termes d'apprentissage profond que nous pouvons produire ces descriptions précises* ». Or l'estimation des expositions environnementales dans les études épidémiologiques en santé environnement utilise très largement les données d'observation des territoires (occupation du sol y compris agricole, topographie, réseaux de transport, espaces verts et bleus, bâtiments industriels, etc.).

3.2 Types de méthodes et cas d'usage pour mobiliser l'IA et les données en santé-environnement

D'une manière générale l'IA devient une assistance aux utilisateurs novices dans l'analyse et la mise en forme des données. Il est désormais possible de préparamétrer des outils permettant à tout un chacun de demander la visualisation d'une analyse ou d'une autre.

La typologie ci-après ne recouvre pas l'ensemble des techniques qui pourraient être qualifiées comme mobilisant l'intelligence artificielle : le travail de référencement se concentre ici sur les méthodes mobilisant intelligence artificielle et données.

Les exemples d'applications des méthodes d'IA dans le champ de la santé-environnement dans la littérature et les différents témoignages recueillis parlent de l'IA comme d'un moyen permettant d'améliorer les analyses statistiques, par l'automatisation et l'augmentation de la vitesse de calcul. L'IA peut apporter une plus-value aux différentes étapes de mobilisation des données :

- Collecte et extraction
- Prétraitement et consolidation
- Analyse
- Valorisation et communication

3.2.1 Collecte et extraction de données

Cela peut recouvrir l'extraction de données de formats non exploitables facilement en formats plus aisément mobilisables (structurés), et à sélectionner et retraiter certaines des informations présentes.

On retrouve notamment :

- Des méthodes d'extraction permettant de reformater les informations non structurées (texte en langage naturel par exemple) sous un format plus aisément exploitable. Par exemple, l'extraction automatisée d'informations cibles dans des fichiers PDF.
- Des méthodes permettant d'extraire et de convertir sous une forme exploitable des données sous forme d'images ou de vidéos (notamment des données satellitaires).

Cas d'usage :

Extraire des données de documents PDF

L'IA est un bon outil pour combiner et structurer les données existantes. Elle va permettre de « consolider l'information » avec des données de mesures locales consolidées dans des fichiers PDF par exemple. Le projet UniBSV mené dans le cadre du Challenge GD4H¹³ visait à récupérer des informations dans les bulletins de santé des végétaux au format PDF pour constituer une base de données exploitable en épidémiologie. Des travaux similaires ont été menés par l'Observatoire français de la biodiversité (OFB) à partir des rapports sur le prix et la qualité de service (RPQS) en utilisant des modèles pré-entraînés type T5 (Google) ou Etalab_ia (DINUM) afin de réduire l'empreinte environnementale et gagner du temps.

Potentiels suggérés :

Participer à alimenter des bases d'entraînement des machines spécifiques à l'épidémiologie environnementale ou aux thématiques environnementales (pesticides, etc.)

Pour la santé environnement nous aurions, par exemple, un intérêt commun avec l'Observatoire français de la biodiversité (OFB) à partager nos expériences et alimenter une base d'annotation dans le cadre de travaux de récupération automatisée de données environnementales à partir de documents au format PDF

Repérage automatisé d'informations spatialisées

À partir de lecture automatisée d'images 3D (LiDAR¹⁴), il serait pertinent, pour les études en santé environnement, de cartographier et recenser les cheminées industrielles ainsi que leur

¹³ Pour plus d'informations sur le Challenge Green Data for Health : <https://gd4h.ecologie.gouv.fr/defis>

¹⁴ La télédétection par laser est une technique de mesure de distance (téléométrie) qui exploite les propriétés de la lumière. En anglais *Light Detection And Ranging* (LiDAR)

hauteur. Ceci constitue une donnée importante pour modéliser les panaches de rejets dans l'atmosphère des substances déclarées dans les bases réglementaires.

Recueillir des données de santé complémentaires

L'IA permettrait d'aller chercher de la donnée sur des pathologies qui n'ont pas de systèmes de surveillance (Système national des données de santé – SNDS, registres, etc.) via les données recueillies par exemple par des réseaux d'associations. Les données captées par les appareils connectés pourraient, de façon encadrée (6.3), enrichir des informations sur l'état de santé des personnes ou de manière indirecte donner un éclairage sur des habitudes de consommation ou sur des pratiques sportives ou de mobilité active. Toutes ces données éparpillées collectées en dehors de tout système de surveillance sont autant d'informations que seules les méthodes d'IA peuvent identifier et prétraiter.

Mobiliser des données d'opinion

Une autre utilisation de l'IA pourrait être de « moissonner » la presse, les réseaux sociaux, les requêtes faites sur Internet par exemple pour recueillir les perceptions ou pour préciser les attentes des parties prenantes sur des problématiques environnementales et/ou sanitaires particulières. L'IA permet la prise en compte et l'identification de « rumeurs » via Internet, les réseaux sociaux, etc.

Recenser de manière automatisée des ressources (par exemple références, rapports)

Certaines IA génératives sont aujourd'hui une véritable aide pour aborder un sujet. C'est notamment le cas de l'outil *Climate Q&A* développé par Ekimetrics, dont les sources sont les rapports du Giec et qui permet d'obtenir une synthèse sur un point spécifique du rapport suivant trois modes d'expression : pour les enfants, pour le grand public et pour les experts¹⁵.

Utiliser les données environnementales captées par le citoyen

Le citoyen est de plus en plus impliqué dans la surveillance de son environnement proche. Les avancées technologiques permettent aujourd'hui de collecter des données de mesures citoyennes sur l'environnement grâce à des capteurs de qualité de l'air intérieur, du bruit ambiant, de températures, etc., comme cela est fait par *Openfood fact*¹⁶. Ces données permettent d'approcher davantage les expositions des personnes.

Optimiser la recherche documentaire

Grâce aux robots conversationnels, il est aujourd'hui possible d'interroger des publications scientifiques, de décrypter des rapports (comme proposé par *Climate Q&A* avec les rapports du Giec, etc.) ou de traiter les corpus documentaires notamment sur les sujets liés au changement climatique où les productions scientifiques sont pléthoriques.

3.2.2 Prétraitement et consolidation des données

- L'IA peut tout d'abord être utilisée pour améliorer la qualité des données. Des algorithmes de *machine learning* peuvent par exemple être utilisés pour l'imputation de valeurs manquantes.
- L'IA permet également de favoriser le croisement de sources de données en permettant une meilleure exploitation croisée de bases de données de nature différente

¹⁵ [ClimateQ&A - a Hugging Face Space by Ekimetrics](#)

¹⁶ Conférence de juin 2023 au Collège de France : <https://www.college-de-france.fr/fr/agenda/seminaire/prevention-nutritionnelle-des-maladies-chroniques-de-la-recherche-action-de-sante-publique>

(indicateurs de santé, trajectoires professionnelles, données de mesures individuelles ou écologiques, etc.). Cf. [Prospective INRS 2022](#).

- Enfin, des méthodes de prétraitement de formats de données non structurées (images, textes) utilisent l'IA, par exemple des techniques d'analyse du langage naturel (*natural language processing*) ou le traitement de données d'imagerie brute grâce au *deep learning*.

Cas d'usage :

Compléter des données manquantes : le projet ExpoCast

À l'international, on retrouve les méthodes de *machine learning* dans l'évaluation des expositions, et notamment pour compléter les données manquantes. C'est le cas du projet ExpoCast (*exposure forecasting*) développé par l'US-EPA (*United States Environmental Protection Agency*) qui a recours aux NAMs (*New Approach Methodologies*) dans la priorisation des risques liés à l'exposition aux substances chimiques [13].

Reconstituer des données d'occupation du sol

Les projets *Gouramic* et *Deep Gouramic* visent la détection automatisée des parcelles de cultures à partir d'images satellites historiques [13]. Faure et coll. ont développé une méthode d'estimation de l'exposition résidentielle agricole sur des périodes anciennes en l'appliquant à 1 155 sujets d'une étude nationale cas témoins, l'étude Testis. La méthode de caractérisation de l'occupation du sol (OCS) s'appuie sur un logiciel de segmentation sémantique semi-automatique, *Gouramic*. Les images photo aériennes utilisées proviennent de l'Institut national de l'information géographique et forestière (IGN), et sont disponibles de 1920 à 1990, permettant de s'intéresser à des expositions anciennes. Les auteurs ont également pris en compte les données météo (vents) dans le calcul de l'exposition.

Dans le projet GeoKPhyto, Andriamanamamonjy et Orazio [14] décrivent des travaux qui visent à améliorer les détections de parcelles de vignes et de vergers à partir de prises de vues aériennes à haute résolution issues de la BD ORTHO® de l'IGN et des techniques d'intelligence artificielle (*deep learning*). La méthode développée, dans les départements de la Gironde et de la Charente-Maritime, a ainsi permis d'accroître sensiblement l'identification des parcelles de vignes et de vergers, pour atteindre un taux d'exhaustivité proche de 100 %. Ces méthodes offrent des perspectives prometteuses pour augmenter l'exhaustivité des données géographiques utilisées pour caractériser un proxy d'exposition des populations aux pesticides.

Potentiels suggérés :

Contrôler la qualité, la cohérence et la complétion automatique de bases de données environnementales utilisées en épidémiologie (Naiade, Basias, etc.) grâce à une IA.

3.2.3 Analyse de données

En dehors des nouvelles façons de travailler que permet déjà l'IA comme super assistant ou l'automatisation de travaux d'état de l'art qui représentent un gain de temps énorme [15], l'IA est aussi mobilisable en épidémiologie notamment pour :

- **Explorer des phénomènes.** Il peut s'agir de caractériser les contributions et liaisons entre variables, mais également la caractérisation et le positionnement des différentes observations.
Exemple de méthode : apprentissage statistique non supervisé (algorithmes de clustering, méthodes de réduction des dimensions comme l'analyse en composante principale – ACP).
- **Expliquer des phénomènes.** L'enjeu consiste à tester l'influence de variables ou de facteurs sur un phénomène d'intérêt.
Exemple de méthode : régression et économétrie.
- **Prévoir un phénomène.** Cette tâche prédictive peut être associée à une étape préalable de sélection de prédicteurs.
Exemples de méthode : algorithmes de sélection de variables (pour sélectionner des prédicteurs), apprentissage statistique supervisé (pour la tâche de prédiction elle-même), modélisation multi-agents...
En santé-environnement, cela peut inclure la modélisation des expositions.

L'IA permet notamment d'étudier des relations non linéaires, ce que les modèles statistiques classiques ont du mal à faire.

De façon plus anecdotique, l'utilisation d'IA générative conversationnelle du type *ChatGPT* permet de générer le code d'un script d'analyse ou d'automatisation de calculs en langage de programmation (Python par exemple) ou encore déboguer un script.

Explorer : exemples de cas d'usage

Construire des typologies d'usagers, de patients, ou de territoires

L'IA pourrait être mobilisée dans la construction d'outils d'aide à la gestion dans les problématiques de *clusters* (agrégats spatio-temporels) de pathologies en permettant de déceler des schémas ou des groupes dans les données de santé qui ne sont pas évidents autrement.

Les méthodes de *clustering* (technique d'apprentissage automatique non supervisée) permettent aussi de détecter des profils de territoires en fonction des nuisances mesurées.

Les méthodes d'IA peuvent faciliter la mobilisation et l'exploitation de données d'observation de la terre pour construire des indicateurs spatialisés de surveillance de l'environnement en quantifiant les modifications de celui-ci et leurs impacts sur la santé humaine, animale et des écosystèmes (espaces verts, modifications des essences d'arbres, érosion de la biodiversité, températures, etc.).

Évaluer des expositions

L'Agence nationale de sécurité sanitaire de l'alimentation, de l'environnement et du travail (Anses) utilise des méthodes de *machine learning* qui s'appuient sur des méthodes de réduction des dimensions, de classification/*clustering* et d'analyse des réseaux afin de prioriser les mélanges et identifier des profils homogènes de co-expositions au sein de la population générale, des travailleurs, ou encore des abeilles [16]. Ces méthodes ont été

implémentées dans le logiciel Anses RSexpo afin de faciliter leur utilisation par des non-statisticiens¹⁷.

Expliquer : exemples de cas d'usage

Évaluer des expositions

En 2019, l'Inserm a mis en place un outil informatique basé sur le traitement automatique du langage permettant d'identifier des effets toxiques du bisphénol S, substitut fréquent du bisphénol A dans les contenants alimentaires, à partir de données déjà publiées dans la littérature scientifique. Plus largement, cet outil peut être déployé sur n'importe quelle substance chimique (ou d'un agent physique) sous réserve qu'elle ait fait l'objet d'études publiées dans la littérature scientifique ou soit présente dans des bases de données [17].

Analyser les facteurs contribuant à un risque afin de détecter des territoires vulnérables

La construction d'algorithmes de détection spatiale et temporelle de pollution ou de cumul d'expositions permettrait d'automatiser la réalisation de diagnostics territoriaux pouvant servir localement lors de l'investigation d'un signalement de *cluster* de pathologies ou pour guider la mise en place de mesures d'atténuation des effets du changement climatique dans certains quartiers prioritaires par exemple (végétalisation, etc.).

Wen Yifan *et al.* ont employé des méthodes de *machine learning* au sein de leur modèle de pollution de l'air dans lequel ont été introduites des données de trafic dynamiques spatialement et temporellement. Ils ont constaté que l'introduction de ces données améliorerait considérablement les modèles spatiaux de simulations du dioxyde d'azote (NO₂) de 47 % et des particules fines (PM_{2.5}) de 15 %. Ils ont ainsi réussi à capturer les niveaux de PM_{2.5} dépassant les limites en raison d'activités de trafic intense et à fournir un outil de « carte hors limites » pour identifier les disparités d'exposition. En revanche, le modèle sans données dynamiques sur le trafic ne permettait pas de saisir les disparités d'exposition induites par le trafic et sous-estimait de manière significative l'exposition des résidents aux PM_{2.5}. Les sous-estimations étaient plus graves pour les communautés défavorisées [18].

Prévoir : exemples de cas d'usage

Prévoir la dynamique des maladies

Prévisions concernant le comportement des maladies à partir de prédicteurs géographiques environnementaux : reproduction de maladies, répartition spatiale des maladies, etc.

Surveiller et prédire des risques associés à la qualité des eaux

La surveillance de la qualité des eaux ou la prédiction de l'étendue des marées noires font aussi appel à des méthodes issues de l'IA.

Les services de gestion de l'eau des collectivités locales se tournent de plus en plus vers l'intelligence artificielle (IA) pour optimiser la gestion de leurs réseaux et anticiper les problèmes de qualité ou de pénurie (Aquasys permet de gérer la ressource en eau et d'anticiper les risques hydriques¹⁸).

Par exemple, le projet *Clean Water AI*¹⁹ se base sur modèle de réseau neuronal, déployé sur des appareils de pointe qui classifient et détectent des bactéries et des particules nocives.

¹⁷ <https://www.anses.fr/fr/system/files/AUTRE2022METH0197Ra.pdf>

¹⁸ <https://www.gouvernement.fr/actualite/lintelligence-artificielle-au-service-de-la-gestion-de-leau>

¹⁹ <https://cleanwaterai.com/#intro>

Les villes peuvent installer des dispositifs IoT²⁰ pour les sources d'eau afin de surveiller la qualité des eaux en temps réel. La régie Eau de Valence a quant à elle mis en place un système de jumeau numérique de son réseau d'eau permettant de définir et d'anticiper les besoins et de détecter toute anomalie.

Aussi, les marées noires peuvent être détectées par traitement d'images satellitaires à l'aide de méthodes de *deep learning*, permettant aussi de prédire leur dispersion et de planifier les stratégies d'intervention et de nettoyage.

Prévoir les catastrophes naturelles et s'adapter au changement climatique

Les conséquences du changement climatique peuvent menacer les espaces urbanisés et leurs populations car les villes n'ont pas été pensées pour faire face à des catastrophes naturelles de plus en plus fréquentes et intenses. La ville d'Hong Kong mobilise l'IA depuis 2018 pour prévoir les typhons et les glissements de terrain en se basant sur des données météorologiques, l'analyse de prises de vues aériennes et des archives des catastrophes précédentes. En Chine, des projets sont en cours sur la prévision de feux de forêt à partir également de données météorologiques passées [19].

Une autre application est liée à la surveillance de la qualité de l'eau dans un contexte de changement climatique. On peut ainsi relever le développement de travaux sur la surveillance automatisée de la qualité des eaux de surface grâce aux données d'imagerie couplées aux méthodes d'IA [20].

Vers des smart cities : mieux surveiller la qualité de l'air et agir pour la réduction de la pollution en ville

Plusieurs projets se développent dans le monde pour diminuer la pollution de l'air liée au trafic routier ou aux activités industrielles au sein des villes. Pittsburgh a mis en place dès 2017 un système de contrôle des feux de circulation par l'IA afin de réduire le temps d'attente des véhicules et l'émission inutile de CO₂ [19]. En diminuant en moyenne de 30 % le temps d'arrêt des automobilistes, ce système a réduit de 20 % les émissions dans l'air. Un système équivalent, basé cette fois sur les données flottantes, a été mis en place par la société Odeven dans l'Allier. En Chine, *Green Horizon IBM* utilise l'IA pour anticiper les pics de pollution dans les villes, permettant aux autorités de mettre en place une stratégie de prévention ciblée en décidant la fermeture ponctuelle d'un site industriel ou en instaurant une politique de circulation alternée par exemple [19].

3.2.4 Valorisation et communication des données et des analyses

L'IA facilite la mobilisation des données massives et la data visualisation : par l'automatisation des traitements, elle permet notamment de créer des tableaux de bord (*dashboards*) et des rapports synthétiques actualisables en temps réel ou en temps proche du réel. Elle permet aussi d'adapter le message au public²¹. Un grand nombre d'outils, davantage orientés marketing et business, existent aujourd'hui pour représenter les données de la manière la plus efficace et dynamique.

²⁰ IoT : l'Internet des objets (IoT) est un réseau d'objets et de terminaux connectés équipés de capteurs (et d'autres technologies) leur permettant de transmettre et de recevoir des données entre eux et avec d'autres systèmes.

²¹ «L'IA, le nouveau game changer de la communication », Alain Garnier (Jamespot) - [Stratégies \(strategies.fr\)](https://strategies.fr)

4. LES MOYENS DONT A BESOIN LA DATA SCIENCE ?

Les liens avec la recherche en France et à l'international sont indispensables pour le conseil en *Data Science* et les compétences techniques spécifiques que nécessite la mise en œuvre des méthodes d'intelligence artificielle.

À l'Institut national de recherche en sciences et technologies du numérique (Inria), les chercheurs viennent de divers horizons dont les États-Unis où les formations en IA sont pointues. L'Université du Colorado Boulder en partenariat avec l'Université d'Arizona et d'Oslo guide les réflexions du laboratoire d'innovation pour la data science environnementale (ESiIL). Le schéma ci-après (figure 3) montre les différentes composantes nécessaires à la mise en place d'une approche IA complète : l'appartenance à une communauté scientifique pour les échanges de pratiques (compétences techniques), les réflexions sur les usages possibles de la data science au vu des données disponibles, les équipements informatiques adaptés aux ambitions d'innovation (espaces de calcul suffisamment puissants, etc.) et une offre de formation aux data sciences qui aille de l'acculturation à l'IA à la maîtrise des méthodes et des outils.

Figure 3 : Structure nécessaire à la mise en place d'une démarche scientifique innovante



(Source : <https://esiil.org/> The Environmental Data Science Innovation & Inclusion Lab (ESiIL) is a NSF-funded data synthesis center led by the University of Colorado Boulder in collaboration with NSF's CyVerse at the University of Arizona and the University of Oslo) – Le laboratoire d'innovation pour la data science environnementale est un centre de synthèse de données financé par la National Science Foundation (NSF) et dirigé par l'Université du Colorado à Boulder en collaboration avec CyVerse de la NSF à l'Université de l'Arizona et à l'Université d'Oslo.

Parmi les compétences spécifiques à l'IA et au big data il y a notamment (tableau 1).

Tableau 1 : compétences spécifiques à l'IA

(source : [Quels sont les métiers et les disciplines de la data science ? \(polytechnique.edu\)](http://polytechnique.edu))

Intitulé	Compétences
Data scientist	Il analyse et exploite la masse de données. Son rôle est de comprendre et de modéliser les données. Il a des compétences poussées en statistiques et dans les différentes méthodes de modélisation et d'analyse (solutions d'analyse et modèles de <i>machine learning</i> sur différents outils).
Data Engineer	Il intervient sur l'infrastructure data. Il est chargé d'extraire, stocker, nettoyer et structurer des données numériques brutes afin de les enregistrer dans des bases de données structurées. Il travaille avec différents langages et outils spécifiques développés dans le big data (SQL, etc.).
Data analyst	Il est chargé du recueil et de l'analyse des données. Il a des compétences clés en analyse (Python et ses bibliothèques, les tests, les méthodes statistiques et le requêtage de bases de données), visualisation des données et possède des connaissances sur le concept du <i>machine learning</i> .

Ces profils sont indispensables pour guider les équipes dans la compréhension des outils développés et le contrôle des résultats en sortie. Internaliser les compétences permet d'augmenter la compréhension des processus de bout en bout, en diminuant l'effet « boîte noire ».

En complément des compétences, les capacités informatiques et technologiques peuvent être exigeantes et nécessiter la mise à disposition de moyens importants en capacité de calcul, mémoire graphique (GPU) notamment.

5. DISCUSSION

Le sujet de l'intelligence artificielle est en plein essor et fait l'objet de nombreux débats en commençant par sa propre définition. Pourtant, elle donne également lieu à la création de communautés de travail autour de l'innovation au sens large (par exemple avec la création du club de l'IA de l'Ecolab). L'intelligence artificielle, bien que parfois controversée, représente un sujet d'avenir, dynamique et de grand intérêt pour la science.

5.1 Critique et perspectives

En même temps que l'IA devient populaire et s'intègre à toutes les étapes du cycle de vie de la donnée, des avertissements se font entendre contre une vision de l'IA providentielle et utopiste. Luc Julia, dans son ouvrage « l'intelligence artificielle n'existe pas » [21] insiste sur le fait qu'« en soit l'intelligence artificielle n'est ni bonne, ni mauvaise. C'est juste un outil, comme un marteau ». Un outil peut engendrer des erreurs ; il est important de garder cela en tête et de ne pas oublier que les capacités de cet outil sont limitées et peuvent comprendre des biais inhérents au modèle de données fournies en apprentissage.

Dans le domaine de l'imagerie médicale par exemple, des études montrent que l'IA se trompe aussi en faisant des erreurs systématiques comme dans la détection des mélanomes sur les peaux noires [22]. Cette limite pointe la nécessaire diversité d'images dans les bases d'entraînement des modèles d'apprentissage.

D'autre part, certains considèrent ces outils comme imparfaits avec un coût important pour le paramétrage, les vérifications de codes, les contrôles des résultats en sortie de modèles, etc.

5.2 Données

Le Rapport Villani 2018 pointe la nécessité de « l'accès aux données, la circulation de celles-ci et leur partage » comme élément *sine qua non* du développement de l'IA en Europe.

La qualité et la complétude des données sont le prérequis pour construire des modèles pertinents et donc d'une IA fiable. En effet à l'instar de Cédric Villani, Olivier Laurent souligne dans son analyse que pour pouvoir automatiser des traitements en entraînant les machines, il faut avoir des données stables et fiables [23]. Or si les données environnementales utilisées pour approcher les expositions des populations aux nuisances environnementales sont de plus en plus nombreuses, pour les besoins des études épidémiologiques, leur qualité et leur complétude à échelles fines ne sont pas toujours satisfaisantes. Dans ces conditions, vouloir mobiliser des outils d'intelligence artificielle est vain.

Par ailleurs, la question de la fiabilité des données d'entraînement collectées par les Gafam pose question, tout comme les enjeux de souveraineté numérique associés. En effet, on peut s'interroger sur la maîtrise du contenu et de la qualité de ces données (les banques de données et d'annotations représentent un volume colossal) et des biais qu'elles peuvent induire. Aujourd'hui, des efforts considérables sont nécessaires pour tenter de s'émanciper et de fonder les méthodes d'IA sur des bases de contenu plus sélectif, adapté aux problèmes posés et maîtrisés.

D'autre part, la question de l'accessibilité des données à l'heure de l'*open data* est primordiale pour permettre aux équipes de chercheurs de progresser. Pour autant les conditions

d'ouverture de la donnée privée posent des questions relatives à la propriété et à la protection des certaines données sensibles (identité des personnes équipées de capteurs, géolocalisation, etc.) qui nécessitent un encadrement juridique et éthique. S'agissant des problématiques environnement et santé, la combinaison *open data* et IA peut laisser craindre des croisements de données « tous azimuts » qui, sans garde-fous scientifiques, pourraient conduire à de fausses associations épidémiologiques.

5.3 Considérations éthiques et juridiques associées à l'IA

Les considérations éthiques peuvent tout d'abord toucher aux usages des données associées à la mobilisation de méthodes d'intelligence artificielle. Pour adresser les enjeux éthiques associés au partage et à l'utilisation des données, Ekitia²² propose une démarche éthique sous la forme d'une charte (<https://www.ekitia.fr/la-charte-ethique-ekitia/>) afin que ces pratiques s'appliquent conformément à des grands principes clés.

Protection des données personnelles et des libertés

Les usages environnementaux « par destination » de certaines données, par exemple des données de mobilité collectées par des applications dédiées, doivent respecter le cadre juridique applicable lors de leur transfert et traitement. Ces usages doivent faire l'objet d'une base légale, *a minima* d'une information des personnes, de l'anonymisation et/ou du respect des droits des personnes. Il reste nécessaire dans tous les cas de garantir la transparence des différents dispositifs et de leurs finalités afin d'engager les citoyens dans le partage de leurs données.

Les données d'entraînement utilisées pour construire les systèmes d'IA peuvent ouvrir sur des possibilités d'extraction des informations pouvant avoir des répercussions conséquentes s'il s'agit de données sensibles, notamment à caractère personnel. Des attaques par exfiltration de modèle, par inversion de modèle, ou encore par inférence d'appartenance, peuvent ainsi mener à extraire des informations sensibles et à une violation de données²³.

Transparence et explicabilité des algorithmes et des systèmes d'IA

Compte tenu de la complexité croissante des modèles de *machine learning*, l'enjeu d'explicabilité et d'interprétabilité des modèles a une pertinence renforcée face à une opacité parfois forte (« effet boîte noire »).

Par explicabilité, on entend la possibilité d'explicitier et de mettre en relation les éléments pris en compte par le modèle d'intelligence artificielle pour mener au résultat produit. La charte éthique des usages des données d'Ekitia mentionnée précédemment distingue ainsi la transparence, qui sous-tend la notion d'explicabilité des algorithmes notamment, comme l'un des principes éthiques clés pour encadrer les usages de données.

Lorsque le système d'IA traite des données à caractère personnel, il convient également de s'assurer que le principe de transparence du RGPD (Règlement général sur la protection des données) est respecté. En effet, selon ce règlement, toute information ou communication relative au traitement de données personnelles doit être « concise, transparente, compréhensible et aisément accessible, en des termes simples et clairs ».

²² Association qui ambitionne de construire un cadre de confiance, éthique et souverain, destiné à permettre aux acteurs de partager et de croiser leurs données tout en respectant les intérêts des individus et des propriétaires des données.
<https://www.ekitia.fr/offre-de-services/>

²³ <https://www.cnil.fr/fr/intelligence-artificielle/ia-comment-etre-en-conformite-avec-le-rgpd>

Par ailleurs, la prise de décision automatisée qui a un effet juridique ou affecte sensiblement une personne est encadrée.

Les biais et discriminations algorithmiques

Les systèmes d'IA peuvent mener à des résultats comportant des biais décisionnels et menant ainsi à des discriminations. Celles-ci peuvent être causées par les données d'apprentissage, non représentatives, qui tendent alors à exacerber ou reproduire un « motif régulier » (« *pattern* » en anglais) ne prévalant pas sur l'ensemble du périmètre d'application du système d'IA, ou bien sur l'algorithme lui-même qui pourrait comporter des failles dans son fonctionnement.

Ces biais se traduisent par des potentielles discriminations sociales et sociétales aggravant les inégalités. À ce propos Mathilde Saliou, dans son ouvrage *Technoféminisme* [24], dénonce le sexisme reproduit dans les modèles entraînés par l'IA lui-même introduit par les humains (une grande majorité d'hommes). Elle évoque également la problématique des données d'entraînement majoritairement masculines en médecine par exemple qui créent un biais dans l'interprétation d'images médicales réalisées sur les femmes.

Bouleversement du travail, santé mentale et exploitation des personnes pour la production de données

En tant qu'agence de santé publique, il faut également considérer l'IA comme nouvel outil de travail susceptible de bouleverser les organisations et d'affecter la santé mentale des travailleurs (sentiment d'inutilité, perte de sens, stress, etc.) [25]. C'est notamment ce que l'IGN pointe comme axe de travail n° 5 (mesure 16 : Analyser systématiquement les enjeux sociaux de l'IA), dans sa feuille de route [6], et qu'elle nomme « Mettre en débat et réguler socialement et écologiquement le déploiement de l'IA ».

Par ailleurs, d'un point de vue éthique, il paraît utile de se questionner sur le système de production des données d'apprentissage. En effet, certaines entreprises exploitent des personnes non ou peu qualifiées qui collectent la donnée via Internet en cliquant sur des images par exemple. L'alimentation des bases d'annotation géantes se fait aussi par l'utilisation gratuite des personnes : lorsque nous cliquons sur les images pour garantir qu'il ne s'agit pas d'un robot ou encore lorsque des applications de location de vélo ou de trottinettes demandent une photo du vélo ou de la trottinette stationnée.

Le Centre scientifique et technique du Bâtiment (CSTB), quant à lui, fait appel à une entreprise éthique d'annotation qui s'assure des bonnes conditions de travail de ses employés.

Enfin des initiatives vertueuses construisent un cadre rassurant autour de l'usage des données et des méthodes propres à l'IA : par exemple, le label Ekitia²⁴ abordé plus haut vise à valoriser les projets respectueux de la charte éthique d'usage des données créée par Ekitia en 2020²⁵. L'objectif est de favoriser le développement des usages des données dans un cadre de confiance éthique, équitable et souverain.

²⁴ [Le Label Ekitia continue son développement ! - Ekitia](#)

²⁵ [La Charte Éthique - Ekitia](#)

5.4 Sobriété numérique : l'IA frugale

Nous devons être conscients de l'impact écologique de l'utilisation des méthodes et technologies liées au développement de l'IA et à ce titre nous interroger sur son apport réel au regard de cet impact. On parle aujourd'hui d'IA frugale à l'heure de la sobriété numérique. L'IA frugale est aussi par ailleurs un domaine d'innovation. En effet, c'est une prouesse technique de limiter les paramètres d'un modèle en restant tout aussi performant. Et être capable de faire tourner des modèles « en local » est aussi un gage d'accessibilité à l'IA pour des petites équipes. On sait aussi que s'agissant des IA génératives, le plus coûteux en ressources énergétiques n'est pas le développement de l'outil en lui-même (par exemple « [Croissant LLM](#) » a nécessité 500 g de CO₂) mais davantage la massification des usages et la génération d'images [26].

Aujourd'hui le numérique est responsable de 10 % du réchauffement climatique (production de CO₂) bien au-delà de l'aérien [21]. Partant de ce constat, l'innovation doit aussi permettre de réduire l'impact écologique du numérique. C'est l'objectif que doivent atteindre les projets retenus dans le cadre de l'appel à projets « Démonstrateurs d'IA frugale » dans les territoires pour la transition écologique lancé par l'Ecolab vague 1 et 2.

Une méthode d'évaluation d'impact environnemental des projets impliquant des méthodes d'IA est proposée par le groupe de recherche EcolInfo dans un document de cadrage. Il s'agit de faire la liste des critères d'évaluation tenant compte des équipements numériques, des spécificités du domaine de l'IA, des impacts des différentes phases de traitement des données (collecte, apprentissage, etc.) [27].

L'Ecolab et l'Afnor s'engagent, en 2024, dans un processus de création d'une norme d'IA frugale.

Protéger les données sensibles, protéger l'environnement ?

Les outils d'intelligence artificielle représentent de grands enjeux dans la protection de la donnée, en particulier en épidémiologie environnementale mêlant données de santé et données environnementales. Soumis, comme tous les outils traitant de la donnée, à des failles et des attaques, ils peuvent également avoir des retombées insoupçonnées sur les individus et la société notamment par le manque d'explicabilité ou de transparence de certains modèles mais aussi par leur forte consommation de données d'entraînement et d'énergie pour les calculs. Ce dernier point, et comme il a été précisé dans ce rapport, soulève les enjeux environnementaux inhérents à l'utilisation de l'IA.

Le récent cahier « Innovation et prospective » de la Cnil « Données, empreintes et libertés » [28] passe en revue ces différents points de questionnement sous l'angle réglementaire. On y apprend que le Règlement général sur la protection de la donnée (RGPD) et la loi informatique et libertés, par leur principe de minimisation des données personnelles et de leur utilisation (droit à la suppression de ces données, droit d'opposition, droit au déréférencement, etc.) participent, de fait, à la réduction de l'empreinte du numérique. Le document met aussi en lumière les éventuelles contradictions ou compromis nécessaires entre protection de la donnée et sobriété énergétique. Par exemple, des systèmes avancés de protection de la donnée comme la cryptographie sont généralement perçus comme augmentant l'empreinte environnementale des traitements informatiques de données. Il convient alors de mettre leur coût carbone réel en regard de leurs bénéfices.

Ce cahier évoque spécifiquement la donnée environnementale, définie par le Conseil national du numérique²⁶ qui est souvent croisée avec des données sensibles de santé²⁷ ou plus largement personnelles. On comprend alors que des outils d'IA et notamment la GeolA qui se basent sur ces types de données doivent particulièrement être encadrés réglementairement en épidémiologie environnementale.

²⁶ Créé en 2011, le Conseil national du numérique est une instance consultative indépendante chargée de conduire une réflexion ouverte sur la relation complexe des humains au numérique. Il est composé de membres bénévoles nommés pour deux ans par le Premier ministre aux domaines de compétences variés (sociologue, économiste, philosophe, psychologue, anthropologue, informaticien, avocat, journaliste...) et de parlementaires désignés par les présidents de l'Assemblée nationale et du Sénat. Il est placé auprès de la secrétaire d'État chargée du numérique. Le CNNum définit les données environnementales ainsi : par nature, elles sont directement produites pour la connaissance et l'analyse du territoire. Les données environnementales par destination sont alors des données collectées et traitées au départ pour des usages qui ne sont pas en lien direct, mais qui peuvent renseigner sur des aspects de l'activité humaine.

²⁷ Données de santé définies par la Cnil pour le RGPD : les données relatives à la santé physique ou mentale, passée, présente ou future, d'une personne physique (y compris la prestation de services de soins de santé) qui révèlent des informations sur l'état de santé de cette personne.

6. PERSPECTIVES

Une étude d'Accenture²⁸ projetait, en France, un gain de productivité supplémentaire de 20 % grâce à l'intelligence artificielle à l'horizon 2035. On mesure alors les enjeux économiques et politiques du développement de l'IA dans tous les domaines. Compte tenu de ces enjeux, la communauté internationale est désormais pleinement consciente de la nécessité de réguler l'usage de l'IA notamment dans le domaine de la santé (IA Act²⁹) au nom de l'éthique[3].

Il existe d'ailleurs des organismes de contrôle de l'usage qui est fait de l'IA comme *AlgorithmWatch*³⁰ qui est une organisation de défense des droits de l'Homme basée à Berlin et à Zurich qui « se bat pour un monde où les algorithmes et l'intelligence artificielle (IA) n'affaiblissent pas la justice, la démocratie et la durabilité, mais les renforcent »³¹. Ils alertent par exemple sur les risques de discrimination des systèmes automatisés de prise de décision.

Les différents scénarios présentés ci-après sont issus et inspirés de la prospective menée par l'INRS « L'intelligence artificielle au service de la santé et sécurité au travail, enjeux et perspectives à l'horizon 2035 » [29]. Ils sont ici adaptés au contexte de cette synthèse portant sur la santé-environnement. Le champ des possibles a volontairement été défini de façon large afin de rester une aide à la réflexion.

- **Le scénario 1** proposé par l'INRS repose essentiellement sur l'absence de contrôle des outils de l'intelligence artificielle par les États et la place prépondérante des géants du numérique comme les AMAMA (Alphabet, Meta, Amazon, Microsoft, Apple) à l'ouest *versus* les BATX (Baidu, Alibaba, Tencent, Xiaomi) en Chine. En santé-environnement, un tel essor technologique pourrait être profitable dans le développement de nouveaux modèles adaptés à la donnée environnementale sans limite de ressources naturelles grâce à des solutions alternatives développées par ces géants pour leurs propres besoins ou de ressources financières du fait de l'attractivité du secteur. Le recours à ces outils se démocratiserait grâce à des formations mises en places par le secteur du numérique. La recherche en santé environnement serait alors totalement familière avec leur utilisation. Cependant, le caractère intrusif de ces technologies questionne, en particulier dans un domaine aussi sensible que celui de la santé, où les données sensibles qui ne seraient plus sous la protection des États seraient sous contrôle d'entreprises privées. Le recours à des capteurs pour mesurer l'exposition environnementale des personnes à travers des objets connectés du quotidien (par exemple smartphone, voiture intelligente, maison connectée etc.), et malgré l'intérêt majeur pour l'épidémiologie environnementale, pourrait poser problème vis-à-vis du recueil constant de ces données géolocalisées et de l'absence de sécurité dans leur utilisation.
- **Le scénario 2** se base, quant à lui, sur la garantie d'un cadre réglementaire des États pour les outils d'intelligence artificielle. Un début de scénario 2 semble se profiler depuis l'adoption très récente (8 décembre 2023) par le Parlement européen de l'IA Act Européen³². Il s'agit d'un règlement provisoire visant à garantir que les droits fondamentaux, la démocratie, l'État de droit et la durabilité environnementale soient protégés contre les risques liés à l'IA (manque de qualité des données d'entrée, absence de signalement des résultats issus de l'IA, non-respect du droit d'auteur,

²⁸ <https://www.accenture.com/fr-fr/insights/artificial-intelligence/ai-maturity-and-transformation>

²⁹ [Textes de la loi | Loi sur l'intelligence artificielle de l'UE \(artificialintelligenceact.eu\)](#)

³⁰ <https://algorithmwatch.org/en/>

³¹ « We fight for a world where algorithms and Artificial Intelligence (AI) do not weaken justice, democracy, and sustainability but strengthen them » (traduction de l'auteur ; source : Algorithmwatch)

³² <https://www.consilium.europa.eu/fr/press/press-releases/2023/12/09/artificial-intelligence-act-council-and-parliament-strike-a-deal-on-the-first-worldwide-rules-for-ai/>

manque de transparence des travaux réalisés avec l'IA, utilisation contraire aux droits fondamentaux, etc.). L'objectif est également d'encourager l'innovation et faire de l'Europe un leader dans ce domaine. Il fixe l'interdiction de certaines applications portant principalement sur la catégorisation biométrique des personnes à partir de caractéristiques sensibles (race, orientation sexuelle, opinions politiques ou religieuses), la reconnaissance faciale et son extraction non ciblée pour la création de bases de données, la reconnaissance d'émotions sur le lieu de travail ou encore la manipulation du comportement humain pour contourner le libre arbitre. Il comprend également des mesures de soutien à l'innovation et notamment aux petites et moyennes entreprises (PME) afin qu'elles puissent développer des solutions d'IA sans pression excessive de la part des géants de l'industrie qui contrôlent la chaîne de valeur. Des « bacs à sable réglementaires » c'est-à-dire des espaces consacrés aux tests et des environnements réels seront mis en place par les autorités nationales pour développer et tester une IA innovante avant sa mise sur le marché. Ce texte, le premier au monde à poser un cadre réglementaire, fait d'ores et déjà débat par rapport à son impact potentiel sur l'innovation et la compétitivité de l'Europe face à des marchés moins régulés. En santé environnement, domaine encore assez confidentiel dans l'utilisation de l'IA, les freins à l'innovation seraient les mêmes que dans les autres domaines. On peut cependant espérer que le soutien promis aux PME pourrait permettre le développement de modèles plus adaptés à la donnée environnementale française. Le texte devra désormais être formellement adopté par le Parlement et le Conseil européens afin d'être intégré à la législation de l'Union européenne.

- **Le scénario 3** imagine un développement démocratique de la technologie IA. Les citoyens sont à la base des processus de contrôle nécessaires au bon développement de l'IA. Les outils sont maîtrisés par tous et leur recours est facilité grâce à l'*open source* et le développement de solutions très accessibles (*low code, no code*). On parle de systèmes d'IA hybrides capables d'être totalement transparents sur les modèles d'apprentissage automatique utilisés et la logique sous-jacente. La confiance collective en l'IA s'instaure du fait du respect du principe d'explicabilité. En santé environnement il ne serait plus alors question de boîte noire, les modèles utilisés seraient également reproductibles et applicables à d'autres données œuvrant ainsi à la mise en place d'une IA frugale. L'ensemble de la recherche en épidémiologie environnementale pourrait disposer des derniers modèles développés et participer à leur amélioration continue de par la formation facilitée des chercheurs et l'accès simplifié et peu onéreux, voire libre, à ces modèles.
- **Le scénario 4** décrit un nouvel hiver de l'IA. L'absence d'explicabilité des modèles IA pourrait conduire à un manque de confiance de plus en plus important du monde de la recherche allant jusqu'à un refus total de les utiliser. Des incidents répétés en particulier dans la sécurité des données de santé limiteraient énormément l'usage de ces outils. Par ailleurs, du fait du manque de modèles adaptés à la donnée environnementale, ceux-ci se multiplient sans synergie entre les équipes amenant à une importante consommation de ressources qui elles-mêmes se raréfient sans solution alternative proposée. Les enjeux sociaux et environnementaux sont jugés prioritaires par les États, l'IA amorce donc son déclin.

7. CONCLUSION

Ce document d'introduction à l'intelligence artificielle, au-delà du domaine spécifique de l'épidémiologie environnementale, doit permettre à ceux qui le souhaitent de projeter des besoins, des utilisations possibles dans des problématiques de santé publique diverses.

Le monde de l'IA est ouvert à tous et le connaître doit permettre de se positionner en tant que scientifiques, techniciens des données, biostatisticiens, géomaticiens, ingénieurs, etc.

Il est aujourd'hui nécessaire de s'emparer de ce champ d'application et ainsi participer humblement à l'ouverture progressive de la « boîte noire » comme le recommande le rapport Villani de 2018.

Pour autant, la présente recherche d'informations sur le domaine spécifique de l'épidémiologie environnementale montre que les possibilités sont grandes mais que les travaux sont encore peu nombreux à mobiliser de l'IA pour « augmenter » les études en santé publique et notamment en santé environnementale. Bien entendu les évolutions dans le domaine étant extrêmement rapides, il est probable que des travaux se développent fortement dans les prochaines années, pour parvenir à des applications plus robustes.

C'est dans ce cadre que Santé publique France organisera un séminaire scientifique à l'automne 2024 avec un panel d'experts de l'IA abordant les utilisations actuelles et futures de l'IA pour les études en santé environnement. L'objectif est de pouvoir échanger des informations sur les pratiques, les méthodes mais aussi d'aborder les questions éthiques et les freins potentiels au développement des méthodes dans ce domaine.

8. RÉFÉRENCES

- [1] Helbert E. Comment réussir le déploiement de l'IA en entreprise? ActulA. 2023.
- [2] La stratégie IA en France. Dossier de presse. 2017.
- [3] Commission européenne, loi européenne sur l'Intelligence Artificielle (AI Act). 2023. Disponible: <https://artificialintelligenceact.eu/fr/>
- [4] Français G. France 2030. Disponible: <https://www.gouvernement.fr/france-2030>
- [5] Ministères transition écologique Cohésion des territoires Mer. Feuille de route : intelligence artificielle et transition écologique du pôle ministériel 2021-2024. Paris : Ministères transition écologique Cohésion des territoires Mer; 2021. 26 p.
- [6] IGN. Feuille de route Intelligence artificielle 2022-2024. 2022.
- [7] United Nations Environment Programme. Définir le rôle de l'intelligence artificielle dans la prévision, l'atténuation et l'adaptation aux impacts du changement climatique. Kenya : UNEP; 2020. 21 p.
- [8] Carolin C. SIG et intelligence artificielle : quels développements et quel futur ? Géomatique Expert. 2017(118).
- [9] VoPham T, Hart JE, Laden F, Chiang YY. Emerging trends in geospatial artificial intelligence (geoAI): potential applications for environmental epidemiology. Environ Health. 2018;17(1):40.
- [10] Bouzeghoub MM, R. Les Big Data à découvert. CNRS Editions éd. 2017.
- [11] Alyssa Imbert NV. Décrire, prendre en compte, imputer et évaluer les valeurs manquantes dans les études statistiques : une revue des approches existantes. Journal de la Societe Française de Statistique. 2018;159 :1-55.
- [12] Rodriguez-Gonzalez A, Zanin M, Menasalvas-Ruiz E. Public Health and Epidemiology Informatics: Can Artificial Intelligence Help Future Global Challenges? An Overview of Antimicrobial Resistance and Impact of Climate Change in Disease Epidemiology. Yearb Med Inform. 2019;28(1):224-31.
- [13] Wambaugh JF, Rager JE. Exposure forecasting – ExpoCast – for data-poor chemicals in commerce and the environment. Journal of Exposure Science & Environmental Epidemiology. 2022;32(6):783-93.
- [14] Angelo A, Sebastien O. Application de méthodes d'intelligence artificielle dans la détection des parcelles agricoles afin de produire des indicateurs spatialisés d'exposition indirecte aux pesticides. Environnement, Risques & Santé. 2023;22(1):27-35.
- [15] Schoene AM, Basinas I, van Tongeren M, Ananiadou S. A Narrative Literature Review of Natural Language Processing Applied to the Occupational Exposome. Int J Environ Res Public Health. 2022;19(14).
- [16] Crépet A, Tressou J, Vanessa G, Béchaux C, Pierlot S, Héraud F, et al. Identification of the main pesticide residue mixtures to which the French population is exposed. Environmental research. 2013;126.
- [17] Carvaillo J-C, Barouki R, Coumoul X, Audouze K. Linking Bisphenol S to Adverse Outcome Pathways Using a Combined Text Mining and Systems Biology Approach. Environmental Health Perspectives. 2019;127(4):047005.

- [18] Wen Y, Zhang S, Wang Y, Yang J, He L, Wu Y, *et al.* Dynamic Traffic Data in Machine-Learning Air Quality Mapping Improves Environmental Justice Assessment. *Environmental Science and Technology*. 2023.
- [19] Arnault Pachot CP. *Intelligence Artificielle & Protection de l'Environnement: le paradoxe d'une technologie énergivore au service des défis écologiques de demain*. 2022.
- [20] Yang L, Driscoll J, Sarigai S, Wu Q, Lippitt CD, Morgan M. Towards Synoptic Water Monitoring Systems: A Review of AI Methods for Automating Water Body Detection and Water Quality Monitoring Using Remote Sensing. *Sensors (Basel)*. 2022;22(6).
- [21] Julia L. *L'intelligence artificielle n'existe pas*. . First éd. 2019.
Disponible: <https://www.cairn.info/revue-paysan-et-societe-2021-1-page-49.htm> ;
https://www.cairn.info/load_pdf.php?ID_ARTICLE=PES_385_0049
- [22] Wen D, Khan SM, Ji Xu A, Ibrahim H, Smith L, Caballero J, *et al.* Characteristics of publicly available skin cancer image datasets: a systematic review. *The Lancet Digital Health*. 2022;4(1):e64-e74.
- [23] Laurent O. Dans la boîte noire de l'IA. Qu'est ce que l'Intelligence Artificielle peut apporter à la recherche en santé environnement ? *ERS*. 2020;19.
- [24] Saliou M. *Technoféminisme, Comment le numérique aggrave les inégalités*. Grasset éd. 2023.
- [25] Stratégie F. *Intelligence artificielle et travail*. 2018.
- [26] Yacine LASJ. Power Hungry Processing: Watts Driving the Cost of AI Deployment? *Arxiv*. 2023.
- [27] Lefèvre L, Ligozat A-L, Trystram D, Bouveret S, Bugeau A, Combaz J, *et al.* Proposition de document de cadrage Évaluation environnementale de projets impliquant des méthodes d'IA. *EcolInfo*; 2022. 1-8 p. Disponible: <https://hal.science/hal-03853135>
- [28] Cnil. *Données, empreinte et libertés*. 2023.
- [29] INRS. *L'intelligence artificielle au service de la santé et sécurité au travail, enjeux et perspectives à l'horizon 2035* ; 2022.

ANNEXES

Annexe 1 : recherche bibliographique

Pour cette synthèse, la recherche bibliographique a porté sur la base Medline (interface Pubmed) sur les cinq dernières années (2018-2023), sur des études écrites en anglais et en français (date d'interrogation de Pubmed : 17 mars 2023). En complément, une recherche de références dans la littérature grise (rapports, synthèses de congrès, numéros thématiques de revues) a été menée selon les mêmes critères.

Les études sur l'intelligence artificielle dans les domaines de l'environnement, l'épidémiologie ou les deux peuvent être de différente nature : études méthodologiques, portant sur l'apport de telle ou telle technologie ou méthode (par exemple des études comparant des modèles) mais aussi des synthèses plus applicatives sur l'apport de ces technologies pour la surveillance environnementale.

Pour délimiter le périmètre de notre recherche, nous avons choisi de restreindre la recherche documentaire aux revues de littérature ou rapports d'organismes portant sur les apports de l'IA (applications, méthodes) pour la santé environnement/la surveillance épidémiologique. Certaines sous-thématiques nous sont apparues prioritaires pour notre question :

- Le traitement massif des données ;
- La modélisation : croisement de données, reconstitution de données historiques à partir d'imagerie, système d'information géographique.

Les articles traitant de l'IA pour la recherche médicale, tels que ceux portant sur l'aide automatisée au diagnostic ont été exclus, s'agissant d'un volet de la littérature extrêmement spécifique et vaste.

Pour la recherche de littérature scientifique, 128 revues de littérature ont été repérées, parmi lesquelles 46 ont été retenues pour leur adéquation avec nos objectifs. Pour la littérature grise, 64 références ont été retenues : 34 articles ou numéros thématiques, 5 synthèses de conférences, 5 initiatives innovantes type *data challenge* et 22 rapports, feuilles de routes, livres blancs. Les requêtes utilisées pour cette recherche documentaires sont disponibles sur demande à l'adresse documentation@santepubliquefrance.fr

Annexe 2 : liste des organismes interrogés

- Cnil : Commissions nationale de l'informatique et des libertés
- HDH : *Health Data Hub*
- IGN : Institut géographique national
- OFB : Observatoire français de la biodiversité
- INRS : Institut national de recherche et de sécurité
- GD4H : *Green Data for Health* (Ecolab)
- Pôle IA du Commissariat général au développement durable du ministère de la transition écologique (Ecolab)
- Ineris : Institut national de l'environnement industriel et des risques
- Inria : Institut national de recherche en informatique et en automatique
- Open Data France
- Medes (Cnes) : Institut de médecine et physiologie spatiales
- ECMWF (Copernicus) : *European Centre for Medium-Range Weather Forecasts*

Annexe 3 : liste non exhaustive de conférences en ligne

- Collège de France :
 - [Bâtir l'Europe de l'intelligence artificielle : par quels artifices ? par quelles intelligences ? | Collège de France \(college-de-france.fr\) – Cédric Villani](#)
 - [L'apprentissage profond : une révolution en intelligence artificielle | Collège de France \(college-de-france.fr\) – Yann LeCun](#)
- Mooc IA gratuits :
 - *Open Class Room* et l'Institut Montaigne : https://openclassrooms.com/fr/courses/6417031-objectif-ia-initiez-vous-a-lintelligence-artificielle?utm_source=pole-emploi&utm_medium=email&utm_campaign=email_students_montaigne&utm_content=objectif-ia
 - L'Inria <https://www.fun-mooc.fr/fr/cours/lintelligence-artificielle-avec-intelligence/>
- Conférence de Luc Julia (créateur de Siri) « L'intelligence artificielle n'existe pas » : https://www.youtube.com/watch?v=yuDBSbng_8o
- « Numérique, intelligence artificielle et développement durable en santé », Faculté de Santé Sorbonne Université : <https://www.youtube.com/watch?v=8Xii6eyNvNc>
- Arte Info plus – KREATUR - « L'IA est-elle sexiste ? » : <https://www.arte.tv/fr/videos/113219-006-A/l-intelligence-artificielle-est-elle-sexiste/>
- Podcast « L'ère des algorithmes » de l'Université Paris Dauphine : <https://shows.acast.com/ex-machina>