



HAL
open science

L'intelligence artificielle des textes. Présentation

Damon Mayaffre, Laurent Vanni

► **To cite this version:**

Damon Mayaffre, Laurent Vanni. L'intelligence artificielle des textes. Présentation. L'intelligence artificielle des textes. Des algorithmes à l'interprétation, 15, Honoré Champion, pp.9-14, 2021, Lettres numériques, 978-2-7453-5640-6. hal-03344917

HAL Id: hal-03344917

<https://hal.science/hal-03344917>

Submitted on 16 Sep 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Encore trop méconnue des Humanités, l'Intelligence artificielle ébranle les pratiques scientifiques et les certitudes. Appliquée aux textes – corpus littéraires, corpus politiques, textes français, textes latins – elle dévoile ses avantages et ses limites.

Sous conditions, la digitalisation de masse de la culture vivifie les sciences et les arts du texte.

La sémantique de corpus tire bénéfice des données massives à condition qu'elles soient circonscrites. L'algorithmique et la puissance de calcul excellent sous condition de se départir des prétentions probatoires au profit d'autres objectifs : explorer, interroger, décrire.

En passant de la fixité du *folio* à la dynamique de l'écran, les textes s'ouvrent sur des parcours de lecture inédits, appuyés sur de nouveaux objets tels que les *motifs*.

Seulement, les programmes de *deep learning* décrits dans cet ouvrage – avant tout, les modèles convolutionnels – ne peuvent rester plus longtemps une boîte noire. Les couches cachées des réseaux de neurones artificiels n'ont d'intérêt que lorsqu'elles nourrissent les parcours interprétatifs et caractérisent des *passages*.

C'est donc la profondeur herméneutique de l'approche numérique qui est exploitée ici. En donnant une autre représentation des corpus, en révélant les structures sous-jacentes des textes, en faisant émerger de nouveaux observables, l'Intelligence artificielle épaulé l'herméneute et le philologue.

Dans cet ouvrage, la machine *apprend* du corpus, *classe* les textes et *décrit* les grandeurs textuelles pertinentes ; reste aux chercheurs, linguistes, informaticiens, historiens, statisticiens le plaisir de comprendre.

Lettres numériques N° 15



LN
15

L'intelligence artificielle des textes

L'intelligence artificielle des textes

Des algorithmes à l'interprétation

Sous la direction de Damon Mayaffre et Laurent Vanni



HONORÉ CHAMPION
PARIS

PRÉSENTATION

L'Intelligence artificielle (IA) reste aujourd'hui une boîte noire tout en étant devenue un mot valise dont la science et la politique – et *a fortiori* nos politiques scientifiques – se sont emparées pour dessiner les projets d'avenir et la société de demain. Des laboratoires au ministère, des villages français au vaste monde, aucun échelon de nos organisations ne semble devoir échapper aujourd'hui, de fait ou en principe, à la révolution des réseaux de neurones artificiels, aux architectures profondes, aux couches cachées, au *deep learning*. Et si la révolution n'est pas nouvelle, puisque Turing aurait aujourd'hui 100 ans et que le Perceptron remonte à l'après-guerre, la puissance actuelle inédite des cartes graphiques et l'ère des *big data*, la rendent effectivement désormais incontournable.

Éclairer autant que possible la boîte noire et entrebâiller le mot valise pour en expliquer la plus-value heuristique en matière de linguistique textuelle, telle est la prétention de cet ouvrage qui réunit linguistes et informaticiens dans une interdisciplinarité non feinte.

C'est en effet par une collaboration effective entre amoureux des textes et amoureux du code que l'on entend procéder : l'expertise linguistique du texte – textes littéraires, textes politiques ; textes français, textes latins – doit servir à comprendre l'algorithme, son fonctionnement, sa portée, ses faiblesses ; l'expertise informatique de l'algorithme – les modèles convolutionnels, la déconvolution, les RNN – doit servir à comprendre le texte, ses structures sous-jacentes, sa complexité sémantique, ses nouveaux observables qu'une lecture naturelle a jusqu'ici pressentis sans réussir à objectiver.

Les chercheurs en SHS qui ont affaire au texte ont tout à gagner de l'Intelligence artificielle à condition d'échapper aux deux écueils symétriques dont la bibliographie a désormais bien balisé les dangers : l'anthropomorphisation de la machine d'un côté, la mécanisation de l'humain de l'autre.

Les contributeurs de cet ouvrage font ainsi un usage métaphorique du mot *intelligence* lorsqu'on l'applique au robot et qu'on le qualifie par *artificielle*. Prêter à l'ordinateur une intuition, une imagination, une créativité – autant d'éléments qui sont partie prenante de l'intelligence humaine – paraît non seulement une thèse mais une hypothèse exorbitante. Cet ouvrage montrera certes, par exemple, les avantages d'une approche néo-connexionniste des textes, pour établir des parcours interprétatifs innovants et complexes, mais le connexionnisme informatique des *networks* aux réseaux de Kohonen, et par delà eux le *programme*, ne sauraient prétendre à l'intelligence humaine, sauf à prêter un esprit aux choses et une âme à l'inanimé.

Pas plus, nous ne pensons scientifiquement possible, ni philosophiquement souhaitable, la mécanisation de l'humain. En l'occurrence, les arts et sciences du texte que l'humanité pratique depuis plusieurs millénaires, et qui établissent dans une posture herméneutique fondamentale l'acte interprétatif comme condition sémantique, ne sauraient envisager que la quête de sens soit entièrement automatisable ; le sème ne peut se laisser réduire au *token* et le raisonnement à l'algorithme. Le trans-humanisme, parfois assumé en post-humanisme, ne saurait trancher la question de la liberté, de l'indécision ou de l'interprétable, sauf à réduire l'homme à la binarité et prêter un déterminisme mécanique à l'indéterminé.

Est-ce à dire, pour autant, que l'Intelligence artificielle est une impasse de la science et de la société, ou ici de la linguistique textuelle et de la sémantique de corpus, là où ailleurs on lui prête le pouvoir vital de décrypter nos radiographies et celui pratique de conduire nos voitures ?

Tout au contraire, cet ouvrage veut montrer les possibilités nouvelles que l'Intelligence artificielle offre aux chercheurs en analyse de corpus en donnant à voir des représentations du texte originales, en objectivant des parcours de lecture heuristiques, en faisant émerger de nouveaux observables linguistiques. En des mots que Bruno Bachimont avait déjà pressentis dans les années 1990 : il ne s'agit certes pas de croire en une machine qui pense mais d'espérer en une machine qui donne à penser [Bachimont 1996].

En effet, si l'outil a prolongé la main sans l'amputer, l'ordinateur peut prolonger l'intelligence – ici l'intelligence des textes – sans l'abolir. Le texte n'est pas un objet naturel mais un artefact culturel qui n'a, dès lors, rien à craindre de l'artifice d'une intelligence numérique

si elle se fixe comme objectif de l’embrasser et pour mieux dire de le représenter : le *présenter à nouveau* à notre savoir-faire analytique ou herméneutique, linguistique ou littéraire ; *en donner une représentation* particulière, réticulaire par exemple plus que linéaire, rhizomique plus qu’orientée, statistique plus que statique.

Car le point commun des contributions qui vont suivre est de poser de manière critique les vices et les vertus de la représentation numérique et du traitement informatique de textes qu’on a longtemps confondus avec leurs supports-papiers traditionnels, et injustement soustraits à la philologie.

Le projet, pour qu’il soit fécond, prend la forme de plusieurs tâches concrètes. Trois seront centrales dans cet ouvrage pour constituer, dans leur succession et complémentarité, un programme de recherche dont nous prétendons ici poser les premières bases : l’apprentissage, la prédiction, la description.

— Apprentissage. Si le numérique modifie déjà depuis plusieurs lustres notre rapport à l’empirie [Rastier 2011], en multipliant les corpus disponibles et en favorisant la manipulation originale (création d’index de lemmes, étiquetage morphosyntaxique des formes, repérage statistiques des cooccurrences, détection de motifs lexico-grammaticaux complexes, etc.), l’Intelligence artificielle surenchérit sur la révolution épistémologique en cours en introduisant la notion d’apprentissage. Sur un jeu de données textuelles identifiées et nécessairement important – appelé le plus souvent *corpus d’apprentissage* – l’algorithme apprend les creux et les reliefs, les normes, les saillances ou les tournures propres à un auteur, à un genre ou à une époque. Aujourd’hui, cet apprentissage empirique – puisqu’il ne repose pas sur des règles linguistiques formelles mais sur des normes endogènes au corpus – est suffisamment performant pour permettre de générer automatiquement des ersatz de textes « à la manière de... » comme il permet de générer automatiquement des ersatz de morceaux de musique « à la manière de Bach ». Et quand bien même les résultats seraient jugés décevants, n’y a-t-il pas moyen, pour nous, d’apprendre de cet apprentissage machinal ?

— Prédiction. Quand la machine a bien appris, l’Intelligence artificielle peut faire valoir des tâches de *prédiction*, que nous connaissons en Analyse statistique de données textuelles sous la forme de classifications. L’apprentissage algorithmique est poussé autant que nécessaire jusqu’à ce qu’il soit estimé satisfaisant : au-delà du corpus

d'apprentissage, un corpus de test composé de textes connus mais un instant anonymisés est constitué pour évaluer les performances, et l'on arrête l'apprentissage aussitôt que le taux d'attribution de ces textes-tests à leur auteur, à leur genre ou à leur époque s'approche des 100%. Dès lors, le corpus d'étude peut être travaillé, et les textes inconnus peuvent être discriminés et classés avec quelques garanties. Ici, les textes d'Émile Ajar seront bien attribués à Romain Gary, et Molière retrouvera la paternité de ses pièces un instant contestée. Avec une précision remarquable, la grande littérature française du XX^e siècle de Proust à Aragon, de Camus à Yourcenar est ordonnée [Brunet 2016]. Cette prédiction de l'IA, déterminante dans la génétique des textes mais aussi dans la modélisation de l'intertextualité – ces discours qui traversent d'autres discours, évidents dans le cadre de reprises plagiaires –, apparaît d'autant plus convaincante qu'elle est éclairée par la statistique textuelle benzécienne [Lebart, Pincemin, Poudat 2019].

— Description. Enfin, si l'IA sait apprendre et prédire, nous lui demandons dans cet ouvrage de décrire. Il s'agit de l'enjeu scientifique principal pour les SHS mais aussi pour l'Informatique : identifier, par *déconvolution* ou tout autre procédé, dans les couches cachées du traitement et les profondeurs du texte, les observables linguistiques sur lesquels la machine s'est appuyée pour reconnaître un auteur ou établir un classement. Ces observables identifiés, à des paliers de descriptions linguistiques croisés (lexique, grammaire, syntaxe), comme le sont les *motifs* [Longrée et Mellet 2013] seront alors considérés comme des *passages* [Rastier 2007] au fondement de notre interprétation, c'est-à-dire de notre intelligence des textes.

*

Cet ouvrage est composé de cinq chapitres dans l'objectif de rendre compte de la pluralité des points de vue – plus au moins critiques sur l'Intelligence artificielle selon les auteurs – et de la pluralité des approches – épistémologique, technique, linguistique, informatique, historique.

Le premier chapitre rédigé par Laurent Vanni et Frédéric Precioso, *Deep learning et description des textes. Architecture méthodologique*, présente par le menu les tenants et les aboutissants du traitement informatique des textes par le *deep learning*. L'attention est avant tout accordée aux modèles convolutionnels (CNN). Les informaticiens répondent alors à l'exigence descriptive du modèle pour les linguistes :

ouvrir la boîte noire algorithmique en éclairant les observables linguistiques responsables de l'apprentissage et de la classification.

Le deuxième chapitre écrit par Étienne Brunet, Ludovic Lebart et Laurent Vanni, *Littérature et Intelligence artificielle* revient sur la génétique des textes, leurs origines parfois remises en cause, ou plus précisément le calcul de la distance intertextuelle dans le cadre de corpus contrastifs. La contribution discute pas à pas l'apport de l'Intelligence artificielle notamment vis-à-vis de méthodes éprouvées comme l'ACP ou l'AFC. Et si, définitivement, Molière avait bien écrit *Amphitryon* ?

Le troisième chapitre proposé par Magali Guaresi et Damon Mayaffre, *Intelligence artificielle et discours politique. Quelles plus-values interprétatives ?* revient sur la tradition de l'analyse du discours politique assistée par ordinateur. Sur la foi de la matérialité textuelle du corpus, l'Intelligence artificielle propose aujourd'hui une objectivation du discours de gauche et de droite, ou de l'idéologie d'Emmanuel Macron, qui n'est pas sans raviver les espoirs méthodologiques de l'AD originelle. Seulement, pour que ces espoirs ne redeviennent pas chimères, la démarche doit prétendre à l'heuristique et non au probatoire.

Le quatrième chapitre conçu par Dominique Longrée, *Motifs linguistiques et deep learning. Vers une détection automatique de nouveaux observables ?* applique le *deep learning* au latin classique et aux corpus étiquetés du L.A.S.L.A. La détection automatique des motifs grammatico-lexicaux reste un enjeu pour le *text mining*, que les modèles convolutionnels pourraient relever quelle que soit la langue considérée.

Le cinquième chapitre, enfin, proposé par François Rastier, *Data vs Corpora*, sonne comme une mise en garde. La machine n'a pas vocation à remplacer le lecteur, et les données se substituer aux corpus. Si l'Intelligence artificielle semble contribuer à la sémantique de corpus, en considérant le cotexte et en « apprenant » des textes rassemblés, elle ne peut réaliser l'acte subversif de l'interprétation linguistique humaine, qui opère par différenciation et même par exclusion intelligente d'information.

Pour les auteurs, Damon Mayaffre et Laurent Vanni

TABLE DES MATIÈRES

<i>IN MEMORIAM</i>	7
--------------------------	---

PRÉSENTATION	9
---------------------------	---

Laurent VANNI et Frédéric PRECIOSO

DEEP LEARNING ET DESCRIPTION DES TEXTES ARCHITECTURE MÉTHODOLOGIQUE	15
--	----

Introduction.....	15
-------------------	----

1. Prérequis	18
--------------------	----

1.1. Le traitement des données.....	18
-------------------------------------	----

1.2. L'architecture du modèle	21
-------------------------------------	----

1.3. Les (hyper)paramètres du réseau	25
--	----

2. <i>Embedding</i> : représentation des mots	32
---	----

3. Convolution : abstraction des données	39
--	----

3.1. <i>Text Deconvolution Saliency</i> (TDS)	41
---	----

3.2. Pondération du TDS.....	51
------------------------------	----

4. Analyse multi-channels : formes, lemmes, étiquettes morpho- syntaxiques	55
---	----

4.1. Modèle	55
-------------------	----

4.2. Exemple 1	58
----------------------	----

4.3. Exemple 2	59
----------------------	----

4.4. Remarques	61
----------------------	----

5. Les entrées et sorties du réseau.....	62
--	----

5.1. Détection des passages-clés	62
--	----

5.2. Filtrage des données.....	66
--------------------------------	----

Conclusion	70
------------------	----

ÉTIENNE BRUNET, LUDOVIC LEBART ET LAURENT VANNI

LITTÉRATURE ET INTELLIGENCE ARTIFICIELLE	73
Introduction.....	73
1. Le roman au XX ^e siècle. Reconnaître que deux textes sont d'un même auteur.....	75
1.1. Une expérience en double aveugle	75
1.2. Les mesures classiques de la distance intertextuelle.....	77
1.3. Méthodes fondées sur le TLE	80
1.4. La solution du <i>deep learning</i>	89
2. Le théâtre classique.....	94
2.1. Les basses fréquences	96
2.2. Les hautes fréquences	101
2.3. Neutralisation du genre	103
2.4. Le <i>Deep learning</i>	104
3. Essai d'explicitation et d'intégration	107
3.1. Imitation de la procédure du <i>deep learning</i> : les triplets....	108
3.2. Réduction et intégration des deux approches, <i>procédurale et algorithmique</i>	114
3.3. La déconvolution	119
3.4. Analyse non-supervisée des « pages » (fragments de 50 lignes consécutives)	125
Conclusion (Éléments pour une...).....	128

MAGALI GUARESI ET DAMON MAYAFFRE

INTELLIGENCE ARTIFICIELLE ET DISCOURS POLITIQUE.....	131
Introduction.....	131
1. <i>Deep learning</i> et textes politiques profonds. Cadrage	135
1.1. Convolution, co-texte, cooccurrences	136
1.2. Réseaux, couches cachées, complexité.....	137
1.3. Du sens en contexte	138
2. Vers une description des discours gauche / droite. Les professions de foi des députés sous la V ^e République	141
2.1. Gauche et droite décrites par <i>deep learning</i> (1958-2017)..	143
2.2. Situer les professions de foi macronistes sur l'échiquier politique	156

3. Vers une objectivation de l'intertexte. Les emprunts de Macron à ses prédécesseurs à l'élysée	162
3.1. Proposition méthodologique : un usage original du <i>deep learning</i>	165
3.2. Résultats : intertexte pluriel et discours patchwork	167
Conclusion	180

DOMINIQUE LONGRÉE

MOTIFS TEXTUELS ET DEEP LEARNING : VERS UNE DÉTECTION AUTOMATIQUE DE NOUVEAUX OBSERVABLES LINGUISTIQUES ?	183
Introduction.....	183
1. Motifs textuels et unités phraséologiques. Critères d'identification	186
2. Motifs textuels et détections semi-supervisées	189
2.1. Motifs et Hyperbase Web Edition	189
2.2. SDMC-Sequential Data Mining under Constraints	190
3. Motifs textuels et hyperdeep.....	191
4. Hyperdeep : un outil heuristique pour le latin ?.....	200

FRANÇOIS RASTIER

DATA VS CORPORA	203
1. Un bref retour sur le mythe de l'IA	203
2. Rien ne nous est donné.....	208
3. Traitements et maltraitements	219
4. Limites de l'émergentisme.....	222
5. Pourquoi les corpus sont indispensables.....	230
6. Le problème des passages et l'expérimentation connexionniste.	238
7. La <i>Data Science</i> contre les sciences de la culture.....	244
BIBLIOGRAPHIE	247
INDEX DES PERSONNES	257
TABLE DES MATIÈRES	263