

Usages de l'IA à la MSHB

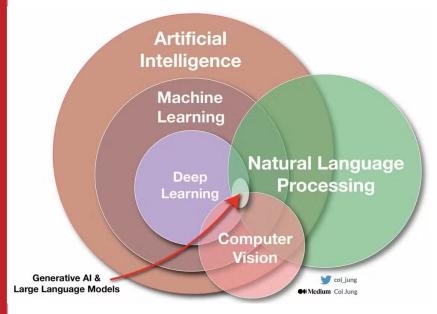
Plateforme Humanités numériques & Plateforme universitaire de données - Bretagne (PUD-B)

Chloé Choquet – Ingénieure en traitement et analyse de corpus textuels – chloe.choquet@mshb.fr

Paul Pinard – Ingénieur-statisticien – <u>paul.pinard@mshb.fr</u>



Intelligence artificielle



https://www.orsys.fr/orsys-lemag/machine-learning-deep-learning-ia-differences/

Intelligence artificielle : technique qui permet à un ordinateur de raisonner, planifier, créer et accomplir des tâches qui nécessitent normalement l'intelligence humaine

Machine Learning (apprentissage automatique): l'ordinateur apprend à partir de données pour faire des prédictions

Deep Learning (apprentissage profond) : machine learning avec des réseaux de neurones pour gérer des données complexes

Computer Vision (vision par ordinateur): l'ordinateur apprend à comprendre des images ou des vidéos

Traitement automatique des langues (NLP/TAL): techniques qui permettent à un ordinateur de comprendre, analyser et produire du texte ou de la parole, en combinant informatique, linguistique et intelligence artificielle

LLM (grand modèle de langage) : modèle d'IA capable de comprendre et produire du texte à grande échelle

IA générative : type d'IA qui crée du contenu nouveau, comme du texte, des images ou de la musique





Traitement automatique des langues & Python

Bibliothèques Python : boîtes à outils de fonctions prêtes à l'emploi pour programmer plus vite

Permettent de : tokeniser, lemmatiser, raciniser, attribuer des étiquettes grammaticales, reconnaître des entités nommées, analyser des sentiments, générer des textes, résumer des textes, traduire automatiquement, etc.

- → règles statistiques simples rapide (pour petites tâches), capacités basiques
- → modèles probabilistes (maching learning) très rapide, bonnes capacités
- → LLM (deep learning) lent (pour grandes tâches), très hautes capacités









```
# Analyse d'une critique
      import spacy
     from transformers import pipeline
     # 1 : Analyse linguistique avec spaCy
     nlp = spacy.load("fr core news md") # modèle français de taille moyenne
     texte = "Une création sans âme ! Dès les premières minutes, ca sonne faux et j'ai su tout de suite que ca n'allait pas le faire. Les décors alternent, sans la
     moindre harmonie, entre fonds verts criards. CGI déqueulasses — je préfère ne pas insister sur les loups (apparaissant lors de l'épisode avec le vieil aveugle !).
     avec des mouvements saccadés dignes d'un jeu de PlayStation 2 - et quelques extérieurs réels, créant ainsi un patchwork visuel sans cohérence. Le bateau pris dans
     les glaces du prologue - avec ses effets numériques soulignant que Hollywood a fortement régressé dans ce domaine depuis de trop nombreuses années - n'aide pas à
     entrer dans le film . Pas plus que cette étrange idée de faire parler Victor jeune et sa mère en français - une langue qu'aucun des deux acteurs les interprétant
     ne maîtrise un tant soit peu. La musique ne contribue pas non plus à élever le niveau. Elle tente à plusieurs reprises d'imiter Danny Elfman, mais sans en avoir le
     talent, ni l'originalité, ni la fantaisie, ni la puissance mélodique. Résultat : un accompagnement fade, sans ampleur, qui se contente de souligner bien platement
      sans jamais sublimer. Qui a composé cette merde ? Ah qui, ce nullos ultra-méga-surestimé d'Alexandre Desplat - m'étonne pas. Une IA aurait pris sa place, elle
     aurait fait beaucoup mieux." # extrait d'une critique de Frankenstein : https://www.senscritique.com/film/frankenstein/critique/331909408
 9
10
     doc = nlp(texte)
11
12
     print("\nEntités nommées (spaCy) : ")
13
     for ent in doc.ents:
14
         print(f"{ent.text} → {ent.label }")
15
16
     # 2 : Analyse de sentiments avec transformers
17
     analyse sentiment = pipeline(
18
         "sentiment-analysis",
         model="nlptown/bert-base-multilingual-uncased-sentiment" # modèle entraîné sur des critiques de produits et films multilingues
19
20
21
22
     resultat = analyse_sentiment(texte)
23
24
     print("\nAnalyse de sentiments (transformers) : ")
     print(resultat)
```

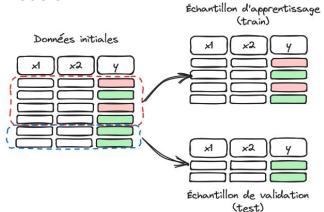
Entités nommées (spaCy):
PlayStation 2 → MISC
Hollywood → LOC
Victor → PER
Danny Elfman → PER
Ah oui → MISC
Alexandre Desplat → PER

Analyse de sentiments (transformers): [{'label': '2 stars', 'score': 0.6387473344802856}]



L'importance des données / Textes :

- Ce que l'on rentre en initialisation impacte grandement le résultat final!
 - Des données qui ne sont pas correctement mises en forme donneront de mauvais résultats
- Méthodes / Codes relativement identique
- Importance surtout sur la structuration des données
- Technique courante en apprentissage automatique :
 - 3 de la base de données servira pour faire apprendre un modèle
 - 1/3 de la base sera pour connaître notre taux d'erreurs du modèle





Exemple 1 : Données brutes / Données retouchées

Données Brutes

- Possède des valeurs aberrantes sur ces variables numériques (Âge, Revenu,...)
- Possède des "typos" sur ces variables qualitatives (CSP, Situation familiale,...)

- Essaie de regrouper des individus en fonction :

- De leurs tailles, leurs sexes, leurs poids, leurs taille de cheveux, leurs situations familiales...

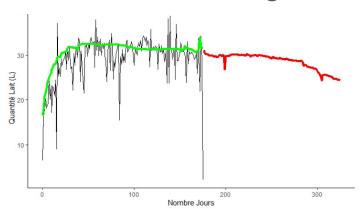
Résultats :

	Données Brutes	Données Retouchées
Précision classification (en %)	~ 90 %	~ 94 %



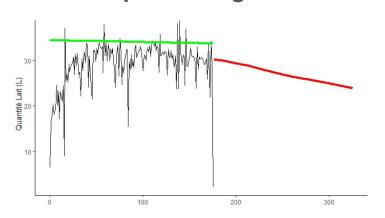
Exemple 2 : Même données - Résultats différents

Machine Learning:



Méthodes	Taux d'erreurs (en moyenne)	
Machine Learning 1	7.96 L	
Machine Learning 2	5.40 L	
Deep Learning	5.69 L	

Deep Learning:

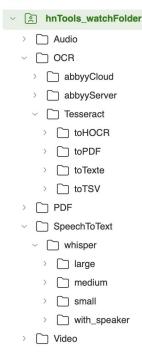


TEMPS	Machine Learning (méthode 1)	Machine Learning (méthode 2)	Deep Learning
Moyenne	2 min 33 sec	10 min 49 sec	1 min 16 sec
Ecart-type	43 sec	16 min 44 sec	59 sec
Maximum	4 min 57 sec	1H 31 min 52 sec	5 min



IR* Huma-Num

H ShareDocs



Outils de traitement de Sharedocs :

- Conversion de la parole en texte (Speech to Text) : Whisper
- Reconnaissance de caractères (OCR) : Abbyy, Tesseract

Ferme de calcul:

L'IR* Huma-num mets à disposition un serveur de LLM dédié sous Ollama.

Ollama est un outil qui permet de faire fonctionner des modèles de langage volumineux (LLM) localement sans utiliser de service de cloud.

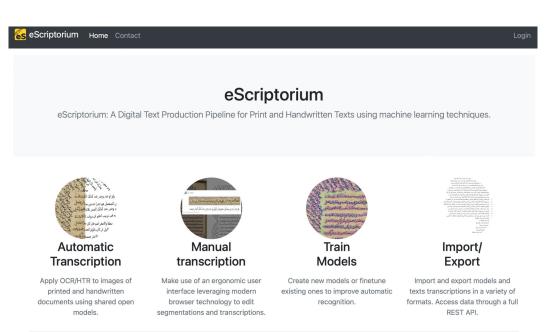


Reconnaissance des textes manuscrits

eScriptorium : plateforme permettant la transcription, l'annotation, la traduction et la publication de documents historiques

Possibilité d'entraîner un modèle

📧 justine.chainiau@mshb.fr (Chargée d'édition de corpus numériques)















Exemple où l'IA n'est pas utile

NVivo : logiciel d'analyse thématique de contenu

Assistant IA pour l'encodage automatique de thèmes : pas encore au point

