

# Apprentissage sur des données de justice

Jean-François COUCHOT  
couchot@femto-st.fr

## 1 Présentation

Le Traitement Automatique des Langues (TAL) repose sur la linguistique et l'intelligence artificielle. Ses principales applications sont la génération automatique de textes (réponses à des questions, résumés, chatbot), l'extraction d'informations (classification, recommandations),... Il commence à peine à être appliqué à des données de justice, domaine critique par excellence.

C'est cependant une volonté de l'état de fournir des données textuelles en entrée de ce genre d'outil avec pour preuve le calendrier de publication des données de justice<sup>1</sup> : chaque année 35000 décisions administratives ou judiciaires sont publiées sur LegiFrance et toutes les décisions de la cours de d'Etat et de cassation seront en ligne courant septembre 2021.

Dans ces décisions mises en ligne, les noms et prénoms des personnes jugées sont systématiquement réduits à leurs initiales. Cela permet d'éviter une association immédiate d'un couple nom prénom (saisi dans un moteur de recherche) à une décision de justice. Cependant cela ne permet pas de protéger la vie privée d'une personne dont le nom apparaîtrait dans une de ces décisions : en y adjoignant des connaissances supplémentaires, rien ne garantit que l'on ne pourra pas ré-identifier cette personne.

Ce projet vise à étudier une partie du vaste domaine d'application du TAL aux décisions de justice en se concentrant sur la protection de la vie privée dans une recherche de similarité entre décisions.

## 2 Plan

Ci dessous une succession d'étapes qui pourraient être suivies.

1. S'ouvrir à la thématique d'informatisation de la décision dans le domaine de la justice [SD19] ;
2. Se documenter sur la protection de la vie privée dans la publication des décisions de justice<sup>2 3</sup>[ABG20] ;
3. Mettre en place un algorithme de détection de similarité par TAL sur des données publiées officiellement [MGGM21, BGPG20] ;
4. Proposer des méthodes de dé-identification de documents de justice.

## Références

- [ABG20] Tristan Allard, Louis Béziaud, and Sébastien Gambis. Online publication of court records : circumventing the privacy-transparency trade-off. *CoRR*, abs/2007.01688, 2020.
- [BGPG20] Paheli Bhattacharya, Kripabandhu Ghosh, Arindam Pal, and Saptarshi Ghosh. Methods for computing legal document similarity : A comparative study. *arXiv preprint arXiv :2004.12307*, 2020.
- [MGGM21] Arpan Mandal, Kripabandhu Ghosh, Saptarshi Ghosh, and Sekhar Mandal. Unsupervised approaches for measuring textual similarity between legal court case reports. *Artificial Intelligence and Law*, pages 1–35, 2021.
- [SD19] Lacour Stéphanie and Piana Daniela. Faites entrer les algorithmes! regards critiques sur la « justice prédictive ». *Cités*, 80 :47–60, April 2019.

---

1. <https://www.justice.gouv.fr/le-ministere-de-la-justice-10017/parution-du-calendrier-de-lopen-data-des-decisions.html>

2. <https://www.nfp77.ch/en/portfolio/court-decisions-in-the-field-of-tension-between-transparency-and-privacy/>

3. <https://www.data.gouv.fr/fr/datasets/conclusion-de-lexperimentation-sur-la-pseudonymisation-des-decisions-de-ju>