

Predilex : NLP appliqué à l'analyse des décisions judiciaires

A.H.KILANI — W.E.K.Marwan

Ecole Normale supérieure Paris-Saclay

March 24, 2020

Outline I

1 Introduction

- Objectif
- Exemple de dates (accident)
- Remarques

2 Pre-processing

- Uniformisation des dates
- Labelisation des contextes

3 Approche classique

- Tf-idf
- Dictionnaires
- Limitations

4 Approche réseaux de neurones

- Transformers - BERT
- Classification des dates
- Classification des dates
- Detection de la non consolidation

Outline II

- Detection de la non consolidation
- Resultats

5 Analyse des features de BERT

6 Limitations

7 Conclusions

Introduction - Objectif :

A partir de documents de jurisprudence, extraire 3 features :

- **Sexe** de la victime
- **Date de l'accident** (AAAA-MM-JJ ou "N.C." si absente)
- **Date de consolidation** (AAAA-MM-JJ, "N.C." si absente, "N.A." si décès avant consolidation)

Accuracy moyenne

Introduction - Exemple de dates (accident)

- « Laetitia X. [...] à l'occasion d'une soirée mousse, **le 1 août 1998**, [...] a été gravement blessée après avoir glissé sur la piste de danse de cette discothèque. »
- « **Le 27 mai 2009**, [...], Alexandre X. a perdu le contrôle de son véhicule à la suite d'un malaise. »
- « **Le 6 mars 2009**, au cours d'un match d'entraînement de football [...], Jean-Ulrich X., né **le 9 mars 1994**, a subi une blessure. »

Introduction - Remarques

- Différents format des dates
- Emplacement des informations
- Des cas particuliers difficiles à gérer (2 accidents dans le même texte)

Pre-processing - Uniformisation des dates

```
split_rows_multiple_dates(rows_in_text_of_rows_with_dates(X_train.dates_uniformed.iloc[10]))
```

```
[ ' 10 avril 1943 sarlat (24200)',  
'(beneficie aide juridictionnelle totale numero 2006 /000004 13 janvier 2006 accordee bureau aide juridictionnelle agen)',  
'demandeur renvoi cassation ordonne arret rendu 15 novembre 2005 cassant annulant arret rendu cour appel bordeaux date',  
'cassant annulant arret rendu cour appel bordeaux date 19 decembre 2002 appel jugement rendu tribunal grande instance bordeaux date',  
'appel jugement rendu tribunal grande instance bordeaux date 13 fevrier 2001 ',  
'arret contradictoire suivant apres cause ete debattue plaidee audience publique 06 juin 2007 devant rene salomon premier president bernard boutin',  
'retour algerie mois 1er octobre 1983 jean-pierre x... ressentait fièvre importante malaises grande fatigue. etait hospitalise',  
'malaises grande fatigue. etait hospitalise abord centre hospitalier agen novembre 1er decembre 1983 puis chu bordeaux',  
'puis chu bordeaux 1er decembre 1983 fevrier 1984. cours hospitalisations recevait quatre transfusions agen puis 61',  
'puis 61 produits sanguins bordeaux. contamination virus hepatite etait diagnostique 22 juin 1994 puis confirme',  
'puis confirme 27 mars 1996 .',
```

Figure 1: Données créées après uniformisation des dates — $T = 10$

Pre-processing - Labelisation des contextes

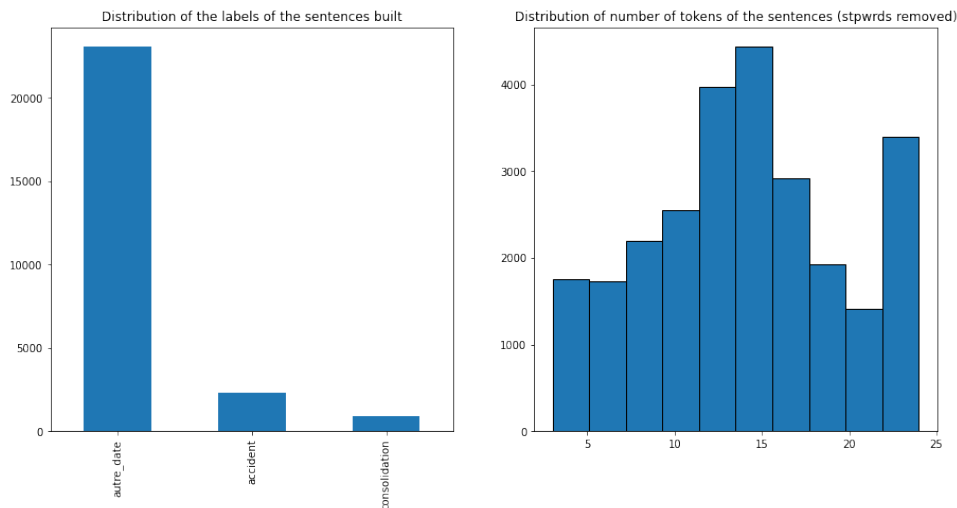


Figure 2: Gauche : Distribution des labels créés — Droite : distribution des nombres de tokens de chaque ligne — $T=10$

- Entraîner un tf-idf sur ces contextes : Matrice Tf-idf
- Représentation d'un contexte : vecteur Tf-idf
- Entraîner un classifieur (SVM, logReg)

- Deux idées clés inspirées de tf-idf: fréquence (à l'aide de dictionnaires) et robustesse
- Dictionnaire du contexte de l'accident
- Prédiction du sexe: chercher le sexe dans le contexte de l'accident
- Dictionnaires pour les sexes basés sur deux corpus
- Pour la consolidation: dictionnaire de décès + dictionnaire de contexte

- Besoin d'éliminer certaines dates lors de l'inférence
- Pour les dictionnaires: aucune exploitation de l'ordre des mots ou de leurs relations
- Décès: Et si ce n'est pas l'appelant(e) qui est décédé(e)?
(Héritage, appelant(e) représentant une personne décédée, etc.)

Approche réseaux de neurones - Transformers - BERT

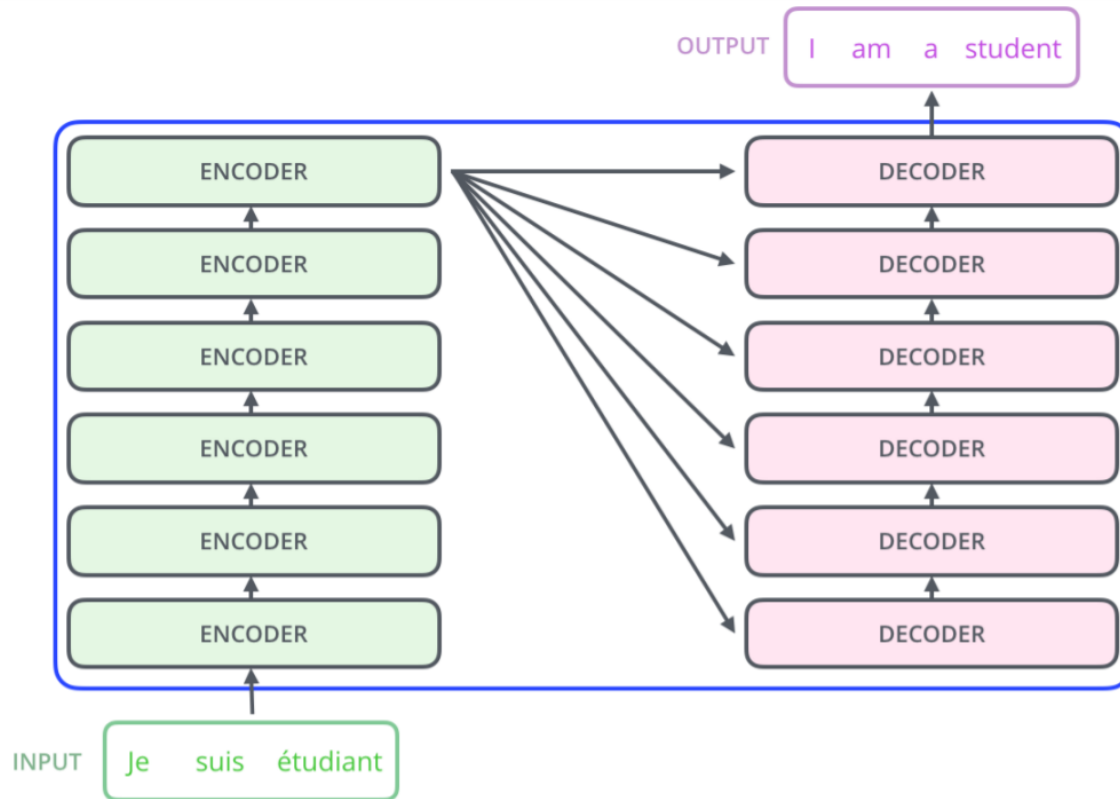


Figure 3: Architecture encodeurs-décodeurs d'un transformer

Approche réseaux de neurones - Transformers - BERT

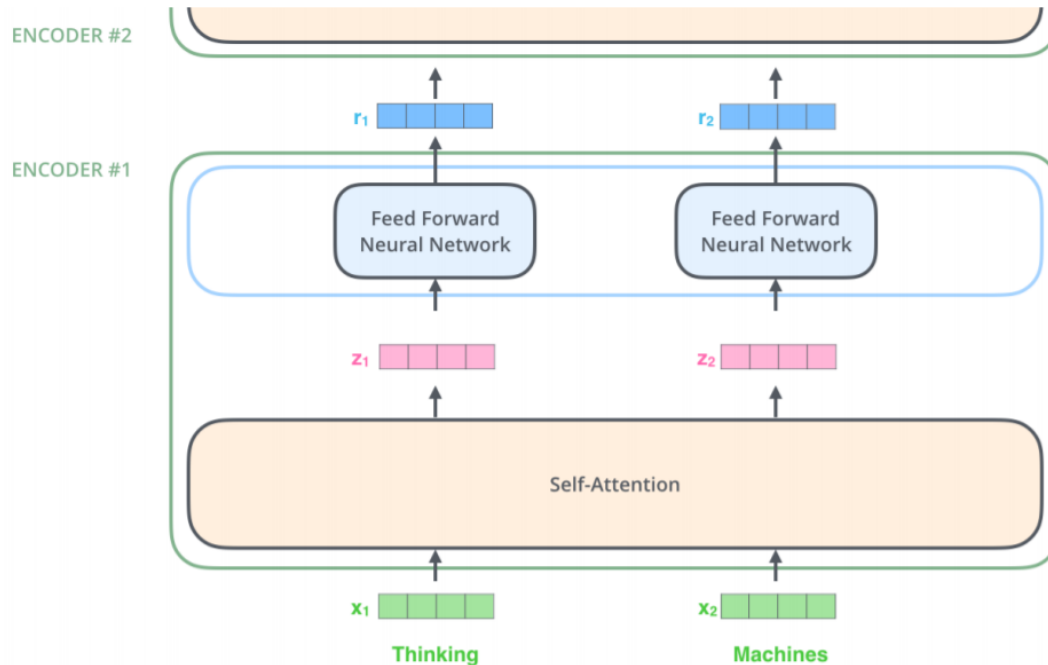


Figure 4: Architecture interne d'un encodeur

Approche réseaux de neurones - Transformers - BERT

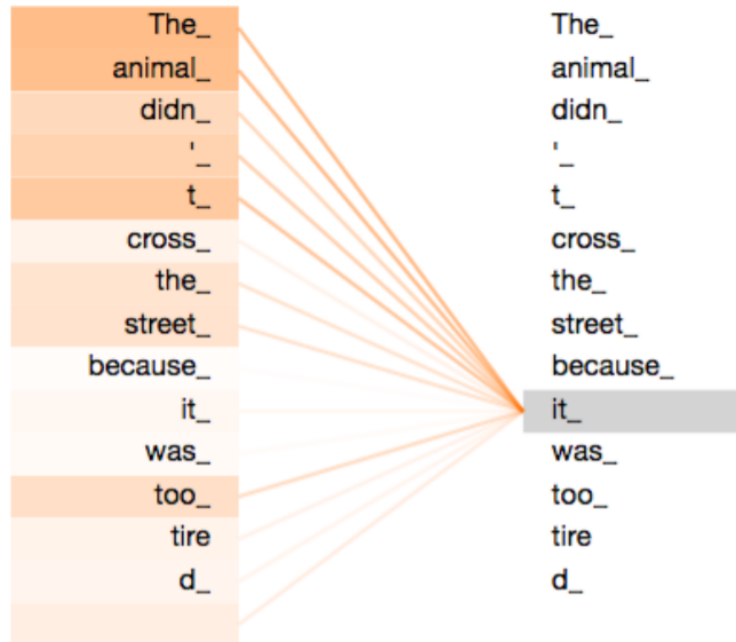


Figure 5: Illustration du principe d'attention

Approche réseaux de neurones - Transformers - BERT

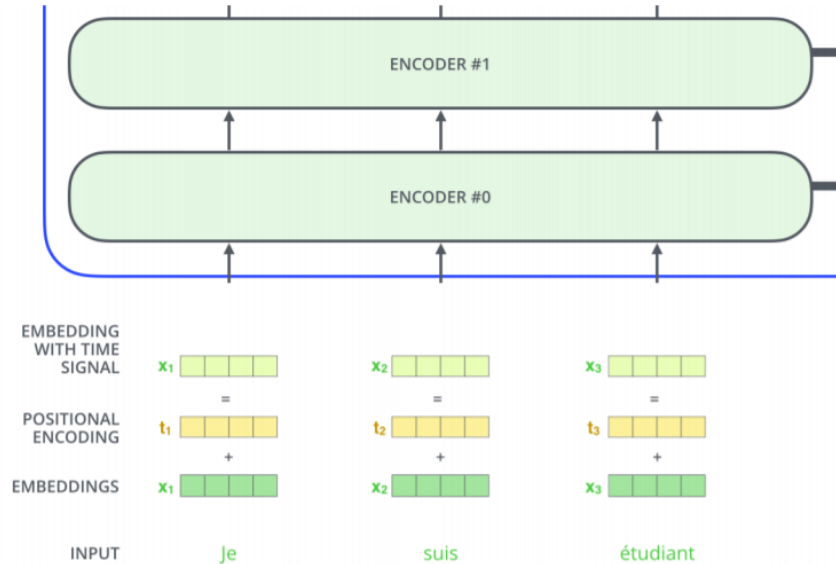


Figure 6: Les vecteurs d'embedding

Approche réseaux de neurones - Transformers - BERT

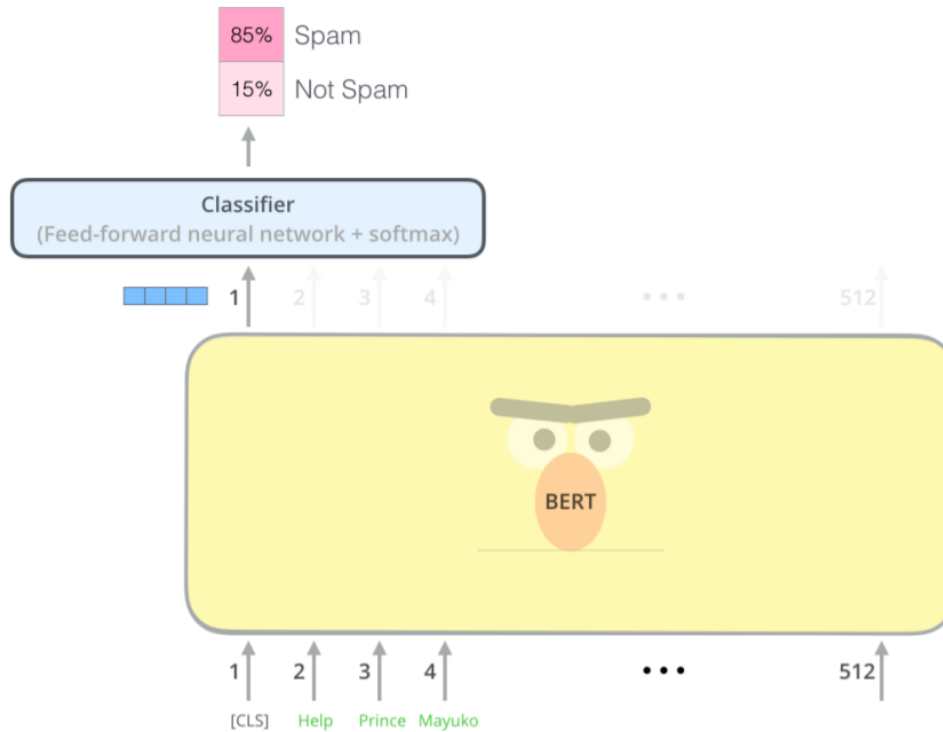


Figure 7: BERT pour la classification : Sortie de l'encodeur du transformer

Approche réseaux de neurones - Classification des dates

- Contexte autour d'une date : Vecteur BERT SOS 768 dimensions

Layer (type)	Output Shape	Param #
input_1 (InputLayer)	(None, 768)	0
dense_1 (Dense)	(None, 768)	590592
dense_2 (Dense)	(None, 3)	2307
Total params: 592,899		
Trainable params: 592,899		
Non-trainable params: 0		

Figure 8: Classifieur

Approche réseaux de neurones - Classification des dates

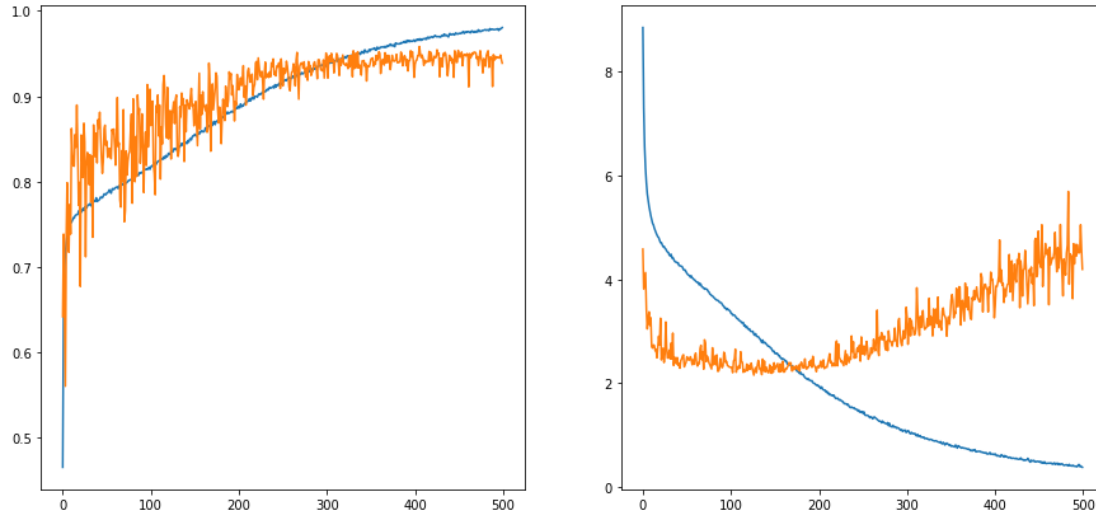


Figure 9: Historique d'entraînement (Gauche: Accuracy — Droite : Loss)

Approche réseaux de neurones - Detection de la non consolidation

Layer (type)	Output Shape	Param #
input_2 (InputLayer)	(None, 20, 768)	0
conv1d_1 (Conv1D)	(None, 18, 400)	922000
global_max_pooling1d_1 (Glob	(None, 400)	0
dense_3 (Dense)	(None, 2)	802
Total params: 922,802		
Trainable params: 922,802		
Non-trainable params: 0		

Figure 10: 1D convolutional network

Approche réseaux de neurones - Detection de la non consolidation

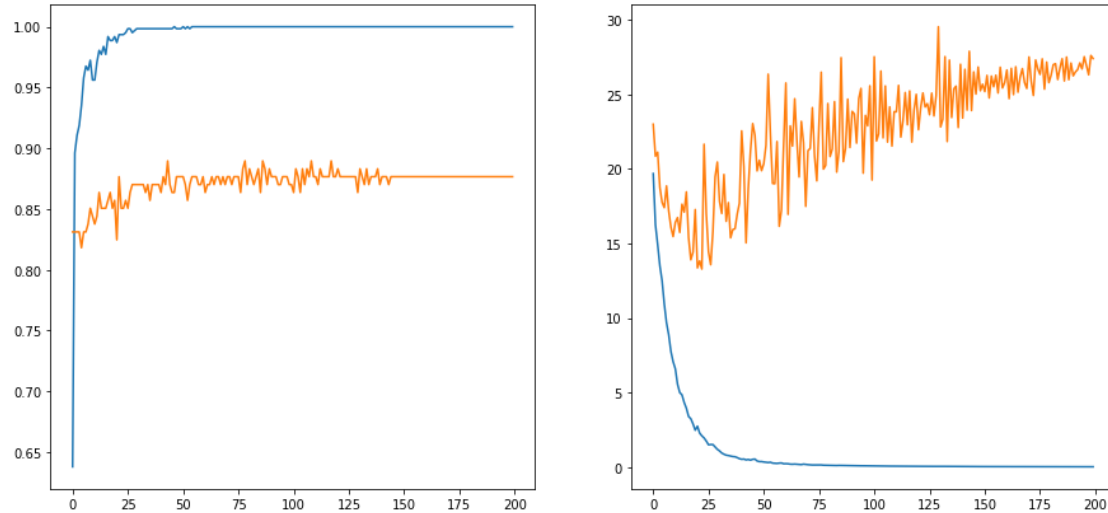


Figure 11: Historique d'entraînement (Gauche: Accuracy — Droite : Loss)

Approach	Accident	Consolidation
2 classes - avg Word vectors	49	41
2 classes - tf idf **	59	50
3 classes - avg Word vectors	54	38
3 classes - tf idf **	65	44
3 classes - BERT - NN*	78	67
Dictionary	72	66

Table 1: Validation accuracy of dates classifiers — * : using the constraint on the date consolidation — ** filtering some dates using a dictionary

Analyse des features de BERT

- Embedding du SOS token tient compte des termes importants du contexte
- Les mots clés ressortent (poids d'attention importants)

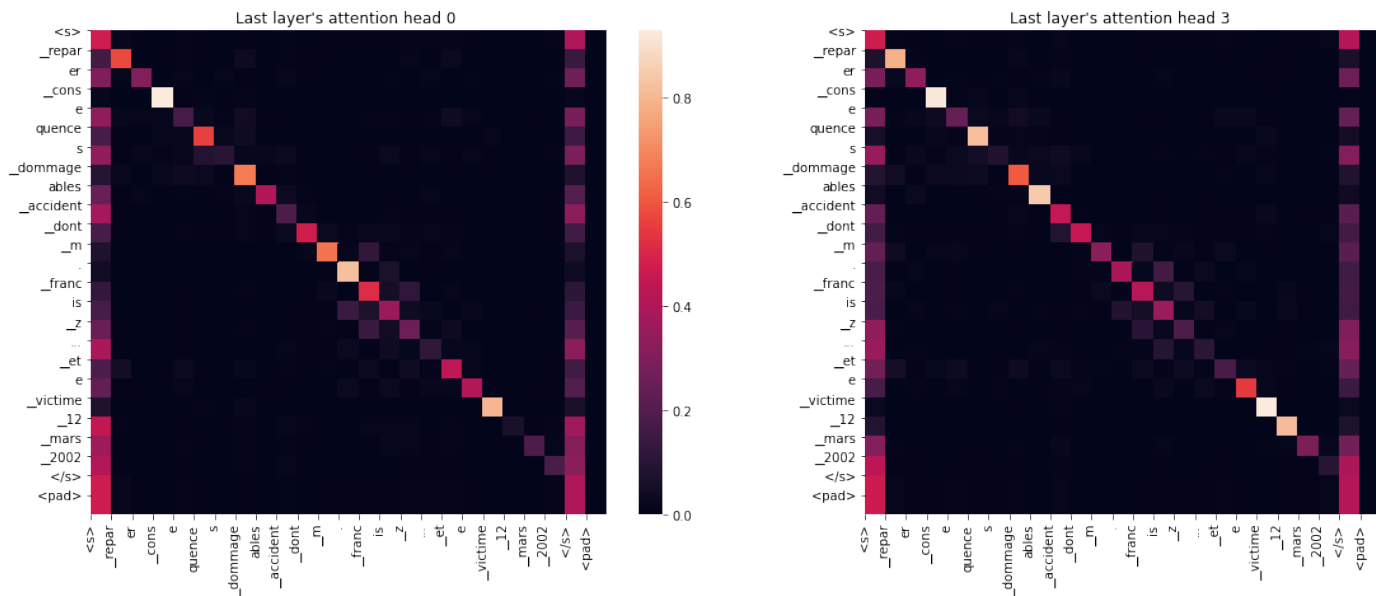


Figure 12: Dernier encodeur — 2 têtes d'attentions

Limitations

- Sur-apprentissage : Pas de couche de Dropouts dans les NN, seuils de probabilités non adéquats
- Pre-processing : On rate les dates dans les intervalles Du 03.06 au 07.06.2003 a été hospitalisé due à....
- Dépendence entre dates : se tromper sur la date d'accident entraine une erreur sur la consolidation
- Features non spécifiques au milieu juridique (BERT pré-entraîné sur OSCAR (Open Super-large Crawled ALMAnaCH coRpus))

- NN n'est pas LA solution pour tout (Classification sexe)
- Information a priori modélisée par BERT : synthèse des contextes, l'ordre des termes, relation sémantiques et syntaxiques
- Rajouter de la régularisation (Dropout NN)