

# Évaluation des apports d'un reseau de neurones antagoniste pour une classification semi-supervisé de textes français issue d'une base de données faiblement labellisés

Dylan Baptiste, Stéphane Cormier et Shufan Jiang

Université de Reims Champagne-Ardenne  
dylan.baptiste@etudiant.univ-reims.fr  
stephane.cormier@univ-reims.fr  
shufan.jiang@univ-reims.fr

**Résumé.** Les tâches de classifications de textes se heurtent souvent au problème de bases de données peu voir pas labellisé. Ici nous nous intéressons au apports de l'ajout d'un GAN afin d'artificiellement créer des données labellisés. Le but est d'améliorer les performances d'un classifieur qui ne se serait entraîné que sur des données labellisé issue d'une petite base de données.

**Mots-clé:** réseau de neurones, réseaux antagonistes génératifs, traitement du langage, classification, semi-supervisé

# Table des matières

1	Contexte.....	2
2	État de l'art.....	2
2.1	Représentations d'encodeur bidirectionnel à partir de transformateurs (BERT) .....	2
2.2	Réseaux antagonistes génératifs .....	3
3	Données .....	3
3.1	Extraction de caractéristiques .....	3
4	Architecture.....	4
4.1	Générateur .....	4
4.2	Discriminateur .....	4
4.3	GAN .....	5
5	Entraînement.....	7
5.1	Fonctions d'erreurs .....	7
5.2	Supervisé.....	7
5.3	Semi-supervisé .....	7
	Compétition ou coopération ?.....	7
	Stabilisation du Générateur et Discriminateur .....	7
6	Evaluation .....	7
6.1	Mesures utilisées .....	7
6.2	Comparaison des mesures .....	7
7	Conclusion et perspectives .....	7

## Liste des figures

## Liste des tables

### 1 Contexte

Ce travail à été réalisé dans le cadre d'un TER (Travail d'Étude et de Recherche) par Dylan Baptiste, étudiant en 2 année de Master, supervisé par Shufan Jiang, doctorante, et Stéphane Cormier, enseignant-chercheur, à l'Université de Reims Champagne-Ardenne.

    presentation des données...

### 2 État de l'art

#### 2.1 Représentations d'encodeur bidirectionnel à partir de transformateurs (BERT)

Les modeles BERT (de l'anglais Bidirectional Encoder Representations from Transformers) sont des modeles de traintement du langage introduit en 2018

par Google (TODO citer <https://arxiv.org/abs/1810.04805>) qui ont permis des améliorations significatives dans ce domaine. BERT a été entraîné sur un corpus de plusieurs langues. L'entraînement des modèles BERT s'effectue en deux étapes, tout d'abord une tâche semi supervisée où des mots caviardés dans un texte doivent être retrouvés par le modèle et une seconde supervisée où le modèle doit déterminer si une phrase B est la suite d'une phrase A ou non (Next-Sentence Prediction). CamemBERT est un modèle basé sur la même architecture mais entraîné uniquement sur un corpus français où la tâche (Next-Sentence Prediction) n'a pas été réalisée. De tels modèles prennent en entrée des phrases encodées, pour chaque token dans la phrase encodée le modèle produit en sortie un vecteur (de dimension 768 pour BERT et CamemBERT) et un autre vecteur (toujours de taille 768) pour l'ensemble de la séquence, ce token est appelé [CLS] et sa représentation en sortie est une représentation dans l'espace, selon le modèle, de la séquence encodée. C'est sur le token [CLS] que le modèle est entraîné pour la tâche de classification supervisée. (TODO image de BERT)

## 2.2 Réseaux antagonistes génératifs

Les GAN (de l'anglais Generative Adversarial Networks) sont une famille de réseaux de neurones communément séparables en deux parties antagonistes: un générateur (G) et un discriminateur (D) qui s'affrontent lors de l'entraînement. Le générateur a pour but de générer des données semblables à des données réelles. Les applications les plus connues concernent la génération d'images (TODO: CITER ). Pour effectuer une telle tâche ce modèle est adjoint à un autre modèle appelé Discriminateur qui doit déterminer si une donnée est réelle (issue d'une base de données) ou générée (c-à-d produite par le générateur), il effectue donc une classification binaire. Il existe une multitude de variantes de réseaux antagonistes génératifs notamment les Conditional GAN (CGAN) dont la génération du G est conditionnée par un ou des labels. Dans ce travail c'est une sorte de CGAN qui sera utilisée.

## 3 Données

Les données qui serviront d'entraînement pour ce projet sont des messages français issus du réseau social twitter. Certains messages sont des observations d'agriculteurs tandis que d'autres non. Les données sont donc labellisées binaires sur cette caractéristique. Nous disposons d'une base de données de XXX messages labellisés à la main: (TODO: donner les stats) XXX observations et XXX non observations, ainsi que XXX messages non labellisés.

### 3.1 Extraction de caractéristiques

Dans ce travail il n'est pas question d'effectuer un entraînement de la partie du modèle CamemBERT. Il y a donc une phase d'extraction des caractéristiques afin de limiter le temps de calcul lors des entraînements. C'est dans cette phase

que l'on doit choisir quelles sorties de CamemBERT utiliser pour la représentation des messages. Pour ce travail seul l'encodage du token [CLS] est utilisé. Ce choix arbitraire ouvre des perspectives de travail qui seront discutées dans la dernière partie. Chaque message est donc à présent "vectorisé" dans un espace de 768 dimensions. Les labels associés aux messages annotés manuellement sont reportés aux vecteurs associés au message et les vecteurs des messages non labellisés restent non labellisés

(TDOD mettre image d'un vecteur)

Les données sont normalisées par une mise à échelle entre -1 et 1

(TODO mettre image après scaling)

## 4 Architecture

L'architecture étudiée ici se compose de deux modèles, le premier est appelé le Générateur (G) et le second le Discriminateur (D)

### 4.1 Générateur

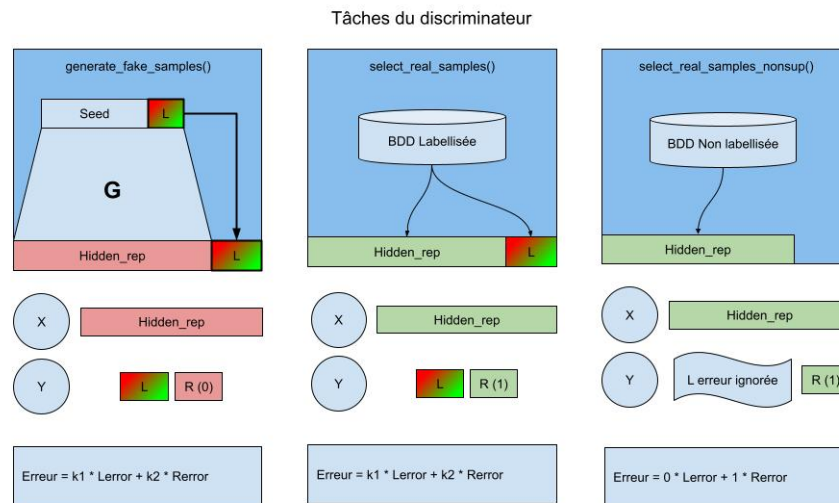
Comme expliqué précédemment le générateur a pour but de générer des données semblables à des données réelles. Ici les données en question sont les vecteurs de représentation de [CLS] du modèle CamemBERT comme présenté plus tôt.

Le Générateur est un simple réseau de neurones artificiel composé de couches denses (TODO préciser le modèle)

Pour effectuer sa tâche de génération le Générateur prendra en entrée un vecteur bruit nommé espace latent qui agira comme une graine de génération, et un label qui conditionnera si le vecteur en sortie doit encoder une observation ou pas. La sortie du modèle est un vecteur de taille 768 activé en tangente hyperbolique

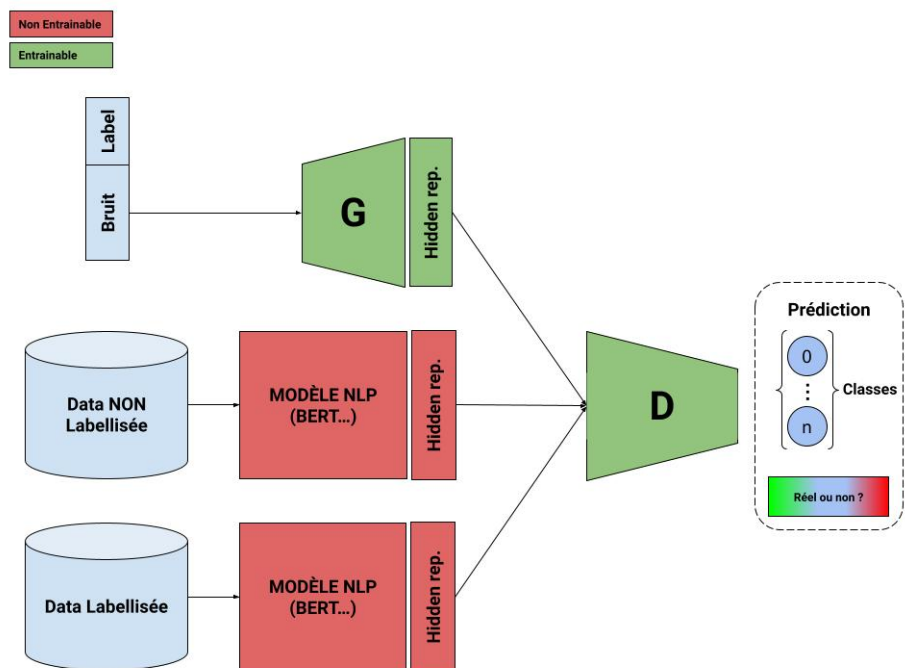
### 4.2 Discriminateur

Le discriminateur a pour tâche de discriminer les données qu'il reçoit en deux catégories : réel ou généré, et en même temps en deux autres catégories : observation ou non observation. Ce modèle prend donc en entrée un vecteur de taille 768 et en sortie deux neurones activés avec la fonction sigmoïde, le premier neurone effectue la première tâche de classification (réel / généré) et le deuxième détermine le label associé au vecteur d'entrée (observation / non observation). Comme pour G les couches cachées sont de neurones artificiels. (TODO parler des branches)



### 4.3 GAN

Le model advertial est donc l'assemblage de G avec D. Afin d'entrainer l'architecture il faut alterné l'entainement de G puis de D c'est a dire ne pas effectuer de retropopagation dans D lors de l'entrainement de G et ne pas retropopager dans G lors de l'entrainement de D. G est recompensé s'il trompe D sur la realité des donné qu'il generer. Similairement D est recompensé s'il determine quel données sont generé et quel donné ne le sont pas. En meme temps la tache de classification du label se fait: (TODO voirsi j'en parle ici)



## 5 Entraînement

### 5.1 Fonctions d'erreurs

Comme il s'agit à chaque fois d'une classification binaire c'est l'erreur d'entropie croisée binaire qui est utilisé:

$$-(y \cdot \log(\hat{y}) + (1 - y) \cdot \log(1 - \hat{y}))$$

Pour le D chacune de ses branches sera évalué avec cette fonction et l'erreur total de la prediction sera la somme des erreurs. Comme le G utilise D pour faire ses predictions son erreur sera calculer comme pour celle du D mais les predictions sont inversé. Ainsi G minimise sa fonction d'erreur s'il arrive à tromper D.

### 5.2 Supervisé

bce seul sur le label

### 5.3 Semi-supervisé

**Compétition ou coopération ?**

**Stabilisation du Générateur et Discriminateur**

## 6 Evaluation

### 6.1 Mesures utilisées

recall, precision, f1 score, APS, AUC au sein d'une cross validation (5folds, équilibre pour l'entraînement)

### 6.2 Comparaison des mesures

## 7 Conclusion et perspectives

perspectives: generation au niveau des token, utilisation d'un modèle NLP entraîné sur le corpus que l'on souhaite classer remplacer la prediction de la réalité par une critique comme cela est fait dans un WGAN

test f zjfhizef un des probleme de cette architecture est que meme lorsque le generateur n'est pas bon, c'est à dire qu'il ne genere pas de données proche de la réalité, le discriminateur doit quand meme donner un avis et etre penalisé, en depit de la pietre qualité du generateur. pour palier a ce probleme, outre le fait de pourvoir attendre une qualité de generation "satisfaisante" un autre type d'architecture aurais pu etre utilisé comme par exemple les WGAN (TODO citer)

Choix de [CLS] : D'autres (ou combinaison) d'autres representations du model auraient pu etre choisies.