

UNIVERSITÉ AL AKHAWAYN À IFRANE

## Mémoire de fin d'études

Titre du projet de fin d'études

Présenté par : Rayane Mestari

Programme : Computer Science – AI & Data Analytics

Encadrant : Dr. Mohammed Mestari

# **Composition du jury**

- Dr. Mohammed Mestari – Encadrant académique
- Dr. Abdellatif Marghich – Membre du jury
- Dr. Imane Nizar – Membre du jury

# **Remerciements**

Je tiens à remercier mon encadrant, les membres du jury ainsi que toute l'équipe pédagogique pour leur soutien précieux.

# Résumé

Ce mémoire présente...

**Mots-clés :** intelligence artificielle, apprentissage automatique, AUI, data analytics

# Abstract

This thesis presents...

**Keywords:** artificial intelligence, machine learning, AUI, data analytics

# Table des matières

<b>Composition du jury</b>	<b>1</b>
<b>Remerciements</b>	<b>2</b>
<b>Résumé</b>	<b>3</b>
<b>Abstract</b>	<b>4</b>
<b>1 Introduction</b>	<b>6</b>
1.1 Tokenisation WordPiece en C++ . . . . .	4
<b>2 Contexte et Problématique</b>	<b>6</b>
<b>3 Méthodologie</b>	<b>7</b>
<b>4 Résultats</b>	<b>8</b>
<b>5 Conclusion et Perspectives</b>	<b>9</b>
<b>A Annexe A : Informations complémentaires</b>	<b>10</b>

# **Chapitre 1**

## **Introduction**

# Fondements Mathématiques des Techniques de NLP

## 1. Tokenization

La tokenization segmente un texte  $T$  en unités appelées *tokens*  $\{t_1, t_2, \dots, t_n\}$ . C'est une fonction de transformation :

$$f : T \rightarrow \{t_i\} \subseteq \Sigma^*$$

avec  $\Sigma$  l'alphabet. Les techniques utilisées incluent :

- Expressions régulières
- Automates finis
- Algorithmes statistiques (ex. pour les langues non segmentées)

## 2. Stemming et Lemmatization

### Stemming

Réduction heuristique basée sur des règles (suffix stripping) :

$$\text{stem}(w) = w' \quad w' \approx \text{racine de } w$$

Exemple : Algorithmes de Porter, Snowball.

### Lemmatization

Forme correcte linguistiquement, dépendant de l'analyse morphosyntaxique :

$$\text{lemma}(w, \text{POS}) = l \quad l \in \text{dictionnaire}$$

Basée sur règles grammaticales et dictionnaires.

## 3. Suppression des Stopwords

Filtrage booléen sur les tokens :

$$f(t) = \begin{cases} 0 & \text{si } t \in \text{stopwords} \\ 1 & \text{sinon} \end{cases}$$

Réduction manuelle de la dimensionnalité de la représentation.

## 4. Bag-of-Words (BoW)

**Modélisation vectorielle** : chaque document devient un vecteur creux :

$$\vec{d} = (c_1, c_2, \dots, c_n) \quad c_i = \text{nb de } w_i \text{ dans } d$$

**Matrice du corpus** :  $D \in \mathbb{N}^{m \times n}$ , avec  $m$  documents et  $n$  mots.

## 5. TF-IDF

**Formule** :

$$\text{TF-IDF}(t, d) = \text{TF}(t, d) \cdot \log \left( \frac{N}{DF(t)} \right)$$

- $\text{TF}(t, d)$  : fréquence du terme dans le document  $d$
- $DF(t)$  : nombre de documents contenant  $t$
- $N$  : nombre total de documents

Pondère l'importance des mots rares et spécifiques.

## 6. Word Embeddings

### 6.1 Word2Vec

Deux variantes :

- CBOW : prédit un mot à partir de son contexte
- Skip-gram : prédit le contexte à partir d'un mot

**Objectif d'apprentissage** :

$$\frac{1}{T} \sum_{t=1}^T \sum_{-k \leq j \leq k, j \neq 0} \log P(w_{t+j} | w_t)$$

$$P(w_O | w_I) = \frac{\exp(v_{w_O}^T v_{w_I})}{\sum_{w \in V} \exp(v_w^T v_{w_I})}$$

Utilise la technique du *negative sampling* pour une approximation efficace.

## 6.2 GloVe

Basé sur les cooccurrences de mots :

$$w_i^T w_j + b_i + b_j = \log(X_{ij})$$

Fonction de coût :

$$J = \sum_{i,j} f(X_{ij}) \left( w_i^T w_j + b_i + b_j - \log(X_{ij}) \right)^2$$

Optimisé par moindres carrés pondérés.

## 6.3 FastText

Utilise les n-grammes de caractères :

$$\vec{w} = \sum_{g \in G(w)} z_g$$

- $G(w)$  : ensemble de n-grammes du mot  $w$
- $z_g$  : vecteur associé à chaque n-gramme

Robuste aux mots rares et aux variantes morphologiques.

## 7. Tableau récapitulatif

Méthode	Type	Représentation	Mathématiques
Tokenization	Symbolique	Séquences	FSA, regex
Stemming	Heuristique	Racine approximée	Substitutions
Lemmatization	Linguistique	Forme canonique	Graphe, dictionnaire
Stopwords	Réduction	Filtres booléens	Indicateur logique
BoW	Statistique	Vecteurs creux	Comptage, matrice $\mathbb{N}^{m \times n}$
TF-IDF	Statistique	Poids pondérés	Log, fréquence inverse
Word2Vec	Neuronale	Vecteurs denses	Skip-gram / CBOW + softmax
GloVe	Factorisation	Cooccurrence	Moindres carrés pondérés
FastText	Neuronale morphologique	N-grammes de caractères	Moyenne vectorielle

## 1.1 Tokenisation WordPiece en C++

```

1 #include <iostream>
2 using namespace std;
3
4 // Fonction de comparaison de cha nes
5 bool compare(const char* a, const char* b) {
6     int i = 0;
7     while (a[i] != '\0' && b[i] != '\0') {
8         if (a[i] != b[i]) return false;
9         i++;
10    }
11    return a[i] == '\0' && b[i] == '\0';
12 }
13
14 // Fonction pour copier une sous-cha ne
15 void substring(const char* word, int start, int end, char* result,
16    bool prefix) {
17     int k = 0;
18     if (prefix) {
19         result[k++] = '#';
20         result[k++] = '#';
21     }
22     for (int i = start; i < end; i++) {
23         result[k++] = word[i];
24     }
25     result[k] = '\0';
26 }
27
28 // Simule un vocabulaire
29 const int VOCAB_SIZE = 5;
30 const char* vocab[VOCAB_SIZE] = {
31     "extra", "##ordin", "##aire", "bon", "##jour"
32 };
33
34 // V rifie si le token est dans le vocabulaire
35 bool in_vocab(const char* token) {
36     for (int i = 0; i < VOCAB_SIZE; i++) {
37         if (compare(token, vocab[i])) return true;
38     }
39     return false;

```

```
39 }
40
41 // Tokeniseur WordPiece
42 void wordpiece_tokenizer(const char* word) {
43     int start = 0;
44     int len = 0;
45     while (word[len] != '\0') len++;
46
47     while (start < len) {
48         bool matched = false;
49         int end = len;
50
51         while (end > start) {
52             char token[50];
53             bool is_prefix = (start > 0);
54             substring(word, start, end, token, is_prefix);
55
56             if (in_vocab(token)) {
57                 cout << token << " ";
58                 start = end;
59                 matched = true;
60                 break;
61             }
62             end--;
63         }
64
65         if (!matched) {
66             cout << "[UNK] ";
67             break;
68         }
69     }
70     cout << endl;
71 }
72
73 int main() {
74     const char* test = "extraordinaire";
75     cout << "Tokenisation de : " << test << endl;
76     wordpiece_tokenizer(test);
77     return 0;
78 }
```

Listing 1.1: Tokeniseur WordPiece en C++

# **Chapitre 2**

## **Contexte et Problématique**

# Chapitre 3

## Méthodologie

# Chapitre 4

## Résultats

# **Chapitre 5**

## **Conclusion et Perspectives**

## **Annexe A**

### **Annexe A : Informations complémentaires**