



ECOLE NATIONALE SUPÉRIEURE D'INFORMATIQUE ET D'ANALYSE DES
SYSTÈMES - RABAT

Rapport de PFA : Healthcare Chatbot design

Réalisé par :

El Mehdi OUDAOUD
Fatima-Ezzahra LAHNINE

Encadré par :

Pr. Azeddine EL HASSOUNY

Année Scolaire 2020/2021

Résumé

Ce Projet propose un design de chatbot pour soins de santé qui répond aux questions les plus fréquemment posées et fait une vérification des symptômes entrées par les clients.

Ce bots permette les utilisateurs de décrire leurs symptômes et leur état de santé en posant une suite de questions. Le bot ensuite comprend les inputs de l'utilisateur en utilisant le traitement du langage naturel (NLP), quelle que soit la variation de l'entrée. Cela est essentiel pour obtenir des réponses précises.

Le bot est intégré dans une application web qui permet l'utilisateur de préciser une partie du corps concernant sa demande, et liste les différents symptômes liés à cette partie. Il peut choisir une des réponses listées qui va être transmise au bot, puis ce dernier répond le client par la description du symptôme choisi et lui donne quelque précautions et conseils à prendre.

Abstract

This project proposes a healthcare chatbot design that answers frequently asked questions and verifies the symptoms entered by customers.

The bot allows users to describe their symptoms and conditions by asking a series of questions. The bot then understands the user's input using natural language processing (NLP), regardless of the variation in the input. This is essential to get accurate answers.

The bot is embedded in a web application that allows the user to specify a body part regarding his query, and lists the different symptoms related to that part. He can choose one of the listed answers which will be transmitted to the bot, then the latter responds to the customer with a description of the chosen symptom and gives him some precautions and advice to take.



Table des matières

1	Architectures générales des chatbots	2
1.1	Conversion parole-texte	2
1.1.1	Reconnaissance vocale à vocabulaire étendu	2
1.1.2	Automatic Speech Recognition (ASR)	2
1.2	Natural Language Processing (NLP)	3
1.2.1	Bag of Words (BoW)	5
1.2.2	Tf-idf Vectorization	6
2	Le traitement de langage naturel avec des réseau de neurones	7
2.1	RNN	7
2.1.1	Architectures des RNN	7
2.1.2	Limitations de RNN	8
2.2	LSTM	8
2.3	CNN	9
2.4	Le modèle BERT	9
2.4.1	L'Architecture de BERT	10
2.4.2	Traitement préalable de texte	10
2.4.3	Les tâches de pré-entraînement	10
2.4.4	Réglage fin(Fine-tuning)	12
3	Concepts et implémentation du NLP	13
3.1	Créer des données d'entraînement	13
3.2	Principes de base du NLP	14
3.3	Implémenter les outils du NLP	16
3.4	Implémentation du Neural Network	16
3.5	Implémentation du Training Pipeline	17
3.6	Implémentation du Chat	19
4	Réalisation finale du Chatbot	22
4.1	Front-end	22

Introduction

Les chatbots sont considérés comme une alternative plus fiable et plus précise aux recherches en ligne que les patients effectuent lorsqu'ils essaient de comprendre la cause de leurs symptômes. Les prestataires de soins de santé pensent que les chatbots pourraient aider les patients qui ne sont pas sûrs de savoir où ils doivent aller pour recevoir des soins.

Les prestataires de soins de santé mettent actuellement en œuvre des bots qui permettent aux utilisateurs de vérifier leurs symptômes et de comprendre leur état de santé depuis le confort de leur domicile. Les chatbots qui utilisent le traitement du langage naturel (NLP) peuvent comprendre les demandes des patients, quelle que soit la variation de l'entrée. Cela est essentiel pour atteindre une grande précision dans les réponses, ce qui est indispensable pour les vérificateurs de symptômes.

Grâce à la connaissance de l'entrée, le bot peut évaluer les informations et aider les utilisateurs à déterminer la cause de leurs symptômes. Avec toutes les données fournies par le robot, les utilisateurs peuvent déterminer si un traitement professionnel est nécessaire ou si des médicaments en vente libre sont suffisants.

Chapitre 1

Architectures générales des chatbots

Un chatbot est un système de dialogue homme-ordinateur via le langage naturel. Il s'agit donc d'un humain ayant une conversation naturelle avec un ordinateur ou un système.

Le chatbot doit être capable de dialoguer et de comprendre l'utilisateur ; on pourrait dire qu'il s'agit d'une fonction de compréhension. Cette compréhension comprend la reconnaissance des intentions et des entités. Les intentions peuvent être considérées comme des verbes et les entités comme des noms. Les robots basés sur le texte ont au minimum un composant de compréhension du langage naturel (NLU).

Toute l'intelligence ne relève pas des capacités de la NLU. Les robots doivent avoir accès à une base externe de connaissances et de bon sens via des API, de sorte qu'ils puissent assurer la fonction de compétence en répondant aux questions des utilisateurs. Enfin, l'agent incarné doit assurer une présence très fonctionnelle. Ironiquement, ces agents numériques n'existaient pas jusqu'à récemment et étaient considérés comme très facultatifs. Aujourd'hui, cette fonction s'avère cruciale dans le cas des utilisateurs ordinaires.

1.1 Conversion parole-texte

1.1.1 Reconnaissance vocale à vocabulaire étendu

La parole est l'un des modes de communication les plus naturels et les plus puissants, et elle est largement acceptée comme l'avenir de l'interaction avec les applications informatiques et mobiles. Elle fait partie de la nature humaine au point que des personnes dont le IQ est inférieur à 50 et dont le cerveau est trois fois plus petit que celui des humains peuvent parler. La recherche neurologique indique que la parole occupe une plus grande partie du cerveau que toute autre fonction de traitement. En intégrant le traitement de la parole, les chatbots seront capables de s'interfacer avec les téléphones et les radios.

La conversion de la parole en texte commence par un processus appelé reconnaissance automatique de la parole (ASR). Les améliorations de la LVCSR peuvent être mesurées selon un certain nombre de critères :

1. **Volume du vocabulaire** : Les vocabulaires, qui étaient à l'origine minuscules et ne comprenaient que des phrases de base (par exemple, oui, non, chiffres, etc.), comprennent aujourd'hui des millions de mots dans de nombreuses langues.
2. **Indépendance du locuteur** : la capacité à reconnaître certains locuteurs. Ceci n'est pertinent que si les chatbots utilisent l'identité du locuteur pour générer des réponses spécifiques à l'utilisateur.
3. **Coordination de l'articulation** : la capacité de traiter un flux continu de mots, qui ne contient pas nécessairement de pauses entre les mots. Il faut une tokénisation et une segmentation appropriées du flux d'entrée, abordées dans la section qui suit.
4. **Gestion du bruit** : la capacité de filtrer le bruit (par exemple, la circulation, la musique de fond ou la parole).
5. **Microphone** : la capacité de traiter la parole à différentes distances du microphone.

1.1.2 Automatic Speech Recognition (ASR)

La reconnaissance vocale ou Speech-To-Text (STT) est un processus de conversion de la parole audio en texte. L'objectif de la ASR est de parvenir à une reconnaissance vocale indépendante du locuteur et portant sur un large vocabulaire. Alors que les chatbots peuvent s'offrir le luxe de s'adresser à un domaine très étroit, le

STT/ASR doit être capable de traiter un large vocabulaire. S'assurer que tout ce qui est dit peut être converti en texte. Le chatbot peut ne pas être en mesure de répondre directement à la requête ou à la demande. Ce qui pourrait tomber en dehors du domaine du chatbot. Mais l'ASR doit au moins présenter un texte précis à la section chatbot/NLU.

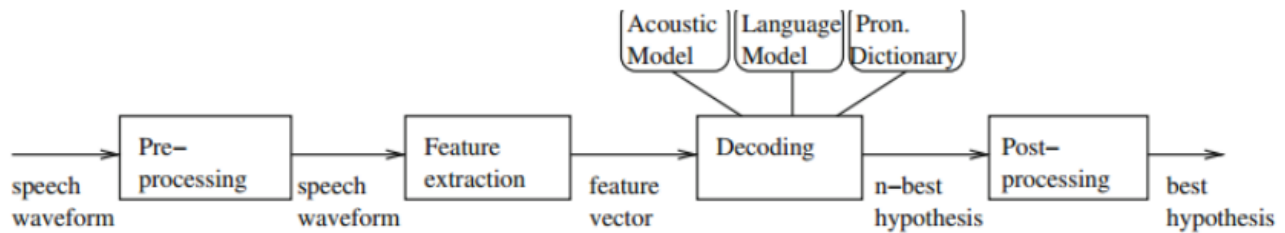


FIGURE 1.1 – Processus de la reconnaissance automatique de la parole

Premièrement, nous avons besoin d'un modèle acoustique qui nous donnera $P(X|W)$ - étant donné un mot, la probabilité que nous entendions un son. Nous trouverons la valeur argmax sur tous les mots de notre langue w , ou le mot ayant la probabilité la plus élevée de représenter ce son. Les modèles acoustiques sont généralement formés sur des enregistrements sonores et les transcriptions correspondantes, qui peuvent être utilisés pour trouver de manière empirique ces probabilités. La représentation statistique de chaque mot (ou phonème) générée par l'analyse du corpus sonore est typiquement représentée par un modèle de Markov caché (HMM). Les HMM sont des processus de Markov où les états sont non observés/cachés (nous ne connaissons pas les phonèmes/mots réels utilisés, et nous nous approximations de cette information).

Deuxièmement, nous avons besoin d'un modèle de langage qui puisse nous indiquer la probabilité d'entendre un mot donné dans notre langue.

Troisièmement, nous avons besoin d'un dictionnaire contenant une liste de mots et de leurs phonèmes. Étant donné un son, nous pouvons décoder le mot en utilisant la règle de Bayes :

$$W = \text{argmax}_{(W \in L)} \frac{P(X|W)P(W)}{P(X)}$$

La quatrième étape est le post-traitement ; la règle de Bayes nous donne une liste des séquences de mots probables avec leur rang, et nous choisissons la plus probable comme étant la séquence de mots que nous avons entendue. Les autres hypothèses les plus probables sont stockées et peuvent être utilisées plus tard dans des algorithmes d'apprentissage par renforcement, où elles seront utilisées pour apprendre et corriger les erreurs de la phase ASR.

1.2 Natural Language Processing (NLP)

La compréhension du langage naturel est à la base des capacités du chatbot. Sans la détection des entités et la reconnaissance des intentions, tous les efforts déployés pour comprendre l'utilisateur sont vains.

La plupart des architectures de chatbot se composent de quatre piliers, à savoir les intentions, les entités, le flux de dialogue (machine à états) et les scripts.

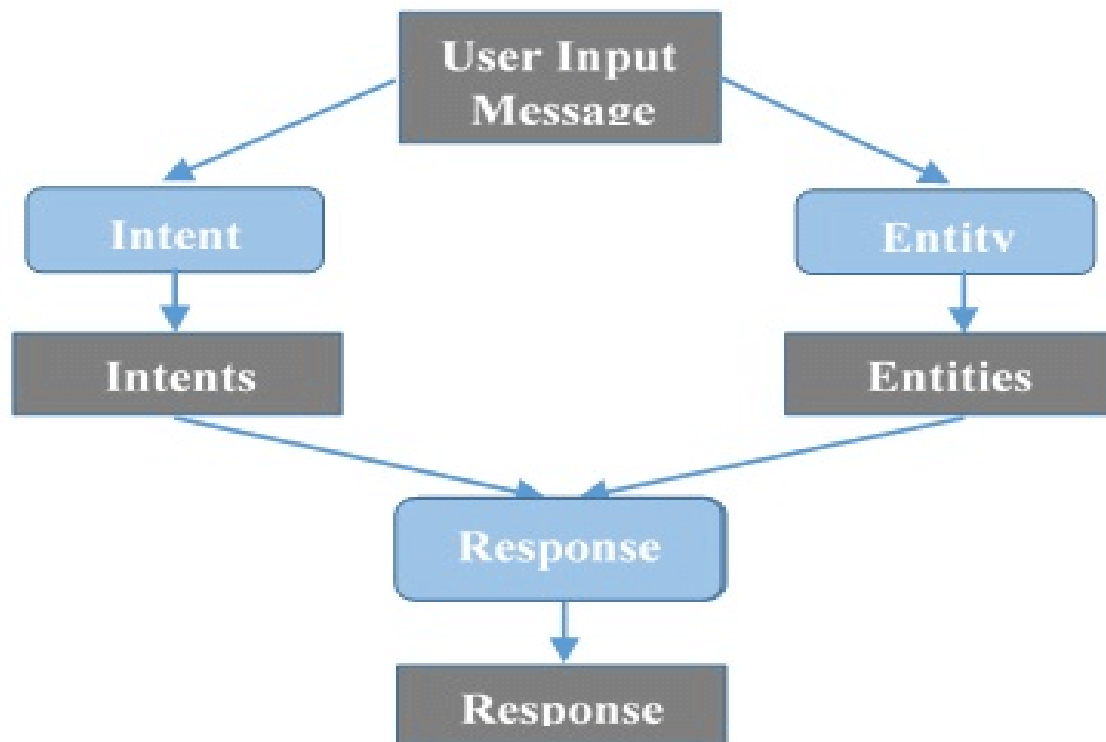


FIGURE 1.2 – Architecture traditionnelle des chatbots

Le dialogue contient les blocs ou états entre lesquels l'utilisateur navigue. Chaque dialogue est associé à une ou plusieurs intentions et/ou entités. Des variables de session peuvent également être utilisées pour décider quels états ou nœuds doivent être visités.

Les intentions et les entités constituent la condition d'accès à ce dialogue. La boîte de dialogue contient la sortie vers le client sous la forme d'un script, d'un message... ou d'une formulation si vous voulez.

C'est l'une des tâches les plus ennuyeuses et les plus laborieuses de la création d'un chatbot. Elle peut devenir complexe et les modifications apportées dans un domaine peuvent avoir un impact sur un autre domaine par inadvertance. Un manque de cohérence peut également conduire à des expériences utilisateur non planifiées. La mise à l'échelle de cet environnement est délicate, en particulier si on souhaite l'étendre à une grande organisation.

Dans la suite, on va introduire les techniques classiques de représentation des données textuelles.

1.2.1 Bag of Words (BoW)

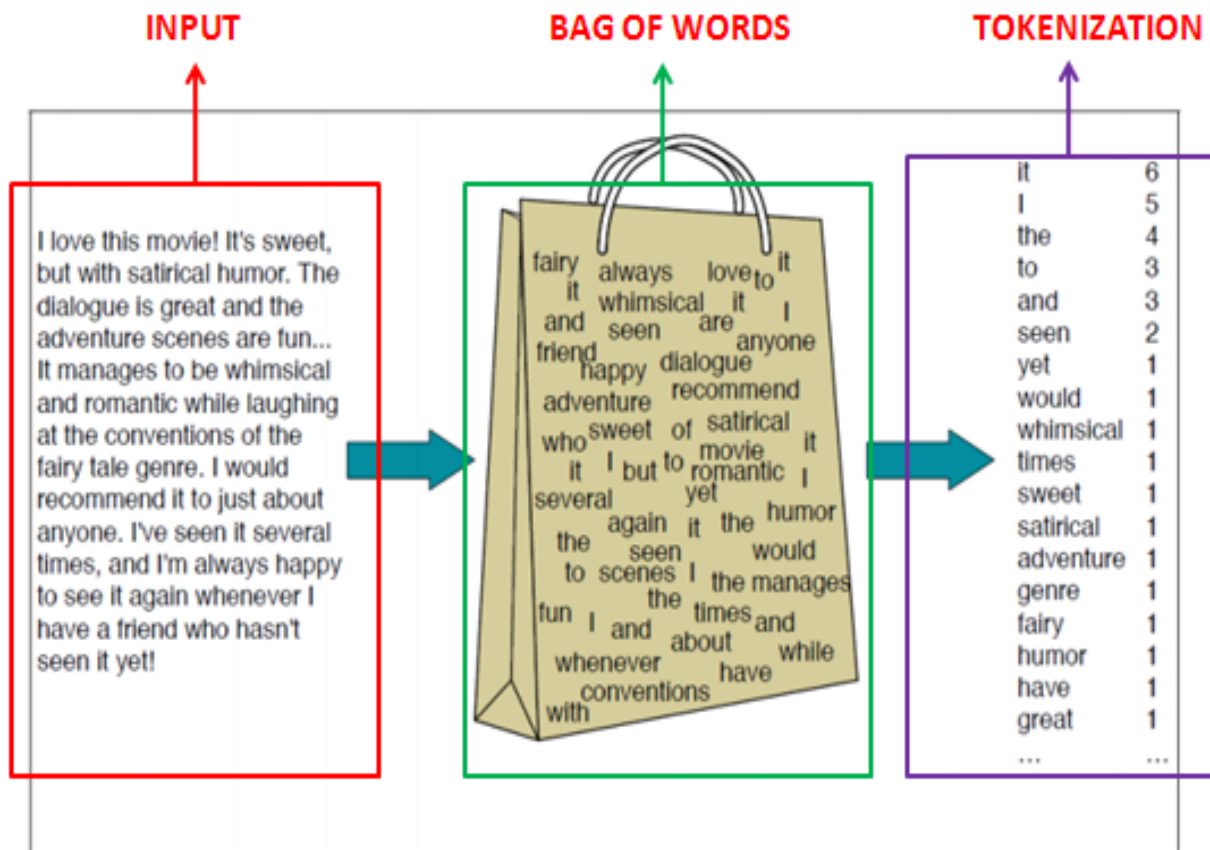


FIGURE 1.3 – Étape de base des algorithmes de classification de texte

Bag of Words (BOW) est un algorithme qui compte le nombre de fois qu'un mot apparaît dans un document. Ces nombres de mots nous permettent de comparer des documents et d'évaluer leurs similitudes pour des applications telles que la recherche, la classification de documents et la modélisation de sujets. Il est nommé ainsi car il ne concerne que l'occurrence du mot et non l'endroit où il est placé (c'est-à-dire l'ordre) dans le sac. De plus, l'intuition derrière une telle approche est que des documents similaires contiennent des mots similaires.

Document	the	cat	sat	in	hat	with
<i>the cat sat</i>	1	1	1	0	0	0
<i>the cat sat in the hat</i>	2	1	1	1	1	0
<i>the cat with the hat</i>	2	1	0	0	1	1

Une meilleure approche sera de créer un vocabulaire de mots groupés. Cela changera la taille du vocabulaire et permettra au sac de mots de capturer un peu plus de sens du document. La création d'un vocabulaire de paires de deux mots s'appelle un modèle bi-gramme. Un vocabulaire de triplés de mots est appelé un modèle de tri-gramme. Par exemple, les bi-grammes de "the cat sat in hat with" sont les suivants :

- "cat sat"
- "sat in"
- "in hat"
- "hat with"
- "the cat"

L'approche générale est appelée le modèle n-gramme, où n est la référence du nombre de mots groupés.

Avantages : simple à comprendre et à implémenter.

Inconvénients :

- Le vocabulaire se compose de tous les mots présents dans la base de données. Ainsi si les données sont trop volumineuses et contiennent de nombreux mots uniques, cela peut créer des matrices très larges et sparse ce qui augmente à la fois la complexité spatiale et temporelle.
- Le modèle BOW ne concerne que l'occurrence du mot et non l'endroit où il est placé (c'est-à-dire l'ordre). Cela conduit à la perte d'informations contextuelles et donc à la signification des mots dans le document (sémantique).

1.2.2 Tf-idf Vectorization

Tf-Idf est l'abréviation du terme fréquence-fréquence de document inverse. Donc, deux choses : la fréquence des termes et la fréquence inverse des documents.

La fréquence du terme (TF) est essentiellement la sortie du modèle BoW. Pour un document spécifique, il détermine l'importance d'un mot en regardant à quelle fréquence il apparaît dans le document. Si un mot apparaît plusieurs fois, le mot doit être important et sa fréquence grande.

Sa formule mathématique est :

$$tf(t, d) = \frac{f_{t,d}}{\sum_{t'} f_{t',d}}$$

avec $f_{t',d}$ le nombre d'occurrences du mot t, dans le document d.

IDF (Inverse Document Frequency) utilisé pour calculer le poids des mots rares dans tous les documents. Les mots qui apparaissent rarement dans le corpus ont un score IDF élevé. Cependant, on sait que certains termes, tels que « je », « a » peuvent apparaître plusieurs fois mais ont peu d'importance. Nous devons donc alourdir les termes fréquents tout en augmentant les rares. Sa formule mathématique est donnée par :

$$idf(t, D) = \log\left(\frac{N}{n_t}\right)$$

avec t le mot dans le corpus D et N la taille du corpus D et n_t le nombre de document contenant t .

Ainsi la formule mathématique pour TF-Idf est :

$$tf - idf(t, d, D) = tf(t, d) \times idf(d, D)$$

Avantages :

- Facile à calculer.
- On dispose d'une métrique de base pour extraire les termes les plus descriptifs d'un document.
- On peut facilement calculer la similitude entre 2 documents en l'utilisant.

Chapitre 2

Le traitement de langage naturel avec des réseaux de neurones

2.1 RNN

Les réseaux de neurones récurrents ou RNN sont une variante très importante des réseaux de neurones largement utilisés dans le traitement du langage naturel.

Conceptuellement, ils diffèrent d'un réseau de neurones standard car l'entrée standard dans un RNN est un mot au lieu de l'échantillon entier comme dans le cas d'un réseau de neurones standard. Cela donne au réseau la flexibilité de travailler avec des longueurs de phrases variables, ce qui ne peut pas être réalisé dans un réseau neuronal standard en raison de sa structure fixe. Il offre également un avantage supplémentaire de partage des fonctionnalités apprises à travers différentes positions de texte qui ne peuvent pas être obtenues dans un réseau neuronal standard.

Un RNN traite chaque mot d'une phrase comme une entrée distincte se produisant au temps « t » et utilise également la valeur d'activation à « $t-1$ » comme une entrée en plus de l'entrée au temps « t ».

2.1.1 Architectures des RNN

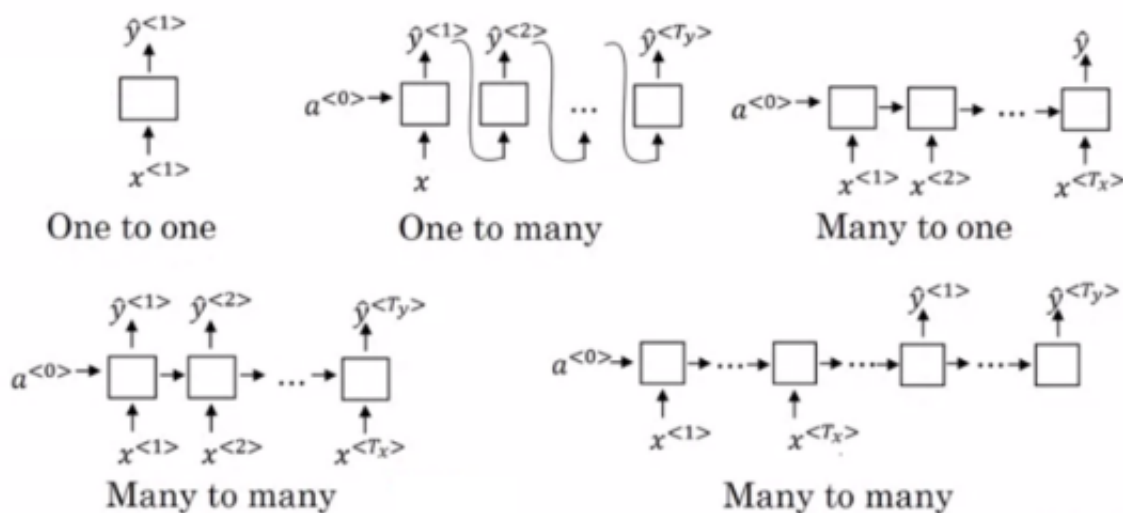


FIGURE 2.1 – Les différentes architectures des RNN

Il existe différentes architectures des RNN :

- **Plusieurs à plusieurs :** c'est-à-dire nombre d'entrées = nombre de sorties. Une telle structure est très utile dans la modélisation de séquence.
- **Plusieurs à un RNN :** l'architecture plusieurs à un fait référence à une architecture RNN où de nombreuses entrées (T_x) sont utilisées pour donner une sortie (T_y). Un exemple approprié d'utilisation d'une telle architecture sera une tâche de classification.
- **Un à plusieurs :** l'architecture un à plusieurs se réfère à une situation dans laquelle un RNN génère une série de valeurs de sortie basées sur une seule valeur d'entrée. Un excellent exemple d'utilisation d'une telle architecture sera une tâche de génération de musique, où une entrée est un jour ou la première note.
- **Plusieurs à plusieurs** (T_x n'est pas égal à T_y) : Cette architecture fait référence à l'endroit où de nombreuses entrées sont lues pour produire de nombreuses sorties, où la longueur des entrées n'est pas égale à la longueur des sorties. Les tâches de traduction automatique constituent un excellent exemple d'utilisation d'une telle architecture. Le codeur fait référence à la partie du réseau qui lit la phrase à traduire, et le décodeur est la partie du réseau qui traduit la phrase dans la langue souhaitée.

2.1.2 Limitations de RNN

Outre toute son utilité, RNN présente certaines limitations dont les principales sont :

Les exemples d'architecture RNN mentionnés ci-dessus sont capables de capturer les dépendances dans une seule direction du langage. Fondamentalement, dans le cas du traitement automatique du langage naturel, il suppose que le mot suivant n'a aucun effet sur la signification du mot précédent. Avec notre expérience des langues, nous savons que ce n'est certainement pas vrai.

Les RNN ne sont pas non plus très bons pour capturer les dépendances à long terme et le problème de la disparition des gradients refait surface dans RNN.

2.2 LSTM

Les réseaux LSTM sont un type de RNN, capable d'apprendre les dépendances à long terme. Les réseaux LSTM représentent la majorité des RNN utilisés aujourd'hui. Face au problème de "vanishing gradient" des RNNs les LSTMs sont présentées comme solution. D'ailleurs la mémoire à long terme est leur comportement par défaut il ne faut pas apprendre à ce réseau de le faire.

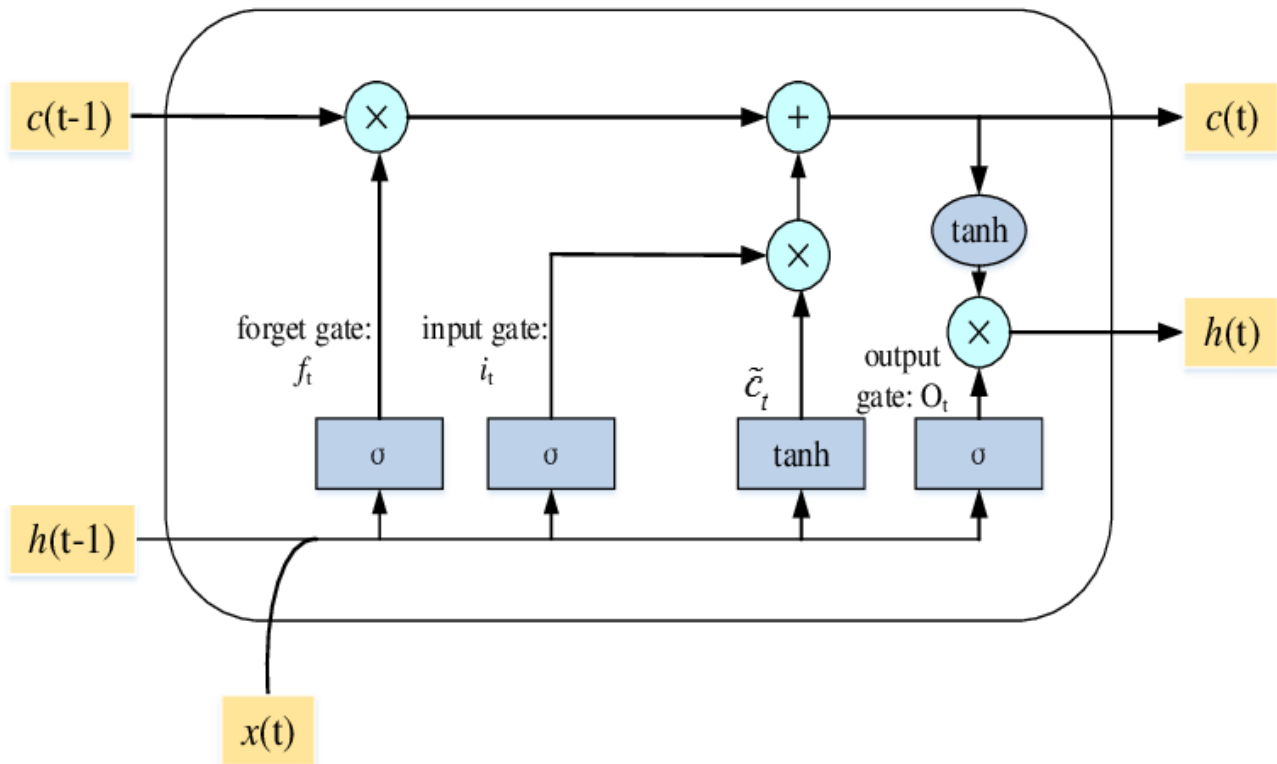


FIGURE 2.2 – La structure de l'unité LSTM

Dans un LSTM il existe trois portes : La porte d'entrée et la porte d'oubli déterminent comment le prochain état est influencé respectivement par l'entrée et le dernier état . La porte de sortie détermine comment la sortie du réseau est influencée par l'état. Parfois, il est utile de se souvenir des détails importants, mais de ne les utiliser que plus tard.

2.3 CNN

Lorsque nous entendons parler de CNN (Convolutional neural networks), nous pensons généralement à la vision par ordinateur. Les CNN ont été responsables de majeures avancées dans la classification des images et sont aujourd'hui au cœur de la plupart des systèmes de vision par ordinateur, du marquage automatique des photos de Facebook aux voitures autonomes. Plus récemment, nous avons également commencé à appliquer les CNN à des problèmes de traitement du langage naturel et obtenu des résultats intéressants pour des problèmes de catégorisation des phrases.

Le fonctionnement des CNN dans NLP ne diffère pas trop de leur fonctionnement sur les images, car ils prennent en entrée une matrice (moralement une image cf la figure ci-dessous) dont les lignes constituent la représentation vectorielle des différents tokens qui composent une phrase du corpus. matrice à utiliser comme entrée. Les modèles n'ont pas besoin d'être complexes pour obtenir des résultats solides. Les filtres utilisés dans ce type d'application du CNN ont une largeur égale à la dimension du modèle utilisé pour faire la vectorisation des tokens.

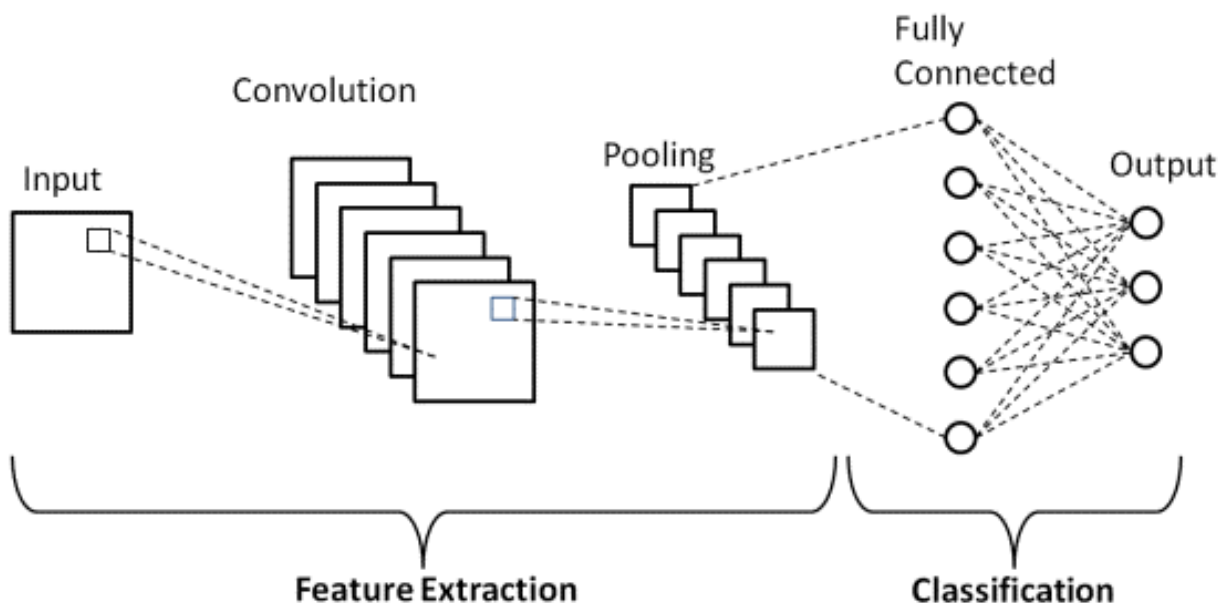


FIGURE 2.3 – Architecture basique du CNN

La figure suivante montre une configuration de base, d'autres configurations proposent l'emploi des canaux multiples en combinant plusieurs représentations vectorielles, ou juste en faisant la concaténation de ces représentations, d'autres proposent d'utiliser un entraînement statique (i.e fixer la matrice d'entrée).

2.4 Le modèle BERT

L'acronyme BERT correspond à "Bidirectional Encoder Representations from Transformers" a été développé par les chercheurs de Google AI . Comme son nom l'indique, il s'agit d'un modèle de représentation du langage qui s'appuie sur un module qu'on appelle "Transformer" pour le traitement naturel du langage (NLP). La principale innovation technique de BERT consiste à pré-entraîner des représentations bidirectionnelles profondes à partir de texte non étiqueté en conditionnant conjointement le contexte gauche et droit. En conséquence, le modèle BERT pré-entraîné peut être affiné avec une seule couche de sortie supplémentaire pour créer des modèles de pointe pour un large éventail de tâches NLP. Cela contraste avec les méthodes précédentes comme une combinaison d'apprentissage de gauche à droite et de droite à gauche.

2.4.1 L'Architecture de BERT

L'architecture BERT s'appuie sur Transformer. Nous avons actuellement deux variantes disponibles :

- **BERT Base** : 12 layers (transformer blocks), 12 attention heads, and 110 million parameters.
- **BERT Large** : 24 layers (transformer blocks), 16 attention heads and, 340 million parameters.

2.4.2 Traitement préalable de texte

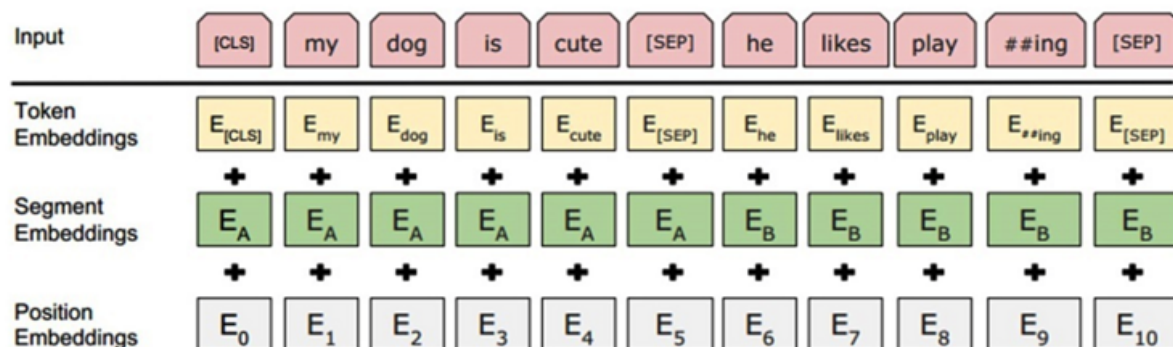


FIGURE 2.4 – Représentation d'entrée BERT

L'entrée est traitée de la manière suivante avant d'entrer dans le modèle :

- Un jeton [CLS] est inséré au début de la première phrase et un jeton [SEP] est inséré à la fin de chaque phrase.
- Incorporation de jetons : il s'agit des incorporations apprises pour le jeton spécifique à partir du vocabulaire de jeton WordPiece.
- Embeddings de segments : BERT peut également prendre des paires de phrases comme entrées pour les tâches (Question-Answering). C'est pourquoi il apprend une incorporation unique pour la première et la deuxième phrase afin d'aider le modèle à les distinguer. Dans l'exemple ci-dessus, tous les jetons marqués comme EA appartiennent à la phrase A (et de même pour EB).
- Embeddings de position : BERT apprend et utilise les incorporations de position pour exprimer la position des mots dans une phrase. Celles-ci sont ajoutées pour surmonter la limitation de Transformer qui, contrairement à un RNN, n'est pas capable de capturer des informations de « séquence » ou « d'ordre ».

2.4.3 Les tâches de pré-entraînement

BERT est pré-entraîné sur deux tâches NLP :

Modélisation de langage masqué(Masked Language Modeling)

Prédiction de la phrase suivante(Next Sentence Prediction)

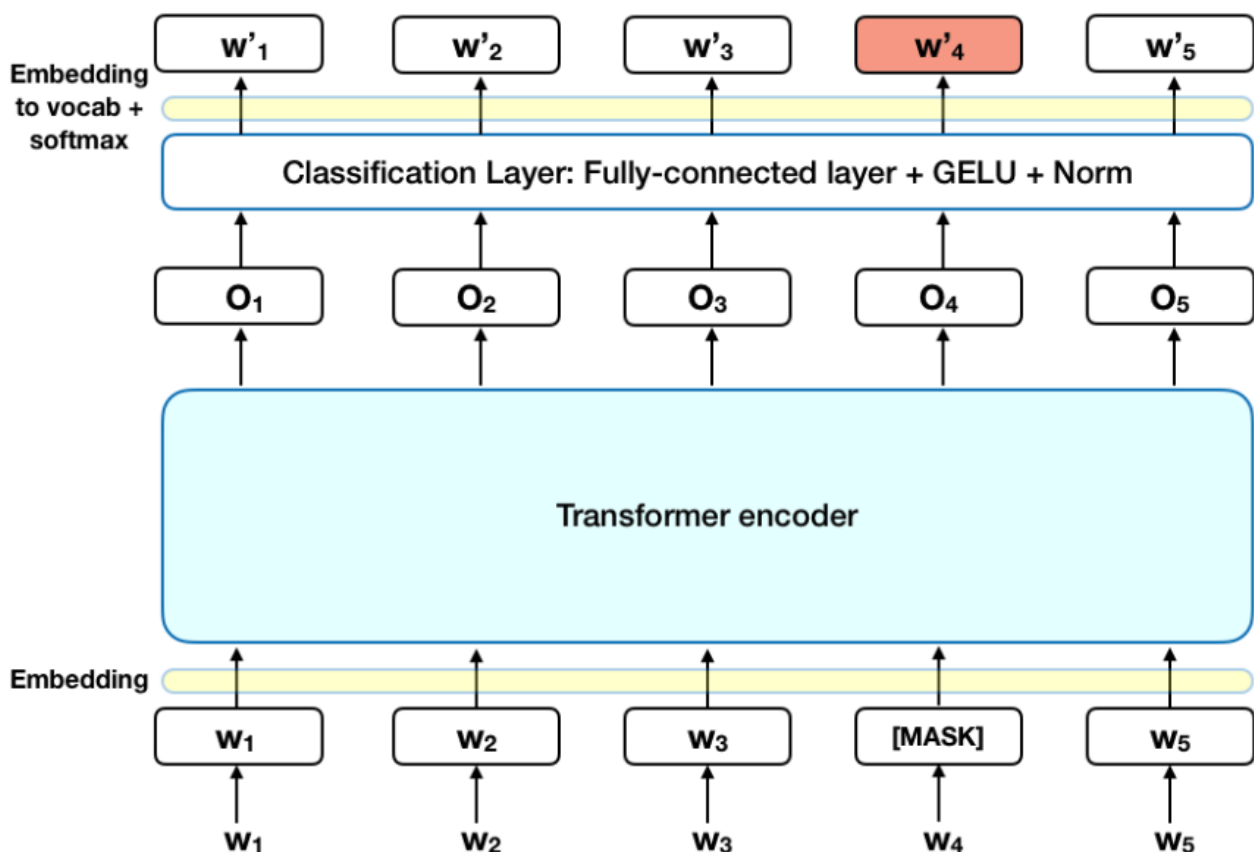


FIGURE 2.5 – Modélisation de langage masqué

Modélisation de langage masqué(Masked Language Modeling) : Avant d'alimenter des séquences de mots dans BERT, 15% des mots de chaque séquence sont remplacés par un jeton [MASK]. Le modèle tente ensuite de prédire la valeur d'origine des mots masqués, sur la base du contexte fourni par les autres mots non masqués de la séquence. En termes techniques, la prédiction des mots de sortie nécessite :

- Ajout d'une couche de classification au-dessus de la sortie de l'encodeur.
- Multiplier les vecteurs de sortie par la matrice d'incorporation, les transformer en dimension de vocabulaire.
- Calcul de la probabilité de chaque mot du vocabulaire avec softmax.

Bien que cela nous permette d'obtenir un modèle pré-entraîné bidirectionnel, un inconvénient est que nous créons un décalage entre le pré-apprentissage et le réglage fin, puisque le jeton [MASK] n'apparaît pas pendant le réglage fin. Pour atténuer cela, nous ne remplaçons pas toujours les mots « masqués » par le jeton [MASK] réel. Le générateur de données d'apprentissage choisit 15% des positions de jeton au hasard pour la prédiction. Si le i -ème jeton est choisi, on remplace le i -ème jeton par (1) le jeton [MASK] 80% du temps (2) un jeton aléatoire 10% du temps (3) le i -ème jeton inchangé 10% du temps. Ensuite, T_i sera utilisé pour prédire le jeton d'origine avec une perte d'entropie croisée.

Prédiction de la phrase suivante(Next Sentence Prediction)NSP : Dans le processus d'apprentissage BERT, le modèle reçoit des paires de phrases en entrée et apprend à prédire si la deuxième phrase de la paire est la phrase suivante dans le document original. Pendant la formation, 50% des entrées sont une paire dans laquelle la deuxième phrase est la phrase suivante dans le document original, tandis que dans les 50% restants, une phrase aléatoire du corpus est choisie comme deuxième phrase. L'hypothèse est que la phrase aléatoire sera déconnectée de la première phrase.

Pour prédire si la deuxième phrase est bien connectée à la première, les étapes suivantes sont effectuées :

- L'ensemble de la séquence d'entrée passe par le modèle Transformer.
- La sortie du jeton [CLS] est transformée en un vecteur en forme de 2×1 , en utilisant une simple couche de classification (matrices apprises de poids et de biais).
- Calcul de la probabilité d'IsNextSequence avec softmax.

2.4.4 Réglage fin(Fine-tuning)

L'utilisation de BERT pour une tâche spécifique est relativement simple : BERT peut être utilisé pour une grande variété de tâches de langage, tout en ajoutant qu'une petite couche au modèle de base :

- Les tâches de classification telles que l'analyse des sentiments sont effectuées de la même manière que la classification de la phrase suivante, en ajoutant une couche de classification au-dessus de la sortie Transformer pour le jeton [CLS].
- Dans les tâches de réponse aux questions (par exemple, SQuAD v1.1), le logiciel reçoit une question concernant une séquence de texte et doit marquer la réponse dans la séquence. En utilisant BERT, un modèle de questions-réponses peut être formé en apprenant deux vecteurs supplémentaires qui marquent le début et la fin de la réponse.
- Dans la reconnaissance d'entités nommées (NER), le logiciel reçoit une séquence de texte et doit marquer les différents types d'entités (personne, organisation, date, etc.) qui apparaissent dans le texte. A l'aide de BERT, un modèle NER peut être formé en alimentant le vecteur de sortie de chaque jeton dans une couche de classification qui prédit l'étiquette NER.

Dans la formation de réglage fin, la plupart des hyper-paramètres restent les mêmes que dans la formation BERT, il reste donc un ensemble des hyper-paramètres qui nécessitent un réglage.

Chapitre 3

Concepts et implémentation du NLP

Un framework de chatbot a besoin d'une structure dans laquelle les intentions de conversation sont définies. Un moyen propre de le faire est d'utiliser un fichier JSON.

Nous allons créer un cadre de chatbot et construire un modèle conversationnel pour un domaine de soins de santé. Le chatbot doit répondre à des questions simples sur les symptômes, les maladies, etc.

Nous allons travailler en 3 étapes :

Nous allons transformer les définitions de l'intention conversationnelle en un modèle Pytorch.

Ensuite, nous allons construire un cadre de chatbot pour traiter les réponses.

Enfin, nous montrerons comment le contexte de base peut être incorporé dans notre processeur de réponses.

Nous utiliserons tflearn, une couche au-dessus de Pytorch, et bien sûr Python.

3.1 Créer des données d'entraînement

Nous avons créé des données d'entraînement dans un fichier Json (intents.json). Il a la structure suivante :

```
intents.json X
intents.json > [ ] intents
2  ✓  "intents": [
3  ✓    {
4      "tag": "greeting",
5  ✓    "patterns": [
6      "Hi",
7      "Hey",
8      "How are you",
9      "Is anyone there?",
10     "Hello",
11     "Good day"
12   ],
13  ✓  "responses": [
14     "Hey :), How are you feeling?",
15     "Hello, thanks for visiting. Where does it hurt?",
16     "Hi there, what seems to be the matter?",
17     "Hi there, how can I help?"
18   ]
19  },
20  ✓  {
21     "tag": "goodbye",
22     "patterns": ["Bye", "See you later", "Goodbye"],
23  ✓    "responses": [
24     "See you later, thanks for visiting",
25     "Have a nice day",
26     "Bye! Come back again soon."
27   ]
28  },
29  ✓  {
30     "tag": "thanks",
31     "patterns": ["Thanks", "Thank you", "That's helpful", "Thank's a lot!"],
32     "responses": ["Happy to help!", "Any time!", "My pleasure"]
33  },
34  ✓  {
35     "tag": "headache",
36  ✓    "patterns": [
```

FIGURE 3.1 – Chatbot intents

3.2 Principes de base du NLP

Nous ne pouvons pas simplement transmettre la phrase d'entrée telle quelle à notre réseau neuronal. Nous devons d'une manière ou d'une autre convertir les chaînes de motifs en nombres que le réseau peut comprendre. Pour cela, nous convertissons chaque phrase en un bag of words (bow). Pour ce faire, nous devons collecter des mots d'entraînement, c'est-à-dire tous les mots que notre bot peut examiner dans les données d'entraînement. Sur la base de tous ces mots, nous pouvons ensuite calculer le sac de mots pour chaque nouvelle phrase. Un sac de mots a la même taille que le tableau de tous les mots, et chaque position contient un 1 si le mot est disponible dans la phrase entrante, ou 0 sinon. Voici un exemple visuel :

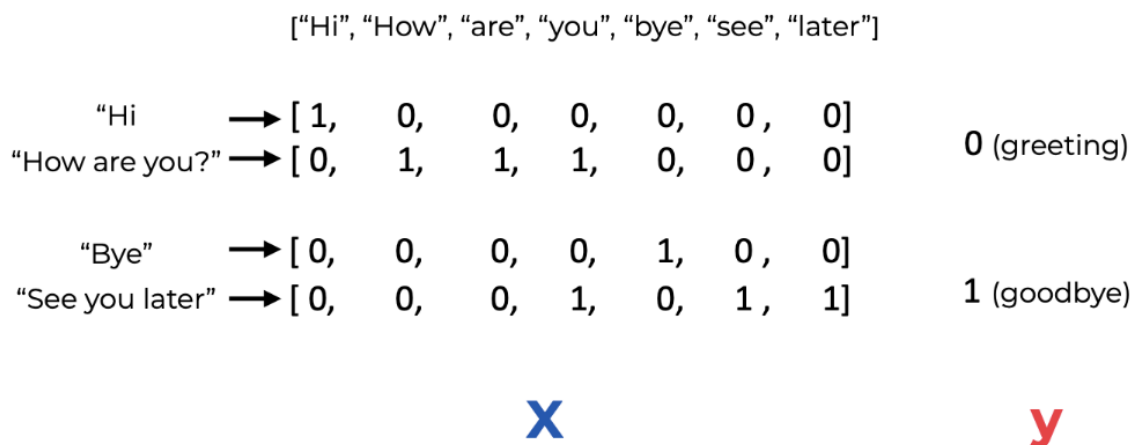


FIGURE 3.2 – Training data bag of words

Avant de pouvoir calculer l'arc, nous appliquons deux autres techniques NLP : la tokenisation et la séparation des mots.

- **Tokenisation** : Découpage d'une chaîne de caractères en unités significatives (par exemple, des mots, des caractères de ponctuation, des chiffres).
- **Stemming** : Un processus de normalisation linguistique, qui réduit les mots à leur racine ou supprime les affixes de dérivation. Par exemple, connexion, connecté, mot de connexion se réduisent à un mot commun "connecter".

Pour nos étiquettes, nous les trions par ordre alphabétique, puis nous utilisons l'index comme label de classe. L'ensemble de notre pipeline de traitement initial ressemble à ceci :

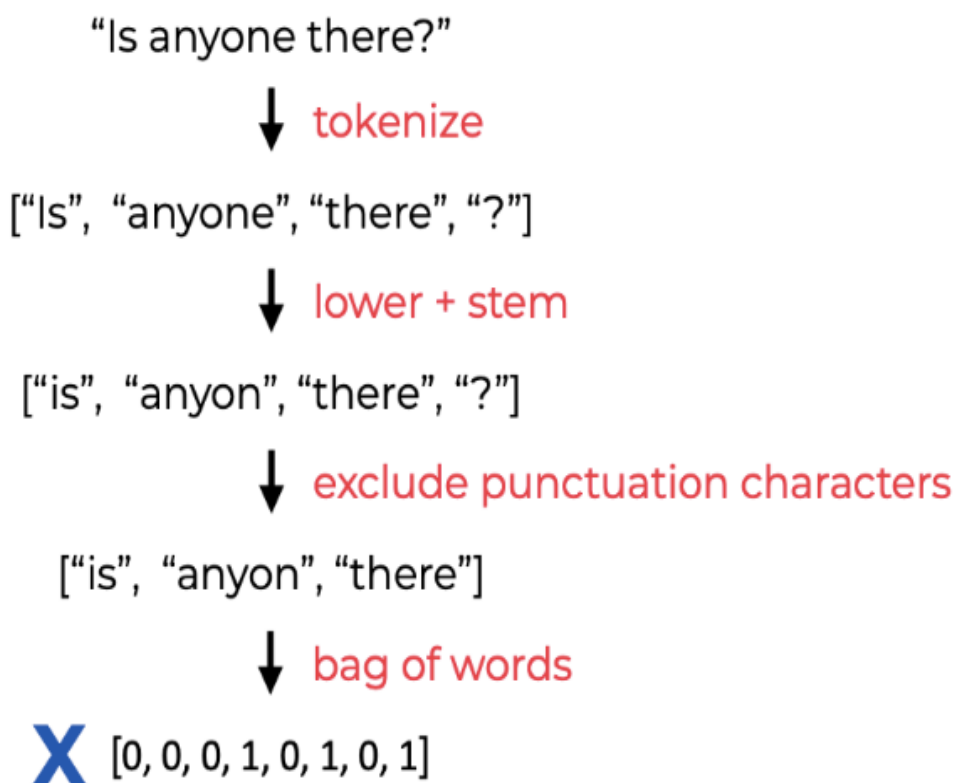


FIGURE 3.3 – NLP preprocessing pipeline

3.3 Implémenter les outils du NLP

Pour cela, nous utilisons le module `nltk`. NLTK (Natural Language Toolkit) est une plateforme de pointe pour la création de programmes Python destinés à travailler avec des données sur le langage humain.

C'est un package Python puissant qui fournit un ensemble d'algorithmes divers pour les langues naturelles. Il est gratuit, opensource, facile à utiliser, avec une grande communauté et bien documenté. NLTK comprend les algorithmes les plus courants tels que la tokenisation, l'étiquetage des parties du langage, la formation de la racine, l'analyse des sentiments, la segmentation des sujets et la reconnaissance des entités nommées. NLTK aide l'ordinateur à analyser, prétraiter et comprendre le texte écrit.

```
nltk_utils.py > ...
1  import numpy as np
2  import nltk
3  from nltk.stem.porter import PorterStemmer
4
5  stemmer = PorterStemmer()
6
7  def tokenize(sentence):
8
9      return nltk.word_tokenize(sentence)
10
11  def stem(word):
12
13      return stemmer.stem(word.lower())
14
15  def bag_of_words(tokenized_sentence, words):
16
17      # stem each word
18      sentence_words = [stem(word) for word in tokenized_sentence]
19      # initialize bag with 0 for each word
20      bag = np.zeros(len(words), dtype=np.float32)
21      for idx, w in enumerate(words):
22          if w in sentence_words:
23              bag[idx] = 1
24
25      return bag
26
```

FIGURE 3.4 – Implémentation des outils NLP

3.4 Implémentation du Neural Network

L'implémentation est simple avec un réseau neuronal Feed Forward à deux couches cachées :

```

model.py > NeuralNet > forward
1  import torch
2  import torch.nn as nn
3
4
5  class NeuralNet(nn.Module):
6      def __init__(self, input_size, hidden_size, num_classes):
7          super(NeuralNet, self).__init__()
8          self.l1 = nn.Linear(input_size, hidden_size)
9          self.l2 = nn.Linear(hidden_size, hidden_size)
10         self.l3 = nn.Linear(hidden_size, num_classes)
11         self.relu = nn.ReLU()
12
13         def forward(self, x):
14             out = self.l1(x)
15             out = self.relu(out)
16             out = self.l2(out)
17             out = self.relu(out)
18             out = self.l3(out)
19             # no activation and no softmax at the end
20         return out

```

FIGURE 3.5 – Implémentation du Neural Network

3.5 Implémentation du Training Pipeline

Dans cette partie, on va expliquer le chargement du modèle formé et les prédictions pour les nouvelles phrases.

On commence par charger notre training data qui existe dans le fichier **intents.json**, puis on passe au traitement de chaque phrase du fichier, en appliquant les méthodes du fichier **nltk_utils.py** : Tokenizing et Stemming :

```

train.py > ...
1
2 with open('intents.json', 'r') as f:
3     intents = json.load(f)
4
5 all_words = []
6 tags = []
7 xy = []
8 # loop through each sentence in our intents patterns
9 for intent in intents['intents']:
10     tag = intent['tag']
11     # add to tag list
12     tags.append(tag)
13     for pattern in intent['patterns']:
14         # tokenize each word in the sentence
15         w = tokenize(pattern)
16         # add to our words list
17         all_words.extend(w)
18         # add to xy pair
19         xy.append((w, tag))
20
21 # stem and lower each word
22 ignore_words = ['?', '.', '!']
23 all_words = [stem(w) for w in all_words if w not in ignore_words]
24 # remove duplicates and sort
25 all_words = sorted(set(all_words))
26 tags = sorted(set(tags))
27
28 print(len(xy), "patterns")
29 print(len(tags), "tags:", tags)
30 print(len(all_words), "unique stemmed words:", all_words)
31

```

FIGURE 3.6 – Chargement et traitement du training data

Ensuite, on crée notre training data d’une manière à avoir les deux vecteurs X et Y, pour avoir à la fin notre Bag of Words list :

```

1 # create training data
2 X_train = []
3 y_train = []
4 for (pattern_sentence, tag) in xy:
5     # X: bag of words for each pattern_sentence
6     bag = bag_of_words(pattern_sentence, all_words)
7     X_train.append(bag)
8     # y: PyTorch CrossEntropyLoss needs only class labels, not one-hot
9     label = tags.index(tag)
10    y_train.append(label)
11
12 X_train = np.array(X_train)
13 y_train = np.array(y_train)
14

```

FIGURE 3.7 – Création du bag of words

Enfin, on applique le training model au dataset :

```
train.py > [?] train_loader
80 dataset = ChatDataset()
81 train_loader = DataLoader(dataset=dataset,
82                             batch_size=batch_size,
83                             shuffle=True,
84                             num_workers=0)
85
86 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
87
88 model = NeuralNet(input_size, hidden_size, output_size).to(device)
89
90 # Loss and optimizer
91 criterion = nn.CrossEntropyLoss()
92 optimizer = torch.optim.Adam(model.parameters(), lr=learning_rate)
93
94 # Train the model
95 for epoch in range(num_epochs):
96     for (words, labels) in train_loader:
97         words = words.to(device)
98         labels = labels.to(dtype=torch.long).to(device)
99
100         # Forward pass
101         outputs = model(words)
102         # if y would be one-hot, we must apply
103         # labels = torch.max(labels, 1)[1]
104         loss = criterion(outputs, labels)
105
106         # Backward and optimize
107         optimizer.zero_grad()
108         loss.backward()
109         optimizer.step()
110
111     if (epoch+1) % 100 == 0:
112         print (f'Epoch [{epoch+1}/{num_epochs}], Loss: {loss.item():.4f}')
113
```

FIGURE 3.8 – Application du training model

Vers la fin du processus du trainig, on sauvegarde notre data résultante dans un fichier **data.pth**.

3.6 Implémentation du Chat

Dans cette partie, on va charger le modèle entraîné et faites des prédictions pour les nouvelles phrases.

On commence d'abord par récupérer les données du fichier **data.pth**.

La fonction `chatbot_response` prend comme paramètre la phrase entrée par l'utilisateur du chatbot et applique les outils du `nlTK`, Tokenizing et Stemming. Elle calcule ensuite la probabilité d'avoir un élément existant parmi les tags de notre training data, et fait un test pour savoir si le tag du pattern entré par l'utilisateur correspond à un tag dans `data.pth`. Enfin elle renvoie une réponse selon le résultat du test.

```

chat.py > ...
37 def chatbot_response(sentence):
38
39     #if sentence == "results":
40     |     #res = listToString(what_are_your_symptoms(filter_disease(intent_tags_list)))
41
42     sentence = tokenize(sentence)
43     X = bag_of_words(sentence, all_words)
44     X = X.reshape(1, X.shape[0])
45     X = torch.from_numpy(X).to(device)
46
47     print("\n *****")
48     print(sentence)
49     print("\n *****")
50
51     output = model(X)
52     _, predicted = torch.max(output, dim=1)
53
54     tag = tags[predicted.item()]
55     print(tag)
56
57     probs = torch.softmax(output, dim=1)
58     prob = probs[0][predicted.item()]
59     if prob.item() > 0.75:
60         for intent in intents['intents']:
61             if tag == intent["tag"]:
62                 res = random.choice(intent['responses'])
63                 intent_tags_list.append(tag)
64     elif "results" in sentence:
65         dic= test.what_are_your_symptoms(filter_disease(intent_tags_list))
66         res=dic
67         print("\n ****//****")
68         print(dic)
69         print("\n ****//****")
70     else:

```

FIGURE 3.9 – Génération des réponses du chatbot

Lorsque l'utilisateur envoie "results" comme phrase, il reçoit comme réponse le symptôme le plus proche aux description donnés. La fonction `what_are_your_symptoms` prend en paramètre la liste de symptômes entrés par l'utilisateur, fait un calcul de probabilité d'avoir une des maladies citées dans le fichier **disease.json** en ce basant sur les symptômes de chaque maladies, et retourne à la fin un dictionnaire qui contient une clé indiquant la maladie et sa valeur, et d'une autre clé qui indique la probabilité et sa valeur.

```

test.py > what_are_your_symptoms
3 def what_are_your_symptoms(patient_symptoms):
4     print(patient_symptoms)
5     """patient_symptoms=[
6         "pain chest",
7         "shortness of breath",
8         "dizziness",
9     ]"""
10
11     potential_diagnosis = []#[[disease,percent]]
12     proba = []#these two are parrallel lists
13     with open('disease.json','r') as f:
14         diseases = json.load(f)
15         for patient_symptom in patient_symptoms:
16             for disease in diseases:
17                 if patient_symptom in diseases[disease]['symptoms']:
18                     total_symp = len(diseases[disease]['symptoms'])
19                     print("symptoms",diseases[disease]['symptoms'])
20                     if diseases[disease]['disease'] in potential_diagnosis:
21                         i = potential_diagnosis.index(diseases[disease]['disease'])
22                         proba[i] += 1/total_symp
23                     else:
24                         potential_diagnosis.append(diseases[disease]['disease'])
25                         proba.append(1/total_symp)
26
27     if proba == []:
28         j = 0
29         potential_diagnosis.append(["you are maybe exhausted",])
30         mx = "?"
31     else:
32         mx = max(proba)
33         mx = "{:.4f}".format(mx)
34         j = proba.index(mx)
35         #print(potential_diagnosis)
36     return {"illness":potential_diagnosis[j][0],"probability":mx}

```

FIGURE 3.10 – Le résultat final du conseil

Chapitre 4

Réalisation finale du Chatbot

4.1 Front-end

Conclusion

Pour conclure, il suffit de dire que le but ultime des chatbots est d'imiter et de s'aligner autant que possible sur une conversation naturelle entre humains. En outre, lorsque nous concevons le flux conversationnel d'un chatbot, nous oublions souvent les éléments qui font partie intégrante d'une véritable conversation humaine.

La digression est une partie importante de la conversation humaine, avec la désambiguïsation bien sûr. La désambiguïsation annule dans une certaine mesure le danger de la prolifération de fallback où le dialogue n'avance pas vraiment.

Avec la désambiguïsation, un bouquet d'options réellement liées et contextuelles est présenté à l'utilisateur, parmi lesquelles il peut choisir, ce qui est sûr de faire avancer la conversation.

Enfin, la pire chose que vous puissiez faire est de présenter un ensemble d'options qui n'est pas lié au contexte actuel ou encore un ensemble d'options prédéfinies et finies qui se répètent constamment. La conscience contextuelle est essentielle dans tous les éléments d'un chatbot.

Bibliographie

- [1] Pytorch documentation. <<https://pytorch.org/docs/stable/index.html>>.
- [2] Jack Cahn. Chatbot : Architecture, design, development.
- [3] NLTK 3.6.2 documentation. <<https://www.nltk.org/>>.
- [4] gk. Contextual chatbots with tensorflow. <<https://chatbotsmagazine.com/contextual-chat-bots-with-tensorflow-4391749d0077>>.
- [5] L. Gutiérrez and B. (2019). Keith. A systematic literature review on word embeddings.
- [6] Salma Jamoussi. Une nouvelle représentation vectorielle pour la classification sémantique. <https://www.researchgate.net/publication/228790563_Une_nouvelle_representation_vectorielle_pour_la_classification_semantique>.
- [7] Y. (2014) Kim. Convolutional neural networks for sentence classification.
- [8] Chen K. Corrado G. Mikolov, T. and J. (2013). Dean. Efficient estimation of word representations in vector space.
- [9] E. G. Nal Kalchbrenner and P. (2014). Blunsom. A convolutional neural network for modelling sentences.
- [10] B. C. W. (2016). Ye Zhang. A sensitivity analysis of (and practitioners' guide to) convolutional neural networks for sentence classification.