

République Algérienne Démocratique et Populaire
Ministère de l'Enseignement supérieur et de la Recherche Scientifique

Université d'Alger 2

Faculté des Lettres et des Langues - Département de Linguistique

MEMOIRE DE MAGISTER

INFORMATIQUE LINGUISTIQUE-

SCIENCES DU LANGAGE ET DE LA COMMUNICATION LINGUISTIQUE

OPTION : TRAITEMENT AUTOMATIQUE DU LANGAGE

*Résolution d'anaphores pronominales de la langue
arabe par apprentissage automatique*

Présenté par : **Leïla NOURI**

Jury

| | | | |
|-----------------|-------------|---------------------------------|--------------|
| M ^r | A. SALMI | Docteur à l'université d'Alger | Président |
| M ^r | A. GUESSOUM | Professeur à l'USTHB | Encadreur |
| M ^{me} | H. ALIANE | Chargée de Recherche au CERIST | Examinatrice |
| M ^{me} | L. MAHDAOUI | Maitre de conférences à l'USTHB | Examinatrice |

Année universitaire 2011/2012

Remerciements

الحمد والشكر لله تعالى

Je tiens à exprimer ma profonde gratitude et mes sincères remerciements au Professeur Ahmed GUESSOUM pour m'avoir proposée ce sujet, pour sa disponibilité et pour ses précieux conseils.

Mes remerciements s'adressent aussi au président du jury Dr Abdelmadjid SALMI et aux membres du jury Dr Hassina ALIANE et Dr Latifa MAHDAOUI, pour leur disponibilité à évaluer le présent mémoire et leur participation au jury de soutenance.

Mes vifs remerciements vont aussi à ma famille, mes amis, et mes collègues du CERIST, qui m'ont toujours soutenue et encouragée au cours de la réalisation de ce mémoire.

ملخص البحث

العائد هو كلمة لها سابقة أو عدة سوابق في النص. حل العائد في اللغة العربية بواسطة التعلم الآلي هو إيجادا بطريقة آلية سابقة العائد باستعمال كم هائل من الأمثلة النصية للعملية.

في الواقع، تواتر تواجد العوائد في نصوص اللغة العربية كبير جدا، الأمر الذي يشكل إشكالية حقيقية للعلاج الآلي للمعطيات النصية، مثل التلخيص الآلي أو استرجاع المعلومات، أو الرد على استعلام في نظم حوار بين الإنسان والآلة، الخ.

ويعتبر التعلم التلقائي للشبكات العصبونية واحدا من الحلول المثيرة للاهتمام لحل العائد .

في هذا العمل قمنا بدراسة مسألة حل العائد في اللغة العربية و قدمنا حلا لمعالجته حاسوبيا عن طريق تحشية مدونة و تشفيرها أليا، ثم استعمالها داخل الشبكات العصبونية باستخدام الشبكة ذات طبقات متعددة بطريقة التغذية الأمامية واستخدام خوارزمية الانتشار الخلفي (للتصحيح داخل الشبكة).

الكلمات الدالة: العائد، حل العائد، اللغة العربية، تحشية المدونة، وسم، الشبكات العصبونية، الشبكة ذات الطبقات المتعددة الأمامية

Résumé

Une anaphore est un mot qui a obligatoirement un ou plusieurs antécédents qui lui réfèrent dans un énoncé. L'objectif de la résolution d'anaphores pronominales de la langue arabe par apprentissage automatique est de trouver de façon informatisée l'antécédent d'une anaphore et ce par apprentissage sur un maximum d'exemples de textes.

En effet, la fréquence d'apparition des anaphores dans un texte arabe est très grande, ce qui pose un problème pour un traitement automatique des données textuelles, tel que l'élaboration d'un résumé de texte ou d'une recherche d'information, ou pour répondre à une requête dans les systèmes de dialogue homme-machine, etc.

L'apprentissage automatique des réseaux de neurones est considéré comme une des solutions intéressantes pour la résolution d'anaphores.

Dans ce travail, nous nous intéressons à la résolution d'anaphores de la langue arabe sur un corpus annoté et codifié automatiquement, soumis à l'apprentissage des réseaux de neurones en utilisant le perceptron multi-couches selon l'algorithme de la rétro-propagation du gradient.

MOTS-CLÉS : anaphores, résolution d'anaphores pronominales, langue arabe, annotation, étiquettes, corpus, réseaux de neurones, perceptron multi-couches.

Abstract

An anaphora is a word that has one or more antecedents that refer to it in a sentence or paragraph. The objective of pronominal anaphora resolution is to find the antecedent of an anaphora in a computerized way based on the use of as large a corpus of examples of manual anaphora resolution as possible.

Indeed, the frequency of occurrence of anaphora in arabic text is very large, which turns out to be an issue that needs to be handled before automatic of textual data, so as to develop quality automatic text summaries information retrieval, or query answers in human-machine dialogue systems.

Automatic training of neural networks is considered as one of the interesting solutions for the resolution of anaphora.

In this work, we are interested in anaphora resolution of the arabic language on an automatically annotated and coded corpora, and a neural network which is trained as a feed-forward multilayer perceptron with backpropagation.

KEYWORDS : *anaphora, anaphora resolution, corpora, arabic language, annotation, tags, neural networks, multilayer perceptron.*

Sommaire

Liste des Figures

Liste des Tableaux

| | |
|---|-----------|
| Introduction Générale..... | i |
| 1. Introduction..... | ii |
| 2. Objectif du mémoire..... | iii |
| 3. Organisation du mémoire..... | iii |
| | |
| CHAPITRE 1 : Traitement Automatique du Langage Naturel | 5 |
| 1. Introduction..... | 5 |
| 2. Les domaines d'application du TALN..... | 6 |
| 2.1 Recherche d'Information | 6 |
| 2.2 Extraction d'Information | 7 |
| 2.3 Systèmes Question-Réponse | 7 |
| 2.4 Systèmes de résumé automatique de textes..... | 8 |
| 2.5 Traduction automatique..... | 8 |
| 2.6 Dialogue homme-machine | 8 |
| 2.7 Génération automatique de textes | 9 |
| 3. Difficultés et ambiguïtés du TALN | 9 |
| 3.1 Ambiguïté morphologique | 9 |
| 3.2 Ambiguïté Syntaxique..... | 9 |
| 3.3 Ambiguïté sémantique et pragmatique..... | 10 |
| 4. Conclusion..... | 10 |
| | |
| CHAPITRE 2 : La Résolution d'anaphores pronominales | 11 |
| 1. Introduction..... | 11 |
| 2. L'anaphore | 12 |
| 3. L'anaphore vs La coréférence | 13 |
| 4. La typologie des anaphores dans la langue arabe | 13 |

| | |
|---|----|
| 4.1 L'anaphore pronominale | 13 |
| 4.1.1 Les pronoms personnels | 13 |
| 4.1.2 Les pronoms relatifs | 15 |
| 4.1.3 Les pronoms démonstratifs | 16 |
| 4.2 L'anaphore lexicale | 16 |
| 4.3 L'anaphore comparative..... | 17 |
| 4.4 L'anaphore verbale..... | 17 |
| 5. Le processus général de la résolution d'anaphores | 17 |
| 5.1 Les connaissances morphologiques et lexicales..... | 18 |
| 5.2 Les connaissances syntaxiques..... | 18 |
| 5.3 Les connaissances sémantiques et pragmatiques | 18 |
| 6. Les mesures d'évaluation dans les systèmes de résolution d'anaphores | 19 |
| 7. Difficulté de la résolution des anaphores | 21 |
| 8. Les différentes approches de la résolution d'anaphores pronominales | 22 |
| 8.1 Présentation | 22 |
| 8.2 L'approche basée sur la syntaxe..... | 23 |
| 8.2.1 L'algorithme de Hobbs..... | 23 |
| 8.2.2 Avantages et inconvénients de l'algorithme de Hobbs | 24 |
| 8.3 L'approche basée sur le discours | 24 |
| 8.3.1 La Théorie du Centrage..... | 24 |
| 8.3.2 L'algorithme BFP..... | 26 |
| 8.3.3 Avantages et inconvénients de la Théorie du Centrage | 28 |
| 8.4 L'approche basée sur le facteur..... | 28 |
| 8.4.1 Lappin et Leass..... | 29 |
| 8.4.2 Kennedy et Boguraev | 30 |
| 8.4.3 GogNIAC de Baldwin..... | 30 |
| 8.4.4 L'approche de Mitkov | 31 |
| 8.5 Les approches basées sur l'apprentissage automatique..... | 35 |
| 8.5.1 Aone et Bennett, McCarthy et Lehnert | 36 |
| 8.5.2 Connolly et. al. | 36 |
| 8.5.3 Cardie et Wagstaff..... | 38 |
| 8.5.4 Soon et. al. | 38 |
| 8.5.5 Strube et. al..... | 39 |
| 8.6 Les approches statistiques | 40 |

| | |
|--|-----------|
| 8.7 Autres approches | 40 |
| 9. Les domaines d'application de la résolution d'anaphores..... | 42 |
| 9.1 La résolution d'anaphores dans l'extraction d'information | 42 |
| 9.2 La résolution d'anaphores dans le résumé automatique du texte | 43 |
| 9.3 La résolution d'anaphores dans la traduction automatique | 44 |
| 9.4 La résolution d'anaphores dans les systèmes question-réponse..... | 44 |
| 10. La résolution des anaphores pronominales dans la langue arabe | 45 |
| 10.1 L'approche multilingue de Mitkov | 45 |
| 10.2 La résolution d'anaphores dans l'Arabe utilisant le web comme corpus..... | 46 |
| 11. Conclusion | 46 |
| | |
| CHAPITRE 3 : Introduction aux réseaux de neurones..... | 47 |
| 1. Introduction..... | 47 |
| 2. Neurone Formel..... | 47 |
| 3. Réseau de Neurones..... | 49 |
| 4. Architecture des réseaux de neurones..... | 49 |
| 4.1 Le perceptron monocouche | 49 |
| 4.2 Le perceptron multi-couches | 50 |
| 4.3 Les réseaux à fonction radiale de base | 51 |
| 5. Apprentissage des réseaux de neurones | 53 |
| 5.1 Apprentissage supervisé..... | 53 |
| 5.2 Apprentissage non-supervisé..... | 55 |
| 5.3 Apprentissage par renforcement..... | 55 |
| 6. Styles d'apprentissage | 55 |
| 6.1 Batch training | 56 |
| 6.2 Incremental training | 56 |
| 7. Avantages et inconvénients des réseaux de neurones | 55 |
| 8. Conclusion..... | 56 |
| | |
| CHAPITRE 4 : Annotation du corpus | 58 |
| 1. Introduction..... | 58 |
| 2. Caractéristiques morphologiques de l'Arabe | 59 |
| 3. Etiquetage morphosyntaxique (POS Tagging)..... | 60 |
| 4. Outil d'étiquetage du corpus..... | 60 |

| | |
|--|-----------|
| 4.1 SAIE POS-Tagger | 61 |
| 4.2 Outil d'annotation et codification d'un corpus en langue Arabe (ANCARTool) | 61 |
| 5. Conclusion..... | 62 |
| CHAPITRE 5 : Approche pour la résolution d'anaphores de la langue arabe | 64 |
| 1. Introduction..... | 64 |
| 2. Codification du corpus..... | 65 |
| 3. Description du Réseau de neurones | 70 |
| 3.1 Le corpus d'apprentissage..... | 71 |
| 3.2 Création du perceptron multi-couches | 71 |
| 3.3 Apprentissage du perceptron multi-couches | 72 |
| 3.3.1 Fonction d'apprentissage..... | 73 |
| 4. Résultats d'apprentissage du perceptron multi-couches | 73 |
| 5. Interprétation des résultats | 75 |
| 5.1 Algorithme d'interprétation des résultats | 76 |
| 5.2 Evaluation des résultats | 79 |
| 6. Conclusion..... | 80 |
| Conclusion & Perspectives | 81 |
| Références Bibliographiques | 81 |
| Annexe | 92 |