

Analyse syntactico-sémantique de la langue amazighe

Siham Boulaknadel

IRCAM, Avenue Allal El Fassi, Madinat Al Irfane,

Rabat-Instituts, Maroc

E-mail: siham_06@yahoo.fr

Meryem TALHA

LRIT – CNRST URAC29, Université Mohammed V-
Agdal Rabat

4, Avenue Ibn Battouta, B.P. 1014 RP, 10006 Rabat,
Maroc

E-mail: meriem.talha@gmail.com

Abstract—L'augmentation des flux de l'information numérique nécessite l'extraction, le filtrage et la classification des informations pertinentes à partir de grands volumes de textes. Toutes ces tâches bénéficient amplement de l'implication de la Reconnaissance des entités nommées (NER) dans l'étape de prétraitement. La tâche NER implique non seulement l'identification des noms propres (entités nommées) dans le texte en langage naturel, mais aussi leur catégorisation en un ensemble de catégories prédéfinies, telles que les noms de personnes, les organisations (entreprises, organismes gouvernementaux, comités, etc), Lieux (villes, pays, fleuves, etc) et divers (titres de films, événements sportifs, etc.) Dans les langues anglaises et françaises, beaucoup de travaux ont été effectués dans ce domaine, où la capitalisation est un indice important pour identifier les entités nommées, alors que la langue amazighe n'a pas cette particularité. Elle souffre d'un manque de disponibilité de ressources et de répertoires toponymiques, se caractérise par sa nature agglutinante, et sa morphologie flexionnelle, ce qui rend la tâche de la reconnaissance plus difficile pour la langue amazighe. Les travaux sur les outils linguistiques automatiques sont une opportunité pour la valorisation de l'amazighe dans la société de l'information. C'est dans ce contexte que s'inscrit cette thèse qui consiste à construire des outils d'analyse syntactico-sémantique automatique indispensables pour l'analyse de textes amazighs. Dans cet article, nous décrivons les caractéristiques de la langue amazighe, nous présentons une étude des différentes approches pour la détection des entités nommées (NE).

Keywords: Reconnaissance des entités nommées; langue amazighe; approches NER;

I. INTRODUCTION

L'extraction d'entités nommées est une sous tâche du domaine d'extraction d'information. Concrètement, la notion de la reconnaissance d'entités nommées (NER) est apparue comme un outil de prétraitement important pour nombreuses applications du traitement de langage naturelle, telles que l'Extraction d'information (IE), la recherche d'information multilingue (IR), la traduction automatique (MT), systèmes de recherche de réponses précises à des questions « question-réponse » (QA), résumé automatique (AS), indexation automatique des documents (AI) entre autres applications de traitement de texte, parce que les entités nommées fournissent principalement des indices importants pour l'identification des informations pertinentes dans le texte. Elle vise à localiser (par exemple, « Le Chef du gouvernement, M. Abdel-Ilah Benkirane »), et classer les éléments de texte

dans des catégories prédéfinies, c'est-à-dire le système va affecter l'étiquette personne au NE « Le Chef du gouvernement, M. Abdel-Ilah Benkirane » en utilisant des méthodes statistiques, fondées sur des règles ou hybrides.

Les entités nommées désignent traditionnellement les noms de personnes, de lieux, d'organisations mais aussi les expressions de dates ou les unités monétaires, les pourcentages, les fonctions et autres mais peuvent aussi se rapporter à des notions plus techniques comme les maladies. (Ehrmann, 2008) analyse la problématique des EN du point de vue théorique des difficultés définitives et catégorielles, proposant la définition suivante : « étant donné un modèle applicatif et un corpus, on appelle EN toute expression linguistique qui réfère à une entité unique du modèle de manière autonome dans le corpus ».

A. Caractéristiques de la langue amazighe :

Au Maroc, la langue berbère ou amazighe (ⵜⴰⴳⴷⵓⴷⴰⵢⵜ), fait partie des langues chamito-sémitiques ou encore appelé afro-asiatiques (Cohen 2007, Chaker 1989) a toujours disposé d'un statut restreint et triplement minoré, bien qu'elle est considérée comme la plus ancienne langue d'Afrique du Nord, et fortement employé sur une aire très vaste par environ 50% (Boukous, 1995), de la population marocaine, mais soumis à une pression conjuguée de l'arabe et du français, autrement dit, elle est régie par l'arabe standard (la langue officielle), l'arabe dialectal marocain et le français.

Sur le plan linguistique, la langue est caractérisée par la prolifération (Outahajala 2012) des dialectes en raison de facteurs historiques, géographiques et sociolinguistiques. Au Maroc, on peut classer la langue berbère en trois grandes zones dialectales capitales dont les trois principales langues vernaculaires peuvent être identifiées: Tarifit (Nord du Maroc), Tamazight (centrale) et Tachelhit (Sud du Maroc). Pour pallier cette situation, vint la création de l'Institut Royal de la Culture Amazighe en 2001 (IRCAM), dans le but était de standardiser, redonner aux culture et langue amazighes la place qu'elles méritent, ainsi que pour uniformiser les structures et à adoucir les pluralités qui représentent des difficultés au niveau de l'intercompréhension.

L'Amazighe est maintenant une langue qui possède tous ses attributs: dotée d'une graphie officielle, un codage propre dans le standard Unicode, une grammaire, une orthographe ainsi qu'un vocabulaire très riche et une littérature orale fabuleusement fortuné. Elle est actuellement une langue nationale et depuis 2003, elle a été introduite dans le système éducatif comme matière

○ Un complément d'objet direct « ⵏⵉⵎⵓ ⵏ ⵓⵎⵓⵏ ⵏ ⵓⵎⵓⵏ » il tient un poisson à la main.

○ Un complément de la particule prédictive Λ = c'est : Λ ⵏⵉⵎⵓ "c'est un poisson".

Etat lié (EA) : qui se manifeste par une variation des noms à initiale vocalique.

○ alternance vocalique ⵏ / ⵙ pour les noms masculins \Rightarrow ⵏⵓⵔⵓⵎ → ⵙⵓⵔⵓⵎ.

○ chute de la voyelle initiale pour les noms féminins \Rightarrow ⵏⵓⵔⵓⵎⵓ → ⵙⵓⵔⵓⵎⵓ

○ addition d'un Λ ou Υ aux noms à voyelle ⵏ ou Υ \Rightarrow ⵏⵓⵔⵓⵎ → ⵏⵓⵔⵓⵎⵓ

○ maintien de la voyelle initiale ⵏ avec apparition de la semi-consonne Λ pour le masculin, les noms féminins restent sans modification. Par exemple :

Masculin \Rightarrow ⵏⵓⵔⵓⵎ "jour" → ⵏⵓⵔⵓⵎⵓ

Féminin \Rightarrow ⵏⵓⵔⵓⵎⵓ "foie" → ⵏⵓⵔⵓⵎⵓ

L'état d'annexion se réalise dans les contextes syntaxiques suivants :

○ Le sujet lexical suit le verbe \Rightarrow ⵙⵓⵔⵓⵎ ⵏ ⵙⵓⵔⵓⵎⵓ → "Le professeur est venu".

○ Après une préposition \Rightarrow ⵏⵓⵔⵓⵎⵓ ⵏ ⵙⵓⵔⵓⵎⵓ → "J'ai parlé au professeur".

○ Après un coordonnant \Rightarrow ⵏⵓⵔⵓⵎⵓ ⵏ ⵙⵓⵔⵓⵎⵓ → "la pluie et le froid".

Verbe :

Le verbe peut prendre deux aspects : une forme simple ou dérivée, qu'il soit simple ou dérivé, il se conjugue à l'un des quatre thèmes: l'aoriste, l'inaccompli, l'accompli positif et l'accompli négatif.

Le verbe simple est composé d'une racine et d'un radical. Par contre Les verbes dérivés sont obtenus à partir des verbes simples par la préfixation de l'un des morphèmes suivants : ⵏⵓⵔⵓⵎ, ++ et Υ / Υ .

○ **Les dérivés en ⵏⵓⵔⵓⵎ :** Ils expriment la causativité, autrement dit l'idée de "faire faire" une action à un patient ou celle de "faire devenir". Par exemple :

ⵏⵓⵔⵓⵎ "accompagner" → ⵏⵓⵔⵓⵎⵓ "faire accompagner", ⵏⵓⵔⵓⵎ "entrer" → ⵏⵓⵔⵓⵎⵓ "faire entrer".

○ **Les dérivés en ++ :** on procède à la préfixation de ++ au verbe simple pour exprimer la forme passive. Ce morphème peut se présenter sous forme de ++ⵓ, ++ⵏⵓ ou ++ⵙⵓ. Par exemple :

ⵏⵓⵔⵓⵎ "attacher" → ++ⵏⵓⵔⵓⵎⵓ / ++ⵙⵓⵔⵓⵎⵓ "être attaché"

ⵏⵓⵔⵓⵎ "tenir" → ++ⵏⵓⵔⵓⵎⵓ / ++ⵙⵓⵔⵓⵎⵓ "être tenu"

○ **Les dérivés en mm :** Un verbe dérivé par la préfixation de Υ ou ses variantes (Υ , Υ) exprime la réciprocité : ⵏⵓⵔⵓⵎ "voir" → Υ ⵏⵓⵔⵓⵎ "se voir".

Les particules :

Les particules sont un ensemble de mots bien évidemment amazighe, en général, assez courts qui jouent le rôle d'indicateurs grammaticaux au sein d'une phrase.

Cet ensemble contient plusieurs formes: les particules d'aspect, d'orientation et de négation; les pronoms indéfinis, démonstratifs, possessifs et interrogatifs; les pronoms personnels autonomes, affixes sujet, affixes d'objet direct et indirect, compléments du nom ordinaire et de parenté, compléments de prépositions; les adverbes de lieu, de temps, de quantité et de manière; les prépositions; les subordonnants et les conjonctions.

Généralement, les particules sont invariables. Or, dans le cas de l'Amazighe, similairement au cas du français, il existe des particules flexionnelles telles que les pronoms possessifs (Υ ⵏⵓⵔⵓⵎ "le sien" → Υ ⵏⵓⵔⵓⵎⵓ "le leur").

B. Difficultés entravant le traitement automatique de la langue amazighe :

Les travaux liés à NER dans la langue amazighe sont encore à l'état naissant en raison des faits potentiels: (i). L'absence de la distinction majuscule/minuscule : c'est un obstacle majeur pour la langue amazighe. En fait, la REN pour certaines langues comme les langues indo-européennes se base principalement sur la présence des lettres majuscules qui est un indicateur très utile pour identifier les noms propres dans les langues utilisant l'alphabet latin. Les lettres majuscules, néanmoins, ne se produisent pas, ni au début ni à l'initiale des noms propres amazighe.

(ii). Il est un fait que la langue amazighe est très agglutinative ayant une morphologie dérivationnelle et flexionnelle assez complexe et riche, les noms peuvent avoir plusieurs formes fléchies et dérivées, la simple suppression des suffixes ne peut suffire à regrouper des familles de mots. En effet, dans la pratique les affixes peuvent altérer le sens d'un mot. Donc, le fait de les éliminer peut engendrer une perte d'information.

(iii). Similairement à d'autres langues naturelles, l'amazighe présentent des incertitudes au niveau des classes grammaticales. En effet, la même forme convient à nombreuses catégories grammaticales, cela dépend du contexte dans la phrase. Par exemple, (Υ ⵏⵓⵔⵓⵎ) peut être considéré comme verbe à l'accompli négatif, il signifie «il existe», ou comme nom de parenté « ma fille ».

(iv). L'amazighe se caractérise par le manque de ressources dictionnaires (de noms etc.), de répertoires toponymiques, de ressources langagières et outils du TAL, à savoir les analyseurs morphologiques, POS taggers.

(v). Manque de corpus annotés qui ont une très grande importance dans diverses tâches du traitement automatique des langues.

(vi). Le nombre de mots fréquemment utilisés comme noms communs et qui peuvent être, également, utilisés comme noms propres est très grand.

(vi). Les substantifs, les adjectifs, les verbes, les pronoms, les adverbes, prépositions, les particules... etc, consistent en un seul mot se produisant entre deux blancs ou des signes de ponctuation (Ameur et al. 2006). Cependant, si une préposition ou un nom de parenté est suivie par un pronom personnel, à la fois la préposition/nom de parenté et le pronom qui suit, forment une chaîne unique délimitée par des espaces ou des signes de ponctuation. Par exemple: (Υ ⵏⵓ) signifiant « pour, au », (Υ) qui signifie « moi » (pronom personnel), donnent (Υ ⵏⵓⵔⵓⵎ ou Υ ⵏⵓⵔⵓⵎⵓ).

(vii). L'état d'annexion :

(viii). Les noms propres au niveau de la langue amazighe sont extrêmement nombreux, ont de nombreuses variantes ainsi qu'ils sont difficiles à détecter sans la présence d'un lexique, et c'est également le cas pour les noms d'organisations ou de produits, les noms de lieux même si ces derniers sont relativement stables par rapport aux

Les chercheurs, de nos jours, penchent populairement vers l'utilisation des méthodes d'apprentissage automatique pour les REN car ces dernières sont facilement adaptables à différents domaines et langues ainsi que leur maintenance est aussi moins coûteuse.

Plusieurs systèmes d'apprentissage automatique ont été traités depuis quelques années. Nous distinguons les méthodes d'apprentissage selon qu'elles soient supervisées, semi-supervisé ou non-supervisé.

• **Méthodes supervisées :**

Il s'agit d'une méthode d'apprentissage qui nécessite une intervention humaine importante, et requiert tout de même d'une manière globale, un volume considérable de données pour les nécessités de l'apprentissage. Elles consistent à faire un apprentissage à partir d'un corpus d'entraînement, préalablement préparé, parmi ces méthodes on peut citer :

▪ **HMM (Hidden Markov Model) :**

Un modèle HMM (Hidden Markov Model) est un modèle stochastique particulier, génératif on modélise comment obtenir x alors que les données x sont disponibles, proposé pour la première fois à la fin des années 60 début des années 70 (Baum. 1970, Baker. 1975), Il représente un objet donné par deux suites de variables aléatoires : l'une dite caché et l'autre observable. Le premier processus génère la suite caché correspond à la suite d'états q_1, q_2, \dots, q_T , le deuxième HMM est responsable de la génération de la suite observable correspond à la suite d'observations o_1, o_2, \dots, o_T à partir de la séquence d'états. Un HMM se définit comme une structure composée d'un ensemble d'états, de graphes de transitions et de distributions de probabilités sur les transitions. En effet, Les HMMs garantissent une fusion cohérente de différents niveaux de modélisation à savoir celle morphologique, lexicale et syntaxique.

Sans s'approfondir dans les détails des formules mathématiques derrière cette méthode, nous présenterons succinctement le principe de cette approche.

Supposant qu'on a choisit au départ un modèle pour chaque catégorie d'entités nommées, sous entendu l'apprentissage se fait à partir d'un corpus d'apprentissage, où chaque élément du texte, inclusivement les entités nommées, a subit une annotation et catégorisation.

La première étape de l'apprentissage consiste à sauvegarder en mémoire un modèle pour chaque mot dans le corpus d'apprentissage. Prenons le cas où un modèle HMM rencontre la phrase suivante dans un corpus de test :

ξΛΛ. ◦◦ ◦◦◦◦◦◦◦◦ = « il est allé jusqu'à Amzmiz. »

Le mot « ◦◦◦◦◦◦◦◦ » qui suit immédiatement la préposition « ◦◦ » n'a pas été reconnu par le modèle HMM, vu que ce dernier n'était pas observé dans le corpus d'apprentissage.

Dans ce cas, le modèle HMM va se servir du contexte antécédent afin d'analyser le mot inconnu.

Le contexte antécédent dans notre cas consiste en la préposition du lieu « ◦◦ ». Supposant que dans notre corpus d'entraînement, la préposition « ◦◦ » est souvent suivie d'un nom de lieu. Dans ce cas, le système va affecter la catégorie nom de lieu au mot « ◦◦◦◦◦◦◦◦ » en se référant sur le modèle probabiliste du contexte

d'apparition de la préposition « ◦◦ » dans le corpus d'apprentissage.

▪ **CRF (Conditional Random Field) :**

Conditional Random Fields (CRF) [3] ou Champs Markoviens Conditionnels sont une famille de modèles statistiques graphiques non orientés non génératif, également connu sous le nom de champ aléatoire, un cas particulier de ce qui est une chaîne linéaire qui correspond à une condition formés de machine à états finis. En se basant sur les valeurs affectées à d'autres nœuds d'entrée assignées, ce modèle est utilisé pour calculer la probabilité conditionnelle des valeurs sur les nœuds de sortie affectés. CRF modèle est très utile dans la reconnaissance d'entités nommées.

Champ aléatoire: Soit $G = (Y, E)$ un graphe où chaque sommet YV est une variable aléatoire. Supposons que $P(Yv | \text{tout autre } Y) = P(Yv | \text{voisins } (Yv))$, alors Y est un champ aléatoire.

L'un des travaux pionniers de la NER pour les langues indiennes a été réalisé par Li and McCallum (2004) [4], où ils ont utilisé le CRF. Dans leur étude, la taille du corpus d'apprentissage était de 340K mots avec 15,063 entités nommées appartenant à trois types, nommés personne, location et organisation. Cette méthode (CRF) a été introduite par d'autres chercheurs (Settles, 2004; Tsai et al., 2006; Leaman and Gonzalez, 2008).

▪ **Maximum Entropy (ME) :**

Introduite par (Berger et al., 1996) communément appelé MaxEnt, Avec les contraintes basées sur les attributs choisis $f_j(x, y)$, la méthode du maximum d'entropie va chercher à maximiser la probabilité conditionnelle de $f_j(x, y)$. Des contraintes sont obtenues à partir des données d'apprentissage, il exprime des relations entre caractéristiques et des résultats de sorties. La probabilité de distribution, qui est satisfaite des propriétés au-dessus, est la plus haute entropie. Il fait l'hypothèse qu'elle suit une loi exponentielle.

▪ **SVM (Support Vector Machine) :**

Les machines à vecteurs de support représentent une famille d'algorithmes d'apprentissage qui permet de remédier à des problèmes de classification, ou de régression, en se référant aux résultats de la théorie de l'apprentissage statistique introduite par Vladimir Vapnik 1995. SVM est une méthode de classification binaire par apprentissage supervisé c'est-à-dire qu'on connaît à l'avance les classes auxquelles appartiennent les échantillons d'apprentissage, et la classification se fait en se basant sur les bases prédéfinies. L'utilisation des SVM est très fréquente pour la tâche de reconnaissance d'entités nommées ainsi que d'autres tâches de traitement du langage naturel. Plusieurs recherches ont été entamés dans ce sens, Takeuchi et Collier (2002) [5] ont utilisé le SVM dans la tâche NER avec une représentation binaire des caractéristiques. Isozaki et Kazawa (2002) [6] ont également utilisé le SVM, ils ont proposé quelques techniques comme la suppression de caractéristiques inutiles pour rendre le classificateur efficace en termes de temps d'apprentissage et de performance.

▪ **Decision Tree (DT) :**

Il s'agit des méthodes les plus célèbres dans le domaine d'apprentissage supervisé, ils représentent graphiquement sous forme d'arbre, un ensemble de règles et sont aisément interprétables. Pour les arbres de grande taille, la procédure globale peut être difficile à appréhender, cependant, la classification d'un élément particulier est toujours compréhensible. Chaque branche de l'arbre en question convient à une décision et dispose d'un certain poids et d'une certaine probabilité pour prendre une décision donnée, elle réclame l'utilisation d'un corpus d'apprentissage. Cette méthode a été utilisée par Sekine 1998 [7] pour la langue japonaise.

Un des avantages de cette méthode c'est que leur structure arborescente permet la génération automatique d'un ensemble de règles facilement lisibles et compréhensibles par un linguiste, contrairement à d'autres approches.

• Méthodes non-supervisées :

Cette méthode se distingue des autres méthodes d'apprentissage supervisé par le fait que l'algorithme procède à un apprentissage avec une intervention humaine minimale. Les techniques se reposent fondamentalement sur des ressources lexicales (par exemple, WordNet), abordé par E. Alfonseca and Manandhar (2002), sur des modèles lexicaux traité par R. Evans (2003) et sur les statistiques obtenues sur un large corpus non annoté. Autres recherches ont été effectués dans le même cadre par Y. Shinyama and Sekine (2004), et O. Etzioni et al. (2005). La méthode non supervisée utilise, entre autres, la technique du « Clustering » qui consiste à regrouper automatiquement les entités similaires. Selon Candillier (2006), l'apprentissage non supervisé « consiste à former différents groupes à partir d'un ensemble de données, de telle manière que les données considérées comme les plus similaires soient associées au même groupe et qu'au contraire les données considérées comme différentes se retrouvent dans des groupes distincts, permettant ainsi d'extraire de la connaissance à partir de ces données ». Les EN seront réunies automatiquement en sous-ensembles selon la ressemblance de leur contexte.

• Apprentissage paresseux :

Cette méthode a été populaire lors de la MUC-6, elle consiste à l'enregistrement des données du corpus d'entraînement dans la mémoire et fait par la suite à la comparaison des nouvelles requêtes avec celles qu'il vient de conserver en mémoire, ensuite, à l'aide d'une formule mathématique, il calcule le degré de similarité des lettres composant les deux requêtes. Finalement le système va affecter à la nouvelle requête la classe de l'exemple d'apprentissage le plus proche.

Prenons un exemple inventé qui un corpus d'apprentissage qui serait restreint juste à deux éléments:

Liste des éléments dans le corpus	Catégorie attribuée d'entraînement manuellement
oCЖCZЖ	Ville
CZCЖ	Personne

Voyons donc les nouvelles requêtes suivantes:

Requête 1 : ЖCЖCZЖ

Requête 2 : oCЖCZЖ

Requête 3 : ⊙ZCЖ

Requête 3 : CЖC

Suite au calcul de la distance entre ces 4 requêtes et les 2 éléments dans le corpus d'apprentissage, le système va affecter aux nouvelles requêtes, les catégories les plus probables en respectant un degré de ressemblance bien précis.

○ Résultat de la première requête : « ЖCЖCZЖ », mot inconnu, touche un degré de ressemblance orthographique avec « oCЖCZЖ ». Il serait affecté dans la catégorie ville.

○ Résultat de la deuxième requête : « oCЖCZЖ », mot existant dans le corpus d'entraînement, il va être attribué dans la classe ville.

○ Résultat de la troisième requête : « ⊙ZCЖ », mot n'est pas connu, et pourtant il est proche d'un point de vue orthographique de « CZCЖ », il reçoit donc la classe personne.

○ Résultat de la quatrième requête : « CЖC », mot n'est pas reconnu, en outre il n'est pas proche d'aucun éléments du corpus, alors il ne reçoit aucune catégorie.

Après l'examen des résultats de sortie, on constate que la méthode de calcul de ressemblance entre les exemples précédents est assez rigoureuse, puisqu'elle ne permet de classer que les mots inconnus qui sont très proches des mots existants dans le corpus.

• Apprentissage semi-supervisé :

Cette méthode est parfois nommée *apprentissage limité* [8]. Parmi les principales techniques employées, nous mentionnons le *bootstrapping* Borthwick (1999). Cette méthode ne nécessite qu'un nombre restreint de données injectées (*appelée semence*) auparavant pour marcher convenablement. Elle peut être sous forme d'une liste de noms de personnes par exemple. Le système procède à l'analyse des phrases qui sont dotées d'un certain type d'entités nommées, par la suite, le système garde les marqueurs lexicaux de ces entités. En répétant ce processus, un grand nombre de noms de personne peut être détecté. Ce principe a été utilisé par S. Brin (1998), Riloff et Jones (1999), Collins et Singer (1999), Yanarber et al. (2000), Nadeau et al. (2006) et M. Pasca et al. (2006).

Les règles abordées sont très facile.

« **Lausanne** » est un lieu; « **Rouen** » est un lieu; « **France** » est un lieu; tout nom qui contient « **Mr.** » est un nom de personne; tout nom qui contient « **Incorporated** » ou « **Inc.** » est un nom d'organisation; « **IBM.** » est une organisation; « **Microsoft** » est une organisation;

Ces règles précédentes, vont permettre de déduire de nouvelles règles. Prenons l'exemple suivant :

Mr. Ahmed.

Le système infère que Ahmed est un nom de personne, car il est précédé par le titre « **Mr.** » qui prédit un nom de personne dans les règles et ainsi et suite.

C. Approches Hybrides

Enfin, il y a les approches dites mixtes, appelées aussi hybrides [9, 10, 11], combinant l'apprentissage automatique (à partir d'un corpus) et l'écriture des règles à la main (création de règles à la main). Les limites de chacune des approches, ont amené les experts du domaine à fusionner les méthodes existantes pour augmenter les

performances de leurs outils. Les approches symboliques, souffrent également du coût de leur développement manuel et de la nécessité d'un expert en linguistique pour pouvoir les modifier et les adapter. Pour remédier à ce problème, les experts se penchent vers des méthodes d'apprentissage automatique de patrons linguistiques. Pourtant, les méthodes statistiques nécessitent, lors de la phase d'entraînement, une grande quantité de textes pré-étiquetés autrement dit un large corpus annoté, vu qu'on ne dispose pas toujours de ces données, cela constitue un réel problème. Ce qui a engendré l'apparition des systèmes hybrides qui tentent de palier problèmes de deux autres approches.

Ces systèmes ont été utilisés par le système suédois Swenam fait par Dalianas et Astrom et bien d'autres.

D. Mesures de performances

Le but de l'extraction des entités nommées et un peu analogue à ce lui de la classification, il s'agit à partir d'un ensemble de textes de *trouver/rechercher* les informations et uniquement celles qui sont *pertinentes* par rapport à une *catégorie considérée*.

Dans le cadre de la détection des entités nommées, deux objets vont être évalués par entité :

La première des choses, **le repérage de l'entité nommée** à travers des frontières fixées par le système (bornes de l'entité). Deuxièmement **la classe attribuée** à cette entité c.à.d. le type de l'entité.

Dans ce cas, nous cherchons à évaluer la capacité du système à trouver les informations pertinentes et uniquement celles là. Pour évaluer les sorties plusieurs critères de mesure des performances ont été proposé le **rappel**, la **précision** et le **F-mesure**. (Van Rijsbergen, 1979) Celles-ci mesurent le nombre d'objets pertinents effectivement récupérés parmi un ensemble d'objets.

E. Conclusion :

Cet article décrit la tâche de la reconnaissance d'entités nommées qui a une double finalité : elle consiste à « **baliser** » des documents textes pour l'identification et la « **catégorisation** » des mots dans des entités nommées correspondantes. Il vise ainsi à accéder aux informations contenues dans des textes dans la perspective de répondre à des questions basiques : Qui? Quoi? Où? Quand? Comment? Pourquoi?

Les conférences MUC, CoNLL ou ACE ont proposé un ensemble limité de catégories (par exemple, une personne, une organisation, l'emplacement et les expressions numériques comme l'argent ou des expressions pour cent), mais il ya beaucoup plus d'étiquettes possibles proposées illustré dans une hiérarchie étendu d'entités nommées[12]

qui contient environ 150 types NE développé par Sekine, et elle atteint encore actuellement 200 types.

REFERENCES

- [1] R. Grishman 1995 The NYU system for MUC-6 or Where's the Syntax. In the proceedings of Sixth Message Understanding Conference (MUC-6) (pp167-195). Fairfax, Virginia.
- [2] S. Sekine, Description of the Japanese NE system used for MET-2, In: MUC-7, Fairfax, Virginia, 1998.
- [3] CRF++:<http://crfpp.sourceforge.net/> Yet Another CRF toolkit (accessed on 3 rd may 2010)
- [4] Li, W., McCallum, A., 2004. Rapid development of Hindi named entity recognition using conditional random fields and feature induction. ACM Tran. on Asian Language Information Processing (TALIP) 2 (3), 290-294.
- [5] Takeuchi, K., Collier, N., 2002. Use of support vector machines in extended named entity. In: Proc. CoNLL-2002.
- [6] Isozaki, H., Kazawa, H., 2002. Efficient support vector classifiers for named entity recognition. In: Proc. COLING-2002, pp. 390-396.
- [7] Nadeau, David and Satoshi Sekine (2007) A survey of named entity recognition and classification. *Linguisticae Investigationes* 30(1):3-26.
- [8] A. Borthwick, J. Sterling, E. Agichtein, R. Grishman, Exploiting diverse knowledge sources via maximum entropy in named entity recognition, in: Proceedings of the 6th Workshop on Very Large Corpora, WVLC-98, Montreal, Canada, 1998.
- [9] A. Mikheev, C. Grover, M. Moens, Description of the LTG System used for MUC-7, In: MUC-7, Fairfax, Virginia, 1998.
- [10] R. Srihari, C. Niu, W. Li, A hybrid approach for named entity and sub-type tagging, In: Proceedings of Sixth Conference on Applied Natural Language Processing (ANLP), 2002, pp. 247-254.
- [11] X. YU, Chinese named entity recognition with cascaded hybrid model, In: Proceedings of NAACL HLT 2007, Prague, 2007, pp. 197-200.
- [12] S. Sekine, K. Sudo, C. Nobata, Extended named entity hierarchy, in: Proceedings of the Third International Conference on Language Resources and Evaluation, LREC2002, Las Palmas, Canary Islands, Spain, 2002.
- [13] Ameur M., Bouhjar A., Boukhris F., Boukouss A., Boumalk A., Elmedlaoui M., Iazzi E.M., Soui_ H., (2004). Initiation à la langue amazighe, IRCAM, Rabat, Maroc.
- [14] Ali Rachidi & Driss Mammass, Informatisation de La Langue Amazighe : Méthodes et Mises En Œuvre. 3rd International Conference: Sciences of Electronic Technologies of Information and Telecommunications, March 27-31, 2005 - TUNISIA