

Corpora based Approach for Arabic/English Word Translation Disambiguation

Farag Ahmed and Andreas Nürnberger

Data and Knowledge Engineering Group, Faculty of Computer Science,
Otto-von-Guericke-University of Magdeburg, Germany
farag.ahmed@ovgu.de
andreas.nuernberger@ovgu.de

ABSTRACT

We are presenting a word sense disambiguation method applied in automatic translation of a query from Arabic into English. The developed machine learning approach is based on statistical models, that can learn from parallel corpora by analysing the relations between the items included in this corpora in order to use them in the word sense disambiguation task. The relations between items in this corpora are obtained by using and developing a purely statistical method from corpora in order to avoid the use of structured linguistic resources like ontology which are not yet available for Arabic in an appropriate quality. The results of this analysis should provide us with some useful semantic information that can help to find the best translation equivalents of the polysemous items.

1. Introduction

Initially, online documents were used predominately by English speakers. Nowadays more than half (50.4%)¹ of web users speak a native language other than English. Therefore, it has become more important that documents of various languages and cultures should be retrieved by web search engines in response to the user's request.

Finding the most effective way to bridge the language barrier between queries and documents is the central challenge in Cross-Language Information Retrieval (CLIR) [48]. Cross Language Information Retrieval (CLIR) allows the user to submit the query in one language and retrieve the results in different languages, providing an important functionality that can help to meet that challenge. Cross-Language Information Retrieval (CLIR) approaches are typically divided into two main categories: approaches that exploit explicit representations of translation knowledge such as bilingual dictionaries or machine translation (MT) and approaches that extract useful translation knowledge from comparable or parallel corpora.

In the last few years, Arabic has become the major focus of many machine translation projects. Many rich resources are now available for Arabic. For example a Giga-Word² Arabic corpora which contains million sentences and Arabic/English Parallel

¹http://www.worldlingo.com/en/resources/languag_statistics.html.

²<http://www ldc.upenn.edu/Catalog/CatalogEntry.jsp?catalogId=LDC2007T40>.

corpus, which contains several thousands, sentence pairs of bilingual text for Arabic and English. The existence of these resources has been a crucial factor in building effective translation tools. Bilingual dictionaries (Arabic with other languages) have been used in several Arabic CLIR experiments. However, bilingual dictionaries sometimes provide multiple translations for the same word, which need to be disambiguated.

This paper proposes a method to disambiguate the user translated query in order to determine the correct word translations of the given query terms by exploiting a large bilingual corpus and statistical co-occurrence. The specific characteristics of Arabic morphology that hinder the correct match are taken into account by bridging the inflectional morphology gap for Arabic. We use one of the well-known Arabic morphological Analyzers [1] that include the araMorph package to translate the user query from Arabic to the English language in order to obtain the sense inventory for each of the ambiguous user query terms.

1.1. Arabic language

Arabic is a Semitic language, consisting of 28 letters, and its basic feature is that most of its words are built up from, and can be analyzed down to common roots. The exceptions to this rule are common nouns and particles. Arabic is a highly inflectional language with 85% of words derived from tri-lateral roots. Nouns and verbs are derived from a closed set of around 10,000 roots [4]. Arabic has three genders, feminine, masculine, and neuter; and three numbers, singular, dual (representing two entities), and plural. The specific characteristics of Arabic morphology make Arabic language particularly difficult for developing natural language processing methods for information retrieval. One of the main problems in retrieving information from Arabic language text is the variation in word forms, for example the Arabic word “kateb” (author) is built up from the root “ktb” (write). Prefixes and suffixes can be added to the words that have been built up from roots to add number or gender, for example adding the Arabic suffix “ن” (an) to the word “kateb” (author) will lead to the word “kateban” (authors) which represents dual masculine. What makes Arabic complicated to process is that Arabic nouns and verbs are heavily prefixed. The definite article “ال” (al) is always attached to nouns, and many conjunctions and prepositions are also attached as prefixes to nouns and verbs, hindering the retrieval of morphological variants of words [5]. Arabic is different from English and other Indo-European languages with respect to a number of important aspects. Words are written from right to left. It is mainly a consonantal language in its written forms, i.e. it excludes vowels. Its two main parts of speech are the verb and the noun in that word order and these consist, for the main part, of trilateral roots (three consonants forming the basis of noun forms that are derived from them). It is a morphologically complex language in that it provides flexibility in word formation: as briefly mentioned above, complex rules govern the creation of morphological variations, making it possible to form hundreds of words from one root [6].

Arabic poses a real translation challenge for many reasons; Arabic sentences are usually long and punctuation has none or little effect on interpretation of an Arabic text. Contextual analysis is important in Arabic in order to understand the exact meaning of some words. Characters are sometimes stretched for justified text (a word is spread over a bigger space than other words), which hinders the exact match for the

same word. In Arabic, synonyms are very common, for example, “year” has three synonyms in Arabic **حول**, **عام** **سنة**, and all are widely used in everyday communication. Despite the previous issues and the complexity of Arabic morphology, which impedes the matching of the Arabic word to the correct stem, e.g. inflectional forms of a word to its basic stem. Another real issue for the Arabic language is the absence of diacritics (sometimes called vowelings). Diacritics can be defined as a symbol over and underscored letters, which are used to indicate the proper pronunciations as well as for disambiguation purposes. The absence of diacritics in Arabic texts poses a real challenge for Arabic natural language processing as well as for translation, leading to high ambiguity. Though the use of diacritics is extremely important for readability and understanding, diacritics is very rarely used in real life situations. Diacritics don't appear in most printed media in Arabic regions nor on Arabic internet web sites. They are visible in religious texts such as the Quran, which is fully diacritized in order to prevent misinterpretation. Furthermore, the diacritics are present in children's books in school for learning purposes. For native speakers, the absence of diacritics is not an issue. They can easily understand the exact meaning of the word from the context, but for inexperienced learners as well as in computer usage, the absence of the diacritics is a real issue. When the texts are unvocalized, it is possible that several words have the same form but different meaning.

1.2. Tim Buckwalter Arabic morphological analyzer (BAMA)

(BAMA) is the most well known tool for analyzing Arabic texts. It consists of a main database of word forms that interact with other concatenation databases. An Arabic word is considered a concatenation of three regions: a prefix region, a stem region and a suffix region. The prefix and suffix regions can be null. Prefix and suffix lexicon entries cover all possible concatenations of Arabic prefixes and suffixes, respectively. Every word form is entered separately. It takes the stem as the base form and also provides information on the root. (BAMA) morphology reconstructs vowel marks and provides an English glossary. It returns all possible compositions of stems and affixes for a word. (BAMA) groups together stems with a similar meaning and associates it with a lemma ID. The (BAMA) contains 38,600 lemmas. For our work, we use the araMorph package. araMorph is a sophisticated java based Buckwalter analyzer.

2. Word Sense Disambiguation

The meaning of a word may vary significantly according to the context in which it occurs. As a result, it is possible that some words can have multiple meanings. This problem is even more complicated when those words are translated from one language into others. Therefore there is a need to disambiguate the ambiguous words that occur during the translations. The word translations disambiguation WTD, or more general Word sense disambiguation (WSD) is the process of determining the correct sense of an ambiguous word given the context in which the ambiguous word occurs. We can define the WSD problem, as the association of an occurrence of an ambiguous word with one of its proper sense. As described in the first section, the absence of the diacritics in most of the Arabic printed media or on the Internet web sites leads to high

ambiguity. This makes the probability that the single word can have multiple meanings a lot higher. For example, the Arabic word *يعد* (yEd) can have the following meanings: “promise”, “prepare”, “count”, “return”, “bring back” in English, the Arabic word *علم* (Elm) can have the following meanings: “flag”, “science”, “he knew”, “it was known”, “he taught”, “he was taught”. The task of disambiguation therefore involves two processes: Firstly, identifying all senses for every word relevant, secondly assigning the appropriate sense each time this word occurs. For the first step, this can be done using a list of senses for each of the ambiguous words existing in everyday dictionaries. The second step can be done by the analysis of the context in which the ambiguous word occurs, or by the use of an external knowledge source, such as lexical resources as well as a hand-devised source, which provides data (eg. grammar rules) useful to assigning the appropriate sense to the ambiguous word. In the WSD task, it is very important to consider the source of the disambiguation information, the way of constructing the rules using this information and the criteria of selecting the proper sense for the ambiguous word, using these rules. WSD is considered an important research problem and is assumed to be helpful for many applications such as machine translation (MT) and information retrieval.

Approaches for WSD can be classified into three categories: supervised learning, unsupervised learning, and corpora based approach. In the following we briefly describe the state of the art for word sense disambiguations.

2.1. Word Sense Disambiguation Approaches

Several methods for word sense disambiguation using a supervised learning technique have been proposed. This include approaches based on Naïve Bayesian [7], Decision List [8], Nearest Neighbor [9], Transformation Based Learning [10], Winnow [11], Boosting [12], and Naïve Bayesian Ensemble [13]. Among all of these methods, the ones using Naïve Bayesian Ensemble are reported to obtain the best performance for word sense disambiguation tasks with respect to the data sets used [13]. The idea behind all these approaches is that it is almost always possible to determine the sense of the ambiguous word by considering its context, and thus all methods attempt to build a classifier, using features that represent the context of the ambiguous word. In addition to supervised approaches for word sense disambiguation, unsupervised approaches and combinations of them have been also proposed for the same purpose. For examples, the authors of [15] proposed an Automatic word sense discrimination which divides the occurrences of a word into a number of classes by determining for any two occurrences whether they belong to the same sense or not, which then used for full word sense disambiguation task. Examples of unsupervised approaches were proposed by [16–21]. In [22] an unsupervised learning method using the Expectation-Maximization (EM) algorithm for text classification problems was proposed which was later improved [23] in order to apply it to WSD tasks. The authors of [24] combine both supervised and unsupervised lexical knowledge methods for word sense disambiguation. In [25] and [26] approaches using rule learning and neural networks were proposed respectively. Hidden Markov Models (HMMs) [31–33] and their extensions are very popular in a variety of fields including computer vision, natural language understanding, and speech recognition and synthesis. HMMs were also proposed for word sense disambiguation tasks [34–36].

Corpora based methods provide an alternative solution for overcoming the lexical acquisition bottleneck by gathering information directly from textual data e.g. bilingual corpora. Due to the expense of manual acquisition of lexical and disambiguation information where all necessary information for disambiguation have to be manually provided, supervised approaches suffer from major limitation in their reliance on predefined knowledge source, which affects their ability to handle large vocabulary in wide variety of contexts. In the last few years amount of parallel corpora available in electronic format have been increased, which helps the WSD researchers to extend the coverage of the existing system or train new system. For example, [27] and [28] used the parallel, aligned Hansard Corpus of Canadian Parliamentary debates for WSD, [29] used monolingual corpora of Hebrew and German. Using bilingual corpus to disambiguate words is leveraged by e.g. [14]. All of these corpora studies are based on the assumption that the mapping between words and word senses is widely different from one language to another. Unlike machine translation dictionaries, parallel corpora usually provide high quality translation equivalents that have been produced by experienced translators, who associate the proper sense of a word based on the context that the ambiguous word used in. However, in order to increase the efficiency of exploiting existing parallel corpora aligned at sentence level, the explicit word-level alignments should be added between sentence pairs in the training corpora. For word alignment two approaches have been proposed, the statistical-based approaches i.e. [37–39] and the lexicon-based approaches i.e. [40]. Several application for word alignment in natural language processing have been studied, i.e. [41, 42]. Some important applications for word alignment methods are, the automatic extraction of bilingual lexica and terminology from corpora [43, 44] and statistical machine translation systems i.e. [45, 46]. For a more detailed overview of word alignment approaches in nature language processing see [47].

In the next section, we describe the proposed algorithm based on Naïve Bayesian classification, explaining the way of solving or at least relaxing the Arabic morphological issues. Afterward, we explain the features used to represent the context in which ambiguous words occur, followed by experimental results, which show the results of disambiguating some ambiguous words using a parallel corpus. The paper closes with a conclusion and future work.

3. Proposed Approach

Our approach is based on exploiting parallel texts, in order to find the correct sense for the translated user query term. The minimum query length that the proposed approach accepts is two and the maximum query length is unlimited. Given the user query, the system begins by translating the query terms using the araMorph package. In case the system suggests more than one translation (senses inventory) for each of the query terms, the system then starts the disambiguation process to select the correct sense for the translated query terms. The disambiguation process starts by exploiting the parallel corpus, in which the Arabic version of the translation sentences matches fragments in the user query. A matched fragment must contain at least one word in the user query besides the ambiguous one. The words could be represented in the surface form or in one of its

variant forms. Therefore, and to detect all word form variants in the translation sentences in the training corpus, special similarity score measures are applied.

3.1. Bridging the Inflectional morphology gap

Languages exhibiting a rich inflectional morphology face a challenge for machine translation systems, as it is not possible to include all word form variants in the dictionaries. Inflected forms of words for those languages contain information that is not relevant for translation. The inflectional morphology difference between high inflectional language and poor inflectional language presents a number of issues for the translation system as well as for disambiguation algorithms. This inflection gap causes a matching challenge when translating between rich inflectional morphology and relatively poor inflectional morphology language. It is possible to have the word in one form in the source language, while having the same word in only a few forms in the target language. This causes several issues for word translation disambiguation, e.g. where more unknown words forms exist in the training data and will not be recognized as being relevant to the searched words. Result, it is possible to have lower matching score for those words even though they have a high occurrence of them in the training data.

The aim of this initial step is to alleviate the Arabic language morphology issues, which has to be done before accessing the Arabic language by the disambiguation algorithm. In order to deal with Arabic morphology issues, we used araMorph package [1].

To describe the problem more clearly, we consider, for simplicity, the Arabic word “دين” as described in section 2. The absence of the diacritics from the Arabic printed media or the Internet web sites causes high ambiguity. The Arabic word “دين” has two translations in English (Religion or debt). We calculate manually the occurrences of this word in the training corpus for both senses. This is done by searching for this word in the corpora and based on its context; we map it to the appropriate sense. As it is shown in Table 1 the word “دين” was found in basic form for the sense (Religion) 49 times and for the sense (Debt) only 10 times.

Table 1. The occurrence of the ambiguous word “دين” in the basic form for both senses

The ambiguous word	Senses	Occurrence of basic form in training data
دين	Religion	49
دين	Debt	10
Total		59

As Table 2 shows, when we consider the inflectional form for the word “دين” we see that the occurrence of the inflectional form for the word “دين” with the sense (Religion) is 1146 and with the sense (Debt) is 240.

Table 3 shows sentence examples from the training corpus where the ambiguous word “دين” appears in basic or inflectional form with both senses. Detecting all word forms variants of the user query terms in the corpus is very essential when computing the score of the synonym sets, as it is shown in Table 2. More than 1386 sentences will

Table 2. The occurrence of the inflectional form for the ambiguous word “دين” for both senses

The ambiguous word	Senses	Occurrence of inflectional form in training data
الدين	The Religion	75
والدين	And the Religion	22
الأديان	The Religions	45
والأديان	The Religions	7
الدينية	The Religious	63
والدينية	And the Religious	28
Total		240
الدين	The debt	860
والدين	And the debt	22
الديون	The debts.	255
والديون	And the debts.	9
Total		1146

Table 3. Sentences examples for the ambiguous word “دين” for both senses in basic and inflectional form

Sense	Form	Arabic sentence	English translation
Religion	Basic	لأن الإسلام الذي هو دين حوار وانفتاح على الناس	because Islam, which is a <i>religion</i> of dialogue and openness to people
Debt	Basic	نضيف إلى ذلك أن الولايات المتحدة الأمريكية أكبر دولة مدينة في العالم، فليها 400 مليار دولار عجزاً في ميزانيتها، يتم تمويلها عن طريق الاقتراض من المؤسسات الدولية والبنوك أو عن طريق تحويل هذا العجز إلى دين في الموازنة	In addition, the USA is the biggest debtor country in the world as it has a budget deficit of \$400 billion which is financed through borrowing from international institutions and banks or through converting such a deficit into a budget <i>debt</i> .

→

Table 3. *Continued*

Sense	Form	Arabic sentence	English translation
Religion	Infl.	ودعوا الوزير إلى التراجع عن قرار افتتاح المدرسة واستبدالها بمركز ثقافي ينشر تعاليم الدين والثقافة العربية	They called on the Minister to backtrack from that decision and to replace that school with a cultural centre promoting tenets of the <i>religion</i> and Arabic culture.
Debt	Infl.	وأكد الوزير أن الدين الخارجي على مصر هو في مستويات آمنة استنادا إلى ترتيبات جدولة الدين في نادي باريس	The minister emphasized that the foreign debt on Egypt was at safe levels due to the arrangements <i>of</i> <i>debt</i> scheduling in Paris Club.

be considered by the WSD algorithm to disambiguate the ambiguous word “دين”. For more details about the word form variant detection and their impact on the retrieval performance, we refer the reader to our previous work [2, 3].

In the following, we describe our approach based on the Naïve Bayesian algorithm, where we learn words and their relationships from a parallel corpus, taking into account that the morphological inflection that differs across the source and target languages.

3.2. Approach based on Naïve Bayesian Classifiers (NB)

The Naïve Bayesian Algorithm was first used for general classification problems. For WSD problems it had been used for the first time in [28]. The approach is based on the assumption that all features representing the problem are conditionally independent giving the value of classification variables. For a word sense disambiguation tasks, giving a word W , candidate classification variables $S = (s_1, s_2, \dots, s_n)$, which represent the senses of the ambiguous word, and the feature $F = (f_1, f_2, \dots, f_l)$ which describe the context in which an ambiguous word occurs, the Naïve Bayesian finds the proper sense s_i for the ambiguous word W by selecting the sense that maximizes the conditional probability of occurring in the given the context. In other words, NB constructs rules that achieve high discrimination level between occurrences of different word-senses

by a probabilistic estimation. The Naïve Bayesian estimation for the proper sense can be defined as follows:

$$P(s_i | f_1, f_2, \dots, f_n) = p(s_i) \prod_{j=0}^m p(f_j | s_i) \quad (1)$$

The sense s_i of a polysemous word w_{amb} in the source language is defined by a synonym set (one or more of its translations) in the target language. The features for WSD, that are useful for identifying the correct sense of the ambiguous words, can be terms such as words or collocations of words. Features are extracted from the parallel corpus in the context of the ambiguous word. The conditional probabilities of the features $F = (f_1, f_2, \dots, f_m)$ with observation of sense s_i , $P(f_j | s_i)$ and the probability of sense s_i , $P(s_i)$ are computed using maximum-likelihood estimates with $P(f_j | s_i) = C(f_j | s_i) / C(s_i)$ and $P(s_i) = C(s_i) / N$. $C(f_j, s_i)$ denoting the number of times feature f_j and sense s_i have been seen together in the training set. $C(s_i)$ denoting the number of occurrences of s_i in the training set and N is the total number of occurrences of the ambiguous word w_{amb} in the training dataset.

3.3. Features Selection

The selection of an effective representation of the context (features) plays an essential role in WSD. The proposed approach is based on building different classifiers from different subset of features and combinations of them. Those features are obtained from the user query terms (not counting the ambiguous terms), topic context and word inflectional form in the topic context and combinations of them.

In our algorithm, query terms are represented as sets of features on which the learning algorithm is trained. Topic context is represented by a bag of surrounding words in a large context of the ambiguous word:

$$F = \{w_{w_{amb-k}}, \dots, w_{w_{amb-2}}, w_{w_{amb-1}}, w_{w_{amb}}, w_{w_{amb+1}}, w_{w_{amb+2}}, \dots, w_{w_{amb+k}}, q_1, q_2, \dots, q_n\},$$

where k is the context size, w_{amb} is the ambiguous word and amb its position. The ambiguous word and the words in the context can be replaced by their inflectional forms. These forms and their contexts can be used as additional features. Thus, we obtain F' which contains in addition to the ambiguous word w_{amb} and its context the inflectional forms w_{inf_i} of the given sense and their context, as it is shown in table 2. Detecting all word form variants of the user query terms in the corpus will make 1386 sentences considered by the WSD algorithm to disambiguate the ambiguous word “دين”. In addition, we count for each context word the number of occurrences of this word and all its inflectional forms, i.e.

$$F' = F \prod_{i=0}^l \{w_{winf_{i-k}}, \dots, w_{winf_{i-2}}, w_{winf_{i-1}}, w_{winf_i}, w_{winf_{i+1}}, \dots, w_{winf_{i+k}}\}.$$

3.4. General Overview of the System

As Figure 1 shows, the system starts by processing the user query. The input is a natural language query Q . The query is then parsed into several words $q_1, q_2, q_3, \dots, q_n$. Each

word is then further processed independent of the other words. Since the dictionary does not contains all word forms of the translated word, only the root form, for each q_m in our query, we find its morphological root using the araMorph tool.³

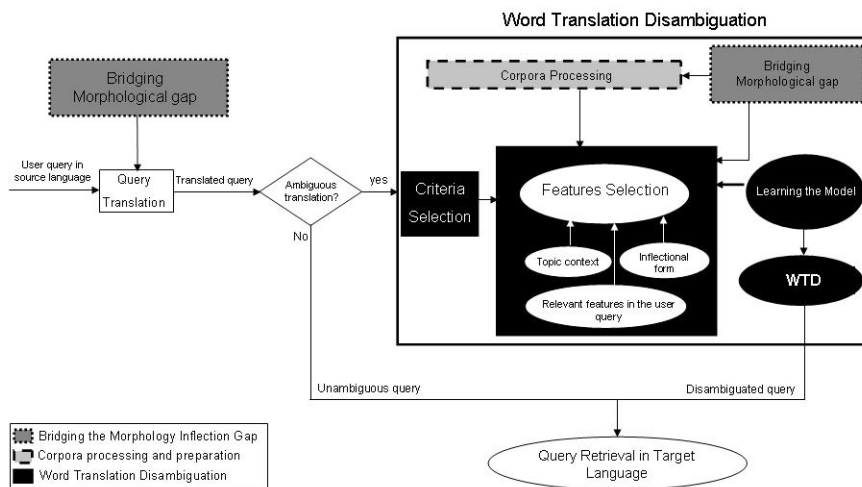


Figure 1. General overview of the system.

After finding the morphological root of each term in the query, the query term is translated. In case the query term has more than one translation, the model provides a list of translations (sense inventory) for each of the ambiguous query terms. Based on the obtained sense inventory for the ambiguous query term, the disambiguation process can be initiated. The algorithm starts by computing the scores of the individual synonym sets. This is done by exploiting the parallel corpora in which the Arabic version of the translated sentences matches words or fragments of the user query, while matched words of the query must map to at least two words that are nearby in the corpus sentence. These words could be represented in the surface form or in one of its inflectional forms. In order to detect all word form variants in the translation sentences in the training corpora, special similarity score measures are applied. Since the Arabic version of the translation sentences in the bilingual corpora matches fragments in the user query, the score of the individual synonym sets can be computed based on the features that represent the context of the ambiguous word. As additional features, the words in the topic context can be replaced by their inflectional form. After we have determined the features, the score of each of the sense sets can be computed. The sense which matches the highest number of features will be considered as the correct sense of the ambiguous query term and then it is assumed to be the best sense that describes the meaning of the ambiguous query term in the context.

³<http://www.nongnu.org/aramorph/>.

3.5. Illustrative examples

To consider how the algorithm performs the disambiguation steps, for simplicity we consider the following query with size 3 however the algorithm work for unlimited query size:

رسم جمركي للسلع

(A customs tax of commodities)

Step 1: The natural language query Q is parsed into several words q_1, q_2, q_3 .

Step 2: For each q_m in the query, we find its morphological root, since the dictionary does not contain all word forms, the algorithm before translation will find the single form of each of the given query terms. For example, the Arabic word قرارهم (their decision) before translation it will be processed and converted to the basic form which is قرار (decision).

Step 3: Translation of the query terms and creation of the sense inventory array in case of any for each of the query term is done. Table 4 shows the sense inventory for each of the ambiguous query terms.

Table 4. Sense inventory for each of the ambiguous query terms

Original Query term	Sense inventory (Possible English Translations)
رسم	[fee, tax, drawing, sketch, illustration, prescribe, trace, sketch, indicate, appoint]
جمركي	[customs, tariff, customs, control]
للسلع	[crack, rift, commodities, commercial, goods]

Step 4: The disambiguation process is initiated. The algorithm starts by computing the scores of the individual synonym sets:

- Number of times feature f_j and sense s_i have been seen together in the training set is computed.
- Number of occurrences of s_i in the training set is computed.
- The total number N of occurrences of the ambiguous word w_{amb} in the training dataset is computed.
- The disambiguation score is computed and the sense which matches the highest number of features is considered as the correct sense of the ambiguous query term.

Table 5 shows the disambiguation scores of the individual synonym sets for each ambiguous query terms with other query terms with 4934 occurrences of the ambiguous word w_{amb} in the training dataset.

As Table 5 shows there are 135 possible translations set for the original query in the source language.

3.6. Training data

The proposed algorithm was developed using Arabic/English parallel corpus.⁴ This corpus contains Arabic news stories and their English translations. It was collected by

⁴<http://www ldc.upenn.edu/Catalog/CatalogEntry.jsp?catalogId=LDC2004T18>.

Table 5. Disambiguation scores for each possible translations sets

S/N	query	score
1	fee AND (customs OR crack)	0
2	fee AND (customs OR rift)	0
3	fee AND (customs OR commodities)	0
4	fee AND (customs OR commercial)	0
5	fee AND (customs OR goods)	0
6	fee AND (control OR crack)	0
7	fee AND (control OR rift)	0
8	fee AND (control OR commodities)	0
9	fee AND (control OR commercial)	0
10	fee AND (control OR goods)	0
11	fee AND (tariff OR crack)	0
12	fee AND (tariff OR rift)	0
13	fee AND (tariff OR commodities)	0
14	fee AND (tariff OR commercial)	0
15	fee AND (tariff OR goods)	0
16	tax AND (customs OR crack)	0,0484
17	tax AND (customs OR rift)	0,0484
18	tax AND (customs OR commodities)	0,05948
19	tax AND (customs OR commercial)	0,05248
20	tax AND (customs OR goods)	0,05539
21	tax AND (control OR crack)	0
22	tax AND (control OR rift)	0
23	tax AND (control OR commodities)	0,01224
24	tax AND (control OR commercial)	0,00525
25	tax AND (control OR goods)	0,01108
26	tax AND (tariff OR crack)	0,00175
27	tax AND (tariff OR rift)	0,00175
28	tax AND (tariff OR commodities)	0,01399
29	tax AND (tariff OR commercial)	0,007
30	tax AND (tariff OR goods)	0,01283
31	prescribe AND (customs OR crack)	0
32	prescribe AND (customs OR rift)	0
33	prescribe AND (customs OR commodities)	0
34	prescribe AND (customs OR commercial)	0
35	prescribe AND (customs OR goods)	0
36	prescribe AND (control OR crack)	0
37	prescribe AND (control OR rift)	0
38	prescribe AND (control OR commodities)	0
39	prescribe AND (control OR commercial)	0
40	prescribe AND (control OR goods)	0
41	prescribe AND (tariff OR crack)	0
42	prescribe AND (tariff OR rift)	0
43	prescribe AND (tariff OR commodities)	0
44	prescribe AND (tariff OR commercial)	0
45	prescribe AND (tariff OR goods)	0
46	indicate AND (customs OR crack)	0
47	indicate AND (customs OR rift)	0
48	indicate AND (customs OR commodities)	0
49	indicate AND (customs OR commercial)	0
50	indicate AND (customs OR goods)	0,00117
51	indicate AND (control OR crack)	0,00058

→

→

Table 5. Continued

S/N	query	score
52	indicate AND (control OR rift)	0,00058
53	indicate AND (control OR commodities)	0,00058
54	indicate AND (control OR commercial)	0,00058
55	indicate AND (control OR goods)	0,00175
56	indicate AND (tariff OR crack)	0
57	indicate AND (tariff OR rift)	0
58	indicate AND (tariff OR commodities)	0
59	indicate AND (tariff OR commercial)	0
60	indicate AND (tariff OR goods)	0,00117
61	appoint AND (customs OR crack)	0
62	appoint AND (customs OR rift)	0
63	appoint AND (customs OR commodities)	0
64	appoint AND (customs OR commercial)	0,00058
65	appoint AND (customs OR goods)	0
66	appoint AND (control OR crack)	0
67	appoint AND (control OR rift)	0
68	appoint AND (control OR commodities)	0
69	appoint AND (control OR commercial)	0,00058
70	appoint AND (control OR goods)	0
71	appoint AND (tariff OR crack)	0
72	appoint AND (tariff OR rift)	0
73	appoint AND (tariff OR commodities)	0
74	appoint AND (tariff OR commercial)	0,00058
75	appoint AND (tariff OR goods)	0
76	trace AND (customs OR crack)	0
77	trace AND (customs OR rift)	0
78	trace AND (customs OR commodities)	0
79	trace AND (customs OR commercial)	0
80	trace AND (customs OR goods)	0
81	trace AND (control OR crack)	0
82	trace AND (control OR rift)	0
83	trace AND (control OR commodities)	0
84	trace AND (control OR commercial)	0
85	trace AND (control OR goods)	0
86	trace AND (tariff OR crack)	0
87	trace AND (tariff OR rift)	0
88	trace AND (tariff OR commodities)	0
89	trace AND (tariff OR commercial)	0
90	trace AND (tariff OR goods)	0
91	sketch AND (customs OR crack)	0
92	sketch AND (customs OR rift)	0
93	sketch AND (customs OR commodities)	0
94	sketch AND (customs OR commercial)	0
95	sketch AND (customs OR goods)	0
96	sketch AND (control OR crack)	0
97	sketch AND (control OR rift)	0
98	sketch AND (control OR commodities)	0
99	sketch AND (control OR commercial)	0
100	sketch AND (control OR goods)	0
101	sketch AND (tariff OR crack)	0

→

→

Table 5. Continued

S/N	query	score
102	sketch AND (tariff OR rift)	0
103	sketch AND (tariff OR commodities)	0
104	sketch AND (tariff OR commercial)	0
105	sketch AND (tariff OR goods)	0
106	drawing AND (customs OR crack)	0,00058
107	drawing AND (customs OR rift)	0,00058
108	drawing AND (customs OR commodities)	0,00117
109	drawing AND (customs OR commercial)	0,0035
110	drawing AND (customs OR goods)	0,00058
111	drawing AND (control OR crack)	0,00058
112	drawing AND (control OR rift)	0,00058
113	drawing AND (control OR commodities)	0,00117
114	drawing AND (control OR commercial)	0,0035
115	drawing AND (control OR goods)	0,00058
116	drawing AND (tariff OR crack)	0,00058
117	drawing AND (tariff OR rift)	0,00058
118	drawing AND (tariff OR commodities)	0,00117
119	drawing AND (tariff OR commercial)	0,0035
120	drawing AND (tariff OR goods)	0,00058
121	illustration AND (customs OR crack)	0
122	illustration AND (customs OR rift)	0
123	illustration AND (customs OR commodities)	0
124	illustration AND (customs OR commercial)	0
125	illustration AND (customs OR goods)	0
126	illustration AND (control OR crack)	0
127	illustration AND (control OR rift)	0
128	illustration AND (control OR commodities)	0
129	illustration AND (control OR commercial)	0
130	illustration AND (control OR goods)	0
131	illustration AND (tariff OR crack)	0
132	illustration AND (tariff OR rift)	0
133	illustration AND (tariff OR commodities)	0
134	illustration AND (tariff OR commercial)	0
135	illustration AND (tariff OR goods)	0

Ummah Press Service from January 2001 to September 2004. It totals 8,439 story pairs (Documents), 68,685 sentence pairs, 93,120 segments pairs, 2 Million Arabic words and 2.5 Million English words. The corpus is aligned at sentence level.

4. Evaluation

We evaluated our approach through an experiment using the Arabic/English parallel corpus aligned at sentence level. We selected 30 Arabic sentences from the corpus as queries to test the approach. These sentences have various lengths starting from two words. These queries had to contain at least one ambiguous word, which has multiple

English translations. In order to enrich the evaluation set, these ambiguous words had to have higher frequencies compared with other words in the training data, ensuring that these words will appear in different contexts in the training data. Furthermore, ambiguous words with high frequency sense were preferred. The senses (multiple translations) of the ambiguous words were obtained from the dictionary. The number of senses per test word ranged from two to nine, and the average was four. For each test word, training data were required by the algorithm to select the proper sense. The algorithm was applied to more than 93,123 parallel sentences. The results of the algorithm were compared with the manually selected sense.

For our evaluation, we built different classifiers from different subsets of features and combinations of them. The first classifier based on features that were obtained from the user query terms and topic context, which was represented by a bag of words in the context of the ambiguous word. The second classifier was based on the topic context and its inflectional form.

In order to evaluate the performance of the different classifiers, we used two measurements: applicability and precision [29]. The applicability is the proportion of the ambiguous words that the algorithm could disambiguate. The precision is the proportion of the correct disambiguated senses for the ambiguous word. The performance of our approach is summarized in Table 6. The sense, which is proposed by the algorithm was compared to the manually selected sense.

As it is expected the approach is better in the case of long query terms which provide more reach features and worse in short queries, especially the one consisting of two words. We consider that the reason for the poor result for the short queries is that, when the query consists of few words it is possible that the features which are extracted from the query terms can appear in the context of different senses. For example, consider the query “الدين الإسلامي” (The Islamic religion). When the algorithm goes through the corpus, the ambiguous word “الدين” (The Religion or The debt) will be found in two different contexts whether in Religion or Debt context. The query term “الإسلامي” (Islamic) can be found in both contexts of the ambiguous word as (Islamic religion) or as a name of bank (Islamic Bank), which is the context of the second sense (Debt). One possible solution for this issue is query expansion. This can be done by exploiting the corpus and suggesting possible term expansion to the user. The user then confirms this term expansion, which will help to disambiguate the ambiguous query term when translating to the target language.

Another reason for the poor performance is that due to the morphological inflectional gap between languages such as Arabic the same word can be found in different forms. In order to increase the performance of the disambiguation process all of these forms need to be detected.

Table 6 shows the overall performance of the algorithm based on building two classifiers from different subsets of features and combinations of them. Those features are user query terms, topic context and word inflectional form in topic context and combinations of them. As is shown in Table 6, the performance of the algorithm is better when using the inflectional forms instead of the basic word form. The reason for that, the Arabic word can be represented not just in its basic form, but in many inflectional forms and so we will have more training sentences that will be visible to the algorithm to disambiguate the ambiguous query terms.

Table 6. The overall performance using Applicability and precision

Classifiers	Applicability	Precision
Query term + Topic context	52%	68%
Query term + feature Inflectional form	82%	93%

5. Conclusion

In this paper we proposed a method for word translation disambiguation using a bilingual parallel corpus together with sense definitions by translations into another language. WTD for each sense of the polysemous word is done by defining a sense of each of the ambiguous words. In order to train the algorithm, a set of features was defined. The algorithm then selects the sense that maximizes the score. Based on the experiments that we performed, using Arabic/English parallel corpus, results could show that our algorithm achieved certain promising results. The applicability and precision using 30 polysemous words were 52% and 68% for the first classifier and 82% and 93% for the second classifier, respectively. Although our algorithm has gained promising results, it still has some problems: the developed approach is based on the premise that the features of each sense is independent from the other sense so that the words that appear in the context of the ambiguous word should appear rarely in the context of the other sense. In future work, we will extend the algorithm so that it can extend and enrich the user query, consequently having more features to describe the context of the ambiguous word.

Furthermore, the proposed algorithm currently estimates the parameters that the algorithm needs for training, according to co-occurrence in the context of the ambiguous word. However, this is not always suitable for all polysemous words. It will be useful to use syntactic co-occurrence as an extra feature that the algorithm can use for training. i.e. the Arabic word “مصر” has the translations (Egypt and insist) in English. Therefore, using some kind of syntactic information will help to disambiguate this kind of word.

BIBLIOGRAPHY

- [1] Tim Buckwalter, Buckwalter. Arabic Morphological Analyzer Version 1.0. Linguistic Data Consortium, University of Pennsylvania, 2002. LDC Catalog No.: LDC2002L49.
- [2] F. Ahmed and A. Nürnberger. N-grams conflation approach for arabic text. In Proceedings of the International Workshop on improving Non English Web Searching (iNEWS 07) in conjunction with 30th Annual International ACM SIGIR Conference, pages 39–46, Amsterdam, Netherlands, 2007.
- [3] F. Ahmed and A. Nürnberger. araSearch: Improving arabic text retrieval via detection of word form variations. In Proc. of the 1st International Conference on Information Systems and Economic Intelligence (SIIE'2008), pages 309–323, Hammamet, Tunisia, 2008.
- [4] Al-Fedaghi and Al-Anzi. A new algorithm to generate arabic root-pattern forms. In Proceedings of the 11th National Computer Conference, King Fahd University of Petroleum & Minerals, pages 04–07, Dhahran, Saudi Arabia, 1989.
- [5] H. Moukdad. Lost in cyberspace: How do search engines handle arabic queries? In Proceedings of the 32nd Annual Conference of the Canadian Association for Information Science, Winnipeg, 2004.

- [6] H. Moukdad and A. Large. Information retrieval from full-text arabic databases: Can search engines designed for english do the job? *International Journal of Libraries and Information Services*, pages 63–74, 2001.
- [7] W. A. Gale, K. W. Church, and D. Yarowsky. A method for disambiguating word senses in a large corpus. 26(5-6):415–439, 1992.
- [8] D. Yarowsky. Decision lists for lexical ambiguity resolution: Application to accent restoration in spanish and french. pages 88–95, 1994.
- [9] T. Ng and H. B. Lee. Integrating multiple knowledge sources to disambiguate word sense: An exemplar-based approach. In *Proceedings of the 34th Annual Meeting of the Association for Computational Linguistics*, pages 40–47, 1996.
- [10] L. Mangu and E. Brill. Automatic rule acquisition for spelling correction. In *Proceedings of the 14th International Conference on Machine Learning*, pages 187–194, 1997.
- [11] R. Golding and D. Roth. A winnow-based approach to context-sensitive spelling correction. *Machine Learning*, 34(1):107–130, 1999.
- [12] L. M. . G. R. Escudero, Gerard. Boosting applied to word sense disambiguation. *Proceedings of the 12th European Conference on Machine Learning (ECML)*, Barcelona, Spain, pp. 129–141, 2000.
- [13] T. Pedersen. A simple approach to building ensembles of Naive Bayesian classifiers for word sense disambiguation. In *Proceedings of the First Annual Meeting of the North American Chapter of the Association for Computational Linguistics*, Seattle, WA, May, pages 63–69, 2000.
- [14] Nancy Ide, N. Parallel translations as sense discriminators. *SIGLEX99: Standardizing Lexical Resources, ACL99 Workshop*, College Park, Maryland, pages 52–61, 1999.
- [15] H. Schütze. Automatic word sense discrimination. *Computational Linguistics*, 24(1):97–124, 1998.
- [16] K. C. Litkowski. Senseval: The cl research experience. *Computers and the Humanities*, 34(1–2): 153–158, 2000.
- [17] Dekang Lin. Word sense disambiguation with a similarity based smoothed l brary. In *Computers and the Humanities: Special Issue on Senseval*, pp. 34:147–152, 2000.
- [18] Philip Resnik. Selectional preference and sense disambiguation. In *Proceedings of ACL Siglex Workshop on Tagging Text with Lexical Semantics, Why, What and How?*, Washington, pp. 4–5, 1997.
- [19] D. Yarowsky. Word-sense disambiguation using statistical models of Ro-get’s categories trained on large corpora. In *Proceedings of COL-ING-92*, Nantes, France, pages 454–460, 1992.
- [20] Indrajit Bhattacharya, Lise Getoor, Yoshua Bengio. Unsupervised Sense Disambiguation Using Bi-lingual Probabilistic Models. *ACL*, 287–294, 2004.
- [21] Hiroyuki Kaji, Yasutsugu Morimoto. Unsupervised Word-Sense Disambiguation Using Bilingual Comparable Corpora. *IEICE Transactions 88-D(2)*, pp. 289–301, 2005.
- [22] K. Nigam, A. McCallum, S. Thrun, and T. Mitchell. Text classification from labeled and unlabeled documents using em. *Machine Learning*, 39(2–3):103–134, 2000.
- [23] H. Shinnou and M. Sasaki. Unsupervised learning of word sense disambiguation rules by estimating an optimum iteration number in the em algorithm. In *Proceedings of the seventh conference on Natural language learning at HLT-NAACL 2003*, pages 41–48, 2003.
- [24] E. Agirre, J. Atserias, L. Padr, and G. Rigau. Combining supervised and unsupervised lexical knowledge methods for word sense disambiguation. In *Computers and the Humanities, Special Double Issue on Senseval*, 34(1):103–108, 2000.
- [25] D. Yarowsky. Unsupervised word sense disambiguation rivaling supervised methods. In *Meeting of the Association for Computational Linguistics* pages 189–196, 1995.
- [26] G. Towell and E. M. Voorhees. Disambiguating highly ambiguous words. *Computational Linguistics*, 24(1):125–146, 1998.

- [27] P. F. Brown, J. C. Lai, and R. L. Mercer. Aligning sentences in parallel corpora. In Proceedings of the 29th Annual Meeting of the Association for Computational Linguistics, pages 169–176, Berkeley, CA, 1991.
- [28] W. A. Gale, K. W. Church, and D. Yarowsky. Using bilingual materials to develop word sense disambiguation methods. In Proceedings of the Fourth International Conference on Theoretical and Methodological Issues in Machine Translation (TMI'92), pages 101–112, 1992.
- [29] I. Dagan and A. Itai. Word sense disambiguation using a second language monolingual corpus. *Computational Linguistics*, 20(4):563–596, 1994.
- [30] Duda, R. O. and Hart, P. E. *Pattern Classification and Scene Analysis*, John Wiley, 1973.
- [31] L. Baum and J. Egon. An inequality with applications to statistical estimation for probabilistic functions of a markov process and to a model for ecology. *Bull. Amer. Meteorol. Soc.*, 73(3): 360–363, 1967.
- [32] L. Baum and T. Petrie. Statistical inference for probabilistic functions of finite state markov chains. *Ann. Math. Stat.*, 37(1):1554–1563, 1966.
- [33] L. Baum, T. Petrie, G. Soules, and N. Weiss. A maximisation technique occurring in the statistical analysis of probabilistic functions of markov chains. *Ann. Math. Stat.*, 41(1):164–171, 1970.
- [34] C. Loupy, M. El-Beze, and P. F. Marteau. Word sense disambiguation using HMM tagger. In Proceedings of the 1st International Conference on Language Resources and Evaluation, LREC, Granada, Spain, pages 1255–1258, 1998.
- [35] A. Molina, F. Pla, and E. Segarra. A hidden markov model approach to word sense disambiguation. In Proceedings of the 8th Ibero-American Conference on AI: Advances in Artificial Intelligence, pages 655–663, 2002.
- [36] A. Molina, F. Pla, and E. Segarra. WSD system based on specialized hidden markov model (upv-shmm-eaw). In Senseval-3: Third International Workshop on the Evaluation of Systems for the Semantic Analysis of Text, Association for Computational Linguistics, pages 171–174, Barcelona, Spain, 2004.
- [37] W. A. Gale and K. W. Church. A program for aligning sentences in bilingual corpora. In Proceedings of the 29th Annual Meeting of the Association for Computational Linguistics (ACLg1), pages 177–184, 1991.
- [38] I. Dagan, K. W. Church, and W. A. Gale. Robust bilingual word alignment for machine aided translation. In Proceedings of the Workshop on Very Large Corpora, pages 1–8, Columbus, Ohio, 1993.
- [39] J. S. Chang and M. H. C. Chert. Using partial aligned parallel text and part-of-speech information in word alignment. In Proceedings of the First Conference of the Association for Machine Translation in the Americas (AMTA94), pages 16–23, 1994.
- [40] M. Ker and J. S. Chang. A class-based approach to word alignment. *Computational Linguistics*, 32(2):313–343, 1997.
- [41] F. J. Och and H. Ney. A comparison of alignment models for statistical machine translation. in coling '00. In Proceedings of 18th International Conference on Computational Linguistics, pages 1086–1090, 2000.
- [42] D. Yarowsky and R. Wicentowski. Minimally supervised morphological analysis by multimodal alignment. In Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics (ACL), pages 207–216, 2000.
- [43] F. Smadja, K. R. McKeown, and V. Hatzivassiloglou. Translating collocations for bilingual lexicons: A statistical approach. *Computational Linguistics*, 22(1):1–38, 1996.
- [44] I. D. Melamed. Models of translational equivalence among words. *Computational Linguistics*, 26(21).

- [45] A. L. Berger, P. F. Brown, S. A. D. Pietra, V. J. Della Pietra, J. R. Gillett, J. D. Lafferty, H. Printz, and L. Ures. The candide system for machine translation. In Proceedings of the ARPA Workshop on Human Language Technology, pages 157–162, Plainsboro, New Jersey, 1994.
- [46] S. NieSSen, S. Vogel, H. Ney, and C. Tillmann. A dp-based search algorithm for statistical machine translation. In Proceedings of COLING-ACL'98: 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics, pages 960–967, Montreal, Canada, 1998.
- [47] F. J. Och and H. Ney. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1):19–51, 2003.
- [48] Yang, Y., & Ma, N. CMU in cross-language information retrieval at NTCIR-3. In K. Oyama, E. Ishida, & N. Kando (Eds.), *NTCIR Workshop3 Proceedings of the third NTCIR workshop on research in information retrieval, automatic text summarization and question answering*. Tokyo, 2002.

