Arabic English Cross-Lingual Plagiarism Detection Based on Keyphrases Extraction, Monolingual and Machine Learning Approach

4	
5	
6	Mokhtar Al-suhaiqi1 ^{1*} , Muneer A.S. Hazaa2 ² , Mohammed Albared3 ³
7	
8 9 10	¹ Master Researcher ,Department of Computer and Information Technology, Yemen Academy For Graduate Studies ² Associate professor, Faculty of Computer and Information Technology, Dhamar
11	University, Yemen
12 13 14	³ Assistant professor, Faculty of Computer and Information Technology, Sana'a University, Yemen
15 16	
17 18	ABSTRACT

Due to rapid growth of research articles in various languages, cross-lingual plagiarism detection problem has received increasing interest in recent years. Cross-lingual plagiarism detection is more challenging task than monolingual plagiarism detection. This paper addresses the problem of cross-lingual plagiarism detection (CLPD) by proposing a method that combines keyphrases extraction, monolingual detection methods and machine learning approach. The research methodology used in this study has facilitated to accomplish the objectives in terms of designing, developing, and implementing an efficient Arabic – English cross lingual plagiarism detection.

This paper empirically evaluates five different monolingual plagiarism detection methods namely i)N-Grams Similarity, ii)Longest Common Subsequence, iii)Dice Coefficient, iv)Fingerprint based Jaccard Similarity and v) Fingerprint based Containment Similarity. In addition, three machine learning approaches namely i) naïve Bayes, ii) Support Vector Machine, and iii) linear logistic regression classifiers are used for Arabic-English Crosslanguage plagiarism detection. Several experiments are conducted to evaluate the performance of the key phrases extraction methods. In addition, Several experiments to investigate the performance of machine learning techniques to find the best method for Arabic-English Cross-language plagiarism detection.

According to the experiments of Arabic-English Cross-language plagiarism detection, the highest result was obtained using SVM classifier with 92% f-measure. In addition, the highest results were obtained by all classifiers are achieved, when most of the monolingual plagiarism detection methods are used.

19

20

Keywords: Cross Language Plagiarism Detection, Mono-Language Plagiarism Detection, Classification, Machine Learning, Key Phrases, Candidate document.

21 22

23 **1. INTRODUCTION**

24

Cross-lingual plagiarism (CLP) happens when texts written in one language are translated 25 26 into another language and used without acknowledging the original sources. Extensive 27 studies have been executed on monolingual plagiarism analysis which content searching for 28 plagiarism in documents of the same language, but CLP still remains a challenge. Previous 29 studies have addressed this problem using methods such as Statistical Machine 30 Translation [1], cross-lingual showed semantic analysis (CL-ESA) [2], syntactic 31 alignment using character N-grams (CL-CNG), dictionaries and thesaurus [3] [4], online 32 machine translators [5, 6], and more recently, semantic networks and word embedding [7] 33 [8], and [9, 10]. Most of the suggested pattern are either limited to bilingual cross-lingual 34 plagiarism detection tasks, when require parallel or comparable corpus which are 35 usually not sufficient or available for low resource languages, while others trust on internet 36 translation services, which are not existing for large scale cross-lingual plagiarism 37 detection.

38 Different methods have been used to solve the cross lingual plagiarism detection. Based on 39 the literature, it could be noticed that the majority of these methods can be classified into 40 machine translation based approaches, parallel corpora based models and hybrid models. The main problems of the existing cross-language plagiarism detection techniques that uses 41 42 machine translation as main method where the quality of the existing machine translation in translating big texts (whole documents) is very low and detecting plagiarism in translated 43 documents is very challenging task because of the lexical and structural changes. In 44 addition, when translated texts are replaced with their synonyms, using online machine 45 translators to detect CLP would result in poor performance. To handle the limitation of these 46 47 methods, this paper aim to design and implement a keyphrases based cross lingual 48 plagiarism detection method. A significant feature of the proposed methodology is that it can be more efficient for detecting mono lingual paraphrased plagiarism where the sentence 49 50 changed and cross lingual translated plagiarism, as it keyphrases based structure is 51 detection method and keyphrases and their translation cannot be paraphrased.

52 This proposed research methodology consists of five phases, denoted as i) documents pre-53 processing phase, ii) Key phrase Extraction, Translation and Fingerprinting phase, iii) 54 Retrieval of Candidate Documents phase, vi) Monolingual plagiarism detection phase and v) 55 Machine Learning phase. The research methodology used in this study has facilitated to 56 accomplish the objectives in terms of designing, developing, and implementing an efficient 57 Arabic – English cross lingual plagiarism detection.

58 The remainder of this paper is structured as follows: Section 2 provides related work of cross-language Arabic – English techniques, as applied to words or sentences. Section 3 is 59 60 proposed methodology, explaining the various proposed algorithms which are used for the pre-processing and framework CLPD; the techniques mentioned in section 3, namely pre-61 62 processing is tokenization and stop word and NLP techniques in section 3.1; in section 3.2, the techniques are the key phrase extraction -based techniques, namely c- value algorithm 63 64 and NC-value and key phrase ranking to find similarity score after that translate Arabic key phrases to English and retrieval candidate document and compare fingerprint for the key 65 66 phrases in section 3.4. Section 3.5 monolingual methods N-Grame and longest common 67 subsequence to compare candidate document and suspicious document by hash table for 68 fingerprint; and section 3.6 Machine Learning phase in this section is plagiarised text or not. 69 in section 4 presents the experimental design, including the tools and packages used in this 70 study, the datasets involving 318 documents from the Arabic and English language 71 benchmark dataset. Section 5 presents the results and discussion of findings and, finally, in 72 section 6, conclusions and recommendations for future research are provided.

73 **2. RELATED WORK**

74

75 In this section, we give an overview of existing research in the area of focused on dataset of 76 document. Specifically focusing on candidate document categorization. In [11], text pre-77 processing techniques, such as stopword removal, and shallow NLP techniques, such as 78 stemming, are applied to documents before counting similarity. Short sentences are also 79 deleted. The degrees of similarity between words are computed by their frequency of cooccurrence and relative distance, as mentioned by a word-correlation matrix generated using 80 81 Wikipedia, A threshold is set to candidate sentences with a low similarity, and the degree of 82 resemblance between two documents is visualized using Dot plot view. Although the results 83 interpreted development over n-gram matching by decreasing the false positives, the 84 approach is still limited to comparison between individual words.

Experiments were created on a domain-specie corpus compounding of English, Arabic,
French, Spanish and Russian texts translated into Italian[12]. The experiment was executed
using an SVM classifier, based on features such as lemmatised words and POS sequences.
The best accuracy was achieved by using a combination of features that includes 1-gram
word with TF-IDF weighting, and 2-grams and 3-grams of POS tags. The experiment
finished that the task biases on the distribution of n-grams of function words and morphosyntactic features.

92 Pouliguen introduced a statistical approach to map multilingual documents for a language-93 independent document representation, which measures similarity between monolingual and 94 cross-lingual documents. A parallel corpus with multilingual interpreted texts was used, and 95 pre-processing techniques including lemmatisation and stopword removal were applied. Parallel texts in various languages are determined by the *tf-idf* of the topic, and the top 100 96 97 words are chosen as descriptors. Each descriptor contains one-to-one interpretations into various languages and is stood for by a vector. The similarity score was computed by 98 99 comparing the vectors between Spanish and English documents[13].

Aljohani and Mohd [14] introduced the first Arabic-English cross-language plagiarism detection using the Winnowing Algorithm to discover the Arabic sentences translated from English sources without indication of the original sources, as well as to diagnosing its main content and processes. The result clarifies that the Winnowing algorithm can be used effectively to discover the Arabic-English cross-language plagiarism with 81% recall, 97% precision and 89% F-measure.

Omar, Alkhatib [15] studied a method for plagiarism detection algorithm in both Arabic and
 English languages. They proved a system to detect plagiarism in both Arabic and English
 languages using "Bing" search machain. The system which bases on plagiarism detection
 algorithm is effective and can supply both Arabic and English languages.

Kent [16] improved a web-based system to discover cross-lingual plagiarism. The system
 decreases candidate document by summarizing. The Summary is interpreted to English.
 Then similar web resources are discovered.

Gottschalk [17]and Demidova improved methods to join text passages written in various languages and consisting of overlapping data. The authors used Named entities and text interpretation to English as features to estimate the similarity between documents. These approaches use text translation as part of the process of obtaining a common comparison space. However, since text translation is a challenging task, it may arrive to high false rate.

Ferrero[9] suggested methods for cross-lingual plagiarism detection using word embeddings.
 These methods require training using decision tree or weights optimization, so here they are
 supervised methods.

España-Bonet and et..al [18] introduced a language autonomous model that measures the semantic similarity between text captures across multiple languages. The system used a Support Vector Machine (SVM) to summarize a number of inter textual features, which contains features divided from embeddings trained using the word2vec model and a multilingual corpora, from lexical similarity measurements, from the internal representation (hidden layer) of a neural network trained using multi-lingual parallel corpora and from CL-ESA. This approach is however best appropriate for low resource languages.

128 3. METHODOLOGY

129

130 This research will study the problem of cross lingual plagiarism detection solution, and 131 proposed solutions for this problem. The primary goal of the research is to design and 132 implement methods for Arabic – English cross lingual plagiarism detection.

This research methodology consists of five main phases, denoted as i) Documents pre processing phase, ii) Key phrase Extraction, Translation and Fingerprinting phase and iii)
 Retrieval of Candidate Documents phase, vi) Monolingual plagiarism detection phase and v)
 Machine Learning phase.



156 In the pre-processing stage, various NLP pre-processing techniques are applied in a first 157 step, each document is spilt into sentences. This work use (.), (;), (:), (!) And (?) Punctuation After splitting documents into sentences, the sentences pre-158 marks as a spilt point. processing consists of three steps: 1) tokenization, 2) normalization, 3) stop word removal. 159 160 All sentences went through a pre-processing stage. In the normalization process, noisy 161 characters are removed. Secondly, in this phase certain stop-words that occur commonly in all documents were removed to avoid plagiarism detection over fitting. After the pre-162 processing stage, each document is represented as a bag of sentences and each sentence 163 164 in its turn is modelled as Bag Of Words.

165

Tokenization								
Input								
قياس_الضوء_الطيفي في الفيزياء , قياس الضوء الطيفي , هو دراسة كمية طيف كهرومغناطيس للطيف الكهرومغناطيسي. وهو أكثر تخصصا من قياس الطيف الكهرومغناطيس , حيث يتعامل فقط مع طيف مرئي , وقريب أشعة فوق البنفسجية وقريب أشعة تحت الحمراء.								
		Out put \	Input Stop wor	d				
•		الفيزياء	في	الطيفي	الضوء	قياس		
كمية	دراسة	ھو	6	الطيفي	الضوء	قياس		
أكثر	I	و هو	الكهرومغناطيسي	للطيف	كهرومغناطيس	طيف		
حيث	4	الكهرومغناطيس	الطيف	قياس	من	تخصصاً		
وقريب	4	مرئي	طيف	مع	فقط	يتعامل		
الحمراء	تحت	اشعة	وقريب	البنفسجية	فوق	اشعة		
		S	top word					
			Out put					
		الفيزياء		الطيفي	الضوء	قياس		
	دراسة			الطيفي	الضوء	قياس		
طيف كهرومغناطيس للطيف الكهرومغناطيسي								
		الكهرومغناطيس	الطيف	قياس		تخصصاً		
		مرئي	طيف			يتعامل		
الحمراء	تحت	اشعة		البنفسجية	فوق	اشعة		

166 167

Fig. 2. Pre-processing tokenization and stop word of Arabic Document

168

169 **3.2 Key phrases Extraction Phase**

170

The main problems of the existing cross-language plagiarism detection techniques that uses machine translation as main method where the quality of the existing machine translation in translating big texts (whole documents) is very low and detecting plagiarism in translated documents is very challenging task because of the lexical and structural changes.

175 Key phrases are the important words/phrases that reflect the subject of the text. The Key 176 phrases describe a document in a coherent and simple way giving the prospective reader a

way to quickly determine whether the document satisfies their information need. According
to that, we index each document by Key phrases and only translate them, if the similarity
score is so high between the Key phrases of two documents, then one of these documents
will be selected as suspicious document. However, the method used here for key phrases
extraction consists of four steps 1) Features Extraction 2) Ranking 3) translation
4) fingerprinting.

183 <u>3.2.1 Features Extraction</u> 184

185 The following features are used for ranking the candidate key phrase:

186 3.2.1.1 Phrase Frequency

Frequency is the number of occurrences of the candidate phrase. Frequency is normalized by the number of all candidate phrases in the document.as [19]

189

190
$$f_{\theta} = tf(kp) = \frac{\#(kp)}{\sum_{n \in \{all_phrases\}} \#(n)}$$
(3.1)

191 3.2.1.2 C-value Approach

The C-Value method is a hybrid domain-independent method combining linguistic and statistical information (with emphasis on the statistical part) for the extraction of key phrases and nested phrases (i.e. phrases that appear within other longer phrases, and may or may not appear by themselves in the corpus). This method takes as input a corpus and produces a list of candidate key phrases, ordered by the likelihood of being valid terms, namely their C-Value measure... C-value is defined as [20]:

198
$$f_{\varphi} = c - value(c) = \log_2 |c| \left\langle f(c) - \frac{1}{p(T_c)} \sum_{P \in T_c} f(b) \right\rangle$$
 (3.2)

199

- 200 Where C is a candidate key phrase, |C| is the number of simple nouns that consist of C, 201 f(C) is its frequency of occurrence in the corpus, T_C is the set of extracted candidate terms 202 that contain C. $P(T_C)$ and is the number of this candidate term.
- 203 3.2.1.3 NC-Value

The NC-Value is used to re-rank and improve the list of the extracted key phrases based on information from the term's <u>neighbourhood</u>. It, therefore, ranks the list of candidate key phrases, trying to bring higher key phrases that <u>are</u> more likely to contain key phrases. The NC-value measure is computed as [19, 21]:

208
$$NC - value(a) = 0.8C - value(a) + 0.2 \sum_{b \in a} f_a(b) w eight(b)$$
 (3.3)

209 210

This main purpose of this phase is to extract the most important Key phrases. To rank each key phrase from the candidate Key phrases.

213 214

215 3.2.3 Translation And Language Normalization

3.2.2 Key phrases Ranking And Filtering

216

In order to overcome the language barrier, all original documents (represented by extracted key phrases) are translated into one language in this case the English language has been chosen as it has bilingual translation between it and most of languages. For this purpose, the present work adopted Google Translate (GT) as it offers API access and is considered the state-of-the-art machine translation system used today.

222 3.2.4 Fingerprinting

223

Document fingerprinting is the process of representing a document as a set of integers resulting from hashing substrings of the document. The comparison is then performed on the fingerprint rather than the whole text. In this work, the process of creating a fingerprint is as follow:

- Key phrasing: key phrases are extracted and each sentence is represented as a
 Bag Of Words.
- Hashing: a hash function is applied to the extracted key phrases to map them to a vector of integers.

232 3.3 Retrieval Of The Candidate Documents Phase

233

The process of candidate documents retrieval is through measuring similarities between the input document and the candidate documents at sentence level. In the fingerprinting method, the amount of similar fingerprints is used as similarity indicator between sentences; measuring similarity between two sentences or subdocuments is calculated by comparing the similarity percentage between a sentence's fingerprint and another sentence's fingerprint. For two sentences A and B, let h(A) h(A) and h(B) be their fingerprints with the corresponding length |h(A)| and |h(B)|. A similarity between A and B based on h(A)

- 241 h(A) and h(B) calculate the percentage of the similar fingerprints as [22, 23]:
- 242

243
$$sim(A,B) = \frac{|h(A) \cap h(B)|}{|h(A)|}$$
 (3.4)

244 If *Stm(A, B)* is greater than a threshold, subdocument B is selected as candidate 245 subdocument.

246 **3.4 Monolingual Plagiarism Detection Techniques**

247

The output of these methods will be used as feature vector that is used to training a machine learning classifier. In this work, several monolingual plagiarism detection techniques have been adopted:

251 3.4.1 N-Grams Similarity

252

253 The number of overlapping n-grams between two documents, d^3 the suspicious document

and d_t^{c} document t from the candidate document, will be counted. the overlapping total is divided by the length of the suspicious subdocument and length of the candidate subdocuments respectively in order to calculate recall and precision.

257 N-gram similarity score is expressed as[23]:

258
$$Score(d^s, d_i^c) = \frac{2^*(R-N)^*(P-N)}{(R-N)+(P-N)}$$
(3.5)

259 3.4.2 Longest Common Subsequence (LCS)

260

Given two documents, LCS is the longest string of matched tokens between these documents. LCS is that unlike n-grams (excluding unigram), LCS allows skip of matched ngrams. LCS score can be expressed as follows[24]:

264
$$ScoreLCS(d^{s}, d^{c}_{i}) = \frac{2*(R - LCS)*(P - LCS)}{(R - LCS)+(P - LCS)}$$

(3.6)

265 3.4.3 Dice Coefficient

266 The Dice similarity between two subdocuments A and B is defined as in[25]:

267
$$Dice(A,B) = \frac{2a}{2a+b+c}$$
 (3.7)

Where (a) refers to the matched key phrases or fingerprints present in both A and B, (b) refers to the key phrases or fingerprints present only in A, and (c) refers to those present only in B.

275 Fig. 3. Dice Coefficient Similarity

276 3.4.4 Fingerprint based Jaccard Similarity

277

Jaccard similarity is a very common set similarity measure that is used in a wide variety of applications. It is defined as

280
$$jaccard(A,B) = \frac{|A \cap B|}{|A \cup B|}$$
 (3.8)

281 Where A is the suspect fingerprints and B is the source fingerprints.

282

283 3.4.5 Fingerprint based Containment Similarity

```
284
```

285 Containment similarity is nearly identical to jaccard similarity, except the denominator is only 286 the number of elements in the suspect fingerprint. Again, let A be the suspect fingerprints 287 and B be the source fingerprints. Due to the size difference in of these fingerprints sets, an 288 asymmetric similarity measure is conducted based on containment similarity as [27]: 289 $C \circ n t a i n m e n t (A, B) = \frac{|A \cap B|}{|A|}$ (3.9)

290

291 3.5 Machine Learning Phase

The main idea is to feed the output of Monolingual plagiarism detection techniques to a machine learning classification framework. As shown in the previous sections, the monolingual plagiarism detection measures are only measure the similarity between suspicious document and candidate documents. However, their scores cannot indicate explicitly whether spacious document is plagiarized or not. To indicate explicitly whether suspicious document is plagiarized or not, we evaluated several classification methods for plagiarism detection.

299 3.5.1 Linear Logistic Regression

Logistic regression predicts the probability of an outcome that can only have binary response, also can handle several predictors (numerical and categorical). The multiple logistic regression model has the form as [24] :

304
$$\log(displag) = b_0 + b_1 X_1 + b_k x_k$$
 (3.10)

$$f(x) = p(displaginasized) = \frac{\exp^{b_0 + b_1 x_1 + \dots + b_k x_k}}{1 + \exp^{b_0 + b_1 x_1 + \dots + b_k x_k}}$$
(3.11)

306 3.5.2 Naive Bayes

305 306 307

300

The major advantage of NB algorithms is that they are easy to implement, often they have a superior performance. Naive Bayes (NB) can be defined as the conditional probability of

310 plagiarized class \mathcal{P}^{c} given monolingual feature vector m_{f}^{f} constructed as follows as[28]:

$$P(pc \mid mf) = p(c \mid s_1, ..., s_j) = p(pc) \prod_j p(s_j \mid pc)$$
(3.12)

311 Thus, the maximum posterior classifier is given as follows:

$$c^* = \arg \max_{c=c} p(c) \prod_{i=1}^{n} p(t_i | c)$$
 (3.13)

312 3.5.3 Linear Discriminate Analysis

313

The basic idea of LDA is to find a one-dimensional projection defined by a vector v that maximizes class separation. This method maximizes the ratio of between-class variance S_{P} to the within-class variance S_{W} in any particular data set thereby guaranteeing maximal separability as[29].

318
$$\max_{v} \frac{v^{'} S_{B} v}{v^{'} S_{w} v}$$
 (3.14)

319 3.5.4 Support Vector Machines

320

SVM is a featured machine learning technique that is developed for the binary classification
 task. SVM proposed to solve two-class problems by finding the optimal separating hyper plane between two classes of data. Suppose that X is set of labelled training points (feature

- 324 vector $(x_1, y_1), \dots, (x_n, y_n)$ where each training point $x_i \in RN$ is given a label $y_i \in RN$
- 325 $\{-1, +1\}$, where i = 1,...,n. The goal in SVM is to estimate a function $\frac{f(x)}{f(x)} = w x + b$
- 326 and to find a classifier y(x) = sign(f(x)) which can be solved through the following
- 327 convex optimization as[18] :

328
$$\min_{w,b} \sum_{i=1}^{n} [1 - y_i (w . x_i + b)] + \frac{\lambda}{2} \|w\| (3.15)$$

329 with λ as a regularization parameter.

330

332

331 4. EXPERIMENTAL RESULTS

333 In this section, several experiments have been conducted in order to evaluate the proposed approaches. First, several experiments have been conducted to evaluate key phrases 334 335 extraction methods. Secondly, Several experiments to empirically compare several monolingual plagiarism detection methods and three classification approaches which are 336 i)Linear Logistic Regression, ii) naïve Bayes, iii) SVM classifiers for Arabic-English Cross-337 338 language plagiarism detection. This research uses the same data set used by ALAA et al 339 2017 [24] for Arabic-English Cross-language plagiarism detection system. The data consists of 318 Arabic files are used for both training and test. All English files were used 340 341 for the comparison of both training and testing stages.

342 Table 4.1 : Detailed description of the experiment dataset

Dataset	Training	Test	Total
Arabic Files	200	118	318
English Files	34	20	54

343

344 **4.1 Experimental Results Of SVM Classifier**

In this experiment, SVM classifier is applied on testing set using 10-fold cross-validation. In
 this work, we used all monolingual plagiarism detection methods namely N-Grams Similarity
 (M1), Longest Common Subsequence (LCS) (M2), Dice Coefficient (M3), Fingerprint based
 Jaccard Similarity (M4), Fingerprint based Containment Similarity(M5) as a features for
 SVM.

350 Table 4.2 shows the performance in terms of the precision, recall, F-measure of Arabic-English Cross-language plagiarism detection by applying the SVM classifier with using 351 different combination set of features. The highest result yield by SVM classifier trained is 352 92% f-measure. As shown in Table 4.2, low performances are obtained when SVM uses 353 only one or two monolingual methods as features and high performances are obtained when 354 SVM uses more than three monolingual methods as features. This means that using all 355 356 monolingual plagiarism detection methods has an obvious positive effect on the quality 357 detection method.

358

^{*} Tel.: +967-771429933; E-mail address: sohaiki1986@gmail.com.

-	M1	M2	M3	M4	M5	PRECISION	F-MEASURE
-	0	1	0	1	0	0.74	0.85
	0	1	0	0	1	0.69	0.82
	0	1	0	0	0	0.59	0.74
	0	1	0	1	0	0.75	0.86
	1	1	0	1	0	0.73	0.84
	0	1	0	1	1	0.67	0.8
	1	1	0	0	0	0.4	0.57
	1	1	0	0	1	0.76	0.86
	0	1	0	0	1	0.71	0.83
	1	0	0	0	1	0.61	0.76
	1	0	0	1	0	0.73	0.84
	1	0	1	0	0	0.79	0.88
	0	1	1	0	0	0.74	0.85
	1	1	1	1	1	0.84	0.91
	0	1	1	1	1	0.85	0.92

360 Table 4.2 the performance of SVM Arabic-English Cross-language plagiarism

361 Detection

362

363

364 4.2 Experimental Results Of NB Classifier

In this experiment, NB classifier is applied on testing set using 10-fold cross-validation. The
idea is to show the best results obtained when the NB classifier is applied. In this work, we
used all monolingual plagiarism detection methods namely N-Grams Similarity (M1),
Longest Common Subsequence (LCS) (M2), Dice Coefficient (M3), Fingerprint based
Jaccard Similarity (M4), Fingerprint based Containment Similarity(M5) as a features for NB.

* Tel.: +967-771429933; E-mail address: sohaiki1986@gmail.com.

359

Table 4.3 shows the performance in terms of the precision, recall, F-measure of Arabic-English Cross-language plagiarism detection by applying the NB classifier using different combination set of features. The highest result yield by NB classifier trained is 89% fmeasure. This means that using all monolingual plagiarism detection methods has an obvious positive effect on the quality detection method. However, the results obtained by NB are lower than that of SVM.

	M1	M2	M3	M4	M5	PRECISION	F-MEASURE
_	0	1	0	1	0	0.53	0.69
	0	1	0	0	1	0.65	0.79
	0	1	0	0	0	0.56	0.72
	0	1	0	1	0	0.68	0.81
	1	1	0	1	0	0.39	0.56
	0	1	0	1	1	0.69	0.82
	1	1	0	0	0	0.61	0.76
	1	1	0	0	1	0.69	0.82
	0	1	0	0	1	0.75	0.86
	1	0	0	0	1	0.77	0.87
	1	0	0	1	0	0.74	0.85
	1	0	1	0	0	0.75	0.86
	0	1	1	0	0	0.7	0.82
	1	1	1	1	1	0.8	0.89
	0	1	1	1	1	0.79	0.88

377 Table 4.3 the performance of NB Arabic-English Cross-language plagiarism detection

378

379 4.3 Experimental Results Of Linear Logistic Regression Classifier

380

In this experiment, linear logistic regression classifier is applied on testing set using 10-fold
 cross-validation. The idea is to show the best results obtained when the linear logistic
 regression classifier is applied. In this work, we used all monolingual plagiarism detection
 methods namely N-Grams Similarity (M1), Longest Common Subsequence (LCS) (M2),

385 Dice Coefficient (M3), Fingerprint based Jaccard Similarity (M4), Fingerprint based 386 Containment Similarity(M5) as a features for NB.

Table 4.4 shows the performance in terms of the precision, recall, F-measure of Arabic-English Cross-language plagiarism detection by applying the linear logistic regression classifier using different combination set of features. The highest result yield by linear logistic regression classifier trained is 86% f-measure. This means that using all monolingual plagiarism detection methods has an obvious positive effect on the quality detection method. However, the results obtained by linear logistic regression are lower than that of SVM and NB.

Table 4.4 The performance of linear logistic regression Arabic-English Crosslanguage plagiarism detection

396

397

M1	M2	M3	M4	M5	PRECISION	F-MEASURE
0	1	0	1	0	0.49	0.66
0	1	0	0	1	0.61	0.76
0	1	0	0	0	0.52	0.68
0	1	0	1	0	0.64	0.78
1	1	0	1	0	0.36	0.53
0	1	0	1	1	0.64	0.78
1	1	0	0	0	0.57	0.73
1	1	0	0	1	0.67	0.8
0	1	0	0	1	0.73	0.84
1	0	0	0	1	0.74	0.85
1	0	0	1	0	0.73	0.84
1	0	1	0	0	0.71	0.83
0	1	1	0	0	0.67	0.8
1	1	1	1	1	0.76	0.86
0	1	1	1	1	0.74	0.85

398 399

400 5. RESULTS DISCUSSION

401

This paper aim to examine the proposed model and observation of the experimental results that have been achieved.

404 In the result tables in the fields (M1, M2, M3, M4, M5) there are values:

405 **1** : indicates that it was used in the experiment.

406 0 : indicates that it was not used in the experiment.

407

According to the experiments of Arabic-English Cross-language plagiarism detection with the
 SVM, NB, linear logistic regression classifiers, the highest result yield by SVM classifier
 with 92% f-measure.

According to the experiments of Arabic-English Cross-language plagiarism detection using SVM, NB, linear logistic regression classifiers with different combination of monolingual plagiarism detection methods namely N-Grams Similarity (M1), Longest Common Subsequence (LCS) (M2), Dice Coefficient (M3), Fingerprint based Jaccard Similarity (M4) and Fingerprint based Containment Similarity(M5) , the highest results obtained by all classifiers are achieved when most of the monolingual plagiarism detection methods used.

417 Furthermore, the obtained results with 92% f-measure were better than the previous 418 work of Aljohani [14]et al. (2014) at 89% and of ALAA [24]et al (2017) with 90%

419



421 Fig. 4. Conclusion of SVM And NB,LLR Result

422 6. CONCLUSION

420

424 Due to rapid growth of research articles in various languages, cross-lingual plagiarism 425 detection problem has received increasing interest in recent years. Cross-lingual plagiarism 426 detection is more challenging task than monolingual plagiarism detection. This paper aims 427 to design and implement a keyphrases based cross lingual plagiarism detection method. 428 This paper empirically investigates five different monolingual plagiarism detection methods 429 with three machine learning approaches namely naïve Bayes, SVM, and linear logistic 430 regression classifiers are used for Arabic-English Cross-language plagiarism detection. 431 Several experiments are conducted to evaluate the performance of the key phrases 432 extraction methods. In addition, several experiments to investigate the performance of 433 machine learning techniques to find the best method for Arabic-English Cross-language 434 plagiarism detection. According to the experiments of Arabic-English Cross-language plagiarism detection, the highest result yield by decision SVM classifier with 92% f-435 measure. In addition, the highest results obtained by all classifiers are achieved when 436 437 most of the monolingual plagiarism detection methods used.

⁴²³

^{*} Tel.: +967-771429933; E-mail address: sohaiki1986@gmail.com.

438 439 440	Future additior langua	work will aim to evaluate the current methodology with different language pairs. In n, future work will studied multilingual plagiarism detection i.e. include more than two ges.
111		
441		
442	ETHIC	
444		
445	CONS	ENT: NA
446	oono	
447		
448		
449		References
450		
451		
452	1.	Barrón-Cedeño, A., P. Gupta, and P. Rosso, Methods for cross-language
453	_	plagiarism detection. Knowledge-Based Systems, 2013. 50: p. 211-217.
454	2.	Potthast, M., et al., Cross-language plagiarism detection. Language
455	•	Resources and Evaluation, 2011. 45(1): p. 45-62.
456	3. 1	Pataki, M., A new approach for searching translated plaglarism. 2012.
457	4.	Gupta, P., A. Barron-Cedeno, and P. Rosso. Cross-language high similarity
458		Cross Language Evaluation Forum for European Languages 2012
409		Springer
461	5	Ebsan N FW Tompa and A Shakery Using a dictionary and n-gram
462	0.	alignment to improve fine-grained cross-language plagiarism detection in
463		Proceedings of the 2016 ACM Symposium on Document Engineering, 2016.
464		ACM.
465	6.	Ferrero, J., Agnes, F., Besacier, L. and Schwab, D., SemEval-2017 Task 1:
466		Cross-Language Plagiarism Detection Methods for Semantic Textual
467		Similarity. arXiv preprint arXiv:1702.03082, 2017c.
468	7.	Franco-Salvador, M., P. Rosso, and M. Montes-y-Gómez, A systematic
469		study of knowledge graph analysis for cross-language plagiarism detection.
470	-	Information Processing & Management, 2016. 52(4): p. 550-570.
471	8.	Speer, R.a.LD., J., ConceptNet at SemEval-2017 Task 2: Extending Word
472		Embeddings with Multilingual Relational Knowledge. arXiv preprint arXiv,
473	0	2017. 1702.(03082): p. 1704-03560.
474	9.	for Crossl anguage Plagiarism Detection arXiv proprint arXiv:1702.03082
475		$2017_2 = 1702 (03082)$
470	10	Glavaš G Franco-Salvador M Ponzetto S.P. and Rosso P. A.
478	10.	Resource-Light Method for
479		Cross-Lingual Semantic Textual Similarity, Knowledge-Based Systems
480		2017.
481	11.	Pera, M.S. and Y.K. Ng, SpamED: A spam E mail detection approach
482		based on phrase similarity. Journal of the American Society for Information
483		Science and Technology, 2009. 60(2): p. 393-409.

12. 484 Baroni, a.B., M. S., A new approach to the study of translations: Machine 485 learning the difference between original and translated text. Literary and Linguistic Computing, (2006). 21(3): p. 259-274. 486 487 13. Pouliguen, S., and Ignat., Automatic identification of document translations in large multilingual document collections. In Proceedings of the International 488 Conference Recent Advances in Natural Language Processing (RANLP'03), 489 490 2003. number 2002.(408): p. 401 491 14. Aljohani, A. and M. Mohd, Arabic-English Cross-language Plagiarism 492 Detection using Winnowing Algorithm. Information Technology Journal, 493 2014. 13(14): p. 2349. 494 15. Omar, K., B. Alkhatib, and M. Dashash, The Implementation of Plagiarism 495 Detection System in Health Sciences Publications in Arabic and English 496 Languages. International Review on Computers & Software, 2013. 8(4). 16. Kent, C.K. and N. Salim. Web based cross language plagiarism detection. in 497 498 Computational Intelligence, Modelling and Simulation (CIMSiM), 2010 Second International Conference on. 2010. IEEE. 499 Gottschalk, S. and E. Demidova, MultiWiki: interlingual text passage 500 17. 501 alignment in Wikipedia. ACM Transactions on the Web (TWEB), 2017. 11(1): 502 p. 6. 503 18. España-Bonet, C. and A. Barrón-Cedeño. Lump at SemEval-2017 Task 1: Towards an Interlingua Semantic Similarity. in Proceedings of the 11th 504 International Workshop on Semantic Evaluation (SemEval-2017). 2017. 505 506 19. Guan, J., A study of the use of keyword and keyphrase extraction techniques for answering biomedical questions. 2016. 507 508 20. Lossio-Ventura, J.A., et al. Combining c-value and keyword extraction methods for biomedical terms extraction. in LBM: Languages in Biology and 509 Medicine, 2013. 510 511 21. Frantzi, K., S. Ananiadou, and H. Mima, Automatic recognition of multi-word terms:. the c-value/nc-value method. International journal on digital libraries, 512 2000. 3(2): p. 115-130. 513 514 22. Lane, P.C., C. Lyon, and James A. Malcolm., "Demonstration of the Ferret plagiarism detector." Proceedings of the 2nd International Plagiarism 515 516 Conference., 2006.. 23. Hoad, T.C. and J. Zobel, Methods for identifying versioned and plagiarized 517 documents. Journal of the American society for information science and 518 519 technology, 2003. 54(3): p. 203-215. Alaa, Z., S. Tiun, and M. Abdulameer, CROSS-LANGUAGE PLAGIARISM 520 24. ARABIC-ENGLISH DOCUMENTS USING LINEAR 521 OF LOGISTIC REGRESSION. Journal of Theoretical & Applied Information Technology, 522 523 2016.83(1). 25. Lin, D. An information-theoretic definition of similarity. in Icml. 1998. 524 525 Citeseer. Yih, W.-T. and C. Meek. Improving similarity measures for short segments of 526 26. 527 text. in AAAI. 2007. 27. Yang, Y., et al., Gb-kmv: An augmented kmv sketch for approximate 528 containment similarity search. arXiv preprint arXiv:1809.00458, 2018. 529 530 28. Ngu, A.H., et al., Smartwatch-based iot fall detection application. Open Journal of Internet Of Things (OJIOT), 2018. 4(1): p. 87-98. 531

Altman, E.I., G. Marco, and F. Varetto, Corporate distress diagnosis:
Comparisons using linear discriminant analysis and neural networks (the Italian experience). Journal of banking & finance, 1994. 18(3): p. 505-529.