

A study on intelligent authorship methods in Persian language

Zeinab Farahmandpour¹, Hooman Nikmehr², Muharram Mansoorizadeh³, Omid Tabibzadeh⁴

Department of Computer Engineering
Faculty of Engineering
Bu-Ali Sina University
Hamedan, Iran

E-mail: ¹ zeinab.farahmandpoor@gmail.com

Department of Computer Architecture
Faculty of Computer Engineering
University of Isfahan
Isfahan, Iran

E-mail: ² nikmehr@eng.ui.ac.ir

Department of Computer Engineering
Faculty of Engineering
Bu-Ali Sina University
Hamedan, Iran

E-mail: ³ mansoorm@basu.ac.ir

Department of Linguistics
Faculty of Literature and Humanities
Bu-Ali Sina University
Hamedan, Iran

E-mail: ⁴ otabibzadeh@yahoo.com

Abstract

Author identification is an attempt to demonstrate the characteristics of the author of a piece of language information so that in the end, it would be possible to significantly distinguish a difference between various texts written by different people. The rapid development of Internet communication has caused Internet tools with anonymous identity, such as emails and weblog, to become popular communication methods for the perpetrators of illegal acts and has raised some security concerns. Persian language is of interest to a great number of different individuals and organizations for various reasons such as political, social, artistic, cultural and religious issues. In this paper, a number of intelligent writeprint methods which help automatic identification of a Persian writer based on their writing style are studied and compared. For this purpose, after collecting two different data-bases, five feature sets were used for extracting stylometric characteristics including lexical; syntactic, semantic and application-specific features. In this study KNN, Delta, Neural Networks, Decision Tree and Linear Discriminate Analysis classification methods were applied to these data-bases. The results and their comparison show that Linear Discriminate Analysis and KNN methods ranked first and second, respectively, in order of accuracy among the studied methods.

Key words: Authorship attribution, writing style, writeprint, human stylome

1. Introduction

Identification of the author means to find or get close enough to the real author of a text particularly in a collection of nominated authors. To achieve this goal, the distinctive information of the authors has to be recognized and analyzed using some models. Author identification is one of the oldest issues in Stylistics from one hand, and one of the newest one on the other hand. Author identification not only can attract the attention of human scientists but also, in a more practical context it can be a favorite issue for politicians, journalists and lawyers. Since the creation of words and documents, there have been debates about the ownership of words and identity and characteristics of the documents' author. So author identification can be defined as an attempt to demonstrate the characteristics of the producer or the author of a piece of language information. The key assumptions in author identification consist of access to a sample text which is definitely written by a member of the set of authors and also specification of the author of that text.

Many researchers believe that people use certain language patterns in their writings which function as the fingerprints of the author. This pattern is called writeprint. In this regard Van Halteren calls the collection of measurable specific characteristics which can be used for identification of a particular author, "Human Stylome" (2005). Since every human being has unique biometric characteristics and behavior patterns, it can be said that the specific language and writing related characteristics such as using specific patterns of words; layouts, syntactic and structural features that are called stylistic characteristics remain relatively constant in people. Learning and recognition of these characteristics with sufficient and high accuracy is a debating subject in the text author identification area. Theoretically, acceptable reasons can be put forward for the hypothesis of possibility of existence of such qualities. Since everyone learns language individually and their experiences as a language learner are different from other people, the language that different people learn would be different in various but small aspects.

Although using the document analysis science and also expert judgments can conventionally result in good conclusions in author identification, recent advances in statistical techniques and application of machine learning on accessible computer sources and objective and automatic inferences from texts, has created a scientific and more reliable method in author identification. This paper implements author identification using automatic machine learning methods and compares them for the Persian language. In Section 2 of this paper a background on author identification will be presented. The characteristics used for identification of the unique writing style of an author are explained in Section 3. Then in Section 4 the author identification methods are introduced and in Section 5 the recommended method for identification of a Persian author is discussed. Section 6 evaluates the recommended system and then compares different methods of machine learning implemented in Persian language author identification. Section 7 is dedicated to discussion and conclusion and Section 8 recommends some future works.

2. Research Background

The basic idea of statistical (or computational) identification of the author is the quantitative measurement of the text's characteristics so that the differences between texts written by different people can be distinguished. The first attempt to measure writing styles were carried out when Mendenhall did some studies on the Shakespeare's plays (1887). Later, in the first half of the twentieth century, some statistical research

projects were conducted on literary texts by Yule and Zipf (Yule 1938, 1944; Zipf 1932). Then to determine the identity of the author, more detailed investigations were performed on the data collection of “The Federalist Papers” by Mosteller and Wallace (1964) which is, without a doubt the strongest and most effective research ever conducted in this field, since it basically started the non-traditional approach of author identification which worked unlike all the conventional methods of the time (which was identifying the author based on human judgment). From then on and until 1990s, studies regarding author identification, trying to define the text attributes for assessment of the writing style (a branch of research called Stylistic), were conducted at the top which resulted in introduction of many but diverse measures such as the length of the sentences, the length of the word, the number of words’ occurrences, the number of repeating characters and the lexical richness functions (Holmes 1994, 1998). In another study, nearly 1000 different measures were introduced (Rudman 1997). It is worth mentioning that during this period, i.e. up until 1990, since the idea of creating a full automated system of author identification was hardly considered, the presented identification methods were only using computers rather than being computer-based.

Considering the amount of available electronic texts, the applicability of the author identification methods in variety of applications can be easily understood. Some of which are: literary researches (e.g. identification of the anonymous author of a controversial literary product from known authors) (Burrows 2002; Hoover 2004), artificial intelligence (such as identification of terroristic messages or statements from a set of known terrorists) (Abbasi and Chen 2005), criminal law (e.g. identification of authors of threatening messages, identification of authors of suicidal notes), civil law (the infringement of laws used in intellectual property issues) (Chaski 2005; Grant 2007), computer debates (e.g. identification of authors of malwares or computer viruses) (Frantzeskou et al. 2006). In such cases, the author attempts to hide his identity for criminal purposes.

Nowadays, statistical and machine learning methods are very common methods for author identification. Burrows was the first to use Principal Component Analysis in 1987 that attracted attention of researchers due to its high differentiation ability. Other multi-variable methods such as Cluster Analysis and Discriminant Analysis (Holmes 1992; Lowe et al. 1995) have produced good results in author identification as well.

Powerful machine learning techniques have had a lot of applications these years. Tweedie, Lowe, Mathews et al. used Neural Network for this purpose (Lowe et al. 1995; Tweedie 1996). While Diedrich, de Vel et al. have successfully applied SVM method for author identification (De Vel 2000; De Vel et al. 2001).

The expansion of the Internet, its borderless nature and the increase of online communications has resulted in creation of security issues in author identification. For this reason, implementation and use of authors identification in different languages is becoming increasingly more important. So far, studies have been conducted on English, Greek, Chinese and Arabic. Peng et al. performed tests on English texts, Chinese novels and Greek newspapers (Peng et al. 2003). Zheng et al. implemented author identification on Internet messages in English and Chinese, etc. (Tambouratzis et al. 2004; Stamatatos et al. 2001; Stamatatos 2008; Zheng et al. 2006). Feature extraction for author identification does not have the same difficulty in different languages. Most features of writing style have been designed for English and it is possible that these features are not equally applicable for other languages. The structural and linguistic

differences in different languages can possibly make the extraction of the features difficult.

Persian is the official language in Iran, Afghanistan and Tajikistan and tens of millions of people speak in this language. Different groups of people are related to this language for political, social, artistic, historical, cultural and religious reasons. The morphological differences of this language with English have resulted in negligence in implementation and designing of natural language processing tools consistent with Persian. This might have partly caused the negligence of author identification in this language. In this paper, some of the well-known techniques in author identification were implemented in Persian texts.

3. Stylistic Features

Previous studies on author identification have defined and classified some features for recognition of writing styles (Holmes 1994; Zheng et al. 2006; Stamatatos et al. 2000). Writing style features are extracted from the text, show the style of people's writings and simplify the author identification process. The writing styles of individuals are different in various but small aspects. The main classification of these features consists of lexical, characteristic, syntactic, semantic and application-specific features.

3.1. Lexical Features

The lexical features of a text are known as a set of clauses that are divided into components called tokens. Each token can be a word, a figure or a punctuation mark. The main advantage of these features is that they are applicable in any language and with any set of texts by using the tokenizer tool (which divides a text to a set of tokens).

The lexical richness function, models the distribution of words in a text. The size of the lexis is another measure which strongly depends on the length of the text. In the way that at the beginning of the text, the size of the lexis increases rapidly and by reaching the end the increasing trend of the lexis size slows down. In another method, regardless of the content information, the text is considered as a set of words each of which is repeated a certain number of times in the text (Fabrizio 2002).

3.2. Characteristic Features

The characteristic features consider the text as a set of character strings. Some of other features of this category can be the number of characters in the text, numbers in the text, capital and small letters, repetitious and punctuation characters (De Vel et al. 2001; Zheng et al. 2006). It has been proven that this group of information, which can be easily obtained for every natural language, is relatively suitable for determination of the writing style (Grieve 2007). By Cavnar and Trenkle (1994), n -gram is introduced as an n -character piece of a longer character string. This feature can extract subtle style differences', including lexical information and its remarkable noise tolerance (spelling errors) is one of its advantages. The Character n -gram extraction method is very simple in terms of computation. For example Character 4-gram for the current paragraph would be as follows: |The-|, |char|, |acte|, |rist|, |ic-f|, |eatu|.... In a study (Forsyth and Holmes 1996), it was found that in an author identification process, the Character n -gram with different " n "s is a more suitable feature than a lexical feature. A diverse number of authors have shown that the number of Character n -gram repeats might be better than even the grammatical features for obtaining lexical information. They are useful without the need for the language information background. It is

noteworthy that the character n -gram feature is conducted on words related to the content and the subject as well. This feature generates useful information only if all the texts of our data-base are on the same subject.

The important issue in Character n -gram method is defining n as the length of the desired string. A large n not only provides better lexical and background information but also it presents better thematic information. However, a large n would increase the dimensions of the feature representation very much by providing hundreds of thousands features for a text. On the other hand, a small n (2 or 3) would be able to display the information of sub-strings (such as syllable) however, this extracted attribute is not sufficient by itself for displaying the text and background information. The procedure of choosing an optimum “ n ” depends on the natural language specifications. For instance, since Greek and German languages have longer words in comparison with English, using larger “ n ”s for them is more suitable. Defining a constant, would prevent application of n -grams with different lengths (Forsyth and Holmes 1996; Houvardas and Stamatatos 2006).

3.3. Syntactic Features

Since an author unconsciously tends to use the same syntactic patterns, using syntactic features provides a more accurate display of the text (Stamatatos et al. 2001). This provides a more reliable method for identification of the author than the lexical features. Extracting syntactic information requires powerful and accurate tools for processing the natural language (NLP) that be able to syntactically analyze texts. Extraction of syntactic features is a text-related procedure and requires a parser which can analyze a particular natural language with high precision. Moreover, the syntactic attributes generate a set of features with low precision due to inevitable errors of the parsers. Baayen, Halteren and Tweedie (1996) were the first people to use syntactic information for author identification. Using the complete tree decomposition procedure of every sentence, they analyzed the corpus set syntactically and extracted the number of occurrences of the re-writing rules. Each re-writing rule expresses one part of a syntactic analysis. For example an adverb preposition phrase consists of the preposition and a noun phrase followed by a prepositional predicate. This detailed information shows the syntactic group of each word and the composition of words to make phrases and other structures. The reported results suggest that this attribute is better than lexical and lexical richness features. The extracted attribute includes the number of noun phrases and clauses, the length of noun phrases and clauses, etc.

The occurrence of the common words, such as prepositions, articles and pronouns, which are called function words, is one of the very good features for author identification (Burrows 1987; Argamon and Levitan 2005). The main advantage of this attribute is its independence to the text’s subject since these words do not carry any semantic information. Another important advantage of this attribute is that the author uses the function words unconsciously this means that the writing style of the writer of different texts with different subjects can be reflected in the function words used by him/her. The choice of these words is usually based on the criteria that are introduced by linguistics. A simple yet successful method to define function words for author identification is to extract words with high occurrences in the corpus of the nominated authors. In studies conducted by Burrows (1987, 1992) a set of 100 function words have been recognized suitable for displaying the author’s style.

3.4. Semantic Features

Semantic features more precisely analyze of the text and provide more applicable features. This is done by obtaining the semantic dependency graph including the semantic features and modification relations.

The tense and the mode of the verbs used by the author and the semantic similarity between the words in the text are another example of such features. The words and phrases that join clauses together are known as juncture adjuncts or letters. Different patterns of using conjunctions (Argamon et al. 2007) result in remarkable differences in writing styles. According to Hildi and Mason conjunctions are categorized into three groups of Descriptive, Exponential and Elaborative. The descriptive adjuncts deepen the content of the text by illustrating and re-focusing. Using descriptive conjuncts can have a good effect on the text and provide solidarity throughout the whole text. The exponential conjuncts would add more information into the text, even sometimes antithetical to the current concepts. Repetitious usage of these conjunct aggregates the text information but if these conjuncts are not used enough it is possible for the reader to get lost among the multiple and various concepts. The elaborative conjuncts describe the text with details or logical relationships.

McCarthy et al. (2006) defined another approach for semantic measure extraction. According to the WordNet (Fellbaum 1998), they have well recognized the information about the word's synonyms. Moreover, they used the Latent Semantic Analysis for lexical features in order to automatically recognize the semantic similarities between words (Deerwester et al. 1990).

3.5. Application-specific Features

The application-specific features can be used for a more delicate presentation of differences in people's styles in a specific confine. The application-specific features in fields such as electronic messages are the application of structural measures which includes using greetings in messages, different types of signatures, dips, the length of the paragraph, etc. (De Vel et al. 2001; Zheng et al. 2006; Teng et al. 2004; Li et al. 2006). In addition, when the text is in the form of HTML, the HTML Tag distribution-specific features (De Vel et al. 2001) such as the colors and the sizes of the fonts can be used as application-specific features (Abbasi and Chen 2005). These features are very important in short texts, where the stylistic features of the text content are not enough, and accurate tools are required for extracting them.

4. Data classification Methods

There are different methods for data classification Statistic Pattern Recognition, which works based on statistical model, Neural Pattern Recognition, which is based on neural network and Syntactic Pattern Recognition, which works based on element structures (clustering).

4.1. Density Estimation by K Nearest Neighbor method

The KNN method is a probabilistic model in which to identify the author of a text, K texts from the authors that have the most similarity to the desired text are recognized and the author who has the most number of texts (n) among the K texts is selected as the answer (Fukunaga 1990).

In this method, to identify the author of an anonymous text, K samples of the training texts, in the feature space, which have the least distance from the anonymous text, are

taken and the author who has the most number of texts among these K texts is selected as the author of the unknown text. This method estimates the probability of the anonymous text belonging to the i^{th} author through equation 1:

$$\hat{p}(x|i) = \frac{n_i}{k} \quad (1)$$

In equation (1), n_i is the number of texts that belong to the i^{th} author in the set of K texts which have the most similarity to the anonymous text. K is constant in every probability calculation.

Suppose that there are n training texts with their authors known and we are going to classify a text with an anonymous author, x , in one of the above c classes (author). To use the KNN method, first a scale should be selected for measuring the distance (such as the Euclidean distance). Then the following steps should be taken:

- A. The distance of each training text with the text x is calculated; K texts which are the closest to the text x are selected regardless of their class (author).
- B. The number of each author's texts in this K texts is determined and called n_i ($i = 1, 2, \dots, c$). where $\sum_i n_i = k$.
- C. The text x is assigned to the author who has the most number of texts among the K texts, i.e. the largest n_i .

4.2. Delta method

In 2002, a method was presented by Burrows exclusively for author identification namely Delta. It calculates the normalized difference of the number of word occurrences (a set of n desired words, usually the most repetitious words) in the text D of the known author and the number occurrences of the same words in text D' with anonymous author by equation (2). In this equation σ_i is the standard deviation of the i^{th} word in the comparative set and $f_i(D)$ is the number of occurrences of the word w_i (the number of occurrence of the i^{th} desired word) (Argamon 2008). This method calculates the distance between the two texts, D being the test anonymous text and D' the written training text by a nominated author, as follows:

$$\Delta_B^{(n)}(D - D') = \sum_1^n \frac{1}{\sigma_i} |f_i(D) - f_i(D')| \quad (2)$$

This method classifies the nominated authors of the D' texts considering their distances from the D test text, while after each difference, the number of words occurrence reduces by factor $\frac{1}{\sigma_i}$. In this equation the test document is classified in the group closest to the specified training document. This method uses the function word distribution feature for classification.

4.3. Neural Networks

The neural network algorithm is one of the most widely used and most practical modeling methods for large and complex problems. In 1958, Rosenblatt introduced the concept of single-layer Perceptron as a useful tool for classification of a set of data into two classes and he offered a proof for stability of the perceptron learning rule. Each neural network consists of an input layer. Each node in this layer is equal to one of the predictor variables. Each input node attaches to all nodes of the hidden layer. The nodes of the hidden layer can be attached to another hidden layer or the output layer. The

output layer can consist of one or more output variables depending on the problem (Merriam and Matthews 1994; Matthews and Merriam 1993). In this study, a three-layer neural network is used for implementation of the author identification system.

4.4. Decision Tree

The Decision trees are used for displaying different concepts in artificial intelligence such as the structure of sentences, equations, game modes, etc. training of a tree is a method for approximating the objective functions with continuous and discrete values. This method is resistant against the noise in the data and it is able to learn the conjunction and disjunction syntax of the predicates. Due to the high efficiency of the decision tree method in problems with high volume data, this method is used in Data Mining. The decision trees are made by consecutive division of data to separated groups and the purpose of this process is to increase the distances of the groups in each division.

In the decision tree ID3, a statistical value is used which is called Information Gain and which determines to what extent a feature is able to divide training examples based on their classification (Uzuner and Katz 2005). Likewise, the features extracted from texts of each author are given to the decision tree. The decision tree method selects features with high entropy. The feature which is in the root has a higher importance than other features. This method uses entropy and interest rate for classification of selected features and in the end composes a tree of features that are classified as important. So that the features which are located in the root of the tree are more discriminant than other features and as we go farther from the root, the importance and separating quality of the selected features reduces. At the end, the precision of the decision tree is also calculated.

4.5. Linear Discriminant

A common algorithm for classification and data dimension reduction is Linear Discriminant Analysis (LDA). This method easily manages the modes in which the numbers of members of the classes are not equal. This algorithm intends to maximize the ratio of data variance of each class to in-class variance in every data-base. In-class dispersion is calculated as the covariance for each class. This method takes the data set into another space. In LDA, the form and location of the original data set do not change when converted into another space and the endeavor is to increase the separation of the classes (Balakrishnama and Ganapathiraju 1998). This method seeks a series of features among the extracted features of the texts that can best divide the author's classes and since it is a linear method, it ignores all nonlinear features.

5. Recommended method for Persian Language

The automatic author identification has been implemented and used for various languages; however it is not applied in this manner to Persian language. In this Section, the works conducted for implementation of the automatic identification in Persian is explained. First, the stylistic features extracted from Persian texts are discussed. Then methods for selecting features of data volume reduction and elimination of additional information will be investigated. Finally, the collected data-base for this application in Persian would be analyzed.

5.1. Stylistic Features Extraction

Different features have been used in different languages for author identification. In this study the most distinguishing these features are selected. Some of the previously discussed features are discriminant by themselves but some others are applicable and provide better results in combination with other features.

In the following a variety of selected features to be extracted from Persian texts are introduced. Some of the features mentioned in English languages are not used due to lack of similar extraction system in Persian.

5.1.1. Lexical Features

The lexical features used in this paper include lexical richness, n -gram and other features that are explained in detail in the following sections.

5.1.1.1. Lexical Richness

The lexical richness is used to measure the lexical frequency in a text. Various metrics are used for measuring the amount of richness for the purpose of author identification. The most common metric of this group is type-token ratio which is calculated as $\frac{V}{N}$; where V is the number of words and N is the number of tokens in the text. Since these features are affected by the length of the texts, many researchers have presented functions describing these features that are claimed to be length-independent. However our investigations did not show any proof for this claim.

Some scholars to achieve better results have used a set of several lexical richness functions, instead of one, in association with multi-variable statistical techniques (Holmes 1992). These functions are not ponderous in terms of calculation. According to some studies (Tweedie and Baayen 1998), the majority of lexical richness functions is very much dependent to the length of the text and is relatively impermanent for texts shorter than 1000 words.

For measuring the lexical richness, Stamatatos et al. (2000) and Baayen et al. (1996) have used a collection of 5 functions namely K , R , W , S and D for Yule (1944), Honore (1979), Brunet, Sichel (1975) and Simpson (1949), respectively. These functions are defined by expressions (3-7).

$$K = \frac{10^4 (\sum_{i=1}^{\infty} i^2 v_i - N)}{N^2} \quad (3)$$

$$R = \frac{(100 \log N)}{\left(1 - \left(\frac{V_1}{V}\right)\right)} \quad (4)$$

$$W = N^{V^{-\alpha}} \quad (5)$$

$$S = \frac{V_2}{V} \quad (6)$$

$$D = \sum_{i=1}^V V_i \frac{i(i-1)}{N(N-1)} \quad (7)$$

In these equations V_i is the number of words that have repeated i times and α is a constant parameter with the value of 0.17. For each Persian text in this paper, these functions are calculated and a vector is composed by these five features.

5.1.1.2. n -gram

The n -gram feature works well in texts with noisy inputs. Since in this method, each character string is decomposed to smaller pieces, any kind of errors in the word would affect only a few numbers of pieces and other pieces of that word would remain flawless.

To create an n -gram profile, first the input text is read and then the following steps are taken:

- The text is divided to separate tokens. Numbers, control characters and punctuation marks would be spared.
- All possible n -grams of n are generated.
- The n -grams are entered in a hash table and the number of their occurrence in the text is founded.

The largest number of n -gram occurrences would occur in uni-gram ($n=1$) which indicates the alphabetic characters in the language. In this paper, the n -gram feature is used due to its importance in the above mentioned studies. In this study, a list of all 72 possible uni-gram in the Persian language is used:

{ ا, ب, پ, ت, ث, ج, چ, د, ذ, ر, ز, س, ش, ص, ض, ط, ظ, ع, غ, ف, ق, ک, گ, ل, م, ن, و, ه, ی, ر, ی, ه, ج, چ, د, ذ, ر, ز, س, ش, ص, ض, ط, ظ, ع, غ, ف, ق, ک, گ, ل, م, ن, و, ه, ی, ر, ی, ه }

The number of possible bi-grams, tri-grams and forth-grams is very large and obtaining the distribution of their occurrence in the texts would result in features with too many dimensions and too small values close to zero. Thus, in this paper only the bi-grams, tri-grams and forth-grams that the number of their occurrence are more than 10% of the most repetitious possible bi-grams, 5% of the most repetitious possible tri-grams and 10% of the most repetitious possible forth-grams in the data-base introduced in Section are used. These threshold values are calculated considering measures such as the number of authors and the length of the texts and by experience. The number of possible and selected bi-grams, tri-grams and forth-grams are listed in Table 1 (The number of possible applicable n -grams used in each data base).

Table 1. The number of possible applicable n -grams used in each data base.

n -gram \ Number	possible applicable features	selected features in University database	selected features in Writer database
Uni-gram	72	72	72
Bi-gram	$72 \times 72 = 5184$	375	71
tri-gram	$72 \times 72 \times 72 = 373248$	560	117
forth-gram	$72 \times 72 \times 72 \times 72 = 26873856$	257	97
Total	27252360	1264	357

5.1.1.3. Other lexical features

In order to benefit the lexical information as much as possible, the following lexical features are used:

- average number of characters used in words
- average number of characters used in sentences
- average number of words used in sentences
- number of all short words equal or shorter than 3 characters
- distribution of Persian characters in the text
- distribution of paragraphs in the text
- distribution of words with the length of 1 to 30 in the text
- distribution of English characters in the text
- distribution of numeric characters in the text
- distribution of half space character in the text

5.1.2. Syntactic features

These features find the rules of sentence generation. Features that are extracted from this structure include noun, adjective and adverbial phrases. In this paper, a POS tagger (Part-Of-Speech tagger) is used for extraction of the number of noun, adjective and adverbial phrases. This POS tagger relates every word to one of these tags.

5.1.2.1. Designing POS tagger for Persian

POS taggers are language related tools that recognize and label words based on their role in the sentence. Due to lack of a suitable software tool for automatic identification of nouns, adjectives and adverbs in a Persian text, such tool is designed to complete the research.

According to Tabatabaee (2009), the structure of noun, adjective and adverbial phrases can be simple or complex. If the structure is a simple one, there is no other way of identification than listing all possible nouns, adjectives and adverbs and then comparing the words with this list. However, if the structure is complex, the structural methods mentioned By Tabatabaee (2009) would be used. In this paper the method recommended in the Persian Orthography & Spelling Dictionary (ZandiMoghadam and Sadeghi 2008) for dividing the words is followed. For instance, an adjective like “زیبا” (beautiful) has a simple structure and it should be put in a list to be identifies; but a noun like “کارگر” (worker) has a complex structure and it consists of “کار” and “گر” whereas “کار” is a noun and “گر” is a noun suffix. Some of the composing structures of noun, adverb and adjective are as follows:

Noun	: (noun+- گر) / (noun+- بان) / (noun+- کار)...
Adjective	: (noun+-s) / (noun+-ناک) / (noun+-وار) ...
Adverb	: (adjective+-آنه) / (noun+بلا) ...

There are some nouns, adverbs and adjectives that are neither simple nor complex, so they cannot be categorized in these structures. A complete list of those words based on Tabatabaee (2009) is made and used as well. For labeling, when come across a word, first the lists of simple and irregular nouns, adverbs and adjectives are searched for that word. If the word is found in those lists, it is labeled accordingly. Otherwise, its composing structure is studied. If it is consistent with one of the known and collected structures for nouns, adverbs and adjectives, then it will be labeled appropriately.

5.1.2.2. Function Words

In Persian language, a set of 922 words are introduced by Davarpanah et al. (2009) as the function words which include the most repetitious pronouns, verbs, conjunctions, prepositions and articles and they are used in this paper. The number of occurrences of each type of function words, “است” (is), “من” (me/I) and ... in each text is found and used as a feature.

5.1.2.3. Punctuation Characters

The number of punctuation characters in the text is one of the syntactic features that can differentiate people’s writing styles. In Persian, the number of usage of punctuation marks by each author such as “?” and “!” in each text is extracted and utilized.

5.1.3. Semantic Features

The way authors use conjunctions can differentiate in the texts written by them. According to classification of Persian conjunction adjuncts and based on the theoretical framework of Hildi and Mason, the conjunctions are categorized in three groups (Jafari 2008). In this study, based on this classification, the type of adjuncts in each text was determined and the number of occurrence of each was calculated and the results were used as semantic features. In any case, to have more prepositions, each group was expanded by adding of synonyms so that each preposition would fit in one of the categories.

5.1.4. Application-specific Features

Considering the type of texts in the available databases, including books and articles, some of application-specific feature used in this paper, were the distribution of tab, enter, space and line feed characters in the text.

5.2. Summary of the extracted features

More than 1500 features along with the number of their occurrence collected from the databases introduced in Section 3-5, are listed in Table 2 (Features extracted from texts).

Table 2. Features extracted from texts.

Features	Sub Features	Number of used Features
Lexical Features	Lexical richness	7
	<i>n</i> -gram	Depends on length of text and number of authors
	Average length or words	1
	average length of sentences' characters	1
	average number of words in sentences	6
	number of all short words equal or shorter than 3 characters	3
	distribution of alphabet characters in texts	72
	distribution of Persian characters in texts	1
	distribution of the number of paragraphs in texts	1
	distribution of words with the length of 1 to 30 in texts	30
	distribution of English characters in texts	1
	distribution of numerical characters in texts	1
	distribution of half space character in texts	1
	Syntactic Features	distribution occurrence of noun, adjective and adverbial phrases
Function words		922
distribution of punctuation characters		1
Semantic Features	distribution of adjuncts in Persian	36
Application-specific features	distribution of tab character	1
	distribution of enter character	1
	distribution of space character	1
	distribution of empty lines	1

5.2. Pre-processing of information collected from texts

Data pre-processing includes all conversions that are done on the preliminary data and changing them to a form simpler and more effective for the next processes such as classification.

There are different tools and methods for data pre-processing such as normalization, which converts data to new one with suitable changing confine and distribution and dimension reduction, which is used for elimination of the repetitive, overflow or irreverent data for classification (Bao 2002; Takamy 2005). Features with larger values have greater effect on the classification cost function that does not necessarily mean they are more important in, so it is considered a diverse effect. In this study the normalization and correlation reduction methods, which are going to be introduced in the next section, were used.

5.2.1. Normalization of Feature's Values

Normalization maps Data to $[-1, 1]$ confine through a linear conversion. For this purpose, if μ_i is the estimation of the average of the i^{th} feature and σ_i is the variance estimation of that feature on n samples, then the i^{th} feature is being normalized by equation (8) (Theodoridis and Koutroumbas 1999).

$$y_{ij} = \frac{x_{ij} - \mu_i}{\sigma_i} \quad (8)$$

Where the mean of converted data in each dimension (feature) is equal to zero and its variance is equal to one.

In normalization of data, it should be taken into consideration that the data extracted from the training texts (the classification design data) and the data extracted from the test texts must be normalized through the same method. Also, the normalization of test data should be done with the mean and variance resulted from the training data.

The dimension of the data must be reduced after normalization. The reason for this dimension reduction can be considered as simplification of the next analyses, performance improvement of the divider based on better display, elimination of repetitive information or an attempt for detection of basic structure by obtaining the graphical display of the data.

5.2.2. Elimination of high correlated features

If one or more features have a high correlation with another feature, they would be considered repetitive. A subset of feature is considered suitable if they have a low correlation. In this study, in order to reduce the correlation and thus reducing the repetitive features, the correlations were calculated and one of two features that had a correlation equal or more than 95% with each other was eliminated.

5.3. Data-bases

Author identification is often used when the author of an anonymous text is searched in a small set of authors. For example, to identify the poet of a poem that is supposed to belong to Hafez.

A suitable training and test sets for an author identification case shall be studies in regard to type and subject so that the identity of the author is the only or the most important dividing factor of the texts (Stamatatos 2009). Ideally, all the texts of the training set shall be exactly related to one subject; however there are very few sets with

this characteristic. Age, education and nationality are other factors that should be closely studied in preparation of ideal evaluation sets. This would reduce the probability of selection of the style of a wide group of people to the writing style of the desired author. Also, all the texts must be written during the same period of time so that the probability of style changes is eliminated (Can and Patton 2004). Due to the lack of an standard data-base (Stamatatos 2009) for Persian writer that can meet the requirements in this area, a database from Bu-Ali Sina University (*University* for abbreviation) was collected. To prepare the database, texts with the exact length of 2009 words, (1500 words for training and 509 words for test,) were collected from 20 junior undergraduate engineering students from the University on a specific subject.

since the texts of the *University* database were informal, (this would result in inaccuracy of many features such as the distribution of function words,) and also became these text were short (this would reduce the validity of some measuring features (Luyckx and Daelemans 2011),) another texts from books and articles written by eight recent Persian authors (Ali Ashraf Sadeghi, Mohammad Ali Forughi, Mojtaba Minavi, Abul Hassan Naafi, Mohammad Amin Riahi, Ahmad Samiee, Fathollah Mojtabaee and Hussein Masumi Hamedani) were collected on subjects in literary field and literary text analysis as *contemporary writers* data-base (*Writers* for abbreviation). For each writer, two documents were collected, with at most 7750 words from which 5000 words were assigned to training and the remaining for test.

According to the profile-based method, in this paper, for each author 70% of the texts are selected for training and 30% for test. In another evaluation, using a sample-based method, each writer's profile is randomly divided into K sub-samples. For making of each sub-sample, that is a set of words, the value of K was increased to such extent that the pieces of the text which are used for training and testing, are able to demonstrate the stylistic features. The value of K can be obtained by experience. From this number of sub-samples, one for test and $K-1$ sub-sample are given to classification models for training. Then the whole process is repeated K times, so that each one of these K samples are used only one time for the test. Then the average of all K resulted from K tests is calculated. This method is called K -fold cross validation. In this study both methods are used for evaluation.

6. Evaluation and comparison of machine learning methods in Persian author identification

In this section, the results of implementation of author identification system in Persian language are presented and in the end, the results compared with one another.

Table 3 (Comparison of Classification methods on two Data-bases) indicates the precision of 5 classification methods in identification of the author of an anonymous text on each data-base separately.

Table 3. Comparison of Classification methods on two Data-bases.

Method Data-base	KNN	Delta	Decision tree	Neural Network	LDA
University	70%	50%	21.6%	15%	81.6%
Writers	100%	87%	57.5%	37.3%	100%

As mentioned before, *University* database would not provide much precision due to the text's short length and informality of the texts. This result can be observed in all classification methods in the table. As high appeared in Table 3, KNN method, is displaying high precisions on both databases. As expected and shown in Table 3, the statistical method of Delta, which is specifically developed for author identification applications, provides acceptable results. However since it uses function word feature which is a syntactic feature, it generates less precision in comparison with KNN which uses all features. Decision tree method has trims the branches in order to prevent heightening and as a result, as it is indicated in Table 3, due to elimination of features important in the decision making process and utilizing few features in the process, does not provide a suitable precision in comparison with other methods.

As presented in Table 3, the precision of Neural Network method in text author identification is much lower than others. This can be due to the small number of training texts for each writer and large feature dimension of samples which can result in a phenomenon named Curse of Dimensionality. Thus, in applications such as the present study, where the number and length of the training texts are small and the number of extracted features for each text is large, using the neural networking method as the classifying method is not recommended. As it can be seen in Table 3, LDA method which uses all kinds of features is showing high precisions on databases, The reason of which can be the high linear coherency of the extracted features of the texts. The studies conducted in this paper suggest that this method generates better result comparing to other methods if only this kind of features is eliminated from it.

The features in the root of decision tree which was implemented on our data-bases included distribution of adverbs in the texts which is a syntactic feature, the average of character distribution in the sentences and the number of desired bi-gram occurrences in the texts that are the lexical features. The level one features included the feature of distribution of adjective phrases from syntactic features and distribution of desired conjunctions in texts which is a semantic feature. The selected feature for the second level is also distribution of the punctuation characters, which is a syntactic feature. This concludes that the mentioned features have the most separating ability among all features and can be used as the important features in author identification applications.

In order to study the effect and the role of different types of features in author identification more closely, in another implementation, the mean of precision of all kinds of features in author identification in two data-bases of *Writers* and *University* are calculated and categorized based on the classification methods. The results are shown in rows 3 to 5 of Table 4 (Comparison of the average precision of classification methods versus features types on both data-bases). Delta method that only uses function words is not addressed in this table. The last row of the table presents the mean of the precision of all four classification methods on all features. This table expresses that the set of lexical features has the higher precision and greater separating ability than the other sets. Table 3 also indicates that the next methods in the ranking belong to the syntactic, semantic and application-specific features respectively.

The result of this comparison is consistent with the result obtained from the decision tree as well.

Table 4. Comparison of the average precision of classification methods versus features types on both data-bases.

Classification Methods \ Feature	Lexical and Characteristic	Syntactic	Semantic	Application-specific
KNN	80%	66.5%	49.5%	49.5%
Decision tree	25%	21.88%	18.75%	12.5%
Neural network	11%	12.5%	16%	8.75%
LDA	38.12%	29.37%	30%	39.37%
average precision	38.53%	32.56%	28.56%	27.53%

This is also noteworthy that it is possible that one feature out of one “ kind” of features has a better position in terms of separating quality to the kind of feature that it belong to among other type of features. Thus, the precision of separating quality of a type of feature being higher does not mean every feature in that group would have a higher precision. Therefore, it is possible that a feature belonging to the syntactic features group has a higher precision in author identification in comparison to other features.

7. Discussion and conclusion

In this study, in addition to designing and implementing of a system for automatic Persian author identification of an anonymous text, a comparison study is conducted on machine learning methods for identification of a Persian author using two databases.

In this research, five classification methods including Linear, *K*-Nearest Neighbors (KNN), Delta, Discriminant Analysis, Decision Tree and Neural Networks are implemented and compared. The results show that implementation of LDA method on the data-bases, results in more precise results compared to other methods.

The results of implementation of these methods show that the shorter the training and test texts the less accurate the author identification will be. However, for short texts, the extracted features shall be selected carefully in order to present the writing style of the author (Hirst and Feiguina 2007). Thus, texts from *University* database do not generate acceptable precision due to the texts’ short length and being informal.

The most discriminant features in the collected data-base in Persian includes the distribution of adverb and adjective in the texts (the output of the software designed for this paper to recognize noun, adverb and adjective in the text), the average of distribution of characters in sentences, the number of desired bi-gram occurrence in the text, distribution of conjunctions and distribution of punctuation characters in the text.

The success rate of the system designed in this study on Persian literature is 90% in averages. Considering that the current study is the first attempt for presenting an automatic and machine learning method for identification of an anonymous Persian author, the obtained precision seems to be quite promising and encouraging.

Due to lack of accurate available natural processing tools for Persian, we encountered some problems in the author identification and thus we designed and implemented a number of these tools. Such as a POS tagger; however, some of these tools were not available and so we were not able to use them. The usage of those tools would definitely increase the precision of the system.

References

- Abbasi, A., & Chen, H. (2005). Applying authorship analysis to extremist group web forum messages. *IEEE Intell Syst App*, 20, 67–75.
- Argamon, S. (2008). Interpreting Burrows' Delta: Geometric and probabilistic foundations. *J Lit Ling Comput*, 23, 131–147.
- Argamon, S., Chase, P., Dhawle, S., Raj, S., Navendu, H., & Levitan, G. S. (2007). Stylistic text classification using functional lexical features. *J Am Soc Inf Sci Technol*, 58, 802-822.
- Argamon, S., & Levitan, S. (2005). Measuring the usefulness of function words for authorship attribution. In: the Joint Conference of the Association for Computers and the Humanities and the Association for Literary and Linguistic Computing, June, Victoria, BC, Canada.
- Baayen, H., Halteren, H. V., & Tweedie, F. (1996). Outside the cave of shadows: using syntactic annotation to enhance authorship attribution. *J Lit Ling Comput*, 11, 121-132.
- Balakrishnama, S., & Ganapathiraju, A. (1998). Linear discriminant analysis- a brief tutorial. Technical Report MSSTATE-ECE-98, Mississippi State University.
- Bao, H. T. (2002). Introduction to knowledge discovery and data mining. Hanoi, Vietnam: Institute of Information Technology National Center for Natural Science and technology.
- Burrows, J. F. (1987). Word-Patterns and story-shapes: the statistical analysis of narrative style. *J Lit Ling Comput*, 2, 61-70.
- Burrows, J. F. (1992). Not unless you ask nicely: The interpretative nexus between analysis and information. *J Lit Ling Comput*, 7, 91-109.
- Burrows, J. F. (2002). Delta: a measure of stylistic difference and a guide to likely authorship. *J Lit Ling Comput*, 17, 267-287.
- Can, F., & Patton, J. M. (2004). Change of writing style with time. *J Comput Humanities*, 38, 61-82.
- Cavnar, W. B., & Trenkle, J. M. (1994). N-gram-based text categorization. In: SDAIR 1994 3rd Annual Symposium on Document Analysis and Information Retrieval, Las Vegas, US, 161–175.
- Chaski, C. E. (2005) Who's at the keyboard? Authorship attribution in digital evidence investigations. *Int J Digit Evid*, 4, 1-13.
- Davarpanah, M. R., Sanji, M., & Aramideh, M. (2009). Farsi lexical analysis and stop word list. *J LIBR HI TECH*, 27, 435 - 449.
- De Vel, O. (2000). Mining E-mail authorship. In: Workshop on Text Mining, ACM International Conference on Knowledge Discovery and Data Mining; 20 Aug, Boston, MA, US, 1-7.
- De Vel, O., Anderson, A., Corney, M., & Mohay, G. (2001). Mining e-mail content for author identification forensics. *ACM SIGMOD Rec*, 30, 55-64.
- Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., & Harshman, R. (1990). Indexing by latent semantic analysis. *J Am Soc Inf Sci*, 41, 391–407.

- Fabrizio, S. (2002). Machine learning in automated text categorization. *ACM Comput Surv*, 34, 1-47.
- Fellbaum, C. (1998). *WordNet: An electronic lexical database*. Cambridge, Massachusetts, USA: MIT Press.
- Forsyth, R. S., & Holmes, D. I. (1996). Feature-finding for text classification. *J Lit Ling Comput*, 11, 163-174.
- Frantzeskou, G., Stamatatos, E., Gritzalis, S., & Katsikas, S. (2006). Effective identification of source code authors using byte-level information. In: *The Proceedings of the 28th international conference on Software engineering*, Shanghai, China, 893-896
- Fukunaga, K. (1990). *Introduction To Statistical Pattern Recognition*. 2nd ed. New York, USA: Academic Press.
- Grant, T. D. (2007) Quantifying evidence for forensic authorship analysis. *Int J Speech Lang Law*, 14, 1–25.
- Grieve, J. (2007). Quantitative authorship attribution: An evaluation of techniques. *J Lit Ling Comput*, 22, 251-270.
- Halteren, H. V., Baayen, H., Tweedie, F., Haverkort, M., & Neijt, A. (2005). New machine learning methods demonstrate the existence of a human stylome. *J Quant Linguist*, 12, 65-77
- Hirst, G., & Feiguina, O. (2007). Bigrams of syntactic labels for authorship discrimination of short texts. *J Lit Ling Comput*, 22, 405-417
- Houvardas, J., & Stamatatos, E. (2006). N-gram feature selection for authorship identification. In: *12th International Conference on Artificial Intelligence: Methodology, Systems, Applications*; Berlin, Germany: Springer, 77–86.
- Holmes, D. I. (1994). Authorship attribution. *J Comput Humanities*, 28, 87-106.
- Holmes, D. I. (1998). The evolution of stylometry in humanities scholarship. *J Lit Ling Comput*, 13, 111-117.
- Holmes, D. I. (1992). A stylometric analysis of Mormon Scripture and related texts. *J Roy Stat Soc A Sta*, 155, 91-120.
- Honore, A. (1979). Some simple measures of richness of vocabulary. *Association for Lit Ling Comput Bulletin*, 7, 172–177.
- Hoover, D. L. (2004). Delta prime? *J Lit Ling Comput*, 19, 477-495.
- Jafari, A. (2008). The study of Persian adjuncts based on functional and formal approaches. *Dastur; Especial Issue of Name-ye Farhangestan*, 5, 128-156.
- Li, J., Zheng, R., & Chen, H. (2006). From fingerprint to writeprint. *J Commun ACM*, 49, 76-82.
- Lowe, D., & Matthews, R. (1995). Shakespeare vs. Fletcher: a stylometric analysis by radial basis functions. *J Comput Humanities*, 29, 449-461.
- Luyckx, K., & Daelemans, W. (2011). The effect of author set size and data size in authorship attribution. *J Lit Ling Comput*, 26, 35-55.
- Matthews, R. A. J., & Merriam, T. V. N. (1993). Neural computation in stylometry I: An application to the works of Shakespeare and Fletcher. *J Lit Ling Comput*, 8, 203-209.
- McCarthy, P. M., Lewis, G. A., Dufty, D. F., & Mcnamara, D. S. (2006). Analyzing writing styles with coh-metrix. In: *FLAIRS 2006 the Florida Artificial Intelligence Research Society International Conference*, California, USA, 764–769.

- Mendenhall, T. C. (1887). The characteristic curves of composition. *Science*, 9, 237-246.
- Merriam, T. V. N., & Matthews, R. A. J. (1994). Neural computation in stylometry II: An application to the works of Shakespeare and Marlowe. *J Lit Ling Comput*, 9, 1-6.
- Mosteller, F., & Wallace, D. (1964). *Inference and Disputed Authorship: The Federalist*. Massachusetts, USA: Addison-Wesley.
- Peng, F., Schuurmans, D., Keselj, V., & Wang, S. (2003). Language independent authorship attribution using character level language models. In: the Proceedings of the 10th conference on European chapter of the Association for Computational Linguistics; Budapest, Hungary, 267-274.
- Rosenblatt, F. (1958). The Perceptron: A probabilistic model for information storage and organization in the brain. *J Psychol Rev*, 65, 386-408.
- Rudman, J. (1997). The State of Authorship Attribution Studies: Some Problems and Solutions. *J Comput Humanities*, 31, 351-365.
- Sichel, H. S. (1975). On a distribution law for word frequencies. *J Am Stat Assoc*, 70, 542-547.
- Simpson, E. H. (1949). Measurement of diversity. London, UK: *Nature*, 1949.
- Stamatatos, E. (2008). Author identification: Using text sampling to handle the class imbalance problem. *J Inform Process Manag*, 44, 790-799.
- Stamatatos, E. (2009). A survey of modern authorship attribution methods. *J AM SOC INF SCI TEC*, 60, 538-556.
- Stamatatos, E., Fakotakis, N., & Kokkinakis, G. (2001). Computer-based authorship attribution without lexical measures. *J Comput Humanities*, 35, 193-214.
- Stamatatos, E., Fakotakis, N., & Kokkinakis, G. (2000). Automatic text categorization in terms of genre and author. *J Comput Linguist*, 26, 471-495.
- Tabatabaee, A. (2009). *Word Formation and Grammatical Category: The Identification of Grammatical Categories of Persian Words, Based on Morphological Characteristics*. Tehran, Iran, Culture, Art and Communication Research.
- Takamy, S. M. M. (2005). Data preprocessing and pattern recognition methods. Technical Report KNTU-AADCL-05, Khaje Nasir Toosi University.
- Tambouratzis, G., Markantonatou, S., Hairetakis, N., Vassiliou, M., Carayannis, G., & Tambouratzis, D. (2004). Discriminating the registers and styles in the modern Greek language-part 2: Extending the feature vector to optimize author discrimination. *J Lit Ling Comput*, 19, 221-242.
- Teng, G., Lai, M., Ma, J., & Li, Y. (2004). E-mail authorship mining based on SVM for computer forensics. In: *Proceedings of 2004 International Conference on Machine Learning and Cybernetics*, Shanghai, China, 1204-1207.
- Theodoridis, S., & Koutroumbas, K. (1999). *Pattern Recognition*. New York, USA: Academic Press.
- Tweedie, F., & Baayen, R. (1998). How variable may a constant be? Measures of lexical richness in perspective. *J Comput Humanities*, 32, 323-352.
- Tweedie, F. J., Singh, S., & Holmes, D. I. (1996). Neural network applications in stylometry: The Federalist Papers. *J Comput Humanities*, 30, 1-10.
- Uzuner, Ö., & Katz, B. A. (2005). comparative study of language models for book and author recognition. In: *IJCNLP 2005 the Second International Joint Conference on NLP*, Jeju Island, Korea. Berlin, Germany, Springer, 969-980.

- Yule, G. U. (1938). On sentence-length as a statistical characteristic of style in prose, with application to two cases of disputed authorship. *Biometrika*, 30, 363–390.
- Yule, G. U. (1944). *The statistical study of literary vocabulary*. Cambridge, England: University Press.
- ZandiMoghadam, Z., & Sadeghi, A. (2008). *A dictionary of Persian spelling, based on inscriptions of academy of Persian language and literature*. Tehran, Iran, Academy of Persian Language and Literature.
- Zheng, R., Li, J., Chen, H., & Huang, Z. (2006). A framework for authorship identification of online messages: Writing-style features and classification techniques. *J Am Soc Inf Sci Technol*, 57, 378-393.
- Zipf, G. K. (1932). *Selected studies of the principle of relative frequency in language*. Cambridge, Massachusetts, USA: Harvard University Press.

Table1. The number of possible applicable n -grams used in each data base.

n-gram	Number	possible applicable features	selected features in University database	selected features in Writer database
Uni-gram		72	72	72
Bi-gram		$72 \times 72 = 5184$	375	71
tri-gram		$72 \times 72 \times 72 = 373248$	560	117
forth-gram		$72 \times 72 \times 72 \times 72 = 26873856$	257	97
Total		27252360	1264	357

Table 2. Features extracted from texts.

Features	Sub Features	Number of used Features
Lexical Features	Lexical richness	7
	n-gram	Depends on length of text and number of authors
	Average length or words	1
	average length of sentences' characters	1
	average number of words in sentences	6
	number of all short words equal or shorter than 3 characters	3
	distribution of alphabet characters in texts	72
	distribution of Persian characters in texts	1
	distribution of the number of paragraphs in	1

	texts	
	distribution of words with the length of 1 to 30 in texts	30
	distribution of English characters in texts	1
	distribution of numerical characters in texts	1
	distribution of half space character in texts	1
Syntactic Features	distribution occurrence of noun, adjective and adverbial phrases	4
	Function words	922
	distribution of punctuation characters	1
Semantic Features	distribution of adjuncts in Persian	36
Application-specific features	distribution of tab character	1
	distribution of enter character	1
	distribution of space character	1
	distribution of empty lines	1

Table 3. Comparison of Classification methods on two Data-bases.

Method Data-base	KNN	Delta	Decision tree	Neural Network	LDA
<i>University</i>	70%	50%	21.6%	15%	81.6%
<i>Writers</i>	100%	87%	57.5%	37.3%	100%

Table 4. Comparison of the average precision of classification methods versus features types on both data-bases.

Feature Classification Methods	Lexical and Characteristic	Synta ctic	Seman tic	Applicatio n-specific
KNN	80%	66.5%	49.5%	49.5%
Decision tree	25%	21.88 %	18.75 %	12.5%
Neural network	11%	12.5%	16%	8.75%
LDA	38.12%	29.37 %	30%	39.37%
average precision	38.53%	32.56 %	28.56 %	27.53%

