



تحسين عملية استنتاج العلامات النصية من الوثائق اعتمادا على المعنى الدلالي

رسالة مقدمة للحصول على درجة الماجستير في الحاسبات والمعلومات

إعداد

ايمان اسماعيل سيد

بكالوريوس الحاسبات والمعلومات
معيدة بقسم نظم المعلومات
كلية الحاسبات والمعلومات
جامعة عين شمس

تحت إشراف

أ.د. خالد البهنسي

أستاذ بقسم نظم المعلومات
كلية الحاسبات والمعلومات
جامعة عين شمس

د. محمد العليمي

أستاذ مساعد بقسم نظم المعلومات
كلية الحاسبات والمعلومات
جامعة عين شمس

د. ولاء جاد

أستاذ مساعد بقسم نظم المعلومات
كلية الحاسبات والمعلومات
جامعة عين شمس

ملخص الرسالة

في الآونة الأخيرة، كمية البيانات المتوفرة على شبكة الإنترنت في تزايد سريع، لذلك فمن الصعب البحث عن البيانات ذات الصلة في مجموعة بيانات ضخمة. وضع علامات نصية للمستند يوفر حلا لهذا النوع من المشاكل. حيث انها تعطي النص معلومات إضافية في شكل ملاحظات أو تعليقات. التعليق التوضيحي للمستند يسهل مهمة العثور على الموضوعات الرئيسية للوثيقة. علاوة على ذلك، فإنه يساعد القارئ على القاء نظرة عامة على المستند وفهمه.

بسبب انتشار تطبيقات وسائل الاعلام الاجتماعية مثل الفيسبوك، تويتر ... الخ. يتم إنتاج الملايين من المستندات القصيرة يوميا. لذلك، يتم استخدام تصنيف النص لاكتشاف المعرفة من هذه المستندات النصية غير المنظمة. المستندات النصية القصيرة لها خصائص خاصة حيث انها مليئة بالاطء وتحتوى على كلمات قليلة لذلك كلماتهم نادرا ما تتكرر. وتستند الأساليب التقليدية لتصنيف هذه الأنواع من الوثائق إلى طريقة "BOW" التي تصنف الوثائق النصية ككلمات مستقلة. بحيث يعبر عن كل مستند بمجموعه من الكلمات. كل كلمة لها قيمه مساويه لعدد مرات ظهورها في المستند. فاستخدام الأساليب التقليدية للتصنيف كمجموعه من الكلمات لديها العديد من العيوب: المستندات النصية القصيره ليس لديها عدد كاف من الكلمات لتمثيل المستندات باستخدام Bow حيث أن هذه الطريقة تعتمد على عدد تكرار الكلمات في المستند، والكلمات في المستندات النصية القصيره نادره التكرار. ايضا طريقه BOW لم تأخذ في الاعتبار العلاقة بين الكلمات وبعضها. لذا يستخدم المعنى الدلالي في وضع العلامات النصية للتغلب على المشاكل الموجوده في المستندات النصية القصيره. ولاثراء وتزويد المستندات ببعض المعلومات الاضافيه لفهم المستند اكثر وتصنيفها جيدا.

في هذا العمل ، تم اقتراح نموذجان فعالان لاستنتاج العلامات النصية من المستندات اعتمادا على المعنى الدلالي. النموذج الاول CBER معتمد على اثراء علميه التصنيف للمستندات النصية حيث انه مكون من النموذج المقترح SAWN والذى يعتمد على قاعده البيانات (WordNet) لاستنتاج المعنى الدلالي ، والنموذج الاخر WVTF الذى يعتمد على BOW في تمثيل البيانات النصية لاداء عمليه التصنيف. النموذج SAWN يستخرج الكلمات ذات المعنى من المستندات النصية القصيره باستخدام قاعده البيانات (WordNet) ويعطيها قيمه تعبر عن اهميتها بالاستعانه بمرادفات هذه الكلمات. وبالتالي الكلمات التى لها نفس المعنى تزيد من قيمه مرادفاتها. علاوة على ذلك تم استغلال العلاقة بين الكلمات وبعضها في حل المشاكل المتعلقة بالمعنى الدلالي مثل تواجد كلمه لها اكثر من معنى او العكس بان توجد مجموعه من الكلمات تعبر عن نفس المعنى.

النموذج المقترح CBER يعمل على اثراء المستندات النصية القصيره بمعلومات اضافيه دون الحاجة لزياده عدد الكلمات المعبره عن كل مستند. كما انه يأخذ في الاعتبار الكلمات التى لا توجد في قاموس ال (WordNet) فانه يعطى للكلمات قيمه للتعبير عن اهميتها اعتمادا على قيمتها الاساسيه التى تعتمد على عدد مرات تكرار الكلمه و قيمه اخرى اعتمادا على المعنى الدلالي ، لذا تأخذ الكلمات ذات المعنى والمعرفه في القاموس قيمه اكبر عن غيرها من الكلمات الغير معرفه .

النموذج ١ الثانى "Wiki_Spots" يستخدم للتعبير عن المعنى الدلالى باستخدام الويكيبيدا عن طريق ربط المستندات النصيه القصيره بالويكيبيدا لاستخراج الكلمات التى تستخدم للاشاره عن مقالات توضيحيه للمعنى الدلالى لهذه الكلمات ، والتعبير عن المستندات بطريقه جديده باستخدام هذه الكلمات التى تعبر عن الموضوعات المتعلقه بالمستندات ، ومن ثم يتم اعطاء هذه الكلمات (Spots) قيمه اعتمادا على اهميتها فى المستند للتصنيف الجيد. علاوه على ذلك ، فان كل spot يمكن ان يعبر عن كلمه واحده او جمله اسميه قصيره مما يحسن من نتائج عمليه التصنيف.

واظهرت التجارب العمليه ان نموذج CBER لاثراء عمليه التصنيف جيد فى تصنيف المستندات النصيره القصيره الذى وصل الى ٠,٩٣٣ ، ٠,٩٣٦ ، ٠,٩٣٤ و ٩٤% فى الدقه والتذكير ، ونموذج SAWN وصل إلى ٠,٨٧٢ ، ٠,٨٧٣ ، ٠,٨٧٧ و ٨٨%. أيضا، نموذج ويكيبيديا القائم على "Wiki_Spots" وصل إلى ٠,٨٣٢ ، ٠,٨٣٣ ، ٠,٨٣٧ و ٨٤% بالمقارنة مع طريقة BOW التقليدية.



Information Systems Department
Faculty of Computer & Information Sciences
Ain Shams University

An Enhanced Automatic Model Based On Semantic Annotation For Text Documents

Thesis submitted as a partial fulfillment of the requirements for the degree of Master of science
in Computer and Information Sciences.

By
Eman Ismail Sayed

B.Sc. in Computer and Information Sciences (2012),
Demonstrator at Information Systems Department,
Faculty of Computer and Information Sciences,
Ain Shams University.

Under Supervision of
Prof. Dr. Khaled El-Bahnasy

Professor
Information Systems Department,
Faculty of Computer and Information Sciences,
Ain Shams University, Egypt

Dr. Mohamed Eleliemy
Associate Professor
Information Systems Department,
Faculty of Computer and Information Sciences,
Ain Shams University, Egypt

Dr. Walaa Gad
Associate Professor
Information Systems Department,
Faculty of Computer and Information Sciences,
Ain Shams University, Egypt

Cairo - 2017

Abstract

The amount of available data on the web is rapidly increasing, so it is difficult and more complicated to search and identify relevant data in huge data sets. Text document annotation provides one approach to ease such types of problems. It provides the text with additional information in the form of notes or comments. Annotation facilitates the task of finding the main topics of a document. Moreover, annotation helps the reader to overview and understand the document.

Due to the spread of social media applications such as Facebook, Twitter.. etc, millions of short texts are being produced daily. Text classification is used to discover knowledge from these unstructured text data. The short text documents (STDs) have special characteristics as being noisy and sparsity because their words are rarely repeated. The traditional methods of classifying such types of documents are based on the Bag of Words (BOW) method, which indexes text documents as independent features. Each feature is a single term or word in a document. A document is represented as a vector in feature space. A document vector contains the word weights, which are the number of word occurrences in the document. Classification of STDs based on BOW has many drawbacks: STDs do not provide enough co-occurrence of words or shared context. Representation of such documents is almost sparse because of empty weights

when using BOW. Hence, the traditional bag of words (BOW) method fails to achieve high accuracy. Moreover, the BOW method treats synonymous words as different features and does not consider the relations between words and documents. Therefore, semantic knowledge was introduced as a background to focus on the semantic relationships between the documents words or terms.

In this thesis, two effective models for semantic annotation are proposed. The first model is Classification Based on Enrichment Representation (CBER). It is composed of the proposed semantic analysis based on WordNet(SAWN) model and the word vector term frequency (WVTF). WVTF is a BOW representation of text documents. SAWN maps the text documents with WordNet to extract the concepts. Concepts are the terms that are defined in WordNet. SAWN chooses the most suitable synonyms for s document concepts by studying and understanding the surrounding concepts in the same document. Thus, concepts with the same meaning will increase the weight of their synonyms. Furthermore, the semantic relationships between concepts have been exploited in order to solve the disambiguation problems such as polysemy and synonyms.

The CBER model enriches the STDs with semantic weights to solve disambiguation problems without increase the document features. It considers all documents terms. However, some terms may not be defined in WordNet. The terms are then provided with a new weight. These weights depend on the term frequency weight based on WVTF and the semantic weight based on SAWN. Thus, the terms that are defined in WordNet gain more weight.

The second model is the Wiki_Spots model. It identifies and extracts cross-referencing text (spot) from the documents using the Wikipedia knowledge base. Each spot is annotated to a Wikipedia article (page) considering the relationships with other spots in the same document. Wiki_Spots model exploits these spots to represent the documents as vectors of topics rather than vector of words in the traditional BOW method. Moreover, each spot is unigram or a noun phrase that helps increasing the classification accuracy.

The experimental results showed that the proposed CBER model is valuable in annotating short text documents to their best labels. CBER showed significant performance 0.933, 0.936, 0.934 and 0.94 in precision, recall, F-measure and accuracy respectively, and SAWN reached 0.872, 0.873, 0.877 and 0.88. Also, the Wikipedia based model "Wiki_Spots" reached 0.832, 0.833, 0.837 and 0.84 compared to the traditional BOW method.

Acknowledgements

Thanks first and foremost to Allah, who gives me the knowledge and patience to produce this work. Thanks Allah to response to my prayers to achieve my goal.

Firstly, I would like to express my sincere gratitude to my supervisor Prof. Khaled Bahnsy for the continuous support of my M.sc study and related research. For his patience, motivation, and immense knowledge. His guidance helped me in all the time of research and writing of this thesis. I could not have imagined having a better advisor and mentor for my M.Sc study.

Also, I would like to extend my thanks and appreciation to my supervisor Dr. Mohamed Hamdy for his guidance, constant support and useful feedback that helped me finalizing my work and my M.Sc thesis.

I also appreciate Dr. Walaa Gad for her effort with me beginning with preparing the proposed master idea till writing and reviewing the thesis, she always followed up my work weekly and her scientific contributions helped me to advance my work and to get better results.

I would to thank all my family especially my mother who helps me in everything. Also, I am lucky to be married to Ali Gamal, who helped me to finalize writing and reviewing my thesis.

List of publications

- Eman Ismail and Walaa Gad. CBER: An Effective Classification Approach Based on Enrichment Representation for Short Text Documents. *Journal of Intelligent Systems*. Vol.26, pp. 233-241. 2016.
- Eman Ismail, Walaa Gad, Mohamed Hamdy, and Khaled Bahnsy. Text document annotation methods. In 2015 IEEE Seventh International Conference on Intelligent Computing and Information Systems (ICICIS), pp. 634-640. IEEE, 2015.

Contents

| | |
|--|-------------|
| Abstract | i |
| Acknowledgements | iv |
| List of Figures | viii |
| List of Tables | x |
| Abbreviations | xi |
| 1 Introduction | 1 |
| 1.1 Overview | 1 |
| 1.2 Research Objectives | 3 |
| 1.3 Main Contributions of this Thesis | 3 |
| 1.4 Thesis Organization | 4 |
| 2 Related Work | 6 |
| 2.1 Introduction | 6 |
| 2.2 Keyword Based Annotations (KWBA) | 10 |
| 2.3 Ontology Based Text Annotation (OBA) | 12 |
| 2.3.1 Wikipedia based approach (WBA) | 19 |
| 2.3.2 Keyword-Semantic based annotation (KW-SA) | 28 |
| 2.4 Enrichment Based Classification | 37 |
| 2.5 Reduction Based Classification | 49 |
| 2.6 Summary | 52 |
| 3 Classification Based Enrichment Representation (CBER) | 53 |
| 3.1 Introduction | 53 |

| | | |
|----------|---|-----------|
| 3.2 | CBER Model | 55 |
| 3.2.1 | Document Preprocessing | 57 |
| 3.2.2 | Word Vector Term Frequency (WVTF) | 62 |
| 3.2.3 | Semantic Analysis Based On WordNet (SAWN) | 62 |
| 3.2.4 | Hybrid CBER Model | 65 |
| 4 | Annotating Short Text Documents Using Wiki_Spots | 67 |
| 4.1 | Introduction | 67 |
| 4.2 | Wiki_Spots Model | 69 |
| 5 | Experimental Results | 75 |
| 5.1 | Dataset description | 75 |
| 5.2 | Classification Based Enriching Representation (CBER) Model Results | 76 |
| 5.3 | Wiki_Spots Model Results | 81 |
| 6 | Conclusion & Future work | 85 |

List of Figures

| | | |
|------|---|----|
| 2.1 | Content Analysis Toolkit Process | 11 |
| 2.2 | Synset Subclasses Hierarchy | 13 |
| 2.3 | WordSense Subclasses Hierarchy | 14 |
| 2.4 | RDF graph data model | 15 |
| 2.5 | WordNet Classes Resources Statement | 17 |
| 2.6 | WSD: Mapping between keywords and Word Net concepts [25]. | 19 |
| 2.7 | The flow of processes to annotate text documents semantically | 20 |
| 2.8 | Topic detection system framework | 21 |
| 2.9 | A snapshot from WikiCFP showing the topic list in the conference CFP. | 22 |
| 2.10 | Term identification process | 23 |
| 2.11 | Malaysia tourism domain taxonomy | 25 |
| 2.12 | People topic content from Wikipedia as in Ref.[1] | 26 |
| 2.13 | Domain Taxonomy | 26 |
| 2.14 | Suggestion algorithm for annotation | 30 |
| 2.15 | GoNTogle architecture | 32 |
| 2.16 | Integration of three techniques on text document annotation | 35 |
| 2.17 | Short text classifier framework | 40 |
| 2.18 | Annotation using Wikipedia entities proposed method processes | 42 |
| 2.19 | Clustering System Approach | 45 |
| 2.20 | System framework for calculating distances between the tweets using Wikipedia knowledge base. | 47 |
| 2.21 | Distance between Wikipedia pages | 48 |
| 2.22 | Associated Wikipedia pages about one tweet | 49 |
| 2.23 | Short text classification using few words Framework | 51 |

| | | |
|-----|---|----|
| 3.1 | Classification Based Enrichment Representation (CBER) Model | 56 |
| 3.2 | Porter Stemmer Rules Steps | 58 |
| 4.1 | <i>Wiki_Sports</i> System over flow | 70 |
| 5.1 | Cross-folding Validation on the Training Dataset. | 78 |
| 5.2 | Results of the CBER Model. | 79 |
| 5.3 | Accuracy of the CBER Model in Comparison with the Baseline Model [2]. | 80 |
| 5.4 | Evaluation of the CBER Model in Terms of Precision. | 80 |
| 5.5 | Evaluation of the CBER Model in Terms of Recall. | 81 |
| 5.6 | Evaluation of the CBER Model in Terms of F-Measure. | 81 |
| 5.7 | Cross-folding Validation on the Training Data Set. | 82 |
| 5.8 | Results of the <i>Wiki_Sports</i> Model. | 84 |

List of Tables

| | | |
|-----|---|----|
| 2.1 | Properties of WordNet Schema | 16 |
| 2.2 | Summary of techniques in text document annotation . . . | 36 |
| 3.1 | Step 1a Rules | 59 |
| 3.2 | Step 1b Rules | 59 |
| 3.3 | Step 1b1 Rules | 59 |
| 3.4 | Step 1c Rules | 59 |
| 3.5 | Step 2 Rules | 60 |
| 3.6 | Step 3 Rules | 60 |
| 3.7 | Step 4 Rules | 61 |
| 3.8 | Step 5a Rules | 61 |
| 3.9 | Step 5b Rules | 61 |
| 5.1 | Snippets Dataset Classes. | 75 |
| 5.2 | Wiki_Spots Results | 84 |

Abbreviations

| | |
|--------------|---|
| BOC | Bag Of Concepts |
| BOW | Bag Of Words |
| CAT | Content Analysis Toolkit |
| CFP | Call For Papers |
| CBER | Classification Based Enrichment Representation |
| DB | Data Base |
| DW | Descriptive Words |
| IR | Information Retrieval |
| KW_SA | Keyword Semantic Based Annotation |
| KW_OA | Keyword Ontolgy Based Annotation |
| KNN | K Nearest Neighbour |
| KWBA | Keyword Based Annotations |
| LDA | Latent Dirichlet Allocation |
| LSA | Latent Semantic Indexing |
| LDA | Latent Dirichlet Allocation |
| LSA | Latent Semantic Indexing |
| NLP | Natural Language Processing |
| NGD | Normalized Google Distance |
| NBC | Naive Bayes Classifier |
| OBA | Ontology Based Annotation |
| OWL | Online Writing Lap |
| POS | Part Of Speech |