

“The Naming of Cats”: Automated Genre Classification

Yunhyong Kim and Seamus Ross

Digital Curation Centre (DCC)
&
Humanities Advanced Technology Information Institute (HATII)
University of Glasgow
Glasgow, UK

“The Naming of Cats is a difficult matter, It isn’t just one of your holiday games; You may think at first I’m as mad as a hatter, When I tell you, a cat must have three different names.” - T.S. Eliot, The Naming of Cats

Abstract. This paper builds on the work presented at the ECDL 2006 ([29]) in automated genre classification as a step toward automating metadata extraction from digital documents for ingest into digital repositories such as those run by archives, libraries and eprint services. We divide features of the documents into five types: features for visual layout, linguistically modeled syntactic features, stylo-metric features, features for semantic structure, and contextual features as an object linked to previously classified objects and other external sources. Results concerning the first two types have been described elsewhere([29]). The current paper discusses results from testing classifiers based on image and stylo-metric features and shows that genres for which image features fail to cluster are the genres for which stylo-metric features cluster very well.

1 Background and Objective

In [29], we summarised the valuable role of automated metadata extraction in the cost-effective efficient management of digital collections: metadata play a key role in management processes ([43], [23]) and the manual creation of metadata is expensive ([15], [23], [40]). As we pointed out in [29], ERPANET’s ([18]) Packaged Object Ingest Project ([19]) identified automatic extraction tools for technical metadata (e.g. [33], [35]), and substantial work on descriptive metadata extraction within specific domains has been conducted (e.g. [32], [13], [2], [50], [21], [22], [6], [26], [47], [51]) along with other work in information extraction from text (e.g. [3], [9], [49], [48]). However, a general tool has yet to be developed to extract metadata from digital objects of varied types and genres. This paper further develops concepts of genre classification introduced in [29] involving the automatic grouping of documents into distinctive document types followed by focused metadata extraction from single document types as a means of creating a tool capable of extracting metadata across many domains at different semantic levels. To reiterate the argument in [29], identifying the genre

first provides a mechanism to limit the scope of document forms from which to extract other metadata. Within a single genre, metadata such as author, title, keywords, identification numbers or references can be expected to appear in a specific style and region, and independent methods have been developed for genre-specific extraction of such metadata for some classes of documents (e.g. Scientific Papers). Note also that different institutions focus on collecting and managing digital materials in different genres; genre classification will support automating the identification, selection, and acquisition of materials in keeping with local collecting policies.

A review of Biber ([7]), Karlgren et al. ([25]), Kessler et al. ([27]), Rauber et al. ([44]), Bagdanov et al. ([4]), Boese ([8]), Finn et al. ([20]) and Santini ([45]) exemplifies the lack of consensus on the definition of genre. Biber's analysis of document genres employed five functional dimensions (information, narration, elaboration, persuasion, abstraction) to characterise text, while Karlgren et al. and Boese concentrated on more popularly accepted genre classes such as FAQ, Job Description, Editorial or Reportage. Kessler et al. tried to address both types, while Finn et al. studied binary classifications (fact versus opinion, positive versus negative reviews). Santini discussed general genre facets, while Bagdanov limited his task to detecting specific journals and brochures. Others ([44], [5]) attempted the clustering of documents rather than classification. An overview of the various efforts in genre analysis can be found in a technical report by Santini ([46]). A broader review of metadata extraction and genre classification is also being prepared by the DELOS NoE Digital Preservation Cluster and is expected to be completed before the publication of this paper.

The variety of definitions adopted by these researchers illustrates a confused interplay of two notions: one of structure and one of function. Structure is defined by the visual layout and is expected to be distinguishable mostly by measurable features such as amount of white space; the length of the document, sentences, or words; and, the presence or absence and location of headers, delimiters, images, or links. Function, on the other hand, is defined by the intended role of the document and is expected to be characterised mostly by linguistic models and semantic analyses of the documents. The two notions are closely linked together by medium, process or event. For example, a scientific research article is usually *structured* so that a title is present on the first page followed by author, affiliation, a body of text consisting of sections, and finally a list of references. It has the *function* of communicating, arguing or describing research. The interrelationship of structure and function are represented by the formatting requirements of journals or conventions in the community or event for which the document was created. The requirements and conventions evolve to optimise the communicative intentionality within the context; other communities or events may find different structures of documents to optimise the same function. Just as biologists study DNA as the building blocks of living organisms to understand the classes into which they have evolved within their environment, it seems important to identify documents by their structure and their function separately as building blocks to infer their genre class within a standardised schema. We seek to be able to

achieve this by a full analysis of five document feature types: image features, syntactic features, stylistic features, semantic structure, and domain knowledge features. We aim to build a system which models the five feature sets for a schema of approximately seventy genres (Table 1).

The genres in Table 1 are not meant to be static: the schema has been evolving as we develop and incorporate well-structured classification standards and as we become aware of digital genres we had not encountered before or which have just emerged in the digital domain. The experiments in this paper have initially limited the study to the image and stylistic feature sets on the nineteen most prolific genres in our experimental data set. Along with the results in [29], the results here are intended to be another step towards a full analysis.

The experimental data in this paper is from the pool of 570 PDF ([37]) files that were sampled randomly from the Internet as described in [29]. As explained in [29], by confining the work to studying PDF files, we hope to put a boundary on the problem space, while working with a widely used portable format for digital objects ingested into digital repositories.

Table 1. Scope of genres

Groups	Genres
Book	Academic book, Fiction, Poetry, Handbook, Other book
Article	Abstract, Scientific research article, Other research article, Magazine article, News report
Short Composition	Fictional Piece, Poems, Dramatic Script, Essay, Short Biographical Sketch, Review
Serial	Periodicals (Newspaper, Magazine), Journals, Conference Proceedings, Newsletter
Correspondence	Email, Letter, Memo, Telegram
Treatise	Thesis, Business/Operational report, Technical report, Misc report
Information Structure	List, Catalogue, Raw Data, Table Calendar, Menu, Form, Programme, Questionnaire, FAQ
Evidential Document	Minutes, Legal proceedings, Financial Record, Receipt, Slips, Contract
Visually Dominant Document	Artwork, Card, Chart, Graph, Diagram, Sheet Music, Poster, Comics
Other Functional Document	Guideline, Regulations, Manual, Grant/Project Proposal, Legal Appeal/Proposal/Order, Job/Course/Project Description, Product/Application Description, Advertisement, Announcement, Appeal/Propaganda, Exam/Worksheet, Fact sheet, Forum Discussion, Interview, Notice, Resume/CV, Slides, Speech transcript

This paper, along with [29] and [30], is intended to show the promise of combining separate classifiers trained on different types of features for genre

classification. Also note that the bottom-up approach of starting from genre-specific extraction may result in several tools, which are overly dependent on the structures of the documents in the domain, with no obvious means of interoperability: the top-down approach of creating a tool which looks across genres, to be refined further within the domain, will enable us to avoid this problem.

2 Classifiers

The experiments described in this paper involve the use of two classifiers:

Image classifier: this classifier depends on features extracted from the PDF document when handled as an image. It uses the module *pdftoppm* from XPDF ([36]) to extract the first page of the document as an image. The resulting image is divided into a sixty-six by sixty-six grid¹. Then Python's Image Library (PIL) ([41], [39]) is employed to extract pixel values in each region. Each region is given a value of 0 or 1 depending on the amount of non-white pixel values it contains. The result is modeled using Naïve Bayes available with the Weka ([52]) machine learning toolkit.

Stylo-metric classifier: this classifier looks at the frequency of selected words, number of font changes, the difference between the largest font size and smallest font size, length of the document, average length of words, and number of words in the front page of the document. The font information was extracted on the level of words using a modified version of *pdftohtml* ([38]), developed by Volker Heydegger at the University of Cologne. The modified version converts a PDF document to a XML file with all the font information for each word in the document. A word list was automatically constructed containing all words which appear in more than half of the files in any one genre. For each file, the frequency of each word was recorded as a vector then augmented by length and font information. The result was modeled using Naïve Bayes in the Weka ([52]) machine learning toolkit.

This paper expresses the view that the image along with the stylistic features will capture the structural elements of genres while the language model combined with the stylistic and semantic features will help to separate documents of distinct functional categories. Involving the image of a document in the classification also enables the management of documents without violating security, maximises the viability of a language independent tool, frees the process from being solely dependent on text processing tools with encoding requirements and problems relating to special characters², and makes the method applicable to paper documents digitally imaged (i.e. scanned).

¹ The choice of the dimension reflects the fact that it seemed to produce the best results at the time but further analysis may be necessary.

² *pdftohtml* failed to extract information from seventeen percent of the documents. The image processing did not fail on any documents.

3 Experiments

There are three main experiments described in this paper:

Clustering experiment: this experiment compared the cluster resolution for two sets of features: the image features and the stylo-metric features. We grouped the data in nineteen genres into two clusters using the Weka Machine Learning Toolkit's ([52]) Estimation-Maximisation algorithm. The purpose was to see how well the files in each genre group into one cluster. The result is expressed in terms of the percentage of files within each genre which have been grouped into one cluster.

Periodicals versus Thesis: in this experiment, we took documents in the genres Periodicals and Thesis. We used the image classifier to classify the documents by using 10 fold cross validation.

Periodicals versus Non-periodicals: we expanded on the experiment above to group four more genres with the genre Thesis as one group labelled Non-periodicals. The four additional classes are Business or Operational Report, Minutes, Fictional Book, and Academic Book. The four extra classes were chosen from the genres that were grouped in the same image cluster as Thesis.

4 Results

Table 2 shows the results of the clustering experiment. The key finding in this experiment is that the genres for which image features fail to cluster are the genres for which stylo-metric features cluster very well. For instance, note that Scientific Research Articles divide half and half into each cluster with no preference when using the visual features while ninety two percent group into one cluster when using stylo-metric features. The opposite is true of Periodicals.

The results described in Tables 3, 4 and 5 use three standard indices in classification tasks: accuracy, precision and recall. Let N be the total number of documents in the data, N_c the number of documents in the data set which are in class C , T the total number of correctly labelled documents in the data set independent of the class, T_c the number of true positives for class C , and F_c the number of false positives for class C . Accuracy is defined to be $A = \frac{T}{N}$; precision and recall for class C is defined to be $P_c = \frac{T_c}{(T_c+F_c)}$ and $R_c = \frac{T_c}{N_c}$, respectively. Table 3 gives the result when the data set was confined to Periodicals and Theses. The accuracy was surprisingly high. To check if the results actually reflect the distinctiveness of image features in periodicals, the experiment was repeated with four more classes of non-periodical documents added to Thesis to form a group of Non-periodicals (results in Tables 4 and 5). A slight decrease in performance is visible in Table 4 (cf. Table 3), but the accuracy is still quite high. On the other hand the results for the stylistic classifier in Table 5 show that stylistic features do not fare as well in distinguishing Periodicals. For a proper evaluation of the performance, a significance test is in order (pending), but a difference of 17.6% in overall accuracy can not be ignored by the strictest of observers, and a decrease in precision on Periodicals from 73.9% to 47.8% (a difference of 26.1%)

Table 2. A Comparison of Visual and Stylo-metric Clusters (percentage of files in one cluster)

Groups	Genres	Visual	Stylistic
Book	Academic Book	87.5	60
	Fiction	87.5	83.3
	Other Book	75	82.4
Article	Scientific research article	50	92
	Other research article	90	73.7
	Magazine article	62.5	84.6
Serial	Periodicals (Newspaper, Magazine)	94.7	62.5
	Newsletter	74.1	83.3
Treatise	Thesis	100	90
	Business/Operational report	66.7	90.9
	Technical report	68.2	72.2
Information Structure	List	68.4	85.7
	Form	68.8	69.2
Evidential Document	Minutes	94.7	76.9
Other Functional Document	Instruction/Guideline	90.5	50
	Job/Course/Project Description	50	66.7
	Product/Application Description	61.1	68.8
	Fact sheet	53.3	78.6
	Slides	60	91.7

inspires the belief that the visual features are better equipped to distinguish periodicals from the other five genres.

Table 3. Distinguishing Periodicals from Thesis using image features

10 fold Cross Validation with the Image classifier, Overall accuracy: 97.26 %		
Genres	Precision(%)	Recall(%)
Periodicals (19 items)	100	94.7
Thesis (18 items)	94.7	100

5 Conclusion and Further Research

The results in [29] and the results in this paper indicate the promise of using a multi-layered decision tree on many different sets of features to classify genres. The results in Table 2 show definite divisions between genres which have strong image features and genres that have strong stylistic features. The results in Tables 2 and 3, indicate that Periodicals have more clearly distinguishing image features than stylo-metric features, while Table 4 suggests that Thesis shares its image features with four other genres. Previous reports ([29], [30]) indicated that

Table 4. Distinguishing Periodicals from Non-periodicals using image features

10 fold Cross Validation with the Image classifier, Overall accuracy: 92.23 %		
Genres	Precision(%)	Recall(%)
Periodicals (19 items)	73.9	89.5
Non-Periodicals (84 items)	97.5	92.9

Table 5. Distinguishing Periodicals from Non-periodicals using stylistic features

10 fold Cross Validation with the Image classifier, Overall accuracy: 74.63 %		
Genres	Precision(%)	Recall(%)
Periodicals (19 items)	47.8	68.8
Non-Periodicals (84 items)	88.6	76.5

Scientific Research Articles are not easily distinguishable by image features from Product Descriptions but better distinguishable when using syntactic features. If all these features are processed in one classifier, the statistical model can be misled by non-distinguishing features. If we were to train on sufficient data, this is not a problem; the non-distinguishing features will be filtered out as noise. It is, however, very difficult to have *sufficient data* when constructing a tool which is intended to have dynamic and domain-independent properties. In [28] and [31], the CANDC part-of-speech tagger ([10]), reputed to have performed well elsewhere, was employed to tag words in an Astronomy research articles. In Astronomy there is frequent usage of the term *He* to refer to the chemical element Helium. The tagger, which was trained on the *Wall Street Journal* articles, tagged *He* to be a pronoun for all instances, propagating further errors on subsequent words. Separating features into smaller groups will minimise the impact of such artefacts, by trying to exclude the noise from the start, making the most of the differing feature strengths for each genre type. The key seems to lie in identifying which genres belong to which type and to combine the classifiers in a reasonable way to build a general classifier.

Further improvement can also be envisioned by integrating more classifiers into the decision process. In [29] we suggested the following classifiers:

- **Extended image classifier** which looks at more than the first page of the document: we could process the image of pages other than the first page of the document or several pages of the document in parallel. This would however involve several decisions: the optimal number of pages to be used and the best way to combine the information from different pages need to be determined (e.g. will several pages be considered to be one image; if not, how will the classification of synchronised pages be statistically combined to give a global classification).
- **Language model classifier on the level of POS and phrases** built on the part-of-speech tags (tags which denote whether a word is a verb, noun or preposition) of the underlying words and also on partial chunk tags (tags indicating noun phrases, verb phrases or prepositional phrases).

- **Semantic classifier** modeling subjective or objective noun phrases (e.g. using [42]) and latent semantic analysis may be necessary for finer distinctions in document
- **Contextual Classifier** built on source information of the document such as the name of the journal or address of the web page, and anchor texts or subject or domain information.

There are two obvious ways of gauging the performance of a genre classifier: comparing against human performance and measuring the stability of the performance as you transfer it across domains. We are undertaking an experiment to examine human performance. A significant amount of disagreement is expected in labelling genre even between human labellers; we intend to cross check the labelled data in two ways:

1. **Document Retrieval Exercise (DRE):** We plan to employ a cohort of postgraduates in information science who will be assigned genres from Table 1. They will retrieve one hundred PDF documents for each of the genres they have been assigned, and give a brief description of the source of the document and the reasons for including the document in their collection.
2. **Re-labelling Experiment:** We will anonymise the file names of the documents collected in the DRE and randomise the document sequence. This corpus will be presented to two new groups of labellers drawn from different backgrounds for re-classifying. They will not have access to the initial genre classification information.

The first experiment will create a pool of PDF files which have already been classified into genres by established organisations and users; this will serve as a reference point, and help us to index the performance on well-designed classification standards. The re-labelling experiment will enable us to compare the disagreement of the three classes of labellers over the same data set: this will help to determine the maximum level of accuracy at which the automated system can be expected to perform and determine which genres are better defined by looking at percentage of files in agreement within each genre.

The longer term aim, once a genre classifier with performance comparable to an average human labeller has been developed, will be to integrate the method with other tools which extract author, title, date, identifier, keywords, language, summarisations and other compositional properties of files within a single genre, and combine the tool with ingest models developed elsewhere.

Acknowledgements: This research is a part of The Digital Curation Centre’s (DCC) ([12]) research programme. The DCC is supported by a grant from the United Kingdom’s Joint Information Systems Committee (JISC) ([24]) and the e-Science Core Programme of the Engineering and Physical Sciences Research Council (EPSRC) ([17], grant GR/T07374/01) provides the support for the research programme. Additional support comes from the DELOS: Network of Excellence on Digital Libraries (G038-507618) ([14]) funded under the European Commission’s IST 6th Framework Programme. We would also like to thank

Volker Heydegger at the Historisch-Kulturwissenschaftliche Informationsverarbeitung (HKI), University of Cologne ([11]), for his programming expertise. HKI at the University of Cologne is a participant in the DELOS Digital Preservation Cluster led by the University of Glasgow.

Note on website citations: All citations of websites were validated on 28 September 2006.

References

1. Aiello, M., Monz, C., Todoran, L., Worring, M.: Document Understanding for a Broad Class of Documents. *Intl. Journal Document Analysis and Recognition* 5(1) (2002) 1–16.
2. Automatic Metadata Generation: <http://www.cs.kuleuven.ac.be/hmdb/amg>
3. Arens, A., Blaesus, K. H.: Domain oriented information extraction from the Internet. *SPIE Document Recognition and Retrieval Vol 5010* (2003) 286.
4. Bagdanov, A. D., Worring, M.: Fine-Grained Document Genre Classification Using First Order Random Graphs. *Intl. Conf. Document Analysis and Recognition 2001* (2001) 79.
5. Barbu, E., Heroux, P., Adam, S., Trupin, E.: Clustering Document Images Using a Bag of Symbols Representation. *International Conference on Document Analysis and Recognition*, (2005) 1216–1220.
6. Bekkerman, R., McCallum, A., Huang, G.: Automatic Categorization of Email into Folders. Benchmark Experiments on Enron and SRI Corpora', CIIR Tech. Report, IR-418 (2004).
7. Biber, D.: *Dimensions of Register Variation: a Cross-Linguistic Comparison*. Cambridge University (1995).
8. Boese, E. S.: Stereotyping the web: genre classification of web documents. Master's thesis, Colorado State University (2005).
9. Breuel, T. M.: An Algorithm for Finding Maximal Whitespace Rectangles at Arbitrary Orientations for Document Layout Analysis. *7th Intl. Conf. Document Analysis and Recognition (ICDAR)*, 66–70 (2003).
10. Curran, J., Clark, S.: Investigating GIS and Smoothing for Maximum Entropy Taggers. *Proceedings, Annual Meeting, European Chapter of the Assoc. of Computational Linguistics* (2003) 91–98.
11. Historisch-Kulturwissenschaftliche Informationsverarbeitung (HKI), University of Koeln: <http://www.hki.uni-koeln.de/>
12. Digital Curation Centre: <http://www.dcc.ac.uk>
13. DC-dot, UKOLN Dublin Core metadata editor: <http://www.ukoln.ac.uk/metadata/dcdot/>
14. DELOS Network of Excellence on Digital Libraries: <http://www.delos.info/>
15. DELOS/NSF Working Groups: Reference Models for Digital Libraries: Actors and Roles (2003) <http://delos-noe.iei.pi.cnr.it/activities/internationalforum/Joint-WGs/actors-Roles.pdf>
16. Dublin Core Initiative: <http://dublincore.org/tools/#automaticextraction>
17. Engineering and Physical Sciences Research Council (EPSRC): <http://www.epsrc.ac.uk/>
18. Electronic Resources Preservation Access Network (ERPANET): <http://www.erpanet.org>

19. ERPANET: Packaged Object Ingest Project.
http://www.erpanet.org/events/2003/rome/presentations/ross_rusbridge_pres.pdf
20. Finn, A., Kushmerick, N.: Learning to Classify Documents According to Genre. *Journal of American Society for Information Science and Technology* (2006) 57 (11), 1506-1518.
21. Giuffrida, G., Shek, E. Yang, J.: Knowledge-based Metadata Extraction from PostScript File. 5th ACM Intl. Conf. Digital Libraries (2000) 77-84.
22. Han, H., Giles, L., Manavoglu, E., Zha, H., Zhang, Z., Fox, E. A.: Automatic Document Metadata Extraction using Support Vector Machines. 3rd ACM/IEEE-CS Conf. Digital libraries (2000) 37-48.
23. Hedstrom, M., Ross, S., Ashley, K., Christensen-Dalsgaard, B., Duff, W., Gladney, H., Huc, C., Kenney, A. R., Moore, R., Neuhold, E.: Invest to Save: Report and Recommendations of the NSF-DELOS Working Group on Digital Archiving and Preservation. Report of the European Union DELOS and US National Science Foundation Workgroup on Digital Preservation and Archiving (2003) <http://delos-noe.iei.pi.cnr.it/activities/internationalforum/Joint-WGs/digitalarchiving/Digitalarchiving.pdf>
24. Joint Information Systems Committee: <http://www.jisc.ac.uk/>
25. Karlgren, J. and Cutting, D.: Recognizing Text Genres with Simple Metric using Discriminant Analysis. 15th Conf. Comp. Ling. Vol 2 (1994) 1071-1075.
26. Ke, S. W., Bowerman, C. Oakes, M. PERC: A Personal Email Classifier. 28th European Conf. Information Retrieval (ECIR 2006) 460-463.
27. Kessler, B., Nunberg, G., Schuetze, H.: Automatic Detection of Text Genre. 35th Ann. Meeting ACL (1997) 32-38.
28. Kim, Y.: Anaphora Resolution for Automatic Citation Linking. Masters Thesis, MSc. for Speech and Language Processing, University of Edinburgh (2004).
29. Kim, Y., Ross, S.: Genre Classification in Automated Ingest and Appraisal Metadata. European Conf. on advanced technology and research in Digital Libraries (2006) Springer LNCS.
30. Kim, Y., Ross, S.: Automating Metadata Extraction: Genre Classification Poster, UK e-Science All Hands Meeting (2006) Nottingham.
31. Kim, Y., Webber, B.: Implicit Reference to Citations: A study of astronomy papers. (preprint, reference available upon request)
32. MetadataExtractor: <http://pami.uwaterloo.ca/> (follow the link for Text Mining)
33. National Archives UK: DROID (Digital Object Identification). <http://www.nationalarchives.gov.uk/aboutapps/pronom/>
34. National Library of Medicine US: <http://www.nlm.nih.gov/>
35. National Library of New Zealand: Metadata Extraction Tool. <http://www.natlib.govt.nz/en/whatsnew/4initiatives.html#extraction>
36. Noonberg, D., B.: XPDF PDF document viewer. <http://www.foolabs.com/xpdf/>
37. Adobe Acrobat PDF specification:
http://partners.adobe.com/public/developer/pdf/index_reference.html
38. PDFTOHTML, PDF to HTML converter. <http://pdftohtml.sourceforge.net/>
39. Python Imaging Library: <http://www.pythonware.com/products/pil/>
40. PREMIS (PREservation Metadata: Implementation Strategy) Working Group: <http://www.oclc.org/research/projects/pmwg/>
41. Python: <http://www.python.org>
42. Riloff, E., Wiebe, J., and Wilson, T.: Learning Subjective Nouns using Extraction Pattern Bootstrapping. 7th CoNLL, (2003) 25-32.
43. Ross, S., Hedstrom, M.: Preservation Research and Sustainable Digital Libraries. *Intl. Journal of Digital Libraries* (2005) DOI: 10.1007/s00799-004-0099-3.

44. Rauber, A., Müller-Kögler, A.: Integrating Automatic Genre Analysis into Digital Libraries. ACM/IEEE Joint Conf. Digital Libraries (2001) Roanoke, VA 1-10.
45. Santini, M.: A Shallow Approach To Syntactic Feature Extraction For Genre Classification. 7th Ann. Colloq. UK Special Interest Group for Comp. Ling. (2004).
46. Santini, M.: State-of-the-art on Automatic Genre Identification. Tech. Report ITRI-04-03 ITRI University of Brighton, UK (2004).
47. Sebastiani, F.: Machine Learning in Automated Text Categorization. ACM Computing Surveys, Vol. 34 (2002) 1-47
48. Shafait, F., Keysers, D., Breuel, T., M.: Performance Comparison of Six Algorithms for Page Segmentation. 7th IAPR Workshop on Document Analysis Systems (DAS) (2006) 368-379.
49. Shao, M., Futrelle, R.: Graphics Recognition in PDF document. 6th IAPR Intl. Workshop on Graphics Recognition (GREC2005), 218-227.
50. Thoma, G.: Automating the production of bibliographic records. R&D report of the Communications Engineering Branch, Lister Hill National Center for Biomedical Communications, National Library of Medicine, 2001.
51. Witte, R., Krestel, R. and Bergler, S.: ERSS 2005: Coreference-based Summarization Reloaded. DUC 2005 Document Understanding Workshop, Canada
52. Witten, I. H., Frank, E.: Data Mining: Practical machine Learning tools and techniques. 2nd Edition, Morgan Kaufmann, San Francisco (2005).