

A Search Engine for Japanese Academic Papers

Emi Ishita (Surugadai University), Teru Agata (Asia University), Atsushi Ikeuchi (Tsukuba University), Michiko Nozue (Railway Technical Research Institute), Yosuke Miyata and Shuichi Ueda (Keio University)



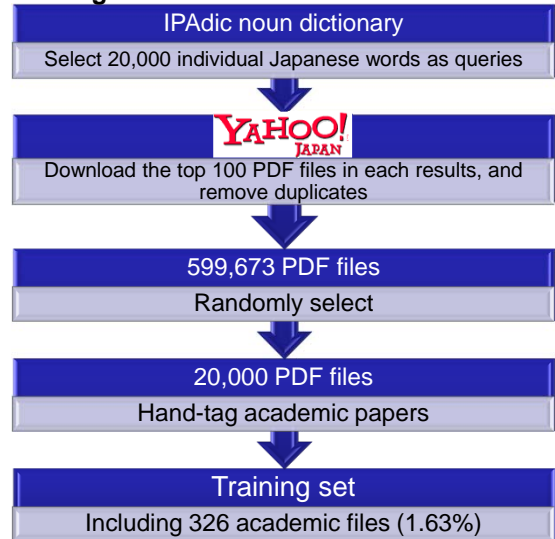
Our Research Purpose

- Development of the "Aletheia" search engine for academic papers written in Japanese using open software

Our Approach

- This system obtains PDF files on a broad range of topics using a commercial search engine, automatically detects Japanese academic papers using classifiers based on text structure, and supports end-user search using Lucene.

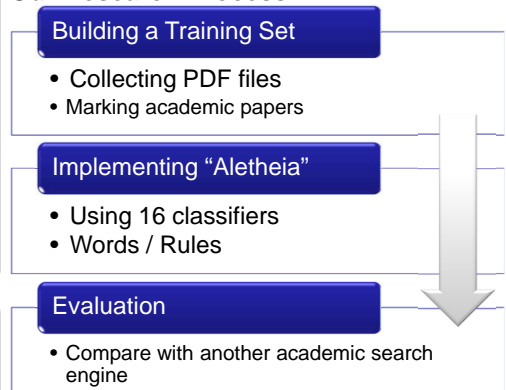
Training Set



The Criteria for Academic Papers

- In the annotator's opinion, had a layout typical of an academic paper
- Included a title, at least one author name, and an affiliation for at least one author
- Included exactly one paper per file
- Included at least one reference
- Included at least two pages.

Our Research Process



Approach, Classifiers, and Token

Approach	classifier	(round)	Line break	Token
Word based	SVM		None	Words
	SVM		None	Bigram
	SVM		Delete	Words
	SVM		Delete	Bigram
	AdaBoost	(10)	None	Words
	AdaBoost	(100)	None	Words
	AdaBoost	(1000)	None	Words
	AdaBoost	(10)	Delete	Words
	AdaBoost	(100)	Delete	Words
	AdaBoost	(1000)	Delete	Words
Rules	AdaBoost	(10)	-	-
	AdaBoost	(100)	-	-
	AdaBoost	(1000)	-	-
	Decision tree(C4.5)		-	-
	NaïveBayes		-	-
	Vote		-	-

Implementation

- When the user enters a query, Lucene is used to identify on-topic PDF files, which are then re-ranked in decreasing order of number of classifiers that classify each file as an academic paper.

Rules + Hand-Built Features

Layout

- File size
- Second level domain of the URL (ac.jp / go.jp)

Style

- Literacy style
- Using Kanji more

Words

- Research, reference, participant, survey, analysis, experiment, research report, research note, figure, table, paper, consideration, conclusion, research institute, research center, University

Example of a Search Result

検索式 = contents:分類器
検索されたPDFファイルは36件です。
スコア URLと要約

160.37 [http://www.ripi.titech.ac.jp/~tamada/2005_suzuki.pdf](#) キーワード: supervised EM アルゴリズムとナイーブベイズ分類器による非線形な高次元データを学習して、02-01 ナイーブベイズ分類器で発音の音素の分類をする際の精度を向上させる。2 次の特徴を導くことは教師付き学習手法としてナイーブベイズ分類器を用いるナイーブベイズ分類器をwしたのEM アルゴリズムと組合せることによりP(x, y) supervised な方法を離れ変数としてL1分布をwしたの事前分布とする対数尤度の確率変数に関する期待値(Q関数)は次のように定義できる: 3.1 ナイーブベイズ分類器による

160.35 [http://chusen.org/~taku/pubs/thesis/02_2001.pdf](#) キーワード: supervised EM アルゴリズムとナイーブベイズ分類器による非線形な高次元データを学習して、02-01 ナイーブベイズ分類器で発音の音素の分類をする際の精度を向上させる。2 次の特徴を導くことは教師付き学習手法としてナイーブベイズ分類器を用いるナイーブベイズ分類器をwしたのEM アルゴリズムと組合せることによりP(x, y) supervised な方法を離れ変数としてL1分布をwしたの事前分布とする対数尤度の確率変数に関する期待値(Q関数)は次のように定義できる: 3.1 ナイーブベイズ分類器による

160.28 [http://research.nii.ac.jp/kenken/ishikawa/reports/H15_A02/A02-21.pdf](#) キーワード: 方法では任意のテキストに対して予め指定された概念体系のどの位置にどの概念がどの程度出現しているかを検出する。この方法は2つのレベルに分けて検出される。まず、1つのレベルでは、複数の概念体系を同時に検出する。次に、2つのレベルでは、検出された概念体系の中から、最も適切な概念体系を選択する。この方法は、従来の概念体系検出方法よりも、検出された概念体系の精度を向上させる。この方法は、従来の概念体系検出方法よりも、検出された概念体系の精度を向上させる。この方法は、従来の概念体系検出方法よりも、検出された概念体系の精度を向上させる。

160.19 [http://www.kielnet.co.jp/kelnet/karuga/paper/2001-04-01.pdf](#) キーワード: 方法では任意のテキストに対して予め指定された概念体系のどの位置にどの概念がどの程度出現しているかを検出する。この方法は2つのレベルに分けて検出される。まず、1つのレベルでは、複数の概念体系を同時に検出する。次に、2つのレベルでは、検出された概念体系の中から、最も適切な概念体系を選択する。この方法は、従来の概念体系検出方法よりも、検出された概念体系の精度を向上させる。この方法は、従来の概念体系検出方法よりも、検出された概念体系の精度を向上させる。この方法は、従来の概念体系検出方法よりも、検出された概念体系の精度を向上させる。

www-lab.sliis.tsukuba.ac.jp/Aletheia

Evaluation

- A measure of topic coverage; How many queries resulted in at least one result.
 - A measure of classifier accuracy; On average, how many of the top 10 documents were (manually determined to be) academic papers by our definition.
- ### Test Queries
- 180 technical terms, randomly selected from titles of doctoral dissertations completed in 2005 at Tokyo, Kyoto and Keio Universities.

	1,110	1,327	201
Number of results			
No-result rate	22.2%	6.7%	79.4%
Number of academic papers	451	386	39
Academic papers rate	40.6%	29.1%	19.4%
Access failures	0	636	36

Conclusion

- Although commercial services such as Google Scholar index a far larger number of papers than Aletheia, we have shown that with suitable attention to language-specific processing and classifier design it is possible to create systems that are somewhat more robust (i.e., that have fewer zero-result sets) and somewhat better focused (i.e., that return fewer non-academic papers).
- In future work we plan to work on focused crawling for academic papers to address the coverage issue.