Multiclass Multi-Level Text Document Classification Yasotha, R.¹, Charles, E.Y.A.⁺ Faculty of Graduate Studies, UoJ *Department of Computer Science, UoJ

Introduction



Level 2

multi level means



Text Document Categorization (TC) - also known as text document classification or text classification or topic spotting, is the task of assigning predefined label to text documents.

Objective

Categorize articles belong to multiclass and multi levels of topics.

- Testing platform for research study: Abstracts of conference papers displayed in ACM digital library.
- Labels: defined by ACM-CCS
- A General and reference
- C Computer systems organization
- E Software and its engineering
- G Mathematics of Computing
- I Security and privacy
- K Computing methodologies
- M Social and professional topics
- 🚹 B Hardware
- D Networks
- F Theory of computation
- H Information systems
- J Human-centered computing
- L Applied computing

N Person nouns people technologie...

Dataset

100 distinct abstracts in each category presented in ACM-CCS in ACM digital library.

Simple Intuition:LDA Model

A document exhibits multiple topics is illustrated in the figure.

> **Text Categorization with Support Vector Machines: Learning with Many Relevant**



assifiers from examples. It analyses the rming methods and behave robustly over Furthermore, they are fully automatic, eliminating the need for manual parame



Topic-per-Document

(TD)

Each **topic** is a distribution over words Each **document** is a mixture of corpus-wide topics Each **word** is drown from one of these topics



Proposed Method

In this approach, we proposed a method to pre- **Result** ... dict a label to any text documents. Here, features are extracted using topic model, LDA (Latent Dirichlet Allocation) proposed by Blei et al.,2003 [1,2].

In LDA, the observed data are the words of each document and the hidden variables represent the latent topical structure, ie., the topics themselves and how each document exhibits them.



Preprocessing

- 1. Extract raw words.
- 2. Remove non-literal characters:{",<,!,',?,etc.}
- 3. Remove stop-words:{if, and, that,etc.}
- 4. Remove words that are generally used in abstracts : {paper, describe, addressing, proposed, framework, etc.}

Result

Gibb sampling algorithm is used to determine the number of clusters in the domain concerned. In this approach, samples posterior distribution at several choice of Topics (T) were computed.

The computation result is shown in the graph shows that, all level I categories under ACM-CCS are divided into four groups according to similarly of documents.

algorithm Κ means divide is to used documents all unall 13 categories der into four clusters and identified categories falling in each cluster.



 χ^2 feature selection method is used for selecting top ten keywords from each cluster.

Clu tim per rea dat soft anr SVS lang algo spe



Cluster	Categories
Cluster 0	A, H, I, J, L & M
Cluster 1	Κ
Cluster 2	F & G
Cluster 3	B, C, D & E

Result ...

ster 0	Cluster 1	Cluster 2	Cluster 3
e	information	performance	problem
formance	data	memory	algorithm
l	user	processor	bound
a	technology	high	log
ware	analysis	power	number
olication	service	hardware	graph
tem	process	architecture	set
guage	work	level	known
orithm	need	application	linear
cific	different	overhead	time

Result ...



References

[1] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent dirichlet allocation," pp. 993–1022, 2003. [2] D. M. Blei, "Probabilistic topic models," vol. 55, pp. 77–84, April 2012.

[3] L. G. Thoman and M. Steyvers, "Finding scientific topics," vol. 101, pp. 5228–5235, April 6 2004.



The table shows the result in terms of F_1 score.

ach	F_1 score(%)
approach with NB	34.62
approach with SVM	34.62
sed Approach	66.66

• Size of the vocabulary

• Total number of topics

• Topics correlation: unified-topics, commontopics, similar-topics, non-overlappingtopics, etc.

• Labeling: Human factor in assigning labels to text documents in the collected dataset.

Instead of common practice of string matching approach, this poster presents an LDA based approach for automatic TC. Proposed model was able to categorize unseen documents with an ac-

Future Direction

• Preliminary step for automatic categorization of text documents in MLMC. This may be extended to a ontology in a polyhierarchical structure.

• In this approach, terms are in single form was considered. Single term cannot provide complete identification of document thematic content.

• Word features in form of bigram and trigram are more influenced in the prediction of topic spotting with more effectiveness.