A Combined Approach to Dynamic web page Classification: merging Structure and Content

Maria Niarou, Sofia Stamou

Department **o**f Archives, Library and Museum Studies Ionian University Corfu, Greece

WEB2020, IARIA Conference Lisbon, Portugal



[l14niar@ionio.gr, stamou@ionio.gr]



Presenter:

Maria (I14niar@ionio.com) is a PhD student at the Ionian University in Corfu. She is a graduate of the Department of Archives, Library Science and Museology (Faculty of Information Science & Informatics). During her Master's studies (Information Service in Digital Environment), she concluded that what interests her most is the data processing via semantic analysis. The research she conducted during the writing of her master's thesis resulted in her first publication: Niarou M., Stamou S. 2012; "Exploring Lexical Ontologies for Hierarchically Organizing the Greek Wikipedia Articles"; International Journal of Digital Information Management, Vol. 10, No. 3, pp. 157-167. The present presentation is the core of the academic research that Maria carries out, supervised by Sofia Stamou (Assistant Professor at the Department of Archives, Library Science and Museology).

Research Interests:

- Natural Language Processing
- Information Retrieval
- Semantics
- Knowledge Extraction from Digital Libraries and Web
- Knowledge Organization Systems

Languages: Greek, English, Italian, French



Presentation Layout

INTRODUCTION

- Motivation
- Related Work
- Contribution

• METHODOLOGY

- Outline
- Algorithm1: Multi-Dimensional Classification (Procedure1)
- Algorithm1: Multi-Dimensional Classification (Procedure2)
- Algorithm2: Re-Classification based on Pages' Change Detection
- Algorithm3: Optimized Re-Classification based on Change's Frequency Detection

• EXPERIMENTAL STUDY

- Framework
- Results

• CONCLUSION and FUTURE WORK

- Conclusion
- Future Work



Intro (1/3)



Intro (2/3)

Related Work (2 main groups)





Contribution





Outline : three autonomous and supplementary algorithms

Algorithm 1 : Multi-Dimensional Page Classification

- Idea: textual & structural classification = complementary
- Goal: to classify the pages structurally and thematically

Algorithm 2: Re-Classification based on Changes' Grade Detection

- Idea: detect pages' dynamic nature via a re-Classification component, in order to maintain data indexes updated
- Goal: to decide which of the modified pages need to be re-Classified

Algorithm 3: Optimized Re-Classification based on Changes' Frequency Detection

- Idea: optimize the re-Classification process
- Goal: save time and resources

Methodology (2/5)

Algorithm 1: Multi-Dimensional Page Classification

Procedure 1: Structure-based classification

Baseline: we reverse the argument that web searches can be classified as either Navigational, Informational or Transactional, and we claim that pages can be classified accordingly.

• <u>Phase1:</u> Page Type Recognition

•



Input elements "Clue" classification

elements

Methodology (3/5)

Algorithm1: Multi-Dimensional Page Classification

Procedure2: Content-based classification

Baseline: having experimented with several textual features, we ended up with **anchor title** and **title** as the most easily extracted and informative of the theme of a page.

• Phase1: Textual elements Extraction

Phase2: Theme Detection

Input elements "Clue" classification elements Methodology (4/5)

Algorithm2: Re-Classification based on Change Detection

Challenge: how to deal with pages' dynamic nature // how to ensure that the classification outcome is up-to-date **Idea**: re-Classification component to our algorithm **Goal**: to detect, measure and identify the *possible* changes

• **<u>Procedure1</u>**: Re-Classification Decision based on Textual Changes

Input elements "Clue" classification elements

Compare (structural and textual) elements of any given page

with their counterparts previously identified!

Methodology (5/5)

Algorithm3: Optimized Re-Classification based on Change's Frequency Detection

Baseline:

capture pages' change frequency => => determine the re-classification policy => optimize the runs of our Re-Classification algorithm =>

=> save time and resources

Experimental Study

Framework (4 points)

Higly changing pages after 3 months

Conclusion & Future Work

Conclusion

Future work

Any Questions...

A Combined Approach to Dynamic web page Classification: merging Structure and Content

Maria Niarou, Sofia Stamou

Department of Archives, Library and Museum Studies Ionian University Corfu, Greece

[l14niar@ionio.gr, stamou@ionio.gr]

Thank you!

