# Automatic Metadata Extraction: The High Energy Physics Use Case

Joseph Boyd

June 24, 2015

#### **Contents**

1	Intr	oduction	4		
	1.1	Motivation	4		
	1.2	Aims	4		
	1.3	Main Results	4		
	1.4	Outline	4		
2	Sup	ervised Sequence Learning	4		
	2.1	Log-linear Models	4		
	2.2	Graphical Models	4		
		2.2.1 Hidden Markov Models	4		
	2.3	Conditional Random Fields	4		
		2.3.1 Feature Engineering	4		
3	Λ+.	omatic Metadata Extraction	,		
J	3.1	Metadata Extraction	7		
	3.2	Related Works	/		
	5.2	Itelated Works	7		
4	Implementation and Data				
	4.1	GROBID	4		
	4.2	Extensions	4		
	4.3	Data Acquisition	4		
5	Res	ults and Analysis	4		
	5.1	Experiment Summary	4		
	5.2	Evaluation Method	4		
	5.3	Baseline	4		
		5.3.1 Header model - Cora dataset	4		
		5.3.2 Header model - Cora dataset appending HEP dataset	4		
		5.3.3 Header model - Cora and HEP combined datasets	4		
		5.3.4 Header model - HEP dataset	4		
		5.3.5 Header model - HEP dataset appending CORA dataset	4		
		5.3.6 Header model - HEP dataset appending $1/3$ CORA dataset	4		
		5.3.7 Header model - HEP dataset appending $2/3$ CORA dataset	4		
		5.3.8 Segmentation model - Cora dataset	4		
		5.3.9 Segmentation model - Cora dataset appending HEP dataset	4		
		5.3.10 Segmentation model - Cora and HEP combined datasets	4		
		5.3.11 Segmentation model - HEP dataset	4		
		5.3.12 Segmentation model - HEP dataset appending CORA dataset	4		

	5.4	Regularisation
		5.4.1 Header model - $L2 = 0$
		5.4.2 Header model - $L2 = 1e^{-6}$
		5.4.3 Header model - $L2 = 1e^{-5}$
		5.4.4 Header model - $L2 = 1e^{-4}$
		5.4.5 Header model - $L2 = 1e^{-3}$
	5.5	Dictionaries
		5.5.1 Header model - HEP dataset
		5.5.2 Header model - HEP dataset appending CORA dataset
		5.5.3 Segmentation model - HEP dataset
		5.5.4 Segmentation model - HEP dataset appending CORA dataset
		5.5.5 Header Model - HEP dataset - 2 <sup>nd</sup> Degree Features
		5.5.6 Header Model - HEP dataset Appending CORA - 2 <sup>nd</sup> Degree Features
		5.5.7 Header Model - HEP dataset - 3 <sup>rd</sup> Degree Features
		5.5.8 Header Model - HEP dataset Appending CORA - 3 <sup>rd</sup> Degree Features
	5.6	Dictionaries + stop words
		5.6.1 Header model - HEP dataset
		5.6.2 Header model - HEP dataset appending CORA dataset
		5.6.3 Segmentation model - HEP dataset
		5.6.4 Segmentation model - HEP dataset appending CORA dataset
		5.6.5 Header Model - HEP dataset - 2 <sup>nd</sup> Degree Features
		5.6.6 Header Model - HEP dataset Appending CORA - 2 <sup>nd</sup> Degree Features
		5.6.7 Header Model - HEP dataset - 3 <sup>rd</sup> Degree Features
		5.6.8 Header Model - HEP dataset Appending CORA - 3 <sup>rd</sup> Degree Features
	5.7	Token Selection
		5.7.1 Segmentation Model - HEP dataset - 5 Tokens
		5.7.2 Segmentation Model - HEP dataset - 10 Tokens
		5.7.3 Segmentation Model - HEP dataset - 15 Tokens
		5.7.4 Segmentation Model - HEP dataset - 20 Tokens
	5.8	Levenshtein
		5.8.1 Segmentation Model - HEP dataset - Binary Threshold (0.05)
		5.8.2 Segmentation Model - HEP dataset - Binary Threshold (0.1)
		5.8.3 Segmentation Model - HEP dataset - Binary Threshold (0.2)
		5.8.4 Segmentation Model - HEP dataset - Binary Threshold (0.4)
		5.8.5 Segmentation Model - HEP dataset - Binary Threshold (0.8)
		5.8.6 Segmentation Model - HEP dataset - Ternary Threshold
		5.8.7 Segmentation Model - HEP dataset - Quaternary Threshold
	5.9	Line Shape
		5.9.1 Segmentation Model - HEP dataset - Binary Threshold
		5.9.2 Segmentation Model - HEP dataset - Ternary Threshold
	5.10	Template Matching
		5.10.1 Segmentation Model - HEP dataset
6	Cone	clusion
Ū	6.1	Summary
	0.1	6.1.1 Key Results
	6.2	Future Work
7	Refe	rences 4
8	App	endices 2

#### 1 Introduction

- 1.1 Motivation
- 1.2 Aims
- 1.3 Main Results
- 1.4 Outline

## 2 Supervised Sequence Learning

- 2.1 Log-linear Models
- 2.2 Graphical Models
- 2.2.1 Hidden Markov Models
- 2.3 Conditional Random Fields
- 2.3.1 Feature Engineering

#### 3 Automatic Metadata Extraction

- 3.1 Metadata Extraction
- 3.2 Related Works

## 4 Implementation and Data

- 4.1 GROBID
- 4.2 Extensions
- 4.3 Data Acquisition

### 5 Results and Analysis

- 5.1 Experiment Summary
- 5.2 Evaluation Method
- 5.3 Baseline
- 5.3.1 Header model Cora dataset
- 5.3.2 Header model Cora dataset appending HEP dataset
- 5.3.3 Header model Cora and HEP combined datasets
- 5.3.4 Header model HEP dataset
- 5.3.5 Header model HEP dataset appending CORA dataset
- F 2.6. Useder model HED detect amonding 1/2 CODA detect