

Automatic Metadata Extraction

The High Energy Physics Use Case

Joseph Boyd

École Polytechnique Fédérale de Lausanne

joseph.boyd@epfl.ch

August 24, 2015

Motivation

- ▶ INSPIRE-HEP digital library at CERN contains over 1 Million documents
- ▶ Manual curation of high energy physics (HEP) papers may be automated with machine learning techniques
- ▶ Custom datasets and specialised features required to model HEP paper characteristics

Aims

Take existing state-of-the-art system for metadata extraction to:

- ▶ demonstrate a qualitative difference between HEP and general papers;
- ▶ propose improvements to model features;
- ▶ run experiments to confirm these improvements, and;
- ▶ draw conclusions about what characterises good feature engineering.

Outline

Introduction

Theory

Automatic Metadata Extraction

Data, Methods, and Implementation

Key Results

Conclusions

Why CRFs?

- ▶ Transition interdependencies implies graphical structure best modelled as a structured sequence
- ▶ Modelling conditional distribution, $p(\mathbf{y}|\mathbf{x})$ sufficient for classification
- ▶ Exploit rich information about observations, \mathbf{x} , without explicitly modelling the underlying probability distribution
- ▶ Classifying metadata may greatly benefit from modelling rich text features (punctuation, font size, layout,...)

Mathematical Formulation

$$p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{y})}{\sum_{\mathbf{y}'} p(\mathbf{x}, \mathbf{y}')} = \frac{1}{Z(\mathbf{x})} \exp \left\{ \sum_k \lambda_k F_k(y_t, y_{t-1}, x_t) \right\}, \quad (1)$$

where $Z(\mathbf{x}) = \sum_{\mathbf{y}'} \exp \left\{ \sum_k \lambda_{ij} F_k(y'_t, y'_{t-1}, x_t) \right\}$ is known as the partition function, ensuring probabilities sum to 1.

$F_k(\mathbf{x}, y) = \sum_t^T f_k(\mathbf{x}, y)$, where f_k is a (typically boolean) function describing one of several features about a token.

The form of the functions themselves, $f(\cdot)$, are known in Wapiti (Section ??) as *templates*. It is in choosing these explicitly that we perform feature engineering.

Solution Approach

- ▶ Formulate convex maximum log likelihood estimator, $l(\Lambda)$, where $\Lambda = \{\lambda_k\}_{k=1}^K$
- ▶ Train (determine Λ) with gradient ascent technique, L-BFGS. Each iteration, l , requires forward-backward algorithm to compute $Z(\mathbf{x}^{(n)})$ for each of N samples – $\mathcal{O}(INT|S|^2)$.
- ▶ Prediction with Viterbi algorithm – $\mathcal{O}(T|S|^2)$.

Outline

Introduction

Theory

Automatic Metadata Extraction

Data, Methods, and Implementation

Key Results

Conclusions

Metadata Extraction

- ▶ *Metadata* refers to content useful to the identification of the document
- ▶ *Extraction* refers to the identification of metadata within the document text
- ▶ Several automatic approaches exist: stylistic analysis, knowledge-base, machine learning, ...

Metadata Extraction (Illustration)

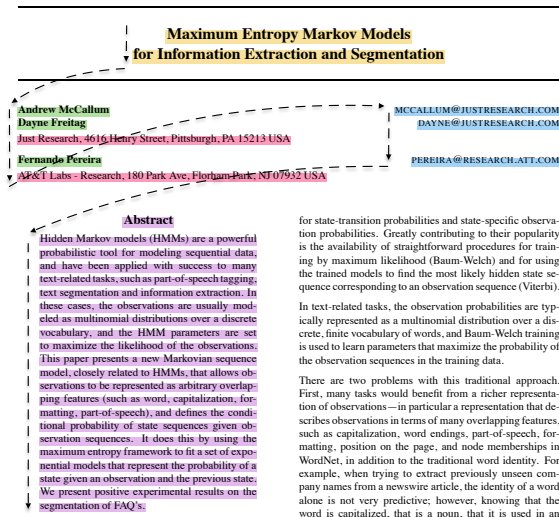


Figure: Tagging of a document header section.

GROBID

- ▶ Selected according to performance in study comparing AME systems [2]
- ▶ Open source Java-based tool developed at INRIA, France
- ▶ Manages *cascade* of CRF models for annotating papers in progressively finer detail
- ▶ Uses C++ library *Wapiti* for back-end calculations (training, prediction)

GROBID - CRF Cascade

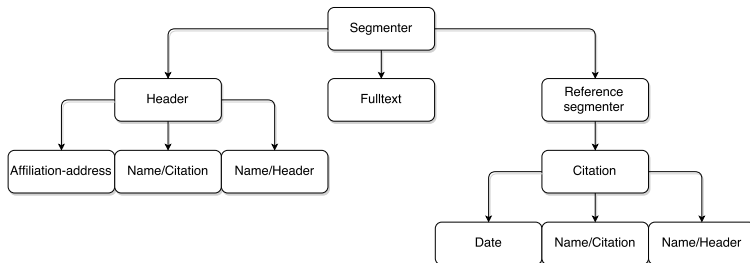


Figure: Cascade of models used by Grobid

Outline

Introduction

Theory

Automatic Metadata Extraction

Data, Methods, and Implementation

Key Results

Conclusions

Identification of beauty and charm quark jets at LHCb

The LHCb collaboration¹

Abstract

Identification of jets originating from beauty and charm quarks is important for measuring Standard Model processes and for searching for new physics. The performance of algorithms developed to select b - and c -quark jets is measured using data recorded by LHCb from proton-proton collisions at $\sqrt{s} = 7$ TeV in 2011 and at $\sqrt{s} = 8$ TeV in 2012. The efficiency for identifying a $b(c)$ jet is about 65%(25%) with a probability for misidentifying a light-parton jet of 0.3% for jets with transverse momentum $p_T > 20$ GeV and pseudorapidity $2.2 < \eta < 4.2$. The dependence of the performance on the p_T and η of the jet is also measured.

Submitted to JINST

© CERN on behalf of the LHCb collaboration, license CC-BY-4.0.

encode different attribute dimensions of an input data space. A good glyph design can enable users to conduct visual search more efficiently during interactive visualization, and facilitate effective learning, memorizing and using the visual encoding scheme. A less effective visual design may suffer from various shortcomings such as being perceptually confusing, semantically ambiguous, difficult to learn and remember, or unable to accommodate low-resolution display devices.

- Eamonn Maguire is with Oxford e-Research Centre and Department of Computer Science, University of Oxford, UK. E-mail: eamonn.maguire@st-annes.ox.ac.uk.
- Philippe Rocca-Serra, Susanna-Assunta Sansone and Min Chen are with Oxford e-Research Centre, University of Oxford, UK. E-mail: {philippe.rocca-serra,susanna-assunta.sansone,min.chen}@oerc.ox.ac.uk.
- Jim Davies is with Department of Computer Science, University of Oxford, UK. E-mail: jim.davies@cs.ox.ac.uk.

Manuscript received 31 March 2012; accepted 1 August 2012; posted online 14 October 2012; mailed on 5 October 2012.

For information on obtaining reprints of this article, please send e-mail to: tvcg@computer.org.

(a) Collaboration field in header section.

(b) Discontinuous header data.

LHCb collaboration

R. Aaij²⁸, B. Adeva³⁷, M. Adinolfi⁴⁶, A. Affolder⁷², Z. Ajaltouni⁵, S. Akar⁶, J. Albrecht⁹, F. Alessio²⁸, M. Alexander³¹, S. Ali⁴¹, G. Alkhazov²⁰, P. Alvarez Cartelle⁵¹, A.A. Alves Jr⁵⁷, S. Amato⁷, S. Amico⁷², Y. Amhis⁷, L. An³, L. Anderlini^{17,4}, J. Anderson⁴⁰, M. Andreotti^{116,1}, J.E. Andrews²⁸, R.B. Appleby⁵⁴, O. Aquines Gutierrez²⁸, F. Archilli³⁸, P. d'Argent⁴¹, A. Artamonov³⁵, M. Artuso³⁹, E. Aslanides⁶, G. Auriemma^{25,48}, M. Baalouch⁵, S. Bachmann¹¹, J.J. Back⁴⁸, A. Badalov³⁶, C. Bacsso⁴⁹, W. Baldini^{16,38}, R.J. Barlow⁵⁴, C. Barschel³⁸, S. Barsuk⁷, W. Barter³⁸, C. Basso⁴⁹, V. Battista³⁹, A. Bay³⁹, L. Beaucourt⁴, J. Beldone²¹, F. Bedeschi¹²³, I. Bediaga⁴, L.J. Be⁴¹, I. Belyaev²¹, E. Ben-Haim⁸, G. Bencivenni¹⁸, S. Benson²⁸, J. Benton⁴⁶, A. Berezhnoy³², R. Bernert⁴⁰, A. Bertolin²², M.-O. Bettler³⁸, M. van Beuzekom⁴⁴, A. Bilen¹¹, S. Bifani⁴⁵, T. Bird⁴⁴, A. Birnkranz⁹, A. Bizzeti^{17,4}, T. Blake⁴⁸, F. Blanc³⁹, J. Blouw¹⁰, S. Blus²⁹, V. Bocci²³, A. Bondar³⁴, N. Bonard^{30,38}, W. Bonivento¹⁵, S. Borghi⁵⁴, M. Borsato⁷, T.J.V. Bowcock⁴², E. Bowen⁴⁶, C. Bozzi¹⁶, S. Braun¹¹, D. Brett⁵⁴, M. Britsch¹⁰, T. Britton²⁹, J. Brodzicka²⁴, N.H. Brook⁴⁶, A. Bursche⁴⁹, J. Buytaert³⁸, S. Cadedotte¹⁵, R. Calabrese^{46,1}, M. Calvi^{20,4}, M. Calvo Gomez^{36,1}, P. Campana³⁹, D. Campora Perez³⁸, L. Capriotti¹⁴, A. Carbone^{14,4}, G. Carbone^{10,4}, R. Cardinale^{10,4}, A. Cardini¹⁵, P. Carutti²⁸, L. Carson⁵⁰, K. Carvalho Akiba^{7,38}, R. Casanova Mohr³⁶, G. Casse³⁷, L. Cassina^{20,4}, L. Castillo Garcia³⁸, M. Cattaneo³⁸, Ch. Cauet⁸, G. Cavallero¹⁹, R. Cenci^{23,4}, M. Charles⁸, Ph. Charpentier³⁸, M. Chefeldt⁴⁶, S. Chen⁵⁴, S.-F. Cheung⁵⁵, N. Chierchia⁴⁰, M. Chruszcz⁴⁰, X. Cid Vidal³⁸, G. Ciezarek⁴¹, P.E.L. Clarke²⁰, M. Clemencic³⁸, H.V. Cliff⁴⁷, J. Clozier³⁸, V. Coco³⁸, J. Cogan⁶, E. Cogneras⁹, V. Cogoni^{15,4}, L. Cojocariu²⁹, G. Collazuol²², P. Collins³⁸, A. Comerma-Montells¹¹, A. Contu^{15,38}, A. Cook⁴⁶, M. Coombes⁴⁶, S. Coquereau⁸, G. Corti³⁸, M. Corvo^{46,1}, I. Counts³⁶, B. Couturier³⁸, G.A. Cowan²⁰, D.C. Craik⁴⁸, A. Crocombe⁴⁸,

¹Laboratori Nazionali dell'INFN di Frascati, Frascati, Italy

²Sezione INFN di Genova, Genova, Italy

³Sezione INFN di Milano Bicocca, Milano, Italy

⁴Sezione INFN di Milano, Milano, Italy

⁵Sezione INFN di Padova, Padova, Italy

⁶Sezione INFN di Pisa, Pisa, Italy

⁷Sezione INFN di Roma Tor Vergata, Roma, Italy

⁸Sezione INFN di Roma La Sapienza, Roma, Italy

⁹Henrik Niewodniczanski Institute of Nuclear Physics Polish Academy of Sciences, Kraków, Poland

¹⁰AGH - University of Science and Technology, Faculty of Physics and Applied Computer Science, Kraków, Poland

¹¹National Center for Nuclear Research (NCBJ), Warsaw, Poland

¹²Horia Hulubei National Institute of Physics and Nuclear Engineering, Bucharest-Magurele, Romania

¹³Petersburg Nuclear Physics Institute (PNPI), Gatchina, Russia

¹⁴Institute of Theoretical and Experimental Physics (ITEP), Moscow, Russia

¹⁵Institute of Nuclear Physics, Moscow State University (SINP MSU), Moscow, Russia

¹⁶Institute for Nuclear Research of the Russian Academy of Sciences (INR RAN), Moscow, Russia

¹⁷Budker Institute of Nuclear Physics (SB RAS) and Novosibirsk State University, Novosibirsk, Russia

¹⁸Institute for High Energy Physics (IHEP), Protvino, Russia

¹⁹Universitat de Barcelona, Barcelona, Spain

²⁰Universidad de Santiago de Compostela, Santiago de Compostela, Spain

²¹European Organization for Nuclear Research (CERN), Geneva, Switzerland

²²École Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland

²³Physik-Institut, Universität Zürich, Zürich, Switzerland

²⁴Nikhef National Institute for Subatomic Physics, Amsterdam, The Netherlands

²⁵Nikhef National Institute for Subatomic Physics and VU University Amsterdam, Amsterdam, The Netherlands

Model	HEP	CORA
Header	157 papers	2506 papers
Segmentation	169 papers	125 papers

Table: Number of training instances for each model from each dataset.

Outline

Introduction

Theory

Automatic Metadata Extraction

Data, Methods, and Implementation

Key Results

Conclusions

Experiment Setup

Confusion matrix - Segmentation (Baseline, HEP)

acknowledgement	544	5	258			10	1		
annex	13	1805	9078			67	1	3	46
body	142	2601	185466		162	303	45	41	590
cover				259	5	26	1		
footnote	6	8	495		1171	47	42	31	18
header		3	906	43	32	12211	11	6	36
headnote		18	372		56	76	1380	52	35
page	2	10	136		18	12	17	2257	12
references	115		392		11	124	18	5	10871
	acknowledgement	annex	body	cover	footnote	header	headnote	page	references

Figure: Baseline confusion segmentation

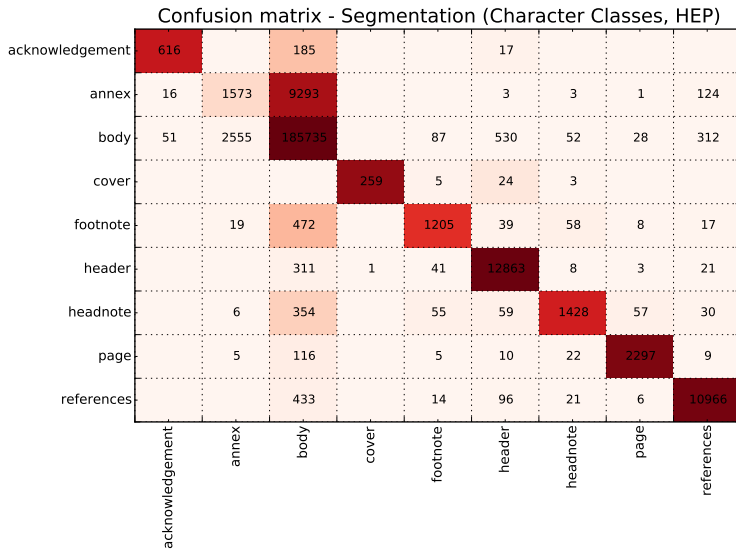


Figure: Classes confusion segmentation

Outline

Introduction

Theory

Automatic Metadata Extraction

Data, Methods, and Implementation

Key Results

Conclusions



R. Aaij, B. Adeva, M. Adinolfi, A. Affolder, Z. Ajaltouni, S. Akar, J. Albrecht, F. Alessio, M. Alexander, S. Ali, et al. Identification of beauty and charm quark jets at LHCb. *arXiv preprint arXiv:1504.07670*, 2015.



M. Lipinski, K. Yao, C. Breitinger, J. Beel, and B. Gipp. Evaluation of header metadata extraction approaches and tools for scientific pdf documents. In *Proceedings of the 13th ACM/IEEE-CS joint conference on Digital libraries*, pages 385–386. ACM, 2013.



E. Maguire, P. Rocca-Serra, S.-A. Sansone, J. Davies, and M. Chen. Taxonomy-based glyph design—with a case study on visualizing workflows of biological experiments. *Visualization and Computer Graphics, IEEE Transactions on*, 18(12):2603–2612, 2012.