# 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

# 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

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

where $Z(\mathbf{x}) = \sum_{y'} \exp\left\{ \sum_k \lambda_{ij} F_k(y'_t, y'_{t-1}, x_t) \right\}$ 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

- Formuluate convex maximum log likelihood estimator, $l(\Lambda)$, where $\Lambda = \{\lambda_k\}_{k=1}^K$
- Train (determine $\Lambda$) with gradient ascent technique, L-BFGS. Each iteration, $I$, 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

# 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

# HEP Paper Characteristics

(a) Collaboration field in header section.

(b) Discontinuous header data.

(c) Collaboration author list.

(d) Collaboration affiliation list.

Figure: Figure (A) shows a collaboration field in a header section. Figure

# Training Data

Two models addressed: *header* and *segmentation*

| Model | HEP | CORA |
|---|---|---|
| Header | 157 papers | **2506 papers** |
| Segmentation | **169 papers** | 125 papers |

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

# Dictionary Features

Dictionaries for INSPIRE-HEP affilations, authors, collaborations, journals, and titles. Stop words modelled as additional dictionary.

$$f_{\text{dict}_i}(x_t) = \mathbb{1}_{\{x_t \in \text{dict}_i\}},$$

# Character Class Features

$$f_{\text{class}_i}(x_t) = \frac{1}{|x_t|} \sum_{n=1}^{|x_t|} \mathbb{1}_{\{x_{ti} \in \text{class}_i\}},$$

for each character class, $\text{class}_i$, where $x_t$ is a token (hence a line for the *segmentation* model), and $x_{ti}$ is the *ith* character in the line. For the decimal (round down) case, we then define,

| Class | Regex |
|---|---|
| Spacing | r'[\s]' |
| Lower case | r'[a-z]' |
| Upper case | r'[A-Z]' |
| Numeric | r'[\d]' |
| Punctuation | r'[.,?:;]' |
| Special character | r'[^\sa-zA-Z d.,?:;]' |

Table: Character classes used as features, along with the regular expressions used to count them.

# Levenshtein Distance Features

$$\text{similarity}(a, b) = 1 - \frac{\text{lev}_{a,b}(|a|, |b|)}{\max(|a|, |b|)}.$$

Due to the constraints on numeric features (see Section ??), we must discretise the result. Thus, for a given line, $x_t$, we define the feature function,

$$f_{lev}(x_t) = \begin{cases} 0 & \text{if } 0 \leq \text{similarity}(x_t, x_{t-1}) \leq T_1 \\ 1 & \text{if } T_1 \leq \text{similarity}(x_t, x_{t-1}) \leq T_2 \\ \vdots & \qquad \vdots \\ \text{N-1} & \text{if } T_{N-1} \leq \text{similarity}(x_t, x_{t-1}) \leq 1 \end{cases}$$

where $T_1, T_2, ..., T_{N-1}$ are thresholds selected to create the $N$ categories. We try several thresholding strategies in our experimentation (see Section ??).
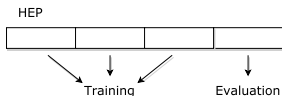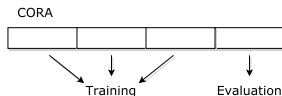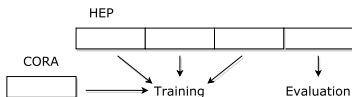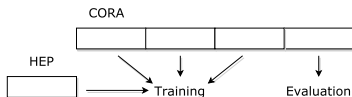
# Outline

# Experiment Setup



(a) CV HEP

(b) CV CORA

(c) CV HEP append CORA

(d) CV CORA append HEP

Figure: The different cross-validation configurations used in our experiments. Figures (A) and (B) show cross-validation on HEP and CORA sets independently. Figures (C) and (D) show cross-validation on the HEP and CORA datasets respectively, appending the other at training time.
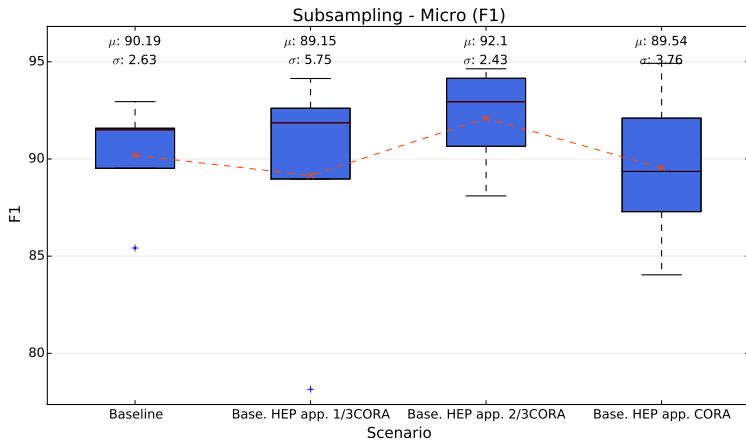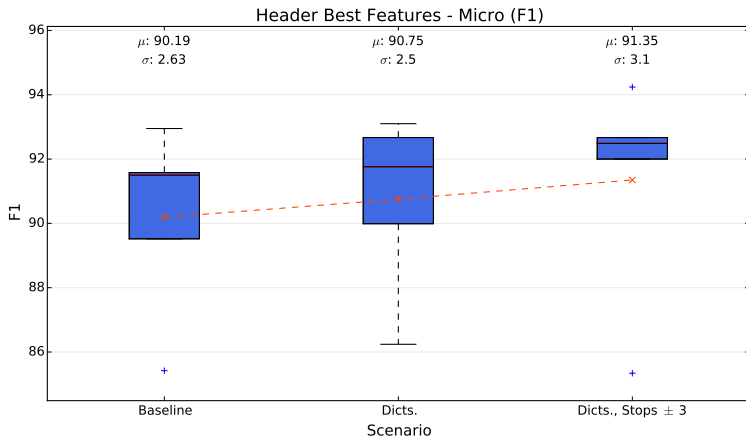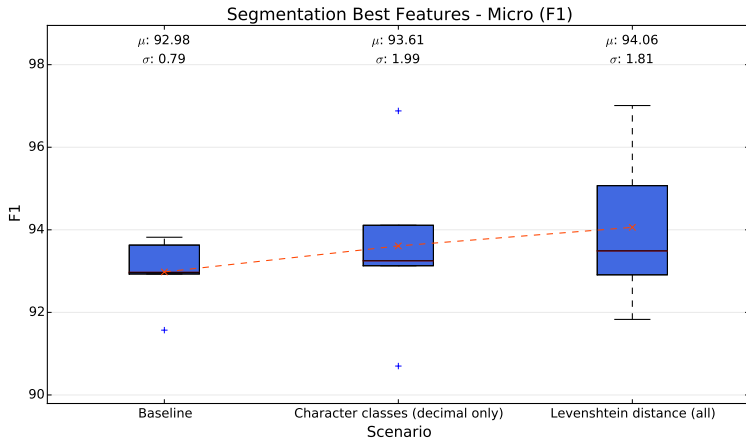
Figure: Baseline confusion segmentation

Figure: Baseline confusion segmentation

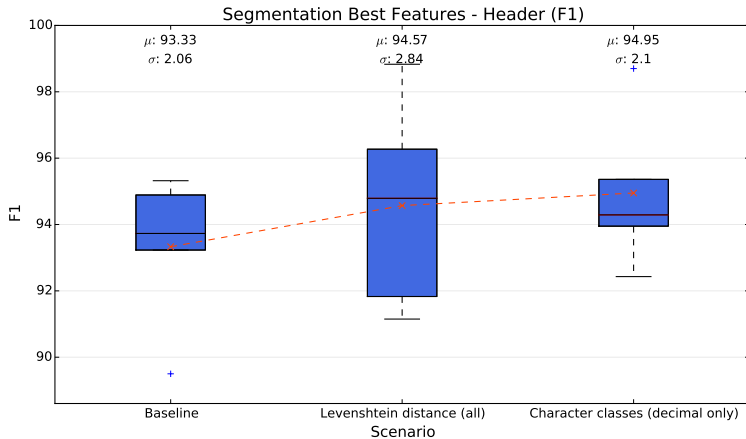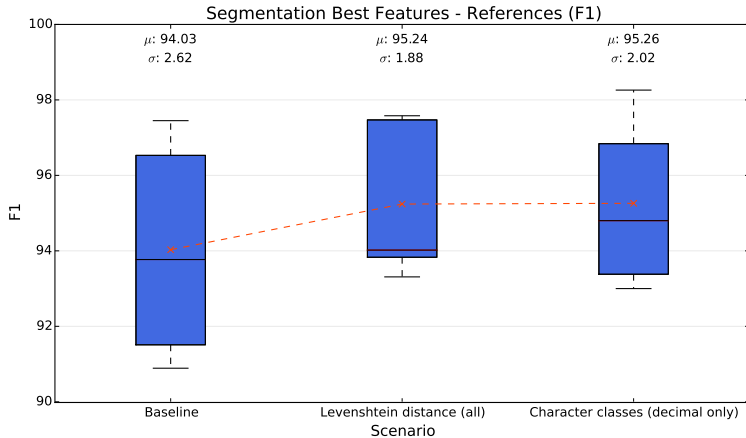Figure: Baseline confusion segmentation
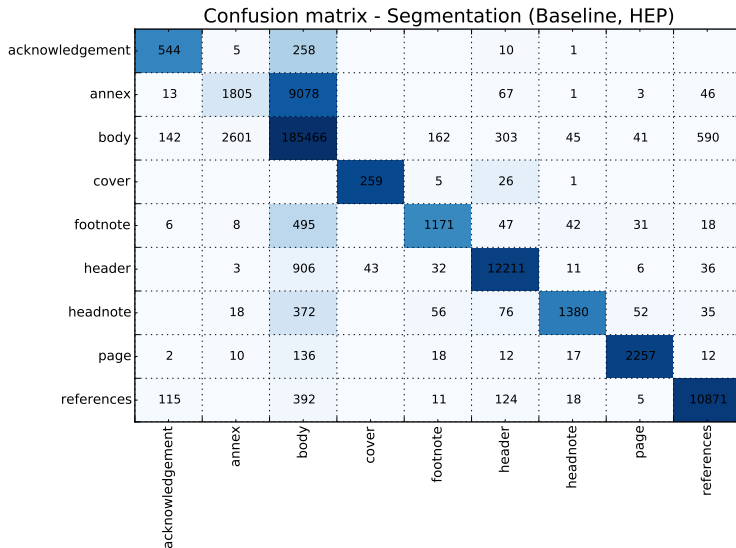
Figure: Baseline confusion segmentation

Figure: Baseline confusion segmentation

Figure: Classes confusion segmentation

Figure: Classes confusion segmentation
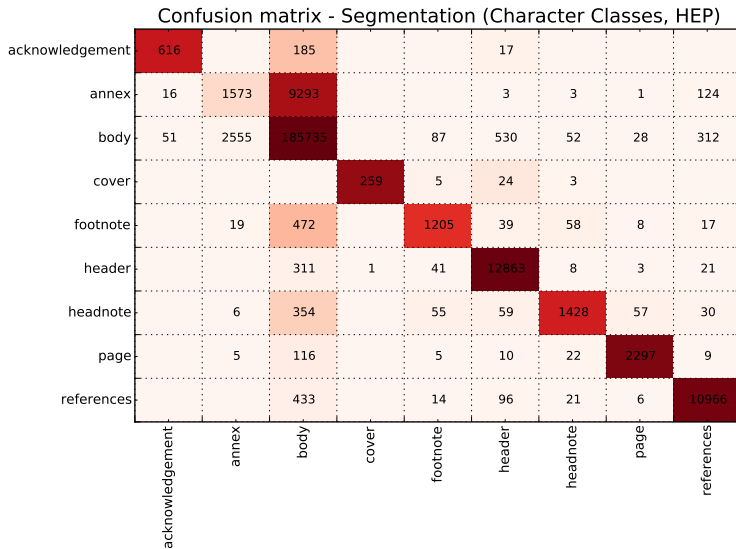
# Outline

# Conclusions

- Qualitative difference between HEP and general papers demonstrated (through inspection, subsampling)
- Most successful features were those offering a dimensionality reduction: dictionaries (12% error reduction), character classes (24% and 21% on header and reference classifications)

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.