School of ACCOUNTING FINANCE 會計及金融學院



# **PolyU\_CBS-AF at FinSBD:**

# **Sentence Boundary Detection of Financial Data with Domain Knowledge Enhancement and Bilingual Training**

Mingyu Wan<sup>1</sup> Rong Xiang<sup>2</sup> Emmanuele Chersoni<sup>1</sup> Natalia Klyueva<sup>1</sup> Kathleen Ahrens<sup>3</sup> Bin Miao<sup>4</sup> David Broadstock<sup>4</sup> Jian Kang<sup>4</sup> Amos Yung<sup>3</sup> Chu-Ren Huang<sup>1</sup> <sup>1</sup>Department of Chinese and Bilingual Studies, The Hong Kong Polytechnic University

<sup>2</sup>Department of Computing, The Hong Kong Polytechnic University <sup>3</sup>Department of English, The Hong Kong Polytechnic University <sup>4</sup>School of Accounting and Finance, The Hong Kong Polytechnic University

12 August 2019

6:

8:

11:

12: end for

15: return keyword\_dict



**RG** Research Centre for **RCE** Professional Communication in English

### **Task Overview**

## **Domain Knowledge Enhancement**

• Purpose:	
------------	--

to detect sentence boundaries (the begins and ends) from noisy texts of financial prospectuses.

Challenges:

documents encoded in non-text formats, such as tables, images, etc. conventional punctuation does not necessarily mark sentence boundary.

Algorithm 1 Keyword Extraction	Algorithm 2 Rule-based validation
nput: dataset	Input: dataset, raw_pred, keyword_dict
Output: keyword_dict	Output: updated_prediction
1: for i in len(dataset) do	1: for keyword in keyword_dict.key do
2: $curword = dataset[i].$	2: for i in len(dataset)-len(keyword) do
3: $nxtword = dataset[i+1].$	3: if word sequence match keyword then

#### **Our Approach:**

Machine Learning: a supervised task of three-class predictions + Feature Engineering: domain-specific and document-specific observation + Rule-based Validation: also know as domain knowledge enhancement.

#### **Performance:**

- English: 0.835 (avg. F-score of BS and ES)
- French: 0.86 (avg. F-score of BS and ES)

## **Feature Engineering**

A fused feature set of 24 dimensions, including:

#### • Two sets of punctuation:

- PUNC SET1 = [`.']
- 'e', '\$', '£', "'', '©', '®']
- **Initially capitalized words**, *e.g.* ". Enter END The BEGIN sales";
- Acronyms, *e.g.*"UBS\_BEGIN" "co" or "kiid";
- **Digital numbers**, ambiguity resolution as in "10.3";
- Letters or Roman numbers, e.g. e.g. A-Z, a-z or I, II, ..., XII;
- **POS tags**, structural cues.

- if "END" in curword then update raw\_pred with NO\_BOUNDARY if "BEGIN" in nxtword then end if  $\operatorname{continue}$ end for else7: end for add keyword to keyword\_dict. 8: return raw\_pred update frequency in keyword\_dict. end if end if
  - The two steps can correct false positive predictions caused by non-text sections, such as the titles, dates, tables, figures, etc., which fail to fall into the traditional category of sentence boundary.

### **Results and Comparison**

13: refine keyword\_dict with *length threshold*.

14: refine keyword\_dict with *frequency threshold*.



- RFC outperforms NN in <sup>3</sup>/<sub>4</sub> tasks
- NN outperforms RFC in the English Test set
- Bilingual training consistently outperforms monolingual training with 1-2 % F gain.

#### **Classifiers**



#### **2. Neural Network**

epoch = 5density = 300 density = 100

$Features \setminus F1$	$BS^a$	$\mathrm{ES}^{\mathrm{b}}$	Mean	$\delta(\%)^{\mathbf{c}}$
Punc1	0.70	0.82	0.76	baseline
+Punc2	0.71	0.83	0.77	$1 \uparrow$
+Cap	0.72	0.86	0.79	$2\uparrow$
+Acro	0.73	0.87	0.80	$1 \uparrow$
$+\mathrm{Dig}$	0.73	0.89	0.81	$1 \uparrow$
+Lett	0.74	0.90	0.82	$1\uparrow$
+POS	0.78	0.92	0.85	$3\uparrow$
+ Enter/All	0.80	0.92	0.86	$1\uparrow$

<sup>a</sup> Beginning boundaries

<sup>b</sup> Ending boundaries

<sup>c</sup> F1 improvement in percentage

Table 1: Performance of feature mining in the English Dev set with RFC

• **POS** vector is the most useful feature.

#### **Adaptation to Test Set**

Test Set	F1-Mean	Rank
Dev_en_rfc1	0.875	
Test_en_rfc1	$0.78^{*}$	$16^{\diamond}$
$Test_en_rfc1_adapted$	$0.835^{\star}$	$10^*$
Dev_fr_rfc1	0.85	
Test_fr_rfc1	$0.86^{*}$	8\$
Test_fr_rfc1_adapted	$0.86^{\star}$	8*

$375$ $78^*$ $16^\diamond$ $335^*$ $10^*$ $35^*$ $10^*$ $85$ $86^*$ $8^\diamond$ $86^*$ $8^\diamond$ $86^*$ $8^*$ $86^*$ $8^*$ French dataset with			
$78^*$ $16^\diamond$ • Our approach $335^*$ $10^*$ • Our significantly be $85$ • Our system is rol $86^*$ $8^\diamond$ • Our system is rol $86^*$ $8^*$ • French dataset with	375		
$335^{\star}$ $10^{*}$ significantlybe adaptation to the Er $85$ adaptation to the Er $86^{*}$ $8^{\diamond}$ • Our system is rol $86^{\star}$ $8^{*}$ French dataset wit	78*	$16^{\diamond}$	• Our approach
$85$ $-$ adaptation to the Er $85$ $ 0$ ur system is rol $86^*$ $8^*$ French dataset wit	$335^{\star}$	$10^*$	significantly be
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	or		adaptation to the Er
$86^*$ $8^\diamond$ • Our system is rol $86^*$ $8^*$ French dataset wit	80		
86* 8* French dataset wit	$86^{*}$	$8^{\diamond}$	• Our system is rol
	86*	8*	French dataset wit

\* Result without adaptation to the test set <sup>\*</sup> Rank without adaptation to the test set \* Result with adaptation to the test set

works with etter nglish test set.

bust for the th or without adaptation to the test set.

• Knowledge enhancement is **significantly** helpful.

Table 2: Performance of RFC in the French Dev and Test

	F1	$\rm NO^a$	$\rm YES^b$	$\delta(\%)$
Dev	BS	0.83	0.83	0
	$\mathbf{ES}$	0.86	0.87	$1\uparrow$
	$\operatorname{Mean}$	0.845	0.85	$0.5\uparrow$
Test	BS	0.81	0.84	$3\uparrow$
	$\mathbf{ES}$	0.88	0.88	0
	Mean	0.845	0.86	$1.5 \uparrow$

<sup>a</sup> Without keyword validation

<sup>b</sup> With keyword validation

sets in terms of keyword validation



\* Rank with adaptation to the test set

Table 3: Performance of RFC w.r.t. knowledge adaptation

#### Conclusions

- NN does not show advantage over statistical classifiers, but it outperforms RFC in terms of unknown features;
- Syntactic information may be helpful for sentence detection;
- Our system is highly effective with minimal training data;
- Future studies will further verify the effectiveness of this feature design and knowledge enrichment.

Contact Person: WAN, Mingyu Clara, Email: mingyu.wan@polyu.edu.hk. Website: https://www.researchgate.net/profile/Mingyu\_Wan/.

Acknowledgements: This work is partially supported by the GRF grant (No. 15608618) and the PDF project (No. 4-ZZKE) at PolyU HK.