Technical report for Events' labels matching computation

Karn Yongsiriwit Telecom SudParis, UMR 5157 CNRS Samovar, France Email: karn.yongsiriwit@telecom-sudparis.eu

I. Events' labels matching

Function $sim.\lambda_p$ computes the similarity of two labels. We use Stanford Part-of-Speech (POS) [1], [2] for stemming string and removing function words. We modify the bag-ofwords similarity with label pruning technique [3] to prunes words from the longer label then measure the similarity of two labels based on pruned words.

The similar label matching between label l_1 and l_2 is defined as follows:

$$\frac{\sin\lambda_{p}(l_{1}, l_{2}) =}{\sum_{i=1}^{|\omega^{1}|} \max_{j=1}^{|\omega^{2}|} (\sin w(pr_{1}^{i}, pr_{2}^{j})) + \sum_{j=1}^{|\omega^{2}|} \max_{i=1}^{|\omega^{2}|} (\sin w(pr_{1}^{i}, pr_{2}^{j}))}{2 \times \min(|\omega^{1}|, |\omega^{2}|)}$$
(1)

Let l_1 and l_2 be the labels of the events e_1 and e_2 , and $\omega^1 = tok(\lambda_1(a_1)), \ \omega^2 = tok(\lambda_2(a_2))$ are tokenized lists of words contained in the labels by POS technique. Further, $pr_1 = pru(\omega^1, \omega^2)$ and $pr_2 = pru(\omega^2, \omega^1)$ are the pruned list of words. Function *sim.w* computes the similarity between two words using existing approaches such as Lin metric ontology matching technique [4]. We define a threshold β_l for $sim.\lambda_p$. Two labels are considered to be functionally similar iff $sim.\lambda_p(l_1, l_2) \ge \beta_l$. Let $pru : P(W) \times P(W) \rightarrow P(W)$ be a generic function. It returns a set of words extracted from its input. $pru(\lambda_p(e_1), \lambda_p(e_2))$ is ω^1 iff $|\omega^1| \le |\omega^2|$, or a subset of ω^1 of size $|\omega^2|$ otherwise. The similarity scores of all word pairs, as well as the maximum score for each word are calculated in $|\omega^1|$. *pru* returns the $|\omega^2|$ -top-scoring words from ω^1 .

In order to illustrate our approach, we present the similarities between labels of the events in the 1^{st} – *zone* of the neighborhood context of *C* and *C2* using $(sim.\lambda_p)$ (see our motivating example). We define the threshold as $\beta_l = 0.5$. The results after descending sort are shown Table II. The details about the example computation are shown in Table I. Note that the word pairs with the similarity value underlined indicates the maximum similarity value among other word pairs.

We create the matching between the most similar events by the ranking the similarity values. The result are event B2 is matched to event B with $sim.\lambda_p(B2, B) = 1.000$ and event D2 is matching to event D with $sim.\lambda_p(D2, D) = 0.543$.

$$N_c^1(C, C2) = \{(B2, B, 1.000), (D2, D, 0.543)\}$$

					-	2,577										
	_			A							В					
		SI	im.w	S	end	e-	mail			si	m.w	ch	eck	cre	dit	
		process		0.	0.000		0.000			process		0.099		0.0	96	
B		che	eck	0.	000	0	.000			check		1.000		0.0	85	
	2 cre rec		dit	0.	0.000 0.576		.000	B2	2	credit request		0.335 0.369		1.0	00	
			uest	0.			.000							0.4	69	
		response		0.	0.000		0.000			response		0.000		0.0	00	
(a) $sim.\lambda_p(l_{B2}, l_A) = 0.288$ (b) $sim.\lambda_p(l_{B2}, A_{label}(B)) = 1.000$)0					
Г					D				E							
	sim.w		cł	check		system			sim.v		w	w acc		1		
		process		0.	0.099		0.000	Г		proces		ss	0.	000	1	
		check		1.	1.000		0.000	0 0 0 3			check		0.	000		
<i>B2</i>		credit		$\overline{0}$	0.335		0.000			B2 credit reques respon		0.		455		
		request		0.	0.369		0.000					st	0.000			
		response		0.	0.000		0.073					nse 0.		000		
(c) $sim \lambda_p(l_{B2}, l_D) = 0.537$ (d) $sim \lambda_p(l_{B2}, l_E) = 0.455$												1				
					I I	4°	4									
r	sin		sım	W	v reject							F		F		
		proces		SS	0.0	000	0			sim.w		send		e-mail		
			check	2	0.0	000				ch	neck	0.0	00	0.0	000	
	Ŀ	\$2	credit		$\frac{0.4}{0.4}$	39		D_{i}	2	pa	iper	0.0	00	0.0	000	
		reques		st	st 0.00					archive		0.000		0.000		
response			nse	$\frac{1}{1}$ 0.000				(1) $sim_{A_p}(l_{D2}, l_A) = 0.000$								
	(e) sin	$n.\Lambda_p(l_B)$	$2, l_F$) = 0.	439										
	Г				В									D		
S		si	m.w	che	check		credit			sim.w		check		system		
D2		check		1.0	1.000		0.335			check		1.000		0.	000	
		paper		0.0	0.000		0.000		2	paper		0.000		0.	085	
		archive		0.000		0.0	0.000			archive		0.000		0.	000	
(g) $sim.\lambda_p(l_{D2}, l_B) = 0.668$ (h) $sim.\lambda_p(l_{D2}, l_D) = 0.543$																
	E				1						I	7				
		sim.		w	v acce		pt				sim	.w	rej	ect		
			chec	heck		0.000					check		0.0	00		
	1	D2	pape	paper		0.000				D2	pape	er	0.0	00		
		arc		ve	0.0						arch	ive	0.0	00		
	(i) $sim.\lambda_p(l_{D2}, l_E) = 0.000$								(j) $sim.\lambda_p(l_{D2}, l_F) = 0.000$							

Table I. Similarity computation of all possible neighbors pairs between events *C* and *C*2 in the 1^{st} – *zone* neighborhood context

Table II. Possible pairs ranked by $sim \lambda_p$.

No.	e_1	e_2	$sim.\lambda_p(e_1,e_2)$
1	<u>B2</u>	B	1.000
2	D2	В	0.668
3	<u>D2</u>	D	0.543
4	D2	В	0.537

Using $sim \lambda_p$, the log-based neighborhood context matching is described in Equation 2. We weight the similar pair

of events based on their labels similarity.

$$M_{sim,\lambda_p}^k(a,b) = \frac{\sum_{i=1}^{z} sim.\lambda_p(l_{r_{a_i}}, l_{r_{b_i}}) \times a_i \times sim.\lambda_p(l_{t_{a_i}}, l_{t_{b_i}}) \times b_i}{|\overrightarrow{e_c(a)}| \times |\overrightarrow{e_c(b)}|}$$
(2)

Using the label similarities from Table II, we have $N_c^1(C, C2) = \{(B2, B, 1.000), (D2, D, 0.543)\}$. So, $e_c(C) = (w(B, C), w(C, D)) = (54, 46), e_c(C2) = ((B2, C2), (C2, D2)) = (70, 70)$ and their matching in the 1st – zone is:

$$\begin{split} M^1_{sim,\lambda_p}(C,C2) &= \frac{1.000\times54\times1.000\times70+1.000\times46\times0.543\times70}{\sqrt{4^2+6^2+54^2+46^2+42^2+13^2+35^2}\times\sqrt{70^2+70^2}}\\ &= 0.615 \end{split}$$

References

- D. Jurafsky and J. H. Martin, Speech and Language Processing (2nd Edition) (Prentice Hall Series in Artificial Intelligence), 2nd ed. Prentice Hall, 2008.
- [2] K. Toutanova and C. D. Manning, "Enriching the knowledge sources used in a maximum entropy part-of-speech tagger," in *In Proceedings of the Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora (EMNLP/VLC-2000)*, 2000, pp. 63–70. [Online]. Available: http://nlp.stanford.edu/~manning/papers/emnlp2000. pdf
- [3] C. Klinkmüller, I. Weber, J. Mendling, H. Leopold, and A. Ludwig, "Increasing recall of process model matching by improved activity label matching," in *BPM*, 2013, pp. 211–218.
- [4] D. Lin, "An information-theoretic definition of similarity," in *ICML*, 1998, pp. 296–304.