



Detekce polarity textu s využitím mezijazyčné transformace

Jakub Šmíd¹

1 Úvod

Detekce polarity textu je úloha zpracování přirozeného jazyka, jejímž cílem je klasifikace vstupního textu dle jeho polarity. Pro tuto úlohu jsou potřeba anotovaná data, kterých může být v některých jazyčích nedostatek. Proto se hledají metody, které umožňují použití anotovaných dat z jednoho jazyka v jazyce druhém. Jednou z možností jsou mezijazyčné transformace.

Cílem této práce bylo vytvořit program pro detekci polarity českého textu a pomocí experimentů ověřit možnosti využití dat z jiného jazyka (konkrétně angličtiny) s použitím mezijazyčných transformací a zhodnotit jejich vliv na úspěšnost detekce polarity.

2 Lineární mezijazyčné transformace

Lineární mezijazyčné transformace transformují předtrénované jednojazyčné slovní vektory do společného prostoru pomocí lineárního zobrazení a dvojjazyčných slovníků, jak shrnul např. Ruder et al. (2019). Lineární zobrazení umožňuje transformaci mezi dvěma vektorovými prostory prostřednictvím affinních transformací (např. rotace, posun a zrcadlení). Tyto přístupy vychází z pozorování, která ukazují, že rozmístění vektorů slov zdrojového jazyka je po provedení vhodné lineární transformace geometricky velmi podobné rozmístění vektorů slov jejich překladů. Cílem je nalézt transformační matici \mathbf{W} , která umožní transformovat vektorový prostor zdrojového jazyka s do vektorového prostoru cílového jazyka t . Vynásobením vektoru \mathbf{x}^s slova v původním prostoru transformační maticí \mathbf{W} je získán příslušný vektor \mathbf{x}^t v prostoru cílovém (viz rovnice 1). V této práci byla pro zisk matice \mathbf{W} použita metoda nejmenších čtverců (MSE), kterou navrhl Mikolov et al. (2013), a ortogonální metoda, kterou popsal Xing et al. (2015). Obě metody ponechávají vždy jeden prostor nezměněný a druhý na něj mapují.

$$\mathbf{x}^t = \mathbf{W}\mathbf{x}^s \quad (1)$$

3 Experimenty a výsledky

Experimenty byly provedeny s použitím tří datových sad obsahujících filmové recenze – jedné české (ČSFD) a dvou anglických (SST a IMDb). Pro detekci polarity textu byly navrženy a implementovány dva modely neuronových sítí – konvoluční neuronová síť (CNN) a rekurrentní neuronová síť (konkrétně obousměrná LSTM) – v kombinaci s natrénovanými slovními vektory *fastText*. Modely byly nejprve natrénovány a vyhodnoceny pouze na českých datech. Dále byly s využitím lineárních mezijazyčných transformací natrénovány na anglických datech a vyhodnoceny na češtině. Transformace proběhla z angličtiny do češtiny i naopak. Pro tvorbu transformační matice bylo vybráno 5 000 nebo 20 000 slov. Tabulka 1 ukazuje dosažené výsledky porovnané s jednojazyčnými experimenty.

¹ student bakalářského studijního programu Inženýrská informatika, obor Informatika, e-mail: biba10@students.zcu.cz

BiLSTM							
Počet tříd	CS	Dataset	Metoda	5 000		20 000	
				EN→CS	CS→EN	EN→CS	CS→EN
3	$84,92 \pm 0,30$	SST	MSE	$46,06 \pm 0,75$	$48,91 \pm 2,59$	$45,86 \pm 0,92$	$47,53 \pm 3,72$
		SST (fráze)	ortog.	$47,25 \pm 3,10$	$48,14 \pm 1,63$	$49,15 \pm 2,10$	$49,71 \pm 3,03$
		SST	MSE	$46,83 \pm 0,60$	$46,52 \pm 2,89$	$46,42 \pm 0,98$	$47,59 \pm 4,29$
		SST (fráze)	ortog.	$49,53 \pm 2,91$	$48,09 \pm 1,39$	$50,09 \pm 2,35$	$49,10 \pm 2,27$
2	$94,29 \pm 0,29$	SST	MSE	$82,82 \pm 0,71$	$84,91 \pm 0,71$	$83,37 \pm 2,49$	$82,54 \pm 2,49$
		SST (fráze)	ortog.	$79,61 \pm 2,69$	$82,41 \pm 2,09$	$81,05 \pm 1,88$	$82,47 \pm 2,49$
		SST	MSE	$82,60 \pm 2,69$	$82,57 \pm 2,71$	$83,33 \pm 0,93$	$84,32 \pm 0,93$
		SST (fráze)	ortog.	$82,42 \pm 1,94$	$83,57 \pm 1,25$	$82,41 \pm 2,46$	$83,14 \pm 2,64$
		IMDb	MSE	$86,20 \pm 1,84$	$87,48 \pm 1,84$	$85,17 \pm 0,89$	$87,98 \pm 0,89$
		IMDb	ortog.	$86,89 \pm 0,64$	$88,25 \pm 1,14$	$87,59 \pm 0,48$	$88,57 \pm 0,36$
CNN							
Počet tříd	CS	Dataset	Metoda	5 000		20 000	
				EN→CS	CS→EN	EN→CS	CS→EN
3	$83,18 \pm 0,14$	SST	MSE	$46,83 \pm 0,09$	$58,87 \pm 1,10$	$47,12 \pm 0,10$	$57,60 \pm 1,61$
		SST (fráze)	ortog.	$50,32 \pm 1,79$	$50,58 \pm 1,10$	$51,53 \pm 1,42$	$50,38 \pm 1,08$
		SST	MSE	$46,81 \pm 0,10$	$57,71 \pm 1,21$	$47,18 \pm 0,09$	$57,63 \pm 0,92$
		SST (fráze)	ortog.	$50,07 \pm 1,27$	$51,50 \pm 0,86$	$51,46 \pm 0,63$	$50,42 \pm 0,75$
2	$93,86 \pm 0,12$	SST	MSE	$86,46 \pm 0,19$	$86,66 \pm 0,46$	$86,99 \pm 0,12$	$86,19 \pm 0,65$
		SST (fráze)	ortog.	$84,25 \pm 1,44$	$85,46 \pm 0,80$	$86,18 \pm 0,30$	$86,82 \pm 0,31$
		SST	MSE	$86,31 \pm 0,37$	$86,47 \pm 0,22$	$86,86 \pm 0,17$	$86,83 \pm 0,40$
		SST (fráze)	ortog.	$84,80 \pm 1,17$	$86,36 \pm 0,40$	$85,93 \pm 0,44$	$86,41 \pm 0,40$
		IMDb	MSE	$88,05 \pm 0,20$	$86,64 \pm 1,35$	$88,36 \pm 0,10$	$88,11 \pm 1,72$
		IMDb	ortog.	$87,18 \pm 0,34$	$87,81 \pm 0,37$	$88,16 \pm 0,17$	$88,99 \pm 0,34$

Tabulka 1: Makro F-míra v procentech pro modely BiLSTM a CNN se slovními vektory *fastText* (natrénované na datech z IMDb a SST datasetů pro angličtinu, z ČSFD datasetu pro češtinu) při použití mezijazyčných transformací s transformační maticí získanou z 5 000 či 20 000 slov metodou nejmenších čtverců nebo ortogonálně a transformací z angličtiny do češtiny (**EN→CS**) i naopak (**CS→EN**). Trénováno jen na datech anglických (sloupec **Dataset**) a vyhodnoceno na českých testovacích. Sloupec **CS** slouží k porovnání s jednojazyčnými výsledky. Pro každý model, počet tříd a směr transformace je nejlepší výsledek zvýrazněný **tučně**, při překrývání intervalu spolehlivosti je navíc **podtržený**.

4 Závěr

Porovnání jednojazyčných experimentů s experimenty s využitím mezijazyčné transformace ukázalo, že s dostatečným množstvím výhradně anglických dat lze dosáhnout velmi dobrých výsledků, které jsou jen o 5 až 6 % horší v porovnání s modely trénovanými jen na českých datech.

Literatura

- Mikolov, T., Le, Q. V., a Sutskever, I. (2013) Exploiting Similarities among Languages for Machine Translation. *CoRR*. Dostupné z: <http://arxiv.org/abs/1309.4168>.
- Ruder, S., Vulic, I., a Søgaard, A. (2019) A survey of cross-lingual word embedding models. *Journal of Artificial Intelligence Research*, s. 569–631.
- Xing, C. et al. (2015) Normalized word embedding and orthogonal transform for bilingual word translation. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, s. 1006-1011.