

Akcelerácia neurónovej siete pre jazykové modelovanie

Dominik Labaš

Abstrakt

Tento dokument sa zaoberá využitím neurónových sietí v jazykovom modelovaní. Naším cieľom bude takýto model vytvoriť a následne implementovať techniky pre urýchlenie výpočtu pravdepodobností. Implementovaná neurónová sieť sa skladá z jednej skrytej vrstvy so zvolenou základnou veľkosťou 150 neurónov. Pri urýchľovaní neurónovej siete sme pri zachovaní presnosti modelu dosiahli 3-násobné zrýchlenie pomocou zmeny aktivačnej funkcie. Predpočítaním matíc pre výpočet skrytej vrstvy sme výpočet skrytej vrstvy urýchlili 2.5krát. Ukladanie videnej histórie model spomalilo o viac ako 40%.

Kľúčové slová: Jazykové modelovanie — Dopredná neurónová sieť — Akcelerácia neurónovej siete

Priložené materiály: [Downloadable Code](#)

xlabas00@stud.fit.vutbr.cz, Faculty of Information Technology, Brno University of Technology

1. Úvod

Neurónová sieť je výpočtový model navrhnutý tak, aby sa tréňovaním naučil vykonávať určitú funkciu. V našom prípade neurónovú sieť využívame pri jazykovom modelovaní, čo je proces rozdelenia pravdepodobnosti nad reťazcami slov, kde samotný jazykový model je funkcia, ktorá tieto pravdepodobnosti dokáže vypočítať ako je znázornené v obrázku (1). Tento spôsob využitia neurónovej siete pri jazykovom modelovaní pre počítanie pravdepodobnosti je používaný od roku 2003 [1].

Pri jazykovom modelovaní, kde je potreba sa učiť komplexné funkcie vzniká takzvaná *Curse of Dimensionality*. Tento problém je dôsledok vysokého počtu rôznych slov nachádzajúcich sa v tréningových dátach. Tieto slová je potrebné od seba rozlíšiť a tréningový algoritmus vyžaduje príklad pre každú dôležitú kombináciu. Neurónové siete sa v jazykovom modelovaní využívajú pre ich schopnosť učiť sa reprezentácie slov, ktoré pomáhajú proti tomuto efektu tým, že pre každé slovo je vytvorený vektor, ktorý toto slovo popisuje. Tieto vektory sa potom nachádzajú v spoločnom priestore, kde je možné popísať podobnosť a vzťahy medzi týmito

slovami.

Skutočná práca jazykového modelu začína až po jeho natrénovaní, kde pri jeho využívaní po ňom vyžadujeme čo najväčšiu rýchlosť. Problému s rýchlosťou doprednej neurónovej siete sa venovali v článku [2], kde navrhli techniky pre urýchlenie výpočtu pre jazykový model. V tomto článku sa pokúsime zreplikovať tri techniky a aplikovať na náš jazykový model.

2. Neurónová sieť ako pravdepodobnostný jazykový model

Jazykové modelovanie pracuje s reťazcami slov v tvare $(w_1, w_2, \dots, w_{N-1}, w_N)$ zkrátene (w_1^N) , pričom jeho cieľ je vypočítať pravdepodobnosť $P(w_1^N)$. Túto pravdepodobnosť je možné faktorizovať na súčin podmienených pravdepodobností ako:

$$P(w_1^T) = P(w_1) \cdot P(w_2|w_1) \cdot \dots \cdot P(w_T|w_1, w_2, \dots, w_{T-1})$$

$$P(w_1^T) = \prod_{t=1}^T P(w_t|w_1^{t-1}) \quad (1)$$

Pravdepodobnosť $P(w_2|w_1)$ vyjadruje pravdepodobnosť výskytu slova w_2 ak predošlé slovo je práve w_1 .

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40 Úlohou jazykového modelovania je teda naučiť sa
41 funkciu rozdelenia pravdepodobností f v tvare:

$$f(w_N, w_{N-1}, \dots, w_1) = P(w_N | w_1^{N-1}) \quad (2)$$

42 Táto úloha bola riešená napríklad za pomoci n -
43 gram modelov [3]. Tie využívajú Markovské pro-
44 cesy n -tého rádu, teda slovo vzdialené n a viac miest
45 od nášho aktuálneho slova nemá na pravdepodobnosť
46 žiaden vplyv:

$$P(w_t | w_1^{t-1}) = P(w_t | w_{t-n+1}^{t-1}) \quad (3)$$

47 V našom prípade využijeme neurónovú sieť, ktorá
48 bude mať za úlohu naučiť sa funkciu pre počítanie
49 pravdepodobnosti:

$$f(w_t, w_{t-1}, \dots, w_{t-n+1}) = P(w_t | w_{t-n+1}^{t-1}) \quad (4)$$

50 2.1 Využitie neurónovej siete

51 Neurónová sieť sa skladá z niekoľkých vrstiev: vstup-
52 nej, výstupnej a jednej alebo viac takzvaných skrytých
53 vrstiev.

54 Na vstupe modelu sa nachádzajú indexy na slová
55 zo slovníka V , kde slovník V je konečná množina slov
56 vyskytujúcich sa vo vstupných dátach. Tieto indexy
57 model mapuje na vektory vlastností $\mathbf{C}(w_i)$, kde $\forall w_i \in$
58 V . Vektory vlastností majú veľkosť m a sú uložené
59 v matic \mathbf{C} o veľkosti $|V| \times m$. Tieto reprezentácie
60 slov sa model učí v priebehu tréningu súčasne s
61 pravdepodobnostnou funkciou (4).

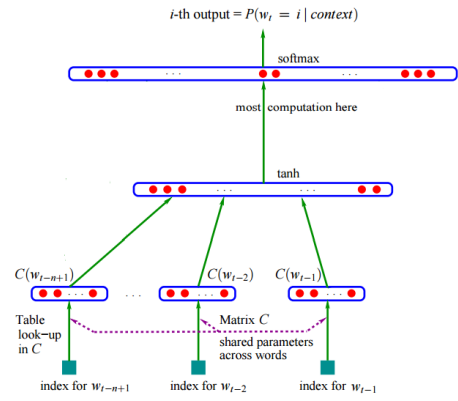
62 Obsah vstupnej vrstvy neurónovej siete je zrežaze-
63 ná sekvencia vektorov $\mathbf{c} = (\mathbf{C}(w_{t-n+1}), \dots, \mathbf{C}(w_{t-1}))$.
64 Túto sekvenciu ďalej model mapuje na rozdelenie
65 pravdepodobnosti pre nasledujúce slovo w_t nad slo-
66 vníkom V . Výsledkom je vektor pravdepodobností
67 o veľkosti $|V|$ vyjadrujúci pravdepodobnosti (4) pre
68 všetky slová slovníka V .

$$f(w_t, w_{t-1}, \dots, w_{t-n+1}) = g(w_t, \mathbf{C}(w_{t-n+1}), \dots, \mathbf{C}(w_{t-1})) \quad (5)$$

69 Funkcia g predstavuje neurónovú sieť, ktorá počíta
70 pravdepodobnosť slova w_t .

71 2.2 Počítanie pravdepodobností pomocou 72 neurónovej siete

73 V našom prípade budeme využívať doprednú neuróno-
74 vú sieť s jednou skrytou vrstvou, ktorá počíta nenor-
75 malizované logaritmické pravdepodobnosti. Po nama-
76 povaní indexov slov na ich reprezentácie a vložení
77 na vstupnú vrstvu siete sa na jednotlivé hodnoty ap-
78 likujú váhy, ktoré predstavujú prechod z jednej vrstvy
79 do druhej.



Obrázok 1. Jazykový model, vstupmi sú indexy na slová zo slovníka V . Vektory vlastností $\mathbf{C}(w_i)$ sú zrežazené na vstupnej vrstve neurónovej siete predstavovanej funkciou g z (5). Výstupom je vektor rozdelenia pravdepodobností nad V . Prevzaté z [1].

$$\mathbf{h} = \tanh(\mathbf{d} + \mathbf{H}\mathbf{c}) \quad (6)$$

Vektor \mathbf{h} predstavuje vektor s výstupom skrytej 80
vrstvy. Vektor \mathbf{d} obsahuje biase pre skrytú vrstvu, 81
matrica \mathbf{H} (o veľkosti $|\mathbf{h}| \times |\mathbf{c}|$) obsahuje váhy pre prechod 82
z vstupnej vrstvy do skrytej a vektor \mathbf{c} obsahuje vs- 83
tupnú vrstvu, kde sú zrežazené sekvencie vektorov 84
($\mathbf{C}(w_{t-n+1}), \dots, \mathbf{C}(w_{t-1})$). Na každý element skrytej 85
vrstvy je potom aplikovaná aktivačná funkcia, ako ak- 86
tivačná funkcia sa používa napríklad logistická funkcia 87
alebo hyperbolický tangens. 88

Následne sú na skrytú vrstvu aplikované váhy pre 89
prechod do výstupnej vrstvy: 90

$$\mathbf{y} = \mathbf{b} + \mathbf{U}\mathbf{h} \quad (7)$$

Vektor \mathbf{y} (o veľkosti $|V|$) predstavuje nenormali- 91
zované logaritmické pravdepodobnosti, vektor biasov 92
 \mathbf{b} pre výstupnú vrstvu a maticou váh \mathbf{U} (o veľkosti 93
 $|V| \times |\mathbf{h}|$) pre prechod zo skrytej do výstupnej vrstvy a 94
 \mathbf{h} predstavuje vektor so skryteou vrstvou. 95

Množina ($\mathbf{b}, \mathbf{d}, \mathbf{U}, \mathbf{H}, \mathbf{C}$) predstavuje všetky voľné 96
parametre, ktoré sa neurónová sieť učí počas tréningu. 97

Na výstupnú vrstvu sa aplikuje softmax, ktorý za- 98
istí že súčet všetkých pravdepodobností je 1: 99

$$P(w_t | w_{t-n+1}^{t-1}) = \text{Softmax}(\mathbf{y}) = \frac{e^{y_{w_t}}}{\sum_{i=1}^{|V|} e^{y_{w_i}}} \quad (8)$$

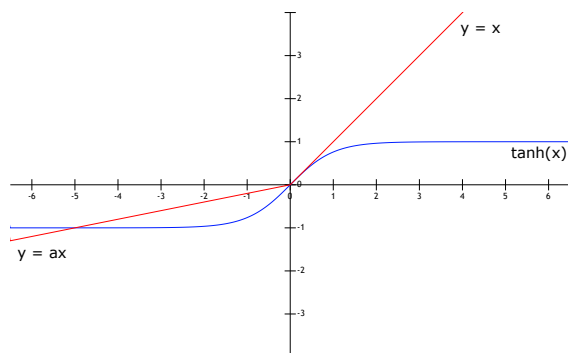
3. Optimalizácia modelu

Pri optimalizácii sa venujeme štyrom technikám po- 101
písaným v [2]. Zmene aktivačnej funkcie za rýchlejšiu 102
šiu, zapamätanie si hodnôt z predošlých priechodov 103
cez model, predpočítaním matíc pre rýchlejší výpočet 104
skrytej vrstvy a odstránenie normalizácie z modelu. 105

106 3.1 Zmena aktivačnej funkcie

107 V tejto časti budeme skúmať zrýchlenie pri nahradení
108 stávajúcej aktivačnej funkcie $\tanh()$ na rýchlejšiu funkciu
109 PReLU() (Parametric Rectified Linear Unit) [4], ktorá
110 počíta výstup skrytej vrstvy ako:

$$\text{PReLU}(x) = \max(0, x) + a \cdot \min(0, x) \quad (9)$$



Obrázok 2. Graf aktivačných funkcií PReLU(x) a tanh(x).

111 Táto zmena nahradí výpočet funkcie $\tanh()$ za
112 rýchlejší výber maxima, s použitím funkcie PReLU()
113 pribúda ďalší voľný parameter a , ktorý sa model učí v
114 priebehu tréningovania.

115 3.2 Ukladanie histórie

116 Pretože pri priechode cez vstupné dáta je šanca, že
117 nám príde rovnaká kombinácia vstupných indexov,
118 budeme si ukladať časť videnej histórie. V tomto
119 prípade ak nám prídu indexy, ktoré máme uložené,
120 cena vypočítania skrytej vrstvy je vynechaná a skrytá
121 vrstva je načítaná z pamäte. Ak čas T je čas potrebný
122 pre výpočet skrytej vrstvy, ukladanie histórie ušetrí p
123 $\cdot T$, kde p značí úspešnosť toho, že sa vstup zhoduje s
124 jedným zo záznamov histórie.

125 3.3 Predpočítanie matíc pre výpočet skrytej 126 vrstvy

127 Pretože po tréningovaní sa výsledky mapovania indexov
128 na reprezentácie slov a váhy pre prechod do skry-
129 tej vrstvy ďalej nemenia, je možné si skrytú vrstvu
130 "predpočítať". Toto docielime tým že si vytvoríme
131 matice M_j o veľkosti $|V| \times |\mathbf{h}|$ pre všetky slová zo
132 slovníku V a všetky možné pozície vo vektore c . Pri
133 počítaní skrytej vrstvy potom stačí sčítať vektory z
134 matíc pre zadané indexy a v modeli pokračovať od
135 aplikačnej funkcie.

$$\mathbf{h} = \tanh\left(\sum_{j=0}^{n-1} M_j[w_j]\right) \quad (10)$$

136 Výraz $M_j[w_j]$ značí indexáciu na pozíciu prvku
137 w_j . Biase pre skrytú vrstvu pripočítavame len v jednej

matici. Týmto spôsobom ušetríme čas potrebný pre 138
výpočet funkcie (6). 139

4. Implementácia 140

Jazykový model je implementovaný v jazyku Python 141
za použitia frameworku PyTorch [5]. 142

Základný model je navrhnutý ako dopredná neuró- 143
nová sieť, ktorá implementuje rovnicu (5). 144

Modely prechádzajú dvomi programami. Prvý na 145
trénovanie modelu, kde sa úspešnosť modelu pozoruje 146
pomocou vypočítavania straty. Druhý pre sledovanie 147
času pri priechode cez model využitím profileru z 148
knížnic PyTorch. 149

4.1 Zmena aktivačnej funkcie 150

Pri tomto kroku sme funkciu $\tanh()$ v programe za- 151
menili za funkciu PReLU() [4], ktorej implementácia 152
je zahrnutá v knižniciach PyTorchu. 153

4.2 Ukladanie histórie 154

Pri ukladaní histórie sa pri evaluačnom využití mod- 155
elu vytvárajú dve polia, kde v jednom sú uložené vs- 156
tupné indexy, ktoré boli nedávno modelom videné a v 157
druhom sú uložené hodnoty skrytej vrstvy po využití 158
aktivačnej funkcie. 159

4.3 Predpočítanie matíc pre výpočet skrytej vrstvy 160 161

V tomto bode sa do modelu vložia vypočítané matice 162
 M_j , do ktorých ukladáme vynásobené matice repre- 163
zentácií slov s váhami pre danú polohu vo vstupnom 164
vektore. Prvá matica sa líši od ostatných pridaním 165
biasov. 166

Pri ďalších evaluačných volaní modelu potom mo- 167
del le indexuje matice M_j , ktorých hodnoty spočíta 168
do skrytej vrstvy [6]. A model pokračuje od využitia 169
aktivačnej funkcie. 170

5. Experimenty 171

Pri tréningovaní modelov sa ich úspešnosť meria v dvoch 172
hodnotách loss a perplexita. Tieto hodnoty sa získa- 173
vajú priebežne z tréningových dát a po každej epoche z 174
validačných dátach. Po dokončení tréningu sa loss 175
a perplexita vyráta pre testovacie dáta. Naše modely 176
porovnávame na základe hodnôt loss a perplexity z 177
testovacích dát, model je presnejší čím je loss alebo 178
perplexita menšia. 179

Rýchlosti modelov sledujeme podľa času, ktorý 180
potrebujú na jeden priechod cez testovacie dáta. Na 181
začiatok sme vytvorili tri štandardné modely s rôznou 182
veľkosťou skrytej vrstvy. Program využíva validačné 183
dáta a po ich spracovaní vypíše najmenší, najväčší čas 184

185 a priemer zo všetkých časov. Hodnoty loss všetkých
 186 modelov mali podobne veľkú hodnotu a ich rozdiel
 187 nebol významný, preto môžeme povedať, že pres-
 188 nosť všetkých modelov je takmer zhodná.

Tabuľka 1. Porovnanie úspešností modelov. Základný model je označený ako $\tanh X$, kde X značí veľkosť skrytej vrstvy a reprezentácií slov.

Model	Perplexita	Loss
$\tanh 50$	173,23	5,15
$\tanh 150$	155,35	5,05
$\tanh 250$	159,39	5,07
PReLU	164,21	5,10
cache	158,01	5,06
projekcia	155,35	5,05

189

Tabuľka 2. Porovnanie rýchlostí základných modelov s rôznou veľkosťou.

Model	Min	Max	Avg
$\tanh 50$	4,72 ms	4,88 ms	4,75 ms
$\tanh 150$	5,55 ms	5,81 ms	5,60 ms
$\tanh 250$	6,27 ms	6,54 ms	6,30 ms

190

Tabuľka 3. Ukážka výstupu z profileru znázorňujúci časovú náročnosť jednotlivých krokov pre baseline model $\tanh 150$.

Funkcia	Čas
embedding	145,2 μs
narrow	4,7 μs
narrow	2,4 μs
cat	27,5 μs
view	3,0 μs
t	3,0 μs
expand	3,3 μs
addmm	77,8 μs
\tanh	147,7 μs
t	2,8 μs
expand	2,6 μs
addmm	1530,9 μs
logsoftmax	4002,0 μs

191 V tabuľke výstupu z profileru môžeme vidieť časové
 192 náročnosti jednotlivých úkonov modelu, kde jednozna-
 193 čne najdlhšie trvajúci je výpočet funkcie pre logmax,

ďalej lineárna transformácia pre výstupnú vrstvu a
 na koniec výpočet aktivačnej funkcie, mapovanie in-
 dexov na reprezentácie slov (embeddings) a lineárna
 transformácia pre skrytú vrstvu.

5.1 Vstupné dáta modelu

Ako vstup pre trénovanie a testovanie modelu sú využívané dáta *Penn Tree Bank*, jedná sa o text s dĺžkou asi 1M slov, počet rôznych slov v texte je 10k. Dáta sú rozdelené na 3 časti: trénovacie, testovacie a validačné. Dĺžka trénovacieho textu je zhruba 10-krát väčšia ako testovacie alebo validačné dáta.

5.2 Zmena aktivačnej funkcie

Pri použití funkcie PReLU() experiment dokázal jednoznačné zrýchlenie pre výpočet aktivačnej funkcie. Výsledkom bolo takmer 3-násobné zrýchlenie výpočtu.

Tabuľka 4. Rýchlosti výpočtu aktivačnej funkcie.

Aktivačná f	Avg
$\tanh()$	146,0 μs
PReLU()	55,7 μs

Tabuľka 5. Celková rýchlosť modelov.

Model	Min	Max	Avg
$\tanh()$	5,55 ms	5,81 ms	5,60 ms
PReLU()	5,37 ms	5,41 ms	5,39 ms

Zrýchlenie celého modelu bolo 3.75%.

5.3 Ukladanie histórie

Pri implementácii tejto techniky sa vyskytol problém, kde pôvodná implementácia umožňuje prácu s 2-rozmerným vstupom, kde potrebný počet zhodných indexov niekoľkonásobne stúpa a tým klesá aj úspešnosť histórie. Preto sme implementáciu pozmenili tak, aby na vstupe modelu bolo vždy len 1-rozmerné pole. Pri využití ukladania videnej histórie pracujeme s úspešnosťou histórie, tá nám značí koľko krát sa v programe hľadaný vstup nachádzal v aktuálne uloženej histórie. Pre porovnanie so štandardným modelom sme pozmenili vstup modelu len na 1-rozmerné pole. Pri experimentovaní sme vytvorili modely s rôznou dĺžkou histórie.

Zo získaných hodnôt sme zistili, že model je asi o 20% rýchlejší ak narazí na hodnotu, ktorú má uloženú v histórii. No ak danú hodnotu v histórii uloženú nemá model je spomalený o viac ako 50%. Pri porovnaní

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

Tabuľka 6. Porovnanie rýchlostí modelov využívajúce ukládanie histórie (cache) s rôznou dĺžkou histórie a základného modelu.

Model	Min	Max	Avg
baseline	0,68 ms	1,04 ms	0,72 ms
cache 150	0,54 ms	1,61 ms	1,01 ms
cache 300	0,54 ms	1,73 ms	1,04 ms
cache 500	0,55 ms	1,96 ms	1,17 ms

230 priemerného času sme dokázali že naše riešenie model
231 nezrýchluje ale naopak zpomaľuje o viac ako 40%.

Tabuľka 7. Úspešnosť histórie podľa jej dĺžky.

Model f	Úspešnosť
cache 150	11,3 %
cache 300	15,2 %
cache 500	18,2 %

232
233 Hoci úspešnosť histórie so zväčšujúcou sa dĺžkou
234 stúpala, stúpala aj čas hľadania v histórii pri neúspeš-
235 ných pokusoch tak, že so zväčšujúcou sa dĺžkou histórie
236 sa model viac a viac spomaľoval.

237 5.4 Predpočítanie matíc pre výpočet skrytej 238 vrstvy

239 Pri využití predpočítaných matíc bola skrytá vrstva
240 vypočítaná 2,5krát rýchlejšie ako so základným mode-
241 lom.

Tabuľka 8. Rýchlosť vypočítania skrytej vrstvy modelu.

Model	Čas na získanie skrytej vrstvy
baseline	254,8 μ s
projekcia	100,2 μ s

Tabuľka 9. Celková rýchlosť modelu podľa typu získavania skrytej vrstvy.

Model	Min	Max	Avg
baseline	5,55 ms	5,81 ms	5,60 ms
projekcia	5,37 ms	5,39 ms	5,38 ms

242
243 Využitie predpočítaných matíc zrýchliło model
244 zhruba o 4%.

6. Záver 245

Táto práca sa venovala zrýchleniu výpočtu pravde- 246
podobností pre jazykový model využívajúci doprednú 247
neurónovú sieť. 248

V tomto dokumente sme videli, ako jednotlivé 249
spôsoby zrýchľovania ovplyvnili náš model. Pri zmene 250
aktivačnej funkcie sme dosiahli 3-násobné zrýchlenie 251
výpočtu aktivačnej funkcie. Uchovávanie histórie nám 252
kvôli cyklickému prehľadávaniu histórie model naopak 253
40% spomalilo. Predpočítanie matíc pre skrytú vrstvu 254
nás vo výpočte skrytej vrstve zrýchliło 2.5-krát. 255

V tomto dokumente sme úspešne implementovali 256
niektoré zrýchlenia pre jazykový model využívajúci 257
doprednú neurónovú sieť. Naďalej chcem v práci 258
pokračovať a implementovať model bez normalizácie, 259
ktorý sa z dôvodov implementačných ťažkostí nepodar- 260
ilo do práce zahrnúť a preskúmať vplyv zrýchľovania 261
pre modely s rôznymi veľkosťami skrytých vrstiev a 262
dĺžkami vstupných reťazcov. 263

Podakovanie 264

Chcel by som poďakovať pánovi Ing. Karlovi Benešovi 265
za jeho rady, skvelú pomoc a nekonečnú trpezlivosť 266
pri vypracovávaní tohto dokumentu. 267

Literatúra 268

- [1] Yoshua Bengio, Réjean Ducharme, Pascal Vin- 269
cent, and Christian Jauvin. A neural probabilistic 270
language model. *Journal of machine learning re-* 271
search, 3(Feb):1137–1155, 2003. 272
- [2] Yinghui Huang, Abhinav Sethy, and Bhuvana 273
Ramabhadran. Fast neural network language 274
model lookups at n-gram speeds. *Proc. Inter-* 275
speech 2017, pages 274–278, 2017. 276
- [3] Peter F Brown, Peter V Desouza, Robert L Mercer, 277
Vincent J Della Pietra, and Jenifer C Lai. Class- 278
based n-gram models of natural language. *Com-* 279
putational linguistics, 18(4):467–479, 1992. 280
- [4] Bing Xu, Naiyan Wang, Tianqi Chen, and 281
Mu Li. Empirical evaluation of rectified activa- 282
tions in convolutional network. *arXiv preprint* 283
arXiv:1505.00853, 2015. 284
- [5] PyTorch. Tensors and dynamic neural networks in 285
python with strong gpu acceleration., 2017. 286
- [6] Jacob Devlin, Chris Quirk, and Arul Menezes. Pre- 287
computable multi-layer neural network language 288
models. In *Proceedings of the 2015 Conference* 289
on Empirical Methods in Natural Language Pro- 290
cessing, pages 256–260, 2015. 291