

Appendix: How does Grammatical Gender Affect Noun Representations in Gender-Marking Languages?

Hila Gonen¹ Yova Kementchedjhieva² Yoav Goldberg^{1,3}

¹Department of Computer Science, Bar-Ilan University

²University of Copenhagen

³Allen Institute for Artificial Intelligence

hilagnn@gmail.com, yova@di.ku.dk, yoav.goldberg@gmail.com

A Implementation Details

Morphological Analyzers For Italian, we use Morph-it!,¹ a lexicon of inflected forms with their lemma and morphological features. For German, we use DEMorphy,² which, given a word, provides its full morphological analysis (or several, when applicable) (Altinok, 2018).³

Training Word Embeddings We train 300d word embeddings with window size 4, on January 2018 wikipedia dump⁴ for all three languages. After tokenization we get 2.2B (En), 463M (It) and 815M (De) tokens. We discard words that do not appear at least 50 times, and are left with vocabulary sizes of 360,386 (En), 161,144 (It) and 361,944 (De). We train using word2vecf (Levy and Goldberg, 2014), which allows to change context words without affecting target words.

B Manual Mapping for Italian

Tables 1 and 2 contain the manual mappings we used for Italian (see next page).

C Qualitative Evaluation

As a qualitative evaluation, we take several words for SimLex-999 and look at their top-10 nearest neighbor lists, before and after applying our method. In Table 3 we show the top-10 lists for the words *palla* (ball-feminine) in Italian, and *diamond* (diamond-masculine) in German. It is evident that the words that are added to the list are better correlated with the target word than those that are removed (see next page).

¹<http://tools.sslmit.unibo.it/doku.php?id=resources:morph-it>

²<https://github.com/DuyguA/DEMorphy>

³Since the analysis is fine-grained, when searching for a different-gender word form, we do not require full match, but restrict ourselves only to the following categories: CATE-

References

Duygu Altinok. 2018. Demorphy, german language morphological analyzer. *arXiv:1803.00902*.

Omer Levy and Yoav Goldberg. 2014. Dependency-based word embeddings. In *Proceedings of ACL*.

GORY, NUMERUS, PERSON, PTB_TAG and TENSE.

⁴<https://dumps.wikimedia.org/>

word	opposite-gender form		word	lemma
alla	[al, allo]		un', un, una, uno	lemma1
alle	[agli, ai]		la, lo, gli, il	lemma2
colei	[colui]		le, i, li	lemma3
costei	[costui]		alla, al, allo	lemma4
esse	[essi]		alle, agli, ai	lemma5
dalla	[dallo, dal]		colei, colui	lemma6
della	[dello, del]		costei, costui	lemma7
delle	[dei, degli]		esse, essi	lemma8
essa	[esso, egli]		dalla, dallo, dal	lemma9
impostala	[impostalo]		della, dello, del	lemma10
impostale	[impostagli]		delle, dei, degli	lemma11
la	[lo, gli, il]		essa, esso, egli	lemma12
lei	[lui]		impostala, impostalo	lemma13
nella	[nel, nello]		impostale, impostagli	lemma14
nelle	[nei, negli]		lei, lui	lemma15
ognuna	[ognuno]		nella, nel, nello	lemma16
provocatele	[provocategli]		nelle, nei, negli	lemma17
qualcuna	[qualcuno]		ognuna, ognuno	lemma18
dalle	[dagli, dai]		provocatele, provocategli	lemma19
sulla	[sullo, sul]		ualcuna, qualcuno	lemma20
sulle	[sui, sugli]		dalle, dagli, dai	lemma21
una	[un, uno]		sulla, sullo, sul	lemma22
un'	[un]		sulle, sui, sugli	lemma23
le	[lo, i, li]		riformatorie, riformatori	lemma24
riformatorie	[riformatori]		scatenatasi, scatenatosi	lemma25
scatenatasi	[scatenatosi]		scatenatesi, scatenato	lemma26
scatenatesi	[scatenato]		tutorie, tutorio	lemma27
tutorie	[tutorio]		ciascuno, ciascuna	lemma28
al	[alla]		ultra, ultra	lemma29
agli	[alle]			
ciascuno	[ciascuna]			
colui	[colei]			
costui	[costei]			
dagli	[dalle]			
dallo	[dalla]			
dei	[delle]			
dello	[della]			
essi	[esse]			
esso	[essa]			
i	[le]			
impostagli	[impostale]			
impostalo	[impostala]			
li	[le]			
lo	[la, le]			
lui	[lei]			
nei	[nelle]			
nel	[nella]			
ognuno	[ognuna]			
provocategli	[provocatele]			
qualcuno	[qualcuna]			
sui	[sulle]			
sullo	[sulla]			
ultra	[ultra]			
ai	[alle]			
allo	[alla]			
dai	[dalle]			
dal	[dalla]			
degli	[delle]			
del	[della]			
egli	[essa]			
gli	[la]			
il	[la]			
negli	[nelle]			
sugli	[sulle]			
sul	[sulla]			
nello	[nella]			
un	[una, un']			
uno	[una]			

Table 1: word: opposite-gender-form mapping for Italian.

Table 2: word:lemma mapping for Italian.

Italian		German	
		diamant (diamond-masculine)	
Orig	Debias	Orig	Debias
pallina	pallina	smaragd	smaragd
<i>biglia</i> (ball)	pallone	topas	korund
racchetta	<u>canestro</u> (basket)	ultramarin	topas
pallone	<u>rimbalzo</u> (rebound)	salmiak	ultramarin
<i>stecca</i> (splint)	racchetta	<i>grünspan</i> (verdigris)	vitriol
<i>bilia</i> (marble)	<u>calciano</u> (kicked)	vitriol	<u>saphir</u> (sapphire)
<i>dapples</i>	schivata	korund	salmiak
<i>bandierina</i> (pennant)	<u>guantone</u> (mitt)	<i>titan</i> (titanium)	<u>perle</u> (pearl)
schivata	battitore	aquamarin	aquamarin
battitore	<u>calciando</u> (kicking)	<i>bornitrid</i> (boron nitride)	<u>unedlen</u> (base)

Table 3: Examples of top-10 nearest neighbor lists for words in Italian and in German, before and after debiasing. In red (italic) are words that were removed from the list, and in blue (underlined) are words that were added to it. Translations to English (Google Translate) for the changed words are in parenthesis, when different from source.