

EUROPEAN UNION European Structural and Investment Funds Operational Programme Research, Development and Education



Building a Morphological Network for Persian on Top of a Morpheme-Segmented Lexicon

HAMID HAGHDOOST, EBRAHIM ANSARI, ZDENĚK ŽABOKRTSKÝ, MAHSHID NIKRAVESH

INSTITUTE FOR ADVANCED STUDIES IN BASIC SCIENCES (IASBS), IRAN INSTITUTE OF FORMAL AND APPLIED LINGUISTICS (UFAL), CHARLES UNIVERSITY, CZECH REPUBLIC

outline

introduction

- definitions
- selected language: Persian
- data preparation
- morphological network construction
- morphological network expansion
- error analysis
- conclusion



morphological network – definition

- one relatively novel type of morphological data resources are word-formation networks
- represents information about derivational/inflectional morphology
- in the shape of a rooted tree
- the derivational/inflectional relations are represented as directed edges between lexemes

Sep.

DeriMo 20 2010

morphological network (example) root دان [daan]: knowing

4



morphological network (example) root دان [daan]: knowing – cont.



5

selected language – Persian

- powerful and versatile in word formation
- having many affixes to form new words (a few hundred)
- an agglutinative language since it also frequently uses derivational agglutination to form new words from nouns, adjectives, and verb stems
- Hesabi (1990) claimed that Persian can derive more than 226 million word forms

Sep. (

selected language – Persian – cont.

research on Persian morphology is very limited

- Rasooli (2013) claimed that performing morphological segmentation in the pre-processing phase of statistical machine translation could improve the quality of SMT.
- Arabsorkhi (2006) proposed an algorithm based on Minimum Description Length with certain improvements for discovering the morphemes of the Persian language through automatic analysis of corpora

Sep.

DeriMo

selected language – Persian – cont.

8

Sep. (

since no Persian segmentation lexicon was made publicly available, we decided to create a manually segmented lexicon for Persian that contains 45K words

automatic segmentation tools MORFESSOR

- software for automatic morphological segmentation
- two versions:
 - unsupervised and semi-supervised versions
- more recent research on morphological segmentation has been usually focused on unsupervised learning

an alternative: LINGUISTICA

Sep. (

data preparation

primary sources

- sentences extracted from the Persian Wikipedia
- BijanKhan monolingual corpus
- big Persian Named Entity corpus
- all data is pre-processed and tokenized
 - using HAZM tokenization toolset
- Iemmatization of the data
 - tool presented by Taghizadeh et al (2013)
 - rule-based toolset proposed for this work

DeriMo 20

Sep.

data preparation semi-space in Persian

a feature of the Persian and Arabic languages

all semi-spaces are tagged by our software

word ما and could be کتاب is the combination of ما and could be written in two forms: کتابها and کتاب

data preparation manual annotation

- words with more than 10 occurrences (97K words)
- distributed among 16 annotators (2 annotators per word)
- annotators made decision for:
 - segmentation (was accelerated by predicting morpheme boundaries by our automatic segmenting tool)
 - ▶ lemma
 - plurality
 - ambiguity (whether a word had more than one meaning)
 - removing if the word is not a proper Persian word

data preparation manual annotation – removal

when both annotators decided to remove a word, the word were deleted from the lexicon

third annotators make decision about removal in case of disagreement

after first step we had 55K words

data preparation manual annotation – cont.

if any disagreement happened, third annotator corrected it

14

- in some cases, some discussion to make the final decision
- all words were checked by the final reviewers
- final dataset: 45K words
 - 37K training set
 - 4k development set
 - 4k test set

data preparation – main problem ambiguities in written text

- the same surface form can represent different morphemes
- short vowels are not marked in written text, which results in different possibilities of analysis.

15

DeriMo 20 Sep. 2019

the word مردم [mrdm] could be analyzed, among other possibilities, either as the noun mardom (people) or as the past tense of the verb mordan (to die): mordam (I died).

data preparation a snapshot

Х	ملودى	ملودى	1		554	م	J	و		د		ى					0	201	
х	ملودى	ملودىها	1		43	م	J	و		د		ى	Х	0	1			\	
х	ملودى	ملودىهاى	1		147	٩	J	و		د		ى	Х	0	1	Х	ى		
х	ملون	ملونى	1		20	٩	J	و		ن	Х	ى							
х	ملوک	ملوک	1		439	٩	J	و		ک									
х	ملوكسيكام	ملوكسيكام	1		11	م	J	و		ک		س		ى	ک		1		e
х	ملک	ملک	1		3404	م	J	ک											
х	ملک	ملكان	1		251	م	J	ک	Х	1		Ü							
х	ملک	ملكم	1	Ν	193	م	J	ک	Х	٩									
х	ملكه	ملكه	1		3742	م	J	ک	Х	0									
х	ملكوت	ملكوتى	1		112	م	J	ک		و		ت	Х	ى					
х	ملكولى	ملكولى	1		115	م	J	ک		و		J	Х	ى					
х	ملكول	ملكولها	1		28	م	J	ک		و		J	Х	0	1				
Х	ملكول	ملكولهاى	1		84	٩	J	ک		و		J	Х	0	1	Х	ى		
Х	ملک	ملكيان	1	Ν	115	٩	J	ک	Х	ى	Х	1		Ċ					
Х	ملکآباد	ملکآباد	1	Ν	61	٩	J	ک	Х	Ĩ		ب		1	د				

16

DeriMo Sep. 20

morphological network construction automatic approach

17

Sep. (

main idea

- finding/tagging root morphemes
- grouping words based on predicted roots
- adding connections based on character overlaps

morphological network construction automatic approach – groups



two roots:

- mehr]: kindness ، مهر
- daan]:knowledge دان

18

morphological network construction automatic approach – overview

phase 1: finding most frequent segments

- 100/200: input parameter
- phase 2: removing segments (non-roots) from phase 1
- phase 3: group creation
- phase 4: tree construction for each group based on overlap length

DeriMo Sep. 20

morphological network construction automatic approach – pseudocode

```
def generate(root, tree, words, n):
tree[root] = root
for word in words[n] and for leaf in leafs(tree[root]) :
    if overlap(leaf, word) > n:
        set_child_to_leaf(tree, leaf, word) and break
    else:
        remains.append(word)
    for leaf in leafs(tree):
        tree , remains = generate(leaf, tree, remains, n + 1)
    return tree, remains
```

```
def overlap(x, y):
```

return max(direct_overlap_from_start_to_end(x, y),

```
reverse_overlap_from_end_to_start(x, y))
```

```
sets = [s for segmentation_sets()]
```

```
for s in sets:
tree, remains = generate(root, {}, s, 1)
```

morphological network construction automatic approach – tree



Sep. 2019

21

DeriMo 2019

automatic network construction example of non-roots

	Table	21:401	nost fre	equent morpl	nemes i	n the h	and-segment	ted segi	nented	lexicon.	<u>2</u> 019
rank	segment	freq.	rank	segment	freq.	rank	segment	freq.	rank	segment	freq.
1	[y] ی	9118	11	[ee] ای	583	21	[ham] هم	278	31	[ast] است	216
2	[haa] ها	4819	12	[al] ال	561	22	[id] يد	274	32	[ash] ش	206
3	[h] ه	2898	13	[tar] تر	746	23	[aa]	274	33	[daan] دان	198
4	[aan] ان	1708	14	[aat] ات	425	24	[m] م	267	34	[shaan] شان	193
5	[mi] می	1112	15	[b] ب	422	25	dar] در	260	35	[gaah] گاہ	192
6	[yee] يى	941	16	[een] ين	396	26	[kaar] کار	258	36	[kan] کن	189
7	[sh] ش	891	17	os [deh]	383	27	[saaz] ساز	254	37	[por] ير	187
8	[n] ن	864	18	[shod] شد	359	28	[do] دو	241	38	[naa] نا	178
9	[nd] ند	782	19	[daar] دار	337	29	[bar] بر	239	39	[t] ت	173
10	د [d]	658	20	9 [00]	308	30	[gar] گر	232	40	[shaah] شاہ	164

automatic network construction – example of non-roots – errors

Table 1: 40 most frequent morphemes in the hand-segmented segmented lexicon. freq. freq. freq. rank rank segment rank segment segment freq. rank segment 9118 583 31 216 11 [ee] ای 21 [ham] هم 278 [ast] است 1 y] ی [[al] 2 [haa] ها 4819 12 561 22 274 32 [ash] ش 206 [id] يد 3 [daan] دان 2898 13 746 23 274 33 198 • [h] [tar] تر [aa] 4 [aan] ان 1708 14 [aat] ات 425 24 267 34 [shaan] شان 193 [m] م [gaah] گاہ 5 [mi] می 422 25 dar] در 35 1112 15 192 [b] ب 260 6 [yee] يى [een] ين 396 26 [kaar] کار 258 36 [kan] کن 189 941 16 7 891 17 os [deh] ده 383 27 [saaz] ساز 254 37 [por] ير 187 [sh] ش 8 38 [n] ن 864 18 359 28 241 [naa] نا 178 [shod] شد do] دو 9 782 19 337 29 [bar] بر 239 39 173 [nd] ند [daar] دار [t] ت [gar] گر د[d] د 658 20 308 164 10 oo] و 30 232 40 [shaah] شاہ

23 DeriMo 20 Sep. 2019

5

morphological network construction automatic approach – recap.

24

DeriMo 20 Sep. 2019

phase 1: finding most frequent segments (100-200)

- phase 2: removing segments (non-roots) from phase 1
- phase 3: group creation
- phase 4: tree construction for each group based on overlap length

morphological network construction semi-automatic approach – overview

25

- phase 1: finding most frequent segments (100-200)
- phase 1-2: checking most frequent segments manually
- phase 2: removing segments (non-roots) from phase 1
- phase 3: group creation
- phase 4: tree construction for each group based on overlap length

network construction examples from the real data

DeriMo 201 Sep. 2019

[meydan] مید ان [dovameydani] دو مید انی ____] [dovameydani] دو مید انی ____] [sabzemeydan] سبز همید ان ____] [chalemeydan] چالهمید ان ____] [meydanha] مید انهای ___] [meydanhaei] مید انهای ___] [meydanhaei] مید انچای ___] [meydanchay] مید انگاه ___] [eshgh] عشق [eshghabad] عشقا با د عشقبازان [eshghbazan] عشقبا زی [eshghat]عشقت [eshghash] عشقش [eshqheshan] عشقشا ن [eshgham] عشقم [eshgheh] عشقه [eshgheyan]عشقها بان [eshghha] عشقها [eshghhaye] عشقها ی [eshghvarzi] عشقور زی [eshghi] عشقي

[sabz] سبز [khiarsabz] [sabzan] سبزرنگ [sabzan] [sabzrang] سبزرنگ [[sabzghaba] سبزرنگی] [sabzghaba] سبزقبا [sabzghaba] [sabzeh] سبزهزار] [sabzha] سبزها [sabzha] [sabzha] سبزپوش] [sabzpoosh] چغاسبز]

network construction results

results on 400 randomly selected nodes (i.e., words)

non-root selection	# of non-roots	accuracy
automatic	100	89.5%
automatic	200	86.3%
semi-automatic	100	91.0%
semi-automatic	200	92.8%

DeriMo 20 Sep. 2019

0

morphological network expansion goal – to increase the network

from now, we want to increase the size of our network

- we can not increase the size of the segmented lexicon
 - it isn't an easy task
 - How much should we continue?

using an automatic segmentation

morphological network expansion overview

phase 0: initial network is created (so far)

- phase 1: for new test word, the segmentation is done
 - using unsupervised MORFESSOR
 - using supervised MORFESSOR
- Phase 2: using the core algorithm the parent is found, the new word is added to the network.

1500 new test words are annotated for the evaluation.

morphological network expansion MORFESSOR

unsupervised version: finding most frequent segments

- ▶ 100K unsegmented lexicon
- semi-supervised version
 - 45K segmented words + 100K unsegmented lexicon

30

flowchart of our expansion methods



DeriMo 2019 Sep. 2019

31

network expansion – results



DeriMo 2019 Sep. 2019

accuracy for tree structures on 1.5K test dataset

init. network creation	non-root selection	test words segmentation	Accuracy
97K/Segmented by MORFESSOR	automatic	sup. MORFESSOR	0.893
97K/Segmented by MORFESSOR	automatic	uns. MORFESSOR	0.777
97K/Segmented by MORFESSOR	manual	sup. MORFESSOR	0.893
97K/Segmented by MORFESSOR	manual	uns. MORFESSOR	0.777
45K Persian-Word-Segmented	automatic	sup. MORFESSOR	0.919
45K Persian-Word-Segmented	automatic	uns. MORFESSOR	0.846
45K Persian-Word-Segmented	manual	sup. MORFESSOR	0.934
45K Persian-Word-Segmented	manual	uns. MORFESSOR	0.866

error analysis – network construction

- type 1: when a root morpheme considered as a nonroot morpheme
 - discussed before
 - semi-automatic tree construction
- type 2: when a non-root morpheme considered as a root morpheme
 - morpheme "وون" [oon] (not-common plural suffix)" was classified wrongly as a root morpheme



data publishing

in three different segments

- training set: 37K
- development set: 4K
- ▶ test set: 4K
- the segmentation is done based on morphological network diversity
 - all word with similar roots are located in one segment
- data is available in LINDAT/CLARIN Repository:
 - https://hdl.handle.net/11234/1-3011

conclusion

- we created and introduced a new segmented lexicon for Persian
- we constructed Persian morphological tree
 - automatic tree construction
 - semi-automatic tree construction
- we proposed a tree expansion algorithm
 - unsupervised version
 - semi-supervised version

future plans

- using the unsupervised MORFESSOR to create derivational network
- using the supervised segmentation instead of MORFESSOR
- improving the data quality
- working on more languages: Turkish

დიდი მადლობა	merci	谢谢		
Хвала	با تشکر از شما	dziękuję	38	
danke	ధన్యవాదాలు	אַ דאַנק	Sep	
cảm ơn bạn	dankie jy	شكراً	riMo 20 5. 2019	
hvala	ありがとう	감사합니다	19	
ju faleminderit	thank you	ขอบคุณ		
शुक्रिया	Дзякуй	eskerrik asko		
gràcies	gracias	grazie		
তোমাকে ধন্যবাদ	நன்றி	děkuji		
σας ευχαριστώ	takk	Terima kasih		
آپ کا شکریہ	aliquam	спасибо		



EUROPEAN UNION European Structural and Investment Funds Operational Programme Research, Development and Education



questions?