

Language Model Based Arabic Word Segmentation

By Saleh Al-Zaid
Software Engineering

Content

- Introduction
- Morpheme Segmentation
- New Stems Acquisition
- Performance Evaluation
- Conclusion

Introduction

- **Paper:** Language Model Based Arabic Word Segmentation, 2003.
- **By:** Young-Suk Lee, Kishore Papineni, Salim Roukos, Ossama Emam and Hany Hassan .
- **Objective:** introduce an algorithm to segment and acquire new Arabic stems form un-segmented Arabic corpus.

Segmentation?

- Divide words into parts (morphemes) using the pattern: prefix*-stem-suffix*
- * = zero or more
- Prefix indicated by #
- Suffix indicated by +
- Example: ال#سيار +ة = السيارة

Algorithm Implementation

1. Language model training on a manually morpheme-segmented small corpus (from 20K to 100K words).
2. Segmentation of input text into a sequence of morphemes using the language model parameters.
3. Acquisition of new stems from a large unsegmented corpus.

Morpheme Segmentation

1. Trigram Language Model
2. Decoder for Morpheme Segmentation
 - **Possible Segmentations of a Word**

Trigram Language Model

- Sample of manually segmented corpus

و# كان ايرفاين الذي حل في ال# مركز
ال# اول في جائز +ة ال# نمسا ال# عام
ال# ماضي علي سيار +ة فبراري شعر ب#
الام في بطن +ه اضطر +ت +ه الي ال#
انسحاب من ال# تجارب و# هو س# ي# عود
الي لندن ل# اجراء ال# فحوص +ات ال#
ضروري +ة حسب ما اشار فريق جاغوار .
و# س# ي# حل سائق ال# تجارب في جاغوار
ال# برازيلي لوسيانو بورتى مكان
ايرفاين في ال# سباق غذا ال# احد
الذي س# ي# كون اولي خطو +ات +ه في
عالم سباق +ات الفورمولا

Trigram Language Model

- Buckwalter equivalence in English:

w# kAn AyrfAyn Al*y Hl fy Al# mrkz Al#
Awl fy jA}z +p Al# nmsA Al# EAm Al#
mADy Ely syAr +p fyrAry \$Er b# AlAm fy
bTn +h ADTr +t +h Aly Al# AnsHAb mn Al#
tjArb w# hw s# y# Ewd Aly lndn l# AjrA' Al#
fHwS +At Al# Drwry +p Hsb mA A\$Ar fryq

- Trigram: $p(m_i \mid m_{i-1}, m_{i-2})$

Table of Segments

Words		Prefixes		Stems		Suffixes	
Arabic	Translit.	Arabic	Translit.	Arabic	Translit.	Arabic	Translit.
الولايات	<i>AlwLAyAt</i>	#ال	<i>Al#</i>	ولاي	<i>wLAy</i>	+ات	<i>+At</i>
حياته	<i>HyAth</i>			حيا	<i>HyA</i>	+ت +ه	<i>+t +h</i>
للحصول	<i>lHSwl</i>	#ال #ل	<i>l# Al#</i>	حصول	<i>HSwl</i>		
الى	<i>AlY</i>			الى	<i>AlY</i>		

Table 1 Segmentation of Arabic Words into Prefix*-Stem-Suffix*

To be used in the algorithm next steps

Decoder for Morpheme Segmentation

- Goal: to find the morpheme sequence which maximizes the trigram probability of the input sentence, as in:

$$\text{SEGMENTATION}_{\text{best}} = \text{Argmax}_{i=1, N} p(m_i | m_{i-1} m_{i-2})$$

N = number of morphemes in the input

Possible Segmentations of a Word

- Depends on the table of Prefix/Suffix

Prefixes		Suffixes	
ال	#ال	ات	+ات
بال	#ال #ب	اتها	+ها +ات
وبال	#ال #ب #و	ونهم	+ون +هم

Table 2 Prefix/Suffix Table

- Each new token is assumed to have: *prefix*-stem-suffix** structure and compared against prefix/suffix table.

Possible Segmentations of a Word

- Steps to find possibilities:
 - i) Identify all of the matching prefixes and suffixes from the table,
 - ii) Further segment each matching prefix/suffix at each character position, and
 - iii) Enumerate all *prefix*-stem-suffix** sequences derivable from (i) and (ii).

Possible Segmentations of a Word

- Example: suppose “واكررها” is new token which is : $wAkrrhA$,
- Using :
 $SEGMENTATION_{best} = \underset{I}{\operatorname{Argmax}} \prod_{i=1, N} p(m_i | m_{i-1} m_{i-2})$,
the possible segmentations are:

Possible Segmentations of a Word

	Prefix	Stem	Suffix	Seg Scores
S1	∅	<i>wAkrrhA</i>	∅	2.6071e-05
S2	∅	<i>wAkrrh</i>	+A	1.36561e-06
S3	∅	<i>wAkrr</i>	+hA	9.45933e-07
S4	w#	<i>AkrrhA</i>	∅	2.72648e-06
S5	w#	<i>Akrrh</i>	+A	5.64843e-07
S6	w#	<i>Akrr</i>	+hA	4.52229e-05
S7	wA#	<i>krrhA</i>	∅	7.58256e-10
S8	wA#	<i>krrh</i>	+A	5.09988e-11
S9	wA#	<i>krr</i>	+hA	1.91774e-08
S10	w# A#	<i>krrhA</i>	∅	7.69038e-07
S11	w# A#	<i>krrh</i>	+A	1.82663e-07
S12	w# A#	<i>krr</i>	+hA	0.000944511

**Table 3 Possible Segmentations of
واكررها (*wAkrrhA*)**

Acquisition of New Stems

- Form large un-segmented corpuses.
- Follow this process:
- **Initialization:** Develop the seed segmenter Segmenter0 trained on the manually segmented corpus Corpus0, using the language model vocabulary, Vocab0, acquired from Corpus0.

Acquisition of New Stems

- **Iteration:** For $i = 1$ to N , $N =$ the number of partitions of the unsegmented corpus:
 - i. Use Segmenter_{i-1} to segment Corpus_i .
 - ii. Acquire new stems from the newly segmented Corpus_i . Add the new stems to Vocab_{i-1} , creating an expanded vocabulary Vocab_i .
 - iii. Develop Segmenter_i trained on Corpus_0 Corpus_i with Vocab_i .

Performance Evaluation

- Using the formula:
- $E = (\text{number of incorrectly segmented tokens} / \text{total number of tokens}) \times 100$
- Example: ليتر : ل# ي# تر
- For the paper, the evaluation is done in test corpus with 28,449 words by using 4 manually segmented seed corpora with 10K, 20K, 40K, and 110K words.

Results

Manually Segmented Training Corpus Size	Unknown Stem	Other Errors	Total # of Errors
10 K Words	1,844 (76.9%)	455 (19.0%)	2,397
20 K Words	1,174 (71.1%)	395 (23.9%)	1,651
40 K Words	1,005 (69.9%)	351 (24.4%)	1,437
110 K Words	333 (39.6%)	426 (50.7%)	841

Table 7 Segmentation Error Analyses

Conclusion

- Language Model Based Arabic Word Segmentation Algorithm can segment and acquire new stems as long we have good training manually segmented corpus
- It give more good results as the training corpus have large number of stems

Thank you

Q & A