

Morphological Analysis for Less-Resourced Languages: Maximum Affix Overlap applied to Zulu

Uwe Quasthoff¹, Sonja Bosch², Dirk Goldhahn¹
¹Natural Language Processing Group, University of Leipzig, Germany
²Dept. of African Languages, University of South Africa

Less-resourced Languages

Lack of text resources:

- Dictionaries are often outdated
- No machine-readable dictionaries
- Limited amount of text in the Web

Some corpora exist (e.g. University of Pretoria, Language Resource Management Agency and Leipzig Corpora Collection) - limited in size, not annotated and often not accessible.

Lack of NLP data and tools:

- Training data for POS-tagger
- Training data for morphology

Complex Morphology of Zulu

Zulu [ISO 639-3: zul]

Nominal classification system: nouns categorized by prefixal morphemes (with class numbers)

Concordial agreement system: noun class prefixes use concordial agreement for linking nouns to verbs, adjectives, pronouns, possessives etc.

abantu	abaningi	bangayichitha	imali	yabo
aba-ntu	aba-ningi	ba-nga-yi-chitha	i-mali	ya-bo
people	many	they-may-it-waste	money	of-them
[Many people may waste their money.]				

The Collaboration Loop

omubi	o<r>mu<n1>	bi<ar>		1364
omubi	o<r>m<o1>	ub<vr>	i<vg>	777
omubi	o<z6 iv>m<n3>	ub<vr>	i<vg>	259
omubi	o<r>m<o1>	ub<vr>	i<hum>	259

Retrain classifier
and
classify unknown instances

No.	Prefix(es)	Root	Suffix	Frequency	Correct
1	ng<p>e<iv>si<n7>	m<vr>	o<in>	402	<input type="checkbox"/>
2	ng<p>e<iv>si<n7>	mo<nr>		17040	<input checked="" type="checkbox"/>
3	ng<p>e<r>si<n7>	mo<ar>		682	<input type="checkbox"/>
4	ng<p>e<iv>s<n7>	imo<nr>		4260	<input type="checkbox"/>

Confirm correct solutions

Machine Learning: Maximum Affix Overlap Algorithm

For all segmentations of the word w into **three segments** w_1 (prefixes), w_2 (root) and w_3 (suffixes), where w_1 and w_3 might be of zero length:

- For each word x in the training set having exactly the **prefix sequence** w_1 we collect the pair (morphological analysis of w_1 , w_2 with the tag of the root of x).
- For each word x in the training set having exactly the **suffix sequence** w_3 we collect the pair (w_2 with the tag of the root of x , morphological analysis of w_3).

From this, form triple by joining on **identical root tags**: (morphological analysis of w_1 , w_2 with the tag of the root of x , morphological analysis of w_3).

Interesting features: **length** of w_2 and **frequency** of identical triples above.

Procedure above allows considering affixes next to stem as part of stem, therefore **shorter stems** should be preferred.

In case of **multiple decompositions** with same stem (or different stems of same length): decompositions ranked according to frequency.

In general, we set a **frequency threshold** of 2 for decompositions to be considered.

Training Data Web Editor

Zulu - Morphological Analysis

Search

Word

Search

No.	Prefix(es)	Root	Suffix	Frequency	Correct
1	y<z9>o<r>	cwaning<vr>	o<in>	201	<input type="checkbox"/>
2	y<z4>o<iv_n11>	cwaningo<nr>		2130	<input checked="" type="checkbox"/>
3	y<z9>o<iv_n3>	cwaningo<nr>		2130	<input type="checkbox"/>
4	y<l9>	ocwaning<vr>	o<in>	2010	<input type="checkbox"/>
5	y<l4>	ocwaning<vr>	o<in>	1206	<input type="checkbox"/>
6	y<z9>o<r>	cwaningo<vr>		214	<input type="checkbox"/>

Confirm

Applications for the Collaboration Loop

- Inflection type / baseform reduction, morphological decomposition, compound decomposition
 - Special problem: Identifying new roots for dictionaries
 - Special problem: Identifying writing variants
- Classification tasks for subject areas or relations (as in WordNet)
- Bilingual translation equivalents

