



# Automatic Morpheme Segmentation and Labeling in Universal Dependencies resources

Miikka Silfverberg  
Mans Hulden

May 22, 2017

# UD & Morphosyntactic info

- UD resources for many languages contain **rich morphological labeling** for lexical and grammatical properties of words

Lexical features	Inflectional features	
	<i>Nominal*</i>	<i>Verbal*</i>
<u>PronType</u>	<u>Gender</u>	<u>VerbForm</u>
<u>NumType</u>	<u>Animacy</u>	<u>Mood</u>
<u>Poss</u>	<u>Number</u>	<u>Tense</u>
<u>Reflex</u>	<u>Case</u>	<u>Aspect</u>
<u>Foreign</u>	<u>Definite</u>	<u>Voice</u>
<u>Abbr</u>	<u>Degree</u>	<u>Evident</u>
		<u>Polarity</u>
		<u>Person</u>
		<u>Polite</u>

<http://universaldependencies.org/u/feat/index.html>


# Segmentation and labeling

## English Data

word form	lemma	POS		morphosyntactic features
busted	bust	VERB	VCN	Tense=Past   VerbForm=Part
authorities	authority	NOUN	NNS	Number=Plur
announced	announce	VERB	VBD	Mood=Ind   Tense=Past   VerbForm=Fin
cells	cell	NOUN	NNS	Number=Plur
cell	cell	NOUN	NN	Number=Sing
mid-nineties	mid-ninety	NOUN	NNS	Number=Plur
tension	tension	NOUN	NN	Number=Sing
tensions	tension	NOUN	NNS	Number=Plur
announcing	announce	VERB	VBG	Tense=Pres   VerbForm=Part
authority	authority	NOUN	NN	Number=Sing
killed	kill	VERB	VBD	Mood=Ind   Tense=Past   VerbForm=Fin

# Segmentation and labeling

## English Data



busted	bust	VERB	VBN	Tense=Past   VerbForm=Part
authorities	authority	NOUN	NNS	Number=Plur
announced	announce	VERB	VBD	Mood=Ind   Tense=Past   VerbForm=Fin
cells	cell	NOUN	NNS	Number=Plur
cell	cell	NOUN	NN	Number=Sing
mid-nineties	mid-ninety	NOUN	NNS	Number=Plur
tension	tension	NOUN	NN	Number=Sing
tensions	tension	NOUN	NNS	Number=Plur
announcing	announce	VERB	VBG	Tense=Pres   VerbForm=Part
authority	authority	NOUN	NN	Number=Sing
killed	kill	VERB	VBD	Mood=Ind   Tense=Past   VerbForm=Fin


# Segmentation and labeling

## English Data

busted	bust	VERB	VRN	Tense=Past   VerbForm=Part
authorities	authority	NOUN	NNS	Number=Plur
announced	announce	VERB	VBD	Mood=Ind   Tense=Past   VerbForm=Fin
cells	cell	NOUN	NNS	Number=Plur
cell	cell	NOUN	NN	Number=Sing
mid-nineties	mid-ninety	NOUN	NNS	Number=Plur
tension	tension	NOUN	NN	Number=Sing
tensions	tension	NOUN	NNS	Number=Plur
announcing	announce	VERB	VBG	Tense=Pres   VerbForm=Part
authority	authority	NOUN	NN	Number=Sing
killed	kill	VERB	VBD	Mood=Ind   Tense=Past   VerbForm=Fin

# Segmentation and labeling

## English Data



busted	bust	VERB	VCN	Tense=Past   VerbForm=Part
authorities	authority	NOUN	NNS	Number=Plur
announced	announce	VERB	VBD	Mood=Ind   Tense=Past   VerbForm=Fin
cells	cell	NOUN	NNS	Number=Plur
cell	cell	NOUN	NN	Number=Sing
mid-nineties	mid-ninety	NOUN	NNS	Number=Plur
tension	tension	NOUN	NN	Number=Sing
tensions	tension	NOUN	NNS	Number=Plur
announcing	announce	VERB	VBG	Tense=Pres   VerbForm=Part
authority	authority	NOUN	NN	Number=Sing
killed	kill	VERB	VBD	Mood=Ind   Tense=Past   VerbForm=Fin

Could labeling/alignments be performed automatically?

# Linguistic intuitions in segmentation

## Data

announced  
announcing  
authorities  
busted  
cell  
cells  
killed  
mid-nineties

## Allomorphs

Number=Plur

Tense=Past

Number=Sing

(Lemmas)

**Raw data**

# Linguistic intuitions in segmentation

## Data

announ ced  
 annou ncing  
 authorit ies  
 buste d  
 ce ll  
 cel ls  
 kil led  
 mid-nine ties

## Allomorphs

Number=Plur	Tense=Past
ies	ced
ties	d
ls	led
(Stems)	Number=Sing
announ	ll
annou	
authorit	VerbForm=Part
buste	
ce	ncing
cel	
kil	
mid-nine	

**Bad Segmentation**



# Linguistic intuitions in segmentation

## Data

announ ced  
 annou ncing  
 authorit ies  
 buste d  
 ce ll  
 cel ls  
 kil led  
 mid-nine ties

**Bad Segmentation**

## Allomorphs

Number=Plur

ies  
 ties  
 ls

(Stems)

announ  
 annou  
 authorit  
 buste  
 ce  
 cel  
 kil  
 mid-nine

Tense=Past

ced  
 d  
 led

Number=Sing

ll

VerbForm=Part

ncing

# Linguistic intuitions in segmentation

## Data

announ ced  
 annou ncing  
 authorit ies  
 buste d  
 ce ll  
 cel ls  
 kil led  
 mid-nine ties

**Bad Segmentation**

## Allomorphs

Number=Plur

ies  
 ties  
 ls

(Stems)

announ  
 annou  
 authorit  
 buste  
 ce  
 cel  
 kil  
 mid-nine

Tense=Past

ced  
 d  
 led

Number=Sing

ll

VerbForm=Part

ncing

# Linguistic intuitions in segmentation

## Data

announc ed  
 announc ing  
 authoriti es  
 bust ed  
 cell ∅  
 cell s  
 kill ed  
 mid-nineti es

## Allomorphs

Number=Plur	Tense=Past
s	ed
es	Number=Sing
	∅
(Stems)	
announc	
authoriti	VerbForm=Part
bus	
cell	ing
kill	
mid-nineti	

**Better Segmentation**

# Linguistic intuitions in segmentation

## Data

announc ed  
 announc ing  
 authoriti es  
 bust ed  
 cell ∅  
 cell s  
 kill ed  
 mid-nineti es

## Allomorphs

Number=Plur

s  
es

Tense=Past

ed

Number=Sing

∅

(Stems)

announc

authoriti

bus

cell

kill

mid-nineti

VerbForm=Part

ing

**Better Segmentation**

# Linguistic intuitions in segmentation

## Data

announc ed  
 announc ing  
 authoriti es  
 bust ed  
 cell ∅  
 cell s  
 kill ed  
 mid-nineti es

## Allomorphs

Number=Plur

s  
es

Tense=Past

ed

Number=Sing

∅

(Stems)

announc

authoriti

bus

cell

kill

mid-nineti

VerbForm=Part

ing

**Better Segmentation**

# Interim

- This is an inference problem with weak supervision
- We can treat this as a search problem in the (huge!) space of all possible segmentations and labelings over a (UD) corpus
- The labels given by UD give us a weak supervision (we know which features are present in each word form, and which aren't)

cells  
cell

cell  
cell

NOUN NNS  
NOUN NN

Number=Plur  
Number=Sing



# Objective function (I)

- A simple objective would be to **minimize the total number of allomorph types** in the corpus

# Minimize allomorphs

## Data

announc ed  
 announc ing  
 authoriti es  
 bust ed  
 cell ∅  
 cell s  
 kill ed  
 mid-nineti es

## Allomorphs

Number=Plur	Tense=Past
s	ed
es	Number=Sing
	∅
(Lemmas)	
announc	
authoriti	VerbForm=Part
bus	
cell	ing
kill	
mid-nineti	

**Better Segmentation**

**allomorphs=11**





# Minimize allomorphs

## Data

announ **ced**  
 annou **ncing**  
 authorit **ies**  
 buste **d**  
 ce **ll**  
 cel **ls**  
 kil **led**  
 mid-nine **ties**

**Bad Segmentation**

## Allomorphs

Number=Plur

**ies**  
**ties**  
**ls**

(Lemmas)

announ  
 annou  
 authorit  
 buste  
 ce  
 cel  
 kil  
 mid-nine

Tense=Past

**ced**  
**d**  
**led**

Number=Sing

**ll**

VerbForm=Part

**ncing**

**allomorphs=16**



# Objective function (I)

- A simple objective would be to **minimize the total number of allomorph types** in the corpus
- Unfortunately, that function has less desirable properties (not continuous, not differentiable, insensitive to small changes in segmentation)

kill led

bus ted

announ ced

talk ed

cal led



# Objective function (I)

- A simple objective would be to **minimize the total number of allomorph types** in the corpus
- Unfortunately, that function has less desirable properties (not continuous, not differentiable, insensitive to small changes in segmentation)

kill ed

bus ted

announ ced

talk ed

cal led



# Objective function (I)

- A simple objective would be to **minimize the total number of allomorph types** in the corpus
- Unfortunately, that function has less desirable properties (not continuous, not differentiable, insensitive to small changes in segmentation)

kill ed

bus ted

announ ced

talk ed

cal led

No change in  
allomorph count!

# Objective function (II)

“Symmetric conditional probability”

$$C(S, F) = \prod_{s \in \Sigma^*, f \in \mathcal{Y}} P(s|f)P(f|s)$$

s is segment (substring)  
 f is feature-value pair

all substrings

set of labels

Intuitively: (1) **substrings** declared allomorphs should be reliable predictors of a feature; (2) **features** should predict a substring

Encourages few different features per allomorph, and few allomorphs per feature

# Objective function (II)

“Symmetric conditional probability”

$$\mathcal{C}(S, F) = \prod_{s \in \Sigma^*, f \in \mathcal{Y}} P(s|f)P(f|s)$$

s is segment (substring)  
f is feature-value pair

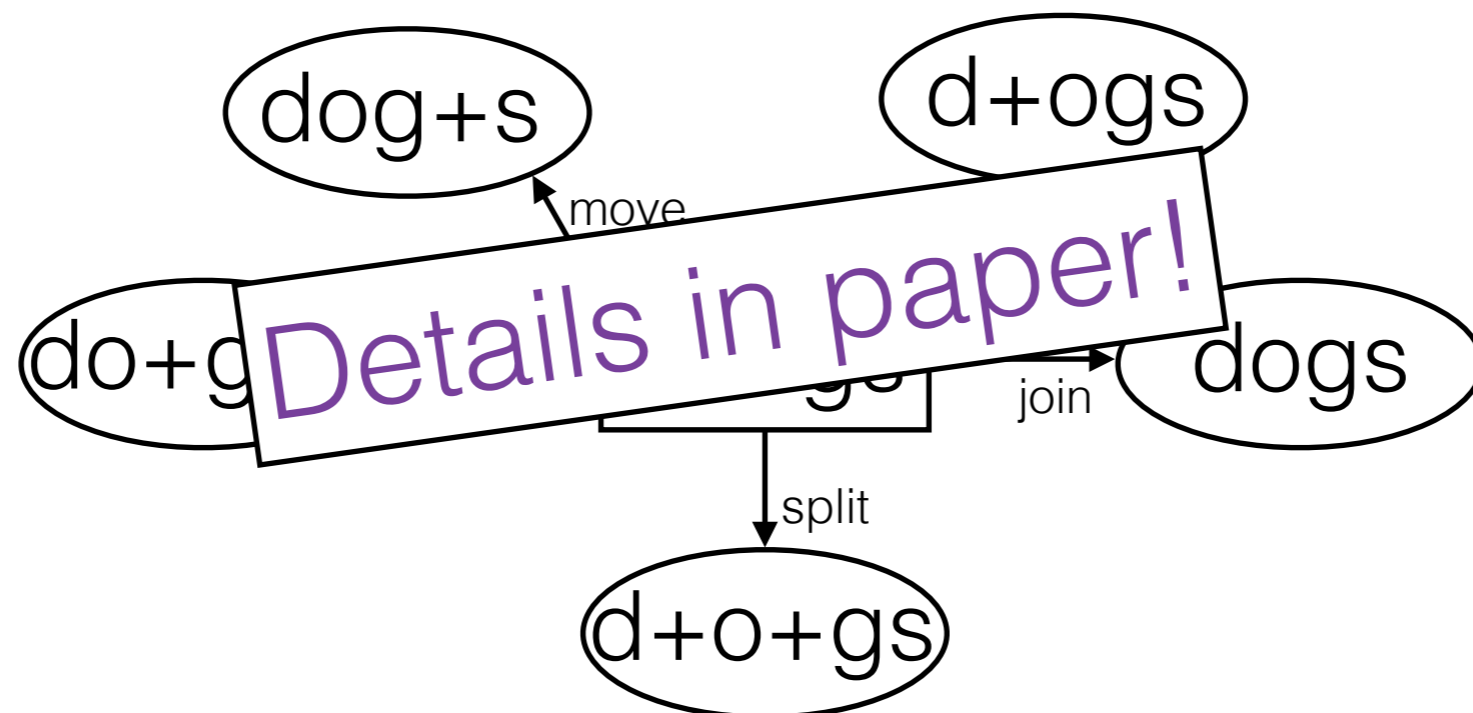
all substrings

set of labels

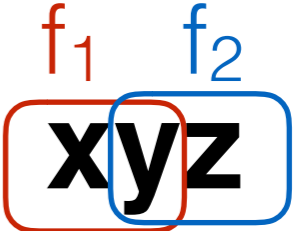
Not strictly true in either direction, but a workable proxy for the allomorph minimization idea (e.g. English has plurals **-s**, **-es**,  $\emptyset$ ; and **-s** can be both pluralizer and present tense 3p allomorph)

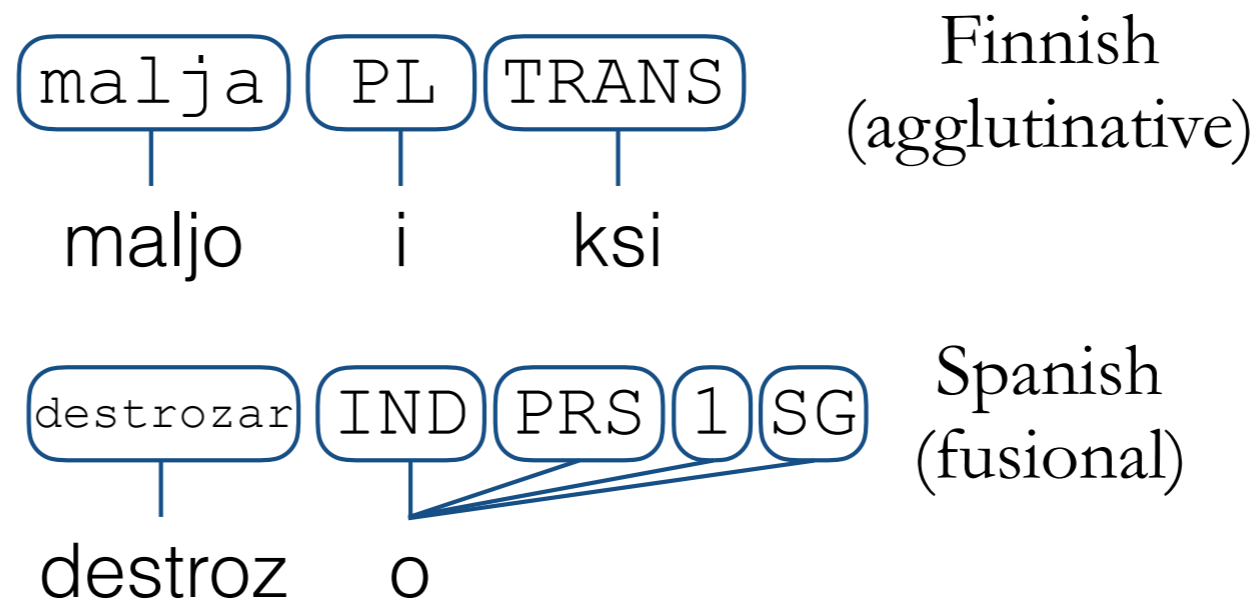
# Sampling

- Search space is still exponential, so we use a **sampling** approach
- Start from a **random segmentation** and labeling
- Make small changes and probabilistically move in the direction that increases the objective function



# Sampling

- Allow many features to correspond to single segments
- Also allow features to correspond to a special NULL ( $\emptyset$ ) segment (similar to IBM alignment models)
- Don't allow overlap in features      forbidden: 
- All actually occurring letters must be associated with a label







# Evaluation

- We compare the method with an unsupervised segmenter **Morfessor** (Creutz and Lagus, 2005)
- Run Morfessor to get a segmentation, then assign labels to maximize objective function with given segmentations (a much easier problem if segmentation is given)
- We use hand-segmented and aligned gold data for Finnish, Swedish, Spanish (a few hundred word forms each) from CoNLL UD shared task this year



# Evaluation

	<b>Finnish</b>	<b>Spanish</b>	<b>Swedish</b>
Recall	87.43	84.38	88.71
Precision	94.63	88.63	94.01
F <sub>1</sub> -score	<b>90.89</b>	<b>86.45</b>	<b>91.28</b>

	Morfessor baseline		
Recall	80.65	81.32	90.82
Precision	76.92	73.64	75.58
F <sub>1</sub> -score	78.74	77.29	82.50

Morpheme boundaries



# Evaluation

	<b>Finnish</b>	<b>Spanish</b>	<b>Swedish</b>
Recall	62.79	50.10	55.87
Precision	71.06	54.22	61.82
F <sub>1</sub> -score	<b>66.67</b>	<b>52.08</b>	<b>58.69</b>
Morfessor baseline			
Recall	30.51	25.93	44.13
Precision	28.45	22.24	32.92
F <sub>1</sub> -score	29.45	23.94	37.71

Unlabeled morphemes




# Evaluation

	<b>Finnish</b>	<b>Spanish</b>	<b>Swedish</b>
Recall	80.07	73.49	88.26
Precision	90.62	79.54	97.66
F <sub>1</sub> -score	<b>85.02</b>	<b>76.39</b>	<b>92.73</b>
Morfessor baseline			
Recall	74.96	48.34	83.10
Precision	69.90	41.47	62.00
F <sub>1</sub> -score	72.34	44.64	71.01

Labeled morphemes

# Wrapup

- Weak supervision helps in inducing segmentation and allomorph labeling
- Can be run on all UD languages
- Gives a consistent segmentation & labeling
- Code at <https://github.com/mpsilfve/ud-segmenter> 
- Errors (differences to linguist-preferred gold standard) remain, some due to objective function, some probably due to well-known differing linguistic notions about gold standard (is it announc+**ed** or announce+**d**?), frequency effects
- Future work: evaluate a range of objective functions, implement raw allomorph minimization

# Thank You