

# Sequence Tagging for Verb Conjugation in Romanian

Liviu P. Dinu   Vlad Niculae   Octavia-Maria Şulea

Center for Computational Linguistics  
University of Bucharest  
<http://nlp.unibuc.ro>

September 2013

# Verbs in Romanian

## Regularity is not black and white

		1 <sup>st</sup>	2 <sup>nd</sup>	3 <sup>rd</sup>
Regular	sg.	merg	mergi	merge
a merge ( <i>to walk</i> )	pl.	mergem	mergeți	merg
Irregular	sg.	sunt	ești	este
a fi ( <i>to be</i> )	pl.	suntem	sunteți	sunt

# Verbs in Romanian

## Regularity is not black and white

		1 <sup>st</sup>	2 <sup>nd</sup>	3 <sup>rd</sup>
Regular	sg.	merg	mergi	merge
a merge ( <i>to walk</i> )	pl.	mergem	mergeți	merg
Irregular	sg.	sunt	ești	este
a fi ( <i>to be</i> )	pl.	suntem	sunteți	sunt
Partially irregular	sg.	port	porți	poartă
a purta ( <i>to wear</i> )	pl.	purțăm	purtați	poartă

Dinu et al, RANLP 2011, EACL 2012

- Hand-crafted sets of regular expressions fully describing conjugation of most verbs
- Predictive model  $h(\text{infinitive}) = \text{regular expression set}$

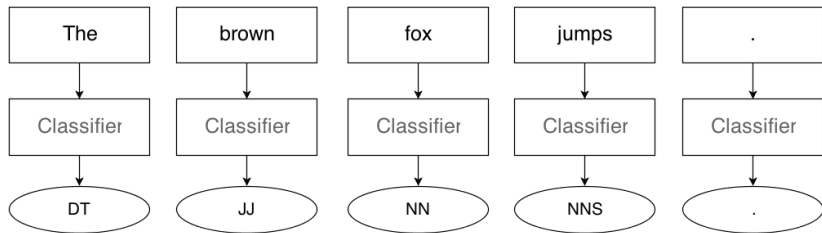
## Running example

	sg.	port	porți	poartă
a purta ( <i>to wear</i> )	pl.	purtăm	purtați	poartă

## Regular expression set

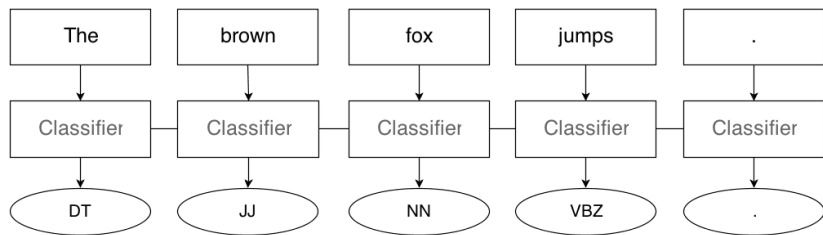
sg.	$\text{^\text{.}(\text{.}^*)\text{o}(\text{.}^*)\text{t}\text{\$}$	$\text{^\text{.}(\text{.}^*)\text{o}(\text{.}^*)\text{\text{ț}}\text{i}\text{\$}$	$\text{^\text{.}(\text{.}^*)\text{o}\text{a}(\text{.}^*)\text{t}\text{ă}\text{\$}$
pl.	$\text{^\text{.}(\text{.}^*)\text{u}(\text{.}^*)\text{t}\text{ă}\text{m}\text{\$}$	$\text{^\text{.}(\text{.}^*)\text{u}(\text{.}^*)\text{t}\text{a}\text{\text{ț}}\text{i}\text{\$}$	$\text{^\text{.}(\text{.}^*)\text{o}\text{a}(\text{.}^*)\text{t}\text{ă}\text{\$}$

# Sequence tagging: POS tagging example



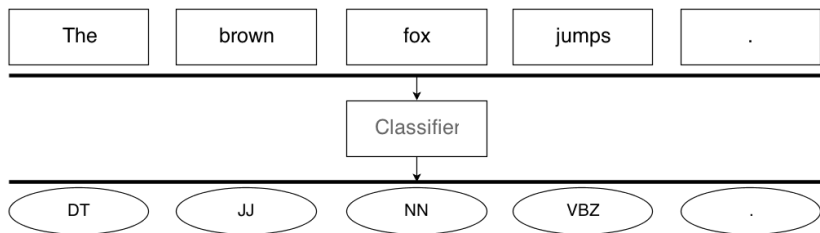
$$\prod \phi(y_i, x_i)$$

# Sequence tagging: POS tagging example (better)



$$\prod \phi_1(y_i, x_i) \phi_2(y_i, y_{i+1})$$

# Sequence tagging: POS tagging example (worse?)



$$\phi(y_1, y_2, \dots, y_n, x_1, x_2, \dots, x_n)$$

# Ignored structure: interaction between classes

a cânta	a deștepta	a deșerta
<i>to sing</i>	<i>to rise</i>	<i>to empty</i>
$\text{^}(\text{.}^*)\text{t}\text{\$}$	$\text{^}(\text{.}^*)\text{e}(\text{.}^*)\text{t}\text{\$}$	$\text{^}(\text{.}^*)\text{e}(\text{.}^*)\text{t}\text{\$}$
$\text{^}(\text{.}^*)\text{ț}\text{i}\text{\$}$	$\text{^}(\text{.}^*)\text{e}(\text{.}^*)\text{ț}\text{i}\text{\$}$	$\text{^}(\text{.}^*)\text{e}(\text{.}^*)\text{ț}\text{i}\text{\$}$
$\text{^}(\text{.}^*)\text{t}\text{ă}\text{\$}$	$\text{^}(\text{.}^*)\text{ea}(\text{.}^*)\text{t}\text{ă}\text{\$}$	$\text{^}(\text{.}^*)\text{a}(\text{.}^*)\text{t}\text{ă}\text{\$}$
$\text{^}(\text{.}^*)\text{t}\text{ă}\text{m}\text{\$}$	$\text{^}(\text{.}^*)\text{e}(\text{.}^*)\text{t}\text{ă}\text{m}\text{\$}$	$\text{^}(\text{.}^*)\text{e}(\text{.}^*)\text{t}\text{ă}\text{m}\text{\$}$
$\text{^}(\text{.}^*)\text{ta}\text{ț}\text{i}\text{\$}$	$\text{^}(\text{.}^*)\text{e}(\text{.}^*)\text{ta}\text{ț}\text{i}\text{\$}$	$\text{^}(\text{.}^*)\text{e}(\text{.}^*)\text{ta}\text{ț}\text{i}\text{\$}$
$\text{^}(\text{.}^*)\text{t}\text{ă}\text{\$}$	$\text{^}(\text{.}^*)\text{ea}(\text{.}^*)\text{t}\text{ă}\text{\$}$	$\text{^}(\text{.}^*)\text{a}(\text{.}^*)\text{t}\text{ă}\text{\$}$



# Conjugation as sequence tagging

## Running example

	sg.	port	porți	poartă
a purta ( <i>to wear</i> )	pl.	purtăm	purtați	poartă

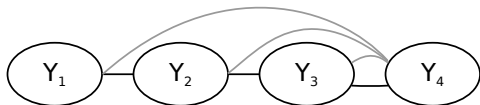
## Variable letters (Moisil)

$\text{form}(u_0 1sg) =$	o	$\text{form}(t_0 1sg) =$	t
$\text{form}(u_0 3sg) =$	oa	$\text{form}(t_0 2sg) =$	ț
$\text{form}(u_0 1pl) =$	u		

## Tagging example

p	u	r	t	a
0	$u_0$	0	$t_0$	$T_4$

- Features: character n-grams to the left and right size up to  $n$
- Dataset: RoMorphoDict (lemmas and forms) labeled using the RegEx sets  
16 ending patterns, 17 variable letters  
4,699 train / 2,257 test / 339 unlabeled
- Grid search, 10-fold cross validation



- An extra factor template allowing the ending to influence all positions
- Inference becomes more complex
- Out-of-the-box sequence tagging no longer appropriate

method	Cross-val. accuracy			Test accuracy		
	word	char	char'	word	char	char'
SVM	0.886	-	-	0.896	-	-
ML	0.924	0.987	0.913	0.914	0.985	0.900
AP	0.923	0.987	0.917	0.912	0.985	0.900
PA	0.925	0.987	0.917	0.912	0.984	0.900
AROW	0.916	0.986	0.912	0.908	0.984	0.895
SKIP	-	0.984	-	0.906	0.983	0.896

Generalization on 105 of the unlabeled verbs:

- many termination patterns are correctly found (30)
- some alternations are found (3)