

Maschinelles Übersetzen natürlicher Sprachen

8. Übungsblatt

2015-01-08

Aufgabe 1 (Extraktion von SCFG-Regeln)

Bestimmen Sie in folgendem alinierten Satzpaar die initial phrase pairs.

	der	Hund	hat	den	Knochen	heute	noch	nicht	gefunden
today						×			
the	×								
dog		×							
has			×						
not								×	
found									×
the				×					
bone					×				
yet							×		

Geben Sie einige der SCFG-Regeln an, die aus diesen Phrasenpaaren extrahiert werden.

Aufgabe 2 (Training von SCFGs)

Betrachten Sie die folgenden beiden alinierten Satzpaare.

$e_1 \setminus f_1$	Kinder	spielen	$e_2 \setminus f_2$	lasst	uns	spielen
children	×		let's	×	×	
play		×	play			×

Ihre Korpushäufigkeiten seien gegeben durch $h(e_1, f_1) = 1$ und $h(e_2, f_2) = 2$ (und für alle Satzpaare $(e, f) \notin \{(e_1, f_1), (e_2, f_2)\}$ soll $h(e, f) = 0$ gelten).

Extrahieren Sie aus dem so gegebenen Korpus SCFG-Regeln und trainieren Sie deren Wahrscheinlichkeiten mittels der in der Vorlesung vorgestellten Methode!

Aufgabe 3

Wir betrachten noch einmal die Übersetzung von arithmetischen Ausdrücken in Infix-Notation in umgekehrte polnische Notation (vgl. 7. Übungsblatt, Aufgabe 3). Entwerfen Sie einen XTT \mathcal{M} , welcher einen Parsebaum $\xi \in PT_G$ eines Infix-Ausdrucks in den entsprechenden Parsebaum einer CFG G' für RPN-Ausdrücke übersetzt!

Zusatzaufgabe: Entwerfen Sie außerdem einen yXTT \mathcal{M}' , welcher direkt den Yield des abgeleiteten Baumes ausgibt!

Aufgabe 4

Gegeben seien die Alphabete

$$\Gamma = \{S, NP, VP, ADJ, NN, VB, DT\},$$

sowie

$$\Sigma = \Gamma \cup \{\text{el, perro, bravo, escucha, coge, hombre}\}$$

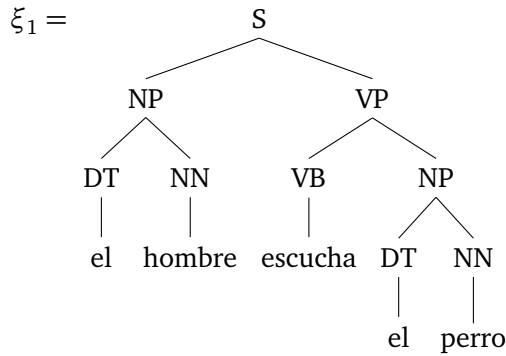
und

$$\Delta = \Gamma \cup \{\text{der, den, Hund, wilde, wilden, Mann, fängt, hört}\}.$$

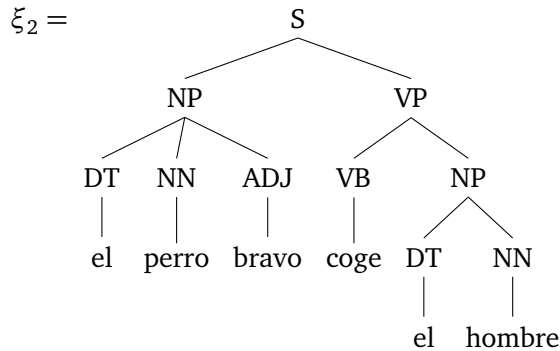
Wir nehmen einen extended tree transducer $\mathcal{M} = (Q, \Sigma, \Delta, q, R)$ an mit der Zustandsmenge $Q = \{q, n, a\}$ und der Regelmenge R wie folgt:

$q(S(x_1 : NP, VP(x_2 : VB, x_3 : NP))) \rightarrow S(n(x_1), VP(q(x_2), a(x_3)))$	$q(NN(\text{perro})) \rightarrow NN(\text{Hund})$
$n(NP(DT(\text{el}), x_1 : NN, x_2 : ADJ)) \rightarrow NP(DT(\text{der}), n(x_2), q(x_1))$	$q(NN(\text{hombre})) \rightarrow NN(\text{Mann})$
$a(NP(DT(\text{el}), x_1 : NN, x_2 : ADJ)) \rightarrow NP(DT(\text{den}), a(x_2), q(x_1))$	$q(VB(\text{escucha})) \rightarrow VB(\text{hört})$
$n(NP(DT(\text{el}), x_1 : NN)) \rightarrow NP(DT(\text{der}), q(x_1))$	$q(VB(\text{coge})) \rightarrow VB(\text{fängt})$
$a(NP(DT(\text{el}), x_1 : NN)) \rightarrow NP(DT(\text{den}), q(x_1))$	$a(ADJ(\text{bravo})) \rightarrow ADJ(\text{wilden})$
	$n(ADJ(\text{bravo})) \rightarrow ADJ(\text{wilde})$

Bestimmen Sie, unter Angabe der entsprechenden Ableitungen, die Übersetzungen der Parsebäume



und



Vollziehen Sie für eine dieser Ableitungen die Arbeitsweise der Projektionsfunktionen $\pi_F: D_{\mathcal{M}} \rightarrow T_{\Sigma}$ sowie $\pi_E: D_{\mathcal{M}} \rightarrow T_{\Delta}$ nach! Wie müssten die Regeln von \mathcal{M} geändert werden, damit der Eingabe- bzw. Ausgabebaum einer Ableitung nicht eindeutig bestimmt ist?