

Einführung in die Computerlinguistik

Hausaufgabe 5, Abgabe 14.05.2012

Laura Kallmeyer

SS 2012, Heinrich-Heine-Universität Düsseldorf

Aufgabe 1 Gegeben sei das folgende POS-getaggte Mini-Korpus:

Satz 1: light the bright fires .
 V Det A N PUNCT

Satz 2: fires are nice .
 N V A PUNCT

Satz 3: the light is bright .
 Det N V A PUNCT

1. Wie sehen Zustands- und Ausgabemenge für ein aus diesem Korpus gelerntes Bigram-HMM aus?
2. Geben Sie für ein aus diesem Korpus gelerntes Bigram-HMM die Übergangswahrscheinlichkeiten und die Ausgabewahrscheinlichkeiten an. (Es reichen die Brüche, die Dezimalzahlen müssen nicht angegeben werden.)

Lösung:

1. $Q = \{q_0, q_f, V, Det, A, N, PUNCT\}$, $O = \{\text{light, the, bright, fires, are, nice, is, .}\}$.

2. Übergangswahrscheinlichkeiten:

$$\begin{array}{llll}
 P(Det|V) = \frac{1}{3} & P(Det|Det) = 0 & P(Det|A) = 0 & P(Det|N) = 0 \\
 P(A|V) = \frac{2}{3} & P(A|Det) = \frac{1}{2} & P(A|A) = 0 & P(A|N) = 0 \\
 P(N|V) = 0 & P(N|Det) = \frac{1}{2} & P(N|A) = \frac{1}{3} & P(N|N) = 0 \\
 P(V|V) = 0 & P(V|Det) = 0 & P(V|A) = 0 & P(V|N) = \frac{2}{3} \\
 P(PUNCT|V) = 0 & P(PUNCT|Det) = 0 & P(PUNCT|A) = \frac{2}{3} & P(PUNCT|N) = \frac{1}{3}
 \end{array}$$

PUNCT steht mit Wahrscheinlichkeit 1 am Satzende, alle anderen Nachfolgetags für PUNCT haben Wahrscheinlichkeit 0.

Für alle anderen Tags ist die Wahrscheinlichkeit, am Satzende zu stehen, 0.

N, V und Det stehen jeweils mit Wahrscheinlichkeit $\frac{1}{3}$ am Satzanfang, alle anderen Tags mit Wahrscheinlichkeit 0.

Ausgabewahrscheinlichkeiten:

	light	the	bright	fires	are	nice	is	.
Det	0	1	0	0	0	0	0	0
N	$\frac{1}{3}$	0	0	$\frac{2}{3}$	0	0	0	0
A	0	0	$\frac{2}{3}$	0	0	$\frac{1}{3}$	0	0
V	$\frac{1}{3}$	0	0	0	$\frac{1}{3}$	0	$\frac{1}{3}$	0
PUNCT	0	0	0	0	0	0	0	1

Aufgabe 2 Nehmen Sie an, Sie haben einen HMM-POS Tagger, unter anderem mit folgenden Wahrscheinlichkeiten:

Emissionswahrscheinlichkeiten:

$$\begin{array}{lll}
 P(\text{the}|Det) = 1 & P(\text{light}|N) = 3 \cdot 10^{-3} & P(\text{fires}|N) = 5 \cdot 10^{-3} \\
 & P(\text{light}|Adj) = 3 \cdot 10^{-3} & P(\text{fires}|V) = 3 \cdot 10^{-3} \\
 & P(\text{light}|V) = 2 \cdot 10^{-3} &
 \end{array}$$

Alle anderen Emissionswahrscheinlichkeiten für light, the und fires seien 0.

Übergangswahrscheinlichkeiten:

$$\begin{aligned}
 P(N|Det) &= 5 \cdot 10^{-1} & P(N|N) &= 1 \cdot 10^{-1} & P(N|Adj) &= 5 \cdot 10^{-1} & P(N|V) &= 2 \cdot 10^{-1} \\
 P(Adj|Det) &= 3 \cdot 10^{-1} & P(V|N) &= 4 \cdot 10^{-1} & P(V|Adj) &= 1 \cdot 10^{-1} & P(V|V) &= 2 \cdot 10^{-1} \\
 P(V|Det) &= 1 \cdot 10^{-1} & & & & & &
 \end{aligned}$$

Angenommen, die Wahrscheinlichkeit, dass ein Det am Satzanfang steht, ist 1, die, dass auf ein N oder V ein Satzende folgt, ist jeweils $0.1 = 1 \cdot 10^{-1}$.

1. Geben Sie die Viterbi Matrix an, die sich bei diesen Wahrscheinlichkeiten für die Eingabe the light fires ergibt. Es reicht, die Einträge anzugeben, die $\neq 0$ sind. Geben Sie für jedes Feld Ihren Rechenweg an.
2. Was ist die beste POS-TAG Sequenz, die sich als Ergebnis für the light fires ergibt?

Lösung:

q_F				225 · 10 ⁻⁹ , N
N		15 · 10 ⁻⁴ , Det	225 · 10 ⁻⁸ , Adj	
V		2 · 10 ⁻⁴ , Det	180 · 10 ⁻⁸ , N	
1. Adj		9 · 10 ⁻⁴ , Det		
Det		1, q_0		
		1	2	3
		the	light	fires

light, Adj: $1 \cdot P(Adj|Det) \cdot P(light|N) = 1 \cdot 3 \cdot 10^{-1} \cdot 3 \cdot 10^{-3}$

light, N: $1 \cdot 5 \cdot 10^{-1} \cdot 3 \cdot 10^{-3}$

light, V: $1 \cdot 1 \cdot 10^{-1} \cdot 2 \cdot 10^{-3}$

fires, N: $\max\{15 \cdot 10^{-4} \cdot 1 \cdot 10^{-1} \cdot 5 \cdot 10^{-3} \text{ Vorgänger N}, 9 \cdot 10^{-4} \cdot 5 \cdot 10^{-1} \cdot 5 \cdot 10^{-3} \text{ Vorgänger Adj}, 2 \cdot 10^{-4} \cdot 2 \cdot 10^{-1} \cdot 5 \cdot 10^{-3} \text{ Vorgänger V}\}$

fires, V: $\max\{15 \cdot 10^{-4} \cdot 4 \cdot 10^{-1} \cdot 3 \cdot 10^{-3} \text{ Vorgänger N}, 9 \cdot 10^{-4} \cdot 1 \cdot 10^{-1} \cdot 3 \cdot 10^{-3} \text{ Vorgänger Adj}, 2 \cdot 10^{-4} \cdot 2 \cdot 10^{-1} \cdot 3 \cdot 10^{-3} \text{ Vorgänger V}\}$

2. Die beste POS-TAG Folge ist demnach Det Adj N.