

# Przetwarzanie Języka Naturalnego

## Lab 5 – Model języka

Aleksander Smywiński-Pohl  
apohl1o@agh.edu.pl

Wydział IEiT  
Katedra Informatyki

28.03.2017



Łańcuchem Markowa nazywamy proces stochastyczny w czasie dyskretnym i o dyskretnym zbiorze stanów bez pamięci.

Brak pamięci oznacza, że kolejny stan osiągany przez system zależy wprost jedynie od stanu bezpośrednio go poprzedzającego (oraz czasu).



Ciąg zmiennych losowych  $X_0, X_1 \dots$  o wartościach całkowitych.

Warunek Markowa:

$$\forall n \in \mathbb{N} \forall i_0, i_1 \dots i_n, i_{n+1} \in \mathbb{Z} :$$

$$P(X_{n+1} = i_{n+1} | X_0 = i_0 \wedge X_1 = i_1 \wedge \dots \wedge X_n = i_n) =$$

$$P(X_{n+1} = i_{n+1} | X_n = i_n) \stackrel{\text{ozn.}}{=} p(i_n, i_{n+1}) \stackrel{\text{ozn.}}{=} p_{i_n i_{n+1}}$$

Łańcuch (czasowo) jednorodny - prawdopodobieństwo nie zależy od  $n$



AGH

## Macierz przejścia

$$P = \begin{bmatrix} p_{0,0} & p_{0,1} & \cdots & p_{0,m} \\ p_{1,0} & p_{1,1} & \cdots & p_{1,m} \\ \vdots & \vdots & \ddots & \vdots \\ p_{m,0} & p_{m,1} & \cdots & p_{m,m} \end{bmatrix}$$

$$X^{(n)} = X^{(n-1)} * P$$

$$X^{(n)} = X^{(0)} * P^n$$



Łańcuch Markowa rzędu  $m$  (z pamięcią  $m$ )

$$P(X_{n+1} = i_{n+1} | X_0 = i_0 \wedge X_1 = i_1 \wedge \dots \wedge X_n = i_n) = \\ P(X_{n+1} = i_{n+1} | X_{n-m+1} = i_{n-m+1} \wedge \dots \wedge X_n = i_n) \quad m \in \mathbb{Z}_+$$

$$Y_n = (X_n, X_{n-1}, X_{n-2}, \dots, X_{n-m+1})$$

- ✠ modelowanie procesów fizycznych (termodynamika, mechanika)
- ✠ chemia
- ✠ modelowanie języka
- ✠ generowanie tekstów, muzyki
- ✠ teoria gier, sztuczna inteligencja
- ✠ ekonomia



$$P(h_i) = P(h_i)_{LM}^\alpha * P(h_i)_{AM}$$

$$P(h_i)_{LM} = \prod_{w_j \in h_i} P(w_j | w_{j-N+1}, \dots, w_{j-1})$$

Np. dla trigramów:

$$P(h_i)_{LM_{trigram}} = \prod_{w_j \in h_i} P(w_j | w_{j-2}, w_{j-1})$$

W praktyce:

$$\log P(h_i)_{LM_{trigram}} = \sum_{w_j \in h_i} \log P(w_j | w_{j-2}, w_{j-1})$$



AGH

## Prawdopodobieństwo n-gramów

-2.387182	abażurze kinkietu	
-1.753636	abażurze lampy	-0.4198931
-2.387182	abażurze lampy.	
-2.387182	abażurze ma	-0.07864241
-2.387182	abażurze marnowanie	
-1.490232	abażurze na	
-2.387182	abażurze nie	
-2.387182	abażurze nocnej	

⊠ logarytm o podstawie 10

⊠ wartość na końcu to tzw. *back-off weight*





AGH

## Użycie *back-off weight* dla n-gramów

$$P(w_j|w_{j-2}, w_{j-1})_{bow} = \begin{cases} P(w_j|w_{j-2}, w_{j-1}) \\ BOW(w_{j-2}, w_{j-1}) * P(w_j|w_{j-1}) \end{cases} \quad P > 0$$

- 1 Zaimplementować algorytm obliczający Word Error Rate, jako odległość Levenshteina dla słów (0,5 pkt.)
- 2 Zaimplementować funkcję obliczającą prawdopodobieństwa zdania na podstawie otrzymanego pliku z modelem językowym (1,5 pkt.)
- 3 Określić optymalną wartość parametr  $\alpha$  na podstawie akustycznych hipotez rozpoznawania mowy oraz modelu językowego (1 pkt.)

Materiały: <http://apohllo.pl/text/lab5.tar.gz>