

dujice z odklínajících textů

$$H = - \sum_{\substack{i \in I_1 \\ j \in J_1}} p(i,j) \cdot \log p(j|i) - \sum_{\substack{i \in I_2 \\ j \in J_2}} p(i,j) \cdot \log p(j|i) - \sum_{\substack{i \in I_1, j \in J_2}} p(i,j) \cdot \log p(i|j)$$

$$- \sum_{i \in I_2, j \in J_1} p(i,j) \cdot \log p(i|j)$$

Sumu lze zjednodušit, a to tak, že jediné dvojice

slov z odklínajících textů se může vyskytovat na rozdílu.

Označme tyto slova  $x_i \in T_1, x_j \in T_2$

$$= - \sum_{\substack{i \in I_1 \\ j \in J_1}} p(i,j) \cdot \log p(j|i) - \sum_{\substack{i \in I_2 \\ j \in J_2}} p(i,j) \cdot \log p(j|i) - p(x_i, x_j) \cdot \log(p(x_j|x_i))$$

$$\hat{p}(i,j) = \frac{c(i,j)}{\# \text{dv. } T_1 + \# \text{dv. } T_2 + 1} \Rightarrow \hat{p}(i,j) \leq p(i,j)$$

$$\hat{p}(j|i) = \hat{p}(j,i) / \hat{p}(i) = \frac{c(j,i)}{\# \text{dv. } T_1 + \# \text{dv. } T_2 + 1} \cdot \frac{\# \text{slv. } T_1 + \# \text{slv. } T_2}{c(i)}$$

$$\#\text{slv} - 1 = \#\text{dvojice} \Rightarrow \#\text{slv} \sim \#\text{dvojice} \Rightarrow \hat{p}(j|i) \doteq \frac{c(j,i)}{c(i)}$$

$\hookrightarrow$  uvažujte ještě vyskytující se dvojice

Nárovní se mohou stát následující situace

( $X_i :=$  poslední slovo v  $T_1, X_j :=$  první slovo v  $T_2$ )

-  $X_i$  se nachází pouze na hranici  $T_1$ :  $\exists i \in I_1 : i \neq x_i$

pak výpočet  $-\sum_{\substack{i \in I_1 \\ j \in J_1}} p(i,j) \cdot \log \hat{p}(j|i)$  není mít mít význam,

dávavši  $p(X_j|X_i) = 1 \Rightarrow \log(p(X_j|X_i)) = 0 \Rightarrow$  zelená část se nezmění,

žlutá část bude 0.

-  $X_i$  se nachází i jinde v textu  $T_1$ :  $\exists i \in I_1 : i = x_i$

pak ve výpočtu  $-\sum_{\substack{i \in I_1 \\ j \in J_1}} p(i,j) \cdot \log \hat{p}(j|i)$  lze ne hodnotu  $\hat{p}(j|i)$  pro  $i = x_i$ ,

tedy  $\log(\hat{p}(j|i))$  se zvýší, protože i hodnota entropy.

Zároveň  $p(X_j|X_i) < 1 \Rightarrow \log(p(X_j|X_i)) > 0 \Rightarrow$  zelená část se zvýší

žlutá část bude mít  $\neq 0$  hodnotu.

$$\Rightarrow = - \sum_{\substack{i \in I_1 \\ j \in J_1}} p(i,j) \cdot \log p(j|i) - \sum_{\substack{i \in I_2 \\ j \in J_2}} p(i,j) \cdot \log p(j|i) - p(x_i, x_j) \cdot \log(p(x_j|x_i))$$

$$P_1(i,j) = \frac{\# dr. T_1}{\# dr. T_1 + \# dr. T_2} \cdot P(i,j)$$

$$P_2(i,j) = \frac{\# dr. T_2}{\# dr. T_1 + \# dr. T_2} \cdot P(i,j)$$

$$= - \sum_{\substack{i \in I_1 \\ j \in J_1}} \frac{\# dr. T_1}{\# dr. T_1 + \# dr. T_2} \cdot P(i,j) \cdot \log p(j|i) - \sum_{\substack{i \in I_2 \\ j \in J_2}} \frac{\# dr. T_2}{\# dr. T_1 + \# dr. T_2} \cdot P(i,j) \cdot \log p(j|i) - p(x_i, x_j) \cdot \log(p(x_j|x_i))$$

$$- \frac{\# dr. T_1}{\# dr. T_1 + \# dr. T_2} \cdot E + \frac{\# dr. T_2}{\# dr. T_1 + \# dr. T_2} \cdot E - p(x_i, x_j) \cdot \log(p(x_j|x_i))$$

$$= E \cdot \left( \frac{\# dr. T_1}{\# dr. T_1 + \# dr. T_2} + \frac{\# dr. T_2}{\# dr. T_1 + \# dr. T_2} \right) - p(x_i, x_j) \cdot \log(p(x_j|x_i))$$

V případě, že žlutá část  $\Rightarrow$  ( $x_i$  se vyskytuje pravděpodobně v  $T_1$ ), je  $H(T_1, T_2) = E$   
a entropie se nezměnila.

V opačném případě v zelené části pro  $i = x_i$ :

$$P(j|i) = \frac{C(j|i)}{C(i)+1} = \frac{C(j|i)}{C(i)} \cdot \frac{C(i)}{C(i)+1} = P(j|i) \cdot \frac{C(i)}{C(i)+1}$$

$$\text{tedy zelená část se zmenší} \Rightarrow \frac{\# dr. v T_1}{\# dr. v T_1 + \# dr. v T_2} \cdot P(i,j) \cdot \log(P(j|i) \cdot \frac{1}{C(i)+1}) = A$$

a do celkové souhytu se přidá i žlutý člen

$$- p(x_i, x_j) \cdot \log(p(x_j|x_i)) = B$$

Proto platí:  $> 0 \Rightarrow$  entropie se zvýší

$(A - B) < 0 \Rightarrow$  entropie se sníží

Cí hýbá měl ještě učebat!

Uniform prob. by méně významných slov má některé  
slova významnější (více v reálném jazyce než v uniformních slovech)

Takhle máme jeho uniform prob.:=  $\frac{1}{|T| + \alpha(|H| - T \cap H)}$   $\alpha \geq 1$

Toto  $\alpha$  je potřeba si  
vytunovat z postr.

$T$  ... fázis voc.

$H$  ... herciovský

$|H| - |T \cap H| \rightarrow$  počet "nových" slov, která  
nebyly v fázisových datotech.