

Wk.. 9.4

$$\text{a) } N(t, c) = N_{t1}$$

Term t kommt in Klasse C vor

$N_{t, c}$   
 $\uparrow \uparrow$   
 Term in Dokument  
 Dokument in Klasse  
 $\bar{N}_{t,c}$  absolute  
 Werte

$N_{t,0}$	... fed, dec
$N_{t,0}$	... fed, dec
$N_{t,1}$	... fed, dec
$N_{t,1}$	... fed, dec

$$N(\text{Brazil}, c) = 51 *$$

$$N(\text{Booted}, c) = 10$$

$$N(\text{producer}, c) = 34 *$$

Fair Feature Selection  
 beide (\*) auswählbar

$$I(U; C) = \sum_{e \in \{0,1\}} \sum_{c \in \{0,1\}} P(U=e, C=c) \cdot \log_2 \frac{P(U=e, C=c)}{P(U=e) \cdot P(C=c)}$$

$$I(U; C) = \frac{N_{t1}}{N} \log_2 \frac{N_{t1}}{N} / \left( \frac{N_{t1}}{N} \cdot \frac{N_{t0}}{N} \right)$$

✓  $N_{t1} + N_{t0}$   $N_{t1} + N_{t0}$

Wk Term in Dokument  
 und Dokument in Klasse

$$= \frac{N_{t1}}{N} \log_2 \frac{N_{t1}}{N} / \left( \frac{N_{t1}}{N} \cdot \frac{N_{t0}}{N} \right)$$

$$+ \frac{N_{t0}}{N} \log_2 \frac{N_{t0}}{N} / \left( \frac{N_{t1}}{N} \cdot \frac{N_{t0}}{N} \right)$$

$$+ \frac{N_{t0}}{N} \log_2 \frac{N_{t0}}{N} / \left( \frac{N_{t1}}{N} \cdot \frac{N_{t0}}{N} \right)$$

$$= \frac{N_{t1}}{N} \log_2 \frac{N_{t1} \cdot N}{N_{t1} \cdot N_{t0}} + \frac{N_{t0}}{N} \log_2 \frac{N_{t0} \cdot N}{N_{t1} \cdot N_{t0}}$$

$$+ \frac{N_{t0}}{N} \log_2 \frac{N_{t1} \cdot N}{N_{t1} \cdot N_{t0}} + \frac{N_{t0}}{N} \log_2 \frac{N_{t0} \cdot N}{N_{t1} \cdot N_{t0}}$$

$$N_{t1} = \frac{0,00748 + 0,000459}{0,00748 + 0,000459} = 0,97847 + 0,000459 = 0,97847$$

$$N_{t0} = 1 - 0,97847 = 0,02152$$

$$N_{t0} = \frac{0,00748 + 0,000459}{0,00748 + 0,000459} = 0,97847$$

$$N_{t0} = \frac{0,00748 + 0,000459}{0,00748 + 0,000459} = 0,97847$$

$$\Rightarrow 1+0$$

$$(L) \quad N_{t1} = N_{t0} + N_{t1}$$

$$N_{t0} = N_{t0} + N_{t0}$$

$$N_{t0} = N_{t0} + N_{t0}$$

$$N_{t1} = N_{t1} + N_{t1}$$

$\Rightarrow I(U; C) = 0,00704$

$\Rightarrow \text{producer, Brazil}$

15.2

$$\hat{P}(f = \text{'the'} | c) \approx 0,08$$

Termfrequenz

tf 'the' =

Dokufrequenz

df 'the' =

15.3

a) China china

$$\text{NB: } T_{\text{NB}}(d) = \underset{G \in C}{\operatorname{argmax}} P(G|d)$$

$$\text{NB: } P(G_j|d) \propto P(G_j) \prod_{t \in d} \text{tf}_{G_j}(t)$$

$$P(G_j) = \frac{q_j}{|D|} \quad \begin{matrix} \text{proportional} \\ \text{with Token} \end{matrix}$$

$$P(\text{China}) = \frac{2}{4} = \frac{1}{2}$$

$$P(\overline{\text{China}}) = \frac{2}{4} = \frac{1}{2}$$

absolute tf Doku in Klasse drin

$$\hat{P}(\text{tf}_{G_j}(t)) = \frac{\text{tf}_{G_j}(t) + 1}{\sum_i \text{tf}_{G_j}(t_i) + V} \quad \begin{matrix} \rightarrow \text{kleinste WSK - Smoothing} \\ V=7 \\ \text{Token mit } P(G_j) \end{matrix}$$

Anteil der TF in gesamte (jedes Mal in Doku c=China)

$$\hat{P}(\text{Taiwan}|\text{China}) = \frac{2+1}{5+7} = \frac{1}{4}$$

$$P(\overline{\text{Taiwan}}|\text{China}) = \frac{0+7}{5+7} = \frac{7}{12} \quad \begin{matrix} \text{(Gesamte Tokens für } c=\text{China)} \\ \text{Sapporo} \end{matrix}$$

$$\hat{P}(\text{Taiwan}|\overline{\text{China}}) = \frac{7+1}{5+7} = \frac{1}{6}$$

$$\hat{P}(\text{Sapporo}|\text{China}) = \frac{2+7}{5+7} = \frac{1}{4} \quad \begin{matrix} \text{TF in } d_5 \\ \text{TF in } d_5 \end{matrix}$$

$$P(\text{China}|d_5) = \hat{P}(\text{China}) \cdot P(\text{Taiwan}|\text{China})^2 \cdot P(\text{Sapporo}|\text{China})^2$$

$$= \frac{1}{2} \cdot \left(\frac{1}{4}\right)^2 \cdot \left(\frac{1}{12}\right)^2 = \frac{1}{384}$$

$$P(\overline{\text{China}}|d_5) = \hat{P}(\overline{\text{China}}) \cdot P(\overline{\text{Taiwan}}|\text{China})^2 \cdot P(\overline{\text{Sapporo}}|\text{China})^2$$

$$= \frac{1}{2} \cdot \left(\frac{1}{6}\right)^2 \cdot \left(\frac{7}{12}\right)^2 = \frac{1}{288}$$

$$P(\text{China}|d_5) > P(\overline{\text{China}}|d_5) \Rightarrow d_5 = \text{China}$$

b) BNB:

$$P(G_j|d) \propto P(G_j) \cdot \prod_{t \in d} P(\text{tf}_{G_j}(t)) \cdot \prod_{t \notin d} P(\overline{\text{tf}}_{G_j}(t))$$

 $\Rightarrow$  Gegenwart, dass TF nicht vorkommen

$$b) \bar{x}_{k,i} \geq \frac{\alpha}{n}$$

$$\sum_{i=1}^n \bar{x}_{k,i} \geq \alpha$$

$$\sum_{i=1}^n \bar{x}_{k,i} \geq 1$$

$$\max d = 1 - \alpha$$

Maximalabstand

$$\lim_{x \rightarrow 1} \max d = 0$$

$$NB: \frac{\alpha}{n} + 1 - \alpha \stackrel{\text{hier:}}{=} \text{Ausgangspkt} = 1$$

$$\alpha = 1$$

$$P = \begin{pmatrix} A & B & C \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{pmatrix}$$

Trick: bei 3 Fälle  
 $\alpha=1$      $\alpha=1$   
 $\frac{1}{n}$      $\frac{1}{n}$

14.3

$$\alpha = 0,1 = \frac{1}{10} \quad n = 3$$

Page Rank

$$P = \begin{pmatrix} \frac{1}{30} & \frac{29}{60} & \frac{29}{60} \\ \frac{1}{30} & \frac{1}{30} & \frac{1}{15} \\ \frac{1}{30} + \frac{1}{2} & \frac{1}{30} & \frac{9}{10} \\ \frac{1}{3} + \frac{1}{10} & \frac{1}{30} & \frac{1}{30} \end{pmatrix}$$

$$\vec{x}_k = P^T \cdot \vec{x}_{k-1} \quad \vec{x}_0 = \begin{pmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{pmatrix} \text{ gleich wahrscheinlich}$$

bis Konvergenz

$$\vec{x}_k = \begin{pmatrix} \frac{1}{30} \\ \frac{29}{60} \\ \frac{60}{60} \end{pmatrix} \quad \text{wie kommen Vom Server zu A?} \\ \text{durch Transportieren } i. (i,j) \notin E \text{ anwenden} \\ \Rightarrow \frac{1}{3} = \frac{1}{30}$$

WSK bei B,C bleibt gleich  $\Rightarrow (1 - \frac{1}{30}) / 2 = \frac{29}{58}$

$$\text{HITS: } A(\vec{i}) = \sum_{\substack{j \in V \setminus \{\vec{i}\} \\ \text{ohne } \vec{i}}} H(j)$$

$$H(\vec{i}) = \sum_{\substack{j \in V \setminus \{\vec{i}\} \\ \text{ohne } \vec{i} \\ \text{ohne selbst}}} A(j)$$

Aufgabentext  
Matrix

$$M = \begin{pmatrix} 0 & 1 & 2 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$$

$$\vec{A} = \vec{x}_k = M^T M \cdot \vec{x}_{k-1} \quad \vec{x}_0 = \begin{pmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{pmatrix} \text{ gleich wahrscheinlich}$$

$$\begin{pmatrix} 0 & 1 & 2 \\ 0 & 0 & 0 \\ 0 & 2 & 1 \\ 0 & 1 & 2 \end{pmatrix} \quad \vec{x}_k \rightarrow \frac{1}{\sqrt{3}} \vec{x}_k \quad \text{theoretisch normalisierung}$$

aber hier braucht man die Normalisierung nicht  
Aufgabenstellung speziell für  $P$

$$g_m = [0, 5; 2; 25; 0; 7, 25]$$

W. 3

a)  $g_m = \alpha g_0 + \beta m_r + \gamma m_{nr}$

$$g_0 = \alpha g_0 + \beta m_r + \gamma m_{nr} - g_0$$

$$0 = (\alpha - 1)g_0 + \beta m_r + \gamma m_{nr}$$

Wenn  
 $\alpha = 1$

$$0 = \beta \cdot m_r + \gamma \cdot m_{nr}$$

$$\gamma m_{nr} = \beta \cdot m_r$$

$$m_{nr} = -\frac{\beta}{\gamma} m_r$$

$m_{nr}$  und  $m_r$  linear abhängig

b) Beispiel  $g_0 = m_r$

$$g_m = \alpha \cdot m_r + \beta m_r - m_{nr}$$

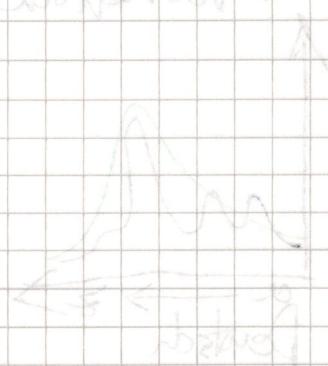
W. 4

$$\alpha = 0, \beta = 1; \gamma = 0$$

W. 5

- Precision → statistische Größe → Tembeschreibung
- Weniger Information bei several wahrscheinlich

ASTRONAUTS  
ARE NOT  
WEIRD



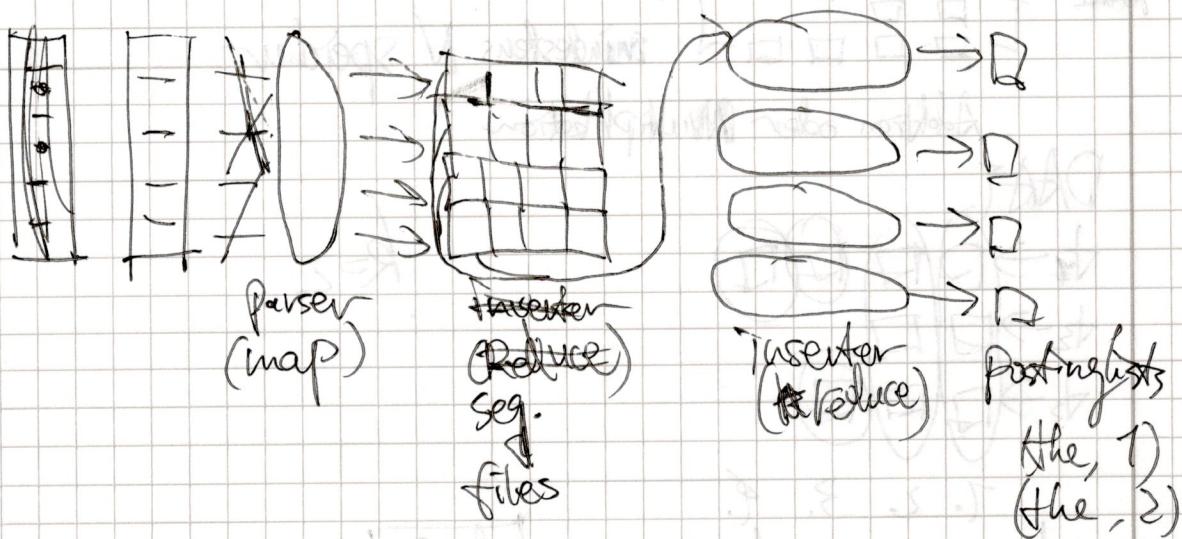
7.2 Pf=1 Gewicht von Term  $t_1$  &  $f_2$  konstant  $\Rightarrow$  unabhängig  
 über: Gewicht von  $t_1$  ED  $D = \{d_1, \dots, d_N\}$

$$\sqrt{t_1, d_1} \Rightarrow \cos\text{-normalisiert}$$

$$t_1 \rightarrow W_{t_1, d_1} > W_{t_1, d_2} > W_{t_1, d_3}$$

z.B.  $k=2$

7.3



Map:  $C \rightarrow \text{list}(\text{Term}, \text{docID})$   
~~Reduce~~:  $\text{list}(\text{Term}, \text{list}(\text{docID})) \rightarrow \text{postingList}_1, \text{postingList}_2, \dots$

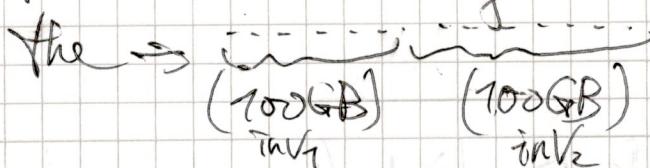
Map:  $C \rightarrow \text{list}(\text{Term}, 1) \rightarrow \text{TypeL}$

bei mehrfache Wörter kommt auch mehrfache TypeL

Reduce:  $\text{list}(\text{word}, (\text{list}(\text{sum}))) \rightarrow (\text{word}, \text{sum})$

7.4 200GB > 100GB

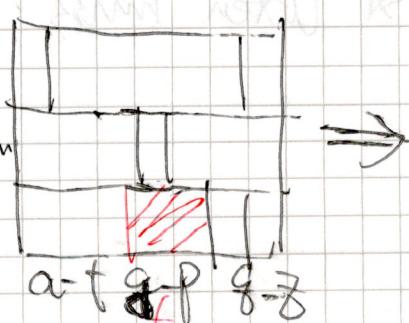
~~Part~~ Partitionierung nach docIDs



7.5

Ian

Schäumstein  
Fall



Before MapReduce  
 Sampling machen  
 a priori