

# Information Engineering 1: Information Retrieval

## Praktikum 4b: Rangierung Relevanz-Rangierung

(Beispiel von Elke Mittendorf)

Es seien die folgenden Dokumente und Anfragen gegeben:

- $q = \text{„Terrorismus, bekämpfen“} = \Phi_1 \Phi_2$
- $D1 = \text{„Gegenmassnahmen gegen Terrorismus“} = \Phi_3 \Phi_1$
- $D2 = \text{„Kampf gegen den Terror“} = \Phi_2 \Phi_1$
- $D3 = \text{„Sicherheit bei asymmetrischer Bedrohung und asymmetrische Sicherheit“} = \Phi_5 \Phi_6 \Phi_7 \Phi_6 \Phi_5$
- $D4 = \text{„Terror bekämpfen“} = \Phi_1 \Phi_2$
- $D5 = \text{„Extremismus und Gewalt“} = \Phi_8 \Phi_9$
- $D6 = \text{„Terrorismus und innere Sicherheit“} = \Phi_1 \Phi_{10} \Phi_5$

Relevant erwiesen sich  $R(q) = \{D1, D2, D3, D6\}$

- Rangieren Sie die Informationsobjekte aus unserem Beispiel mit dem BIR (RSJ-Gewichtung). Diese Rangierung ist a posteriori (weiss zum Zeitpunkt des Vergleichs, welche Dokumente relevant sind).

$$RSV(q, d_j) := \sum_{\phi_i \in \Phi(q) \cap \Phi(d_j)} \log \frac{p_i(1 - q_i)}{q_i(1 - p_i)}$$

(zur Berechnung der  $p_i$  und  $q_i$  siehe BIR.pdf)

- Vergleichen Sie mit probabilistischer idf-Gewichtung

$$RSV(q, d_j) := \sum_{\phi_i \in \Phi(q) \cap \Phi(d_j)} idf(\phi_i), \text{ mit}$$

$$idf(\phi_i) := \log \left( \frac{1 + n}{1 + df(\phi_i)} \right)$$

- Vergleichen Sie mit tf.idf-Cosinus (siehe Folie 72)



Ein paar Anregungen zur Diskussion:

- Bleibt die Reihenfolge gleich?
- Ist immer das gleiche Dokument der beste Match?
- Was ist der fundamentale Unterschied zwischen BIR/RSJ und der probabilistischen idf-Gewichtung/tf.idf-Cosinus?
- Was ist mit D2 gegenüber D4?
- Was ist mit D1 gegenüber D6?
- Bonus-Verständnisfrage: was wäre die optimale Rangierung für das Beispiel?