

Similarity Thesauri and Cross-Language Retrieval

Sebastian Marius Kirsch
skirsch@moebius.inka.de



Back

Close

What is a thesaurus?

the·sau·rus [...]

2a: a book of words or of information about a particular field or set of concepts; *especially*: a book of words and their synonyms b: a list of subject headings or descriptors usually with a cross-reference system for use in the organization of a collection of documents for reference and retrieval

– definition from the Merriam-Webster Online Dictionary



Back

Close

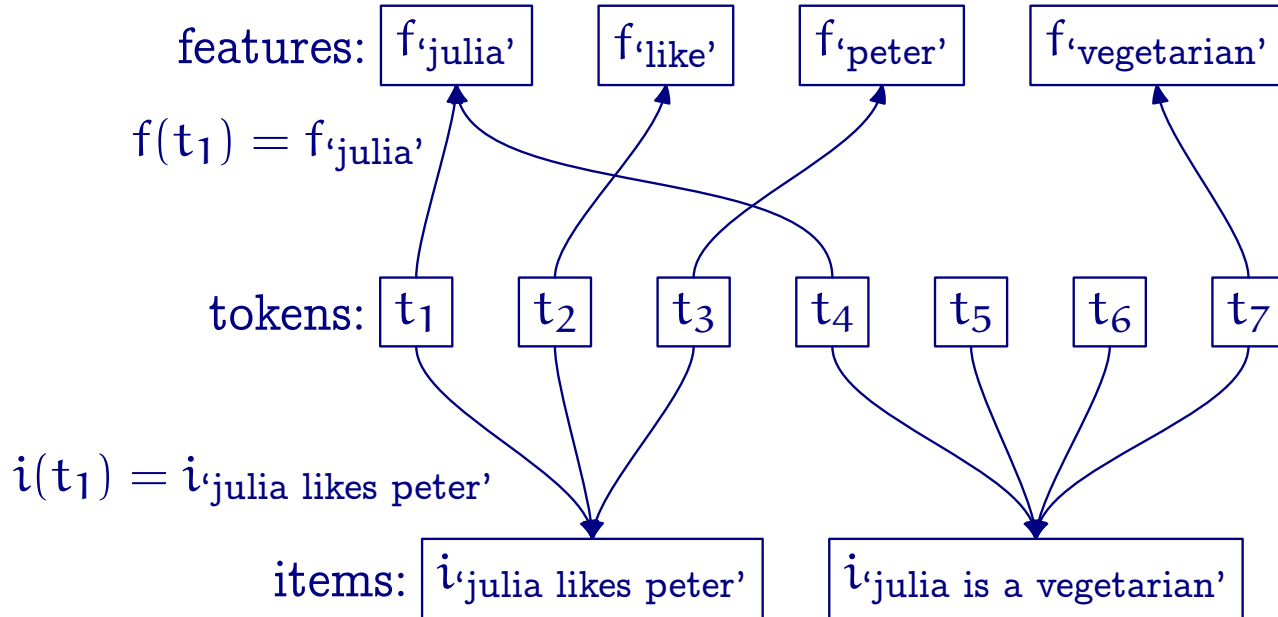
Vector Space Model



Back

Close

Tokens, Items, and Features



Back

Close

Weighting



Back

Close

Weighting

$$\mathbb{f}(f_i, i_j) = |\{\theta \in \Theta | f(\theta) = f_i \wedge i(\theta) = i_j\}|$$



Back

Close

Weighting

$$\text{ff}(f_i, i_j) = |\{\theta \in \Theta \mid f(\theta) = f_i \wedge i(\theta) = i_j\}|$$

$$\text{weight}_{\text{ff}}(i_j, f_i) = \text{ff}(f_i, i_j)$$



Back

Close

Weighting

$$\text{ff}(f_i, i_j) = |\{\theta \in \Theta \mid f(\theta) = f_i \wedge i(\theta) = i_j\}|$$

$$\text{weight}_{\text{ff}}(i_j, f_i) = \text{ff}(f_i, i_j)$$

$$\text{if}(f_i) = |\{i \in I \mid \exists \theta \in \Theta : f(\theta) = f_i \wedge i(\theta) = i\}|$$



Back

Close

Weighting

$$\text{ff}(f_i, i_j) = |\{\theta \in \Theta \mid f(\theta) = f_i \wedge i(\theta) = i_j\}|$$

$$\text{weight}_{\text{ff}}(i_j, f_i) = \text{ff}(f_i, i_j)$$

$$\text{if}(f_i) = |\{i \in I \mid \exists \theta \in \Theta : f(\theta) = f_i \wedge i(\theta) = i\}|$$

$$\text{iif}(f_i) = \frac{1}{\log(\text{if}(f_i) + 1)}$$



Back

Close

Weighting

$$\mathbf{ff}(f_i, i_j) = |\{\theta \in \Theta \mid f(\theta) = f_i \wedge i(\theta) = i_j\}|$$

$$\text{weight}_{\mathbf{ff}}(i_j, f_i) = \mathbf{ff}(f_i, i_j)$$

$$\mathbf{if}(f_i) = |\{i \in I \mid \exists \theta \in \Theta : f(\theta) = f_i \wedge i(\theta) = i\}|$$

$$\mathbf{iif}(f_i) = \frac{1}{\log(\mathbf{if}(f_i) + 1)}$$

$$\text{weight}_{\mathbf{tf} \cdot \mathbf{idf}}(i_j, f_i) = \mathbf{ff}(f_i, i_j) \cdot \mathbf{iif}(f_i)$$



Back

Close

Construction

- normalization: $\text{weight}_{\text{norm.}}(i_j, f_i) = \frac{\text{weight}(i_j, f_i)}{\sqrt{\sum_{f \in F} (\text{weight}(i_j, f))^2}}$
- split similarity computation

$$\text{SIM}(t_i, t_j) = \vec{t}_i \cdot \vec{t}_j$$

$$\begin{aligned}
 & \overbrace{\sum_{d \in D} \text{weight}(t_i, d) \cdot \text{weight}(t_j, d)}^{=\text{sim}(t_i, t_j)} \\
 = & \frac{\sum_{d \in D} \text{weight}(t_i, d) \cdot \text{weight}(t_j, d)}{\sqrt{\underbrace{\left(\sum_{d \in D} (\text{weight}(t_i, d))^2 \right)}_{=c(t_i)} \cdot \underbrace{\left(\sum_{d \in D} (\text{weight}(t_j, d))^2 \right)}_{=c(t_j)}}
 \end{aligned}$$



Back

Close

Monolingual thesaurus

Items big, cabbage, car, drive, julia, ketchup, like, peace, peter, vegetable, vegetarian, war

Features peter drives a big car, julia likes peter, julia is a vegetarian, vegetarians like vegetables, cabbage is a vegetable, big vegetarians who like cabbage do not drive cars, war is peace, ketchup is a vegetable



Back

Close



Back

Close

Query Expansion



Back

Close

Query Expansion

$$\vec{q}_c = \sum_{t \in T} \text{weight}(q, t) \cdot \vec{t}$$



Back

Close

Query Expansion

$$\vec{q}_c = \sum_{t \in T} \text{weight}(q, t) \cdot \vec{t}$$

$$\text{sim}_{qt}(q, t_i) := \vec{q}_c \cdot \vec{t}_i = \sum_{t \in T} \text{weight}(q, t) \cdot \vec{t} \cdot \vec{t}_i$$



Back

Close

Query Expansion

$$\vec{q}_c = \sum_{t \in T} \text{weight}(q, t) \cdot \vec{t}$$

$$\text{simqt}(q, t_i) := \vec{q}_c \cdot \vec{t}_i = \sum_{t \in T} \text{weight}(q, t) \cdot \vec{t} \cdot \vec{t}_i$$

$$\vec{q}_e = \left(\frac{\text{simqt}(q, t_1)}{\sum_{t \in T} \text{weight}(q, t)}, \dots, \frac{\text{simqt}(q, t_n)}{\sum_{t \in T} \text{weight}(q, t)} \right)^T$$



Back

Close

Query Expansion

$$\vec{q}_c = \sum_{t \in T} \text{weight}(q, t) \cdot \vec{t}$$

$$\text{simqt}(q, t_i) := \vec{q}_c \cdot \vec{t}_i = \sum_{t \in T} \text{weight}(q, t) \cdot \vec{t} \cdot \vec{t}_i$$

$$\vec{q}_e = \left(\frac{\text{simqt}(q, t_1)}{\sum_{t \in T} \text{weight}(q, t)}, \dots, \frac{\text{simqt}(q, t_n)}{\sum_{t \in T} \text{weight}(q, t)} \right)^T$$

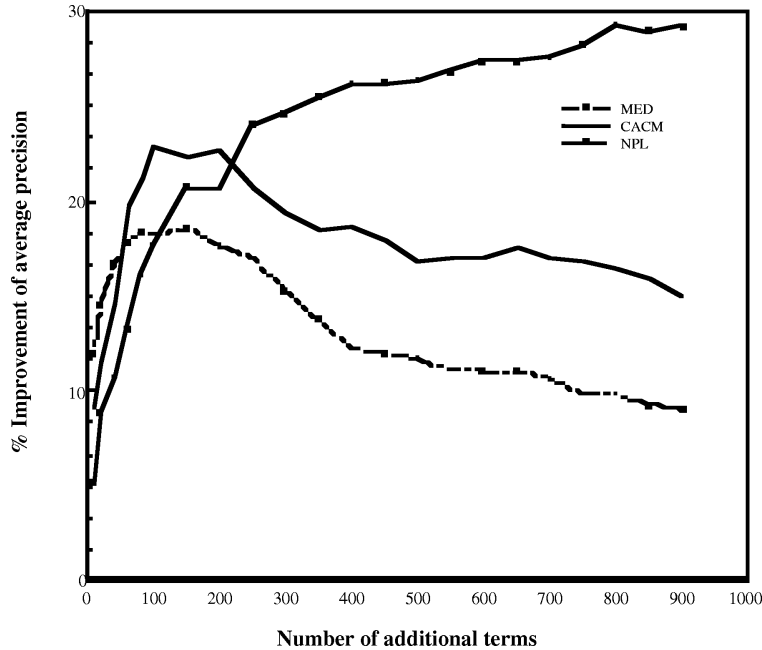
$$\vec{q}_{\text{exp.}} = \vec{q} + \vec{q}_e$$



Back

Close

Evaluation: Query Expansion



Back

Close

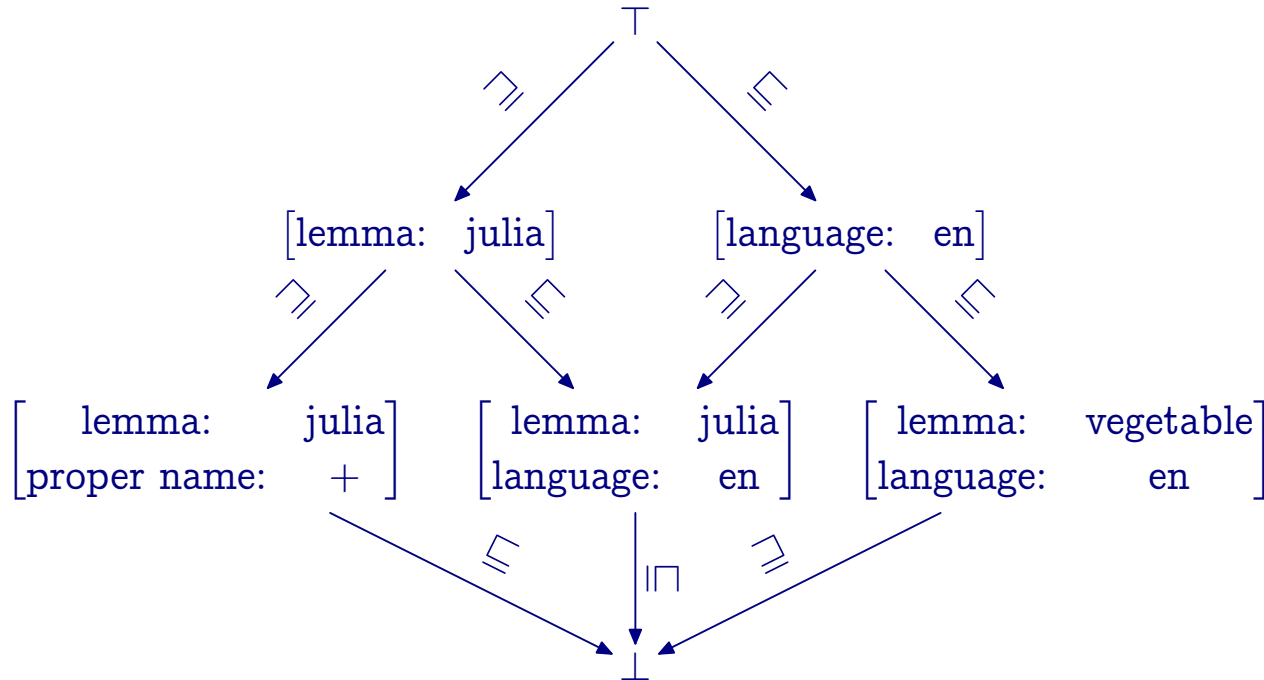
Cross-Language Information Retrieval



Back

Close

Subsumption order



Redefining \mathbb{f} , i

A feature f_i is deemed to occur in an item i_j for every token θ where

1. f_i is equal to or more general than that token's feature $f(\theta)$,
and
2. i_j is equal to or more general than that token's item $i(\theta)$.



Back

Close

Redefining ff , if

A feature f_i is deemed to occur in an item i_j for every token θ where

1. f_i is equal to or more general than that token's feature $f(\theta)$,
and
2. i_j is equal to or more general than that token's item $i(\theta)$.

$$\text{ff}(f_i, i_j) = \{|\theta \in \Theta | f_i \sqsubseteq f(\theta) \wedge i_j \sqsubseteq i(\theta)\}$$



Back

Close

Redefining ff , if

A feature f_i is deemed to occur in an item i_j for every token θ where

1. f_i is equal to or more general than that token's feature $f(\theta)$,
and
2. i_j is equal to or more general than that token's item $i(\theta)$.

$$\text{ff}(f_i, i_j) = \{ \theta \in \Theta \mid f_i \sqsubseteq f(\theta) \wedge i_j \sqsubseteq i(\theta) \}$$

$$\text{if}(f_i) = \{ i \in I \mid \exists \theta \in \Theta : f_i \sqsubseteq f(\theta) \wedge i_j \sqsubseteq i(\theta) \}$$



Back

Close

Cross-lingual thesaurus

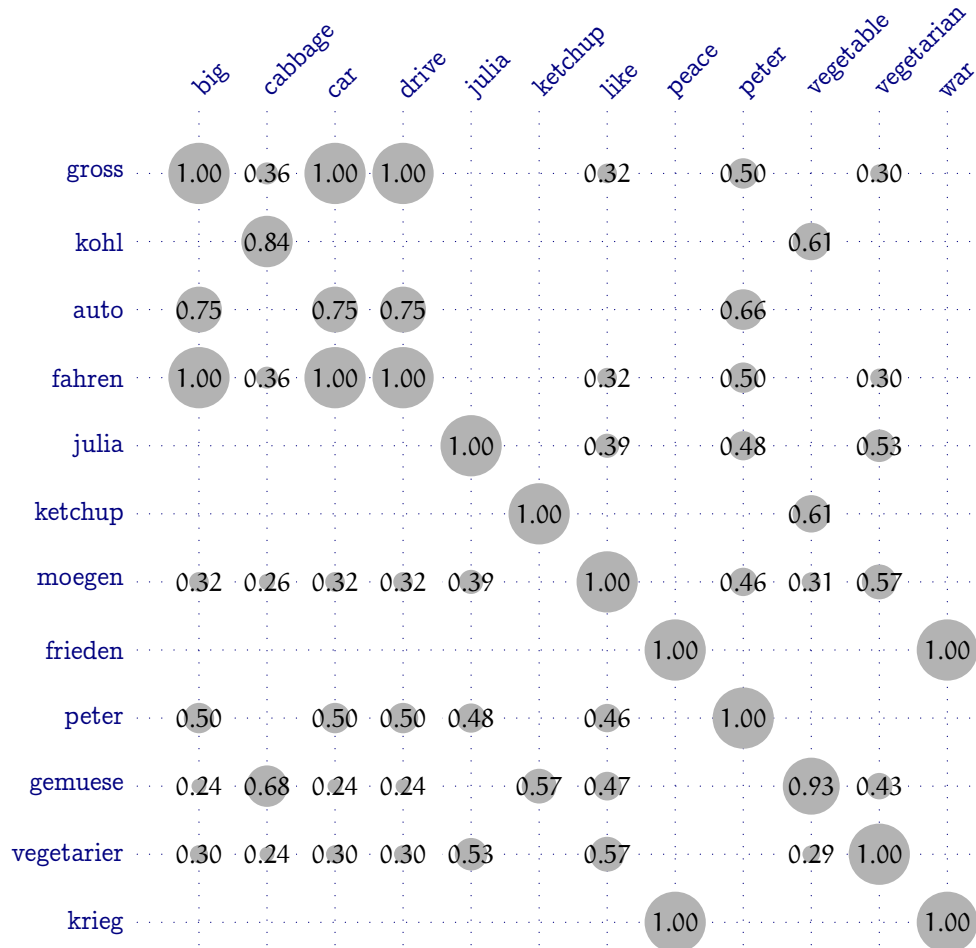
Items big, cabbage, car, drive, julia, ketchup, like, peace, peter, vegetable, vegetarian, war, gross, kohl, auto, fahren, julia, ketchup, moegen, frieden, peter, gemuese, vegetarier, krieg

Features [peter drives a big car, peter faehrt ein grosses auto], [julia likes peter, julia mag peter], [julia is a vegetarian, julia ist vegetarierin], [vegetarians like vegetables, vegetarier moegen gemuese], [cabbage is a vegetable, kohl ist gemuese], [big vegetarians who like cabbage do not drive cars, grosse vegetarier die gemuese moegen fahren keine autos], [war is peace, krieg ist frieden], [ketchup is a vegetable, ketchup ist gemuese]



Back

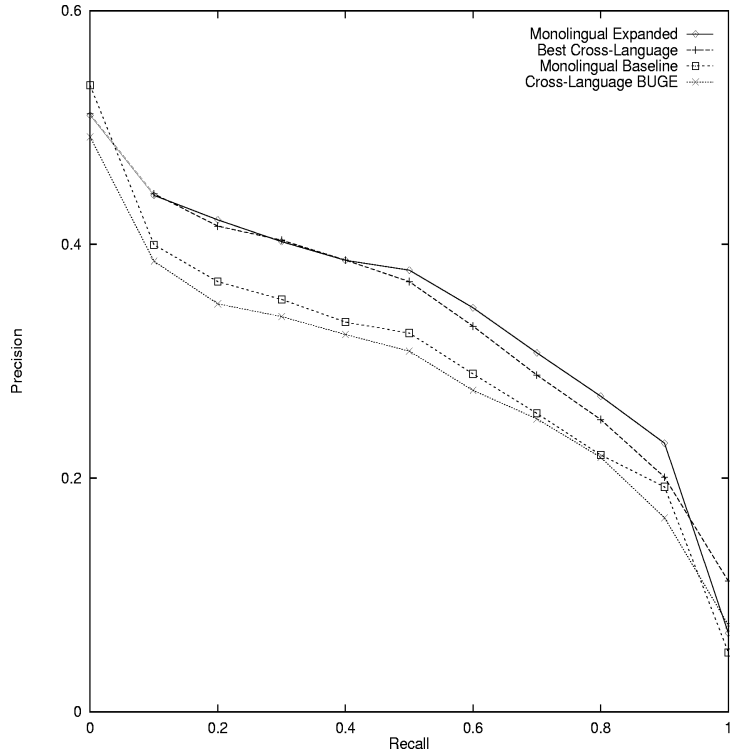
Close



Back

Close

Evaluation: Cross-lingual thesauri



Back

Close

Conclusion

- Theoretical foundations
- Performance
- Improved precision
- Ease of integration
- Multiple use
- Cross-language retrieval



Back

Close

Handout, Sample code

The handout for this talk is available at

<http://sites.inka.de/moebius/docs/simthes-ho.pdf>

Sample code is available at

<http://sites.inka.de/moebius/comp/simthes/>

Thanks very much for your attendance!



Back

Close