# Utilizing Query Facets for Search Result Navigation

Jan Friedrich, Christoph Lindemann, Michael Petrifke

Computer Networks and Distributed Systems
University of Leipzig

01.09.2015

# Table of contents

# Examples I: Microsoft Academic Search



Figure: Example of facets on Microsoft Academic Search

# Examples II: Amazon.com



Figure: Example of facets on Amazon.com

# Examples III: Implemented Results



Figure: Facets generated for the query "hip fractures"

# Introduction to Facets

- A facet is a flat set of terms
- Facets provide selectors / filters for, mostly nominal, object attributes
- Displayed facets traditionally represent existing attributes of the listed objects
- Facets show aspects that help to easily distinguish objects on the level of one attribute $\rightarrow$ it is of no use to show a specific facet if all relevant objects match one and the same value of the corresponding attribute
- Facets provide insight and help to navigate the search result space

# Faceted Web Search Characteristics

- Semi-structured documents
- Some explicit document attributes like in document reference systems (e.g. author, title, publication date, keywords) $\rightarrow$ however, not useful in the context of general web search
- Useful facets are not connected to predefined document attributes (e.g. search results for "IFA Berlin" might benefit from the facets "vendors" or "exhibition hall" $\rightarrow$ this information is hidden in the text)
- Huge number of possible facets and facet terms $\rightarrow$ every existing taxonomy provides many sets of related terms

# Requirements of Faceted Web Search Systems

- Behave similar to Boolean filters $\rightarrow$ learned behavior from other faceted applications
- Terms of one and the same facet should be mutual exclusive $\rightarrow$ only few terms match the same document
- Small number of facets and terms per facet $\rightarrow$ facets distract the user
- Proposal: Use ranking features that characterize the partition properties of the candidate facets

# Faceted Web Search Problems

1. Generation of facets and assignment of facet terms to documents
2. Ranking and selection of relevant facets for the user and query
3. Utilization of user-selected facet terms (user feedback)

# First Work on Facet Generation

- Facetedpedia: Wikipedia provides categories and hyperlinks between articles [5]
- Blogs provide keywords and categories
- External resources like WordNet's hypernym information [1] and other taxonomies
- Above methods not applicable to the general web or require expensive offline computations
- Topic discovery, search result clustering $\rightarrow$ search for labels that fit subsets of the result documents $\rightarrow$ facet generation searches for one-level hierarchies that are representative for the search results

# Facet Extraction from Lists

- Dou et al. [2] introduced the idea to exclusively utilize lists of terms that can be found in the search result documents $\rightarrow$ no external resources required

- Types of lists:
  - Lists in free text
  - Fixed HTML patterns (e.g. ol, ul and tables)
  - Visual repeat regions to extract lists that use CSS and other HTML structures than the fixed patterns above

- Above lists (list candidates) are post-processed, clustered and than ranked to generate the final facets

# HTML Meta Patterns

- Modern web design sometimes utilizes Cascading Style Sheets (CSS) to generate visual lists from general HTML tags like span or p
- Observation I: fixed HTML patterns are not able to extract these lists
- Observation II: visual information is not required to extract most of these lists
- Proposal: HTML Meta Pattern, that finds elements whose children are mostly structurally identical (i.e. same HTML subtree based on the element names)
- ignore comments, script, ...

# HTML Meta Pattern Example



Figure: Example of the requirement of the HTML Meta Pattern

# Candidate List Ranking

- Dou et al. [2] clusters similar lists together and ranks lists high if many result documents contain many terms of the list; they also require lists to appear and different websites
- Kong et al. [3] clusters terms of candidate lists into clusters based on their text and list context; afterwards he uses multiple TF and IDF measures to rank the facets
- Both do not penalize facets whose terms often appear together on each document
- Both do not differentiate between terms in lists and terms occurring on their own

# Navigation Focused Idea

- Binary relevance assessment to decide if a specific facet term $t$ is relevant for a specific search result document $d$: $t$ is relevant for $d$ if $d$ contains $t$ outside of lists $\rightarrow$ in this case $t$ is a valid value for $d$ in each facet that contains $t$

- Each facet term $t$ induces a subset of the search results $D'_t$ where $t$ is relevant

- Idea: Measure the quality of the partition properties of the set of subsets $\{D'_{t1}, D'_{t2}, \ldots, D'_{tn}\}$ of facet $F = \{t_1, t_2, \ldots, t_n\} \rightarrow$ facet extraction algorithm **NAV**

# Search Result Pre-Processing

- Each search result document $d$ is transformed into the bag of words representation $d' = \{t_1, t_2, \ldots, t_n\}$, containing only the terms not contained in lists $\rightarrow$ condensed document representation

- $d'$ is generated at no cost: the candidate list extraction phase removes sub-tress / text snippets that contain the extracted list

$$D'_t = \{d' | d' \cap \{t\} \neq \emptyset\}$$

- We further define $D'_F = \bigcap_{t \in F} D'_t$ as the condensed search result

# Facet Ranking Function

$$R_F = \alpha C_F + \beta S_F + \gamma P_F + \delta T_F$$

# Partition Features I: Subtopic Coverage

- Subtopic coverage $C_F$ recognizes the fact that the original query might have numerous interpretations, but each facet is only relevant for one of these possible search intents
- We approximate the number of sub-intents $\#I$ and calculate a distance measure to the expected number of documents matching at least one of the facet terms of F

$$\#I(D) = log(|D|)$$

$$C_F = \exp\left(-\frac{\left|\frac{|D|}{\#I(D)} - |D'_F|\right|}{10}\right)$$

# Partition Features II: Size Equality

▶ $S_F$ is a measure of the equality of the $D'_t$ document set sizes with $\mu_F^S$ being the mean set size

$$\mu_F^S = \frac{\sum_{t \in F} |D'_t|}{|F|}$$

$$S_F = 1 - \frac{\sum_{t \in F} (\mu_F^S - |D'_t|)^2}{\sum_{t \in F} |D'_t|^2}$$

# Partition Features III: Mean Number Facets

- The reciprocal of the mean number of facet terms per page $P_F$ is used to prefer facets whose facet terms' co-occurrence rate is very low

$$\mu_F^C = \frac{\sum_{d' \in D_F'} |d' \cap F|}{|D_F'|}$$

$$P_F = \frac{1}{\mu_F^C}$$

# Partition Features IV: Facet Size

- $T_F$ is used to prioritize larger facets

$$T_F = \log |F|$$

# Feedback Theory

- The feedback model defines how user selected facet terms are used to improve the web search result in terms of matching the user intent
- $t^u$ represents a user-selected terms (feedback terms)
- $F^u = \{t_1^u, t_2^u, .., t_o^u\}$ is the set of feedback terms of facet F (feedback facet)
- $\mathcal{F}^u = \{F_1^u, F_2^u, ..., F_p^u\}$ is the set of non-empty feedback facets

# Feedback Model

- Kong et al. [4] found Boolean filtering not useful
- They proposed soft ranking $\rightarrow$ original document score is combined with a score that depends on the feedback terms

$$S'_E(d, q, \mathcal{F}^u) = \lambda S(d, q) + (1 - \lambda) S_E(d, \mathcal{F}^u)$$

- Two implementations of $S_E$

$$S_{ST}(d, \mathcal{F}^u) = \frac{1}{N} \sum_{F^u \in \mathcal{F}^u} \sum_{t^u \in F^u} S(d, t^u)$$

$$S_{TT}(d, \mathcal{F}^u) = \sum_{F^u \in \mathcal{F}^u} \sum_{t^u \in F^u} S(d, t^u)$$

# Evaluation Model

- Extrinsic evaluation $\rightarrow$ impact an the search quality (NDCG)
- ClueWeb09 Category B corpus and TREC 2011 relevance measurements of the diversity task
  - Queries with sub-intents $\rightarrow$ relevance judgments for the sub-intents
  - Macro-Averaging: first average over the sub-intents per query, than over the queries
- Per sub-intent, incrementally add the remaining best facet term to the feedback terms
- BM25F ranking function

# Single-Term Feedback I: Top-1 Facet

| Facet Ranking | Candidate List Extraction | Parsed Docs | nDCG @10 | nDCG @20 |
|---|---|---|---|---|
| No facets | | | 0.0672 | 0.0759 |
| QF-I | HTML | 20 | 0.0699 | 0.0805 |
| QF-I | HTML | 50 | 0.0662 | 0.0798 |
| QF-I | HTML + Meta | 20 | 0.0673 | 0.0788 |
| QF-I | HTML + Meta | 50 | 0.0649 | 0.0763 |
| NAV | HTML | 20 | 0.0736 | 0.0877 |
| NAV | HTML | 50 | 0.0704 | 0.0839 |
| NAV | HTML + Meta | 20 | 0.0721 | 0.0858 |
| NAV | HTML + Meta | 50 | 0.0705 | 0.0778 |

Figure: Single term feedback performance using top-1 facet

- ▶ NAV considerable higher scores than QF-I
- ▶ Meta Pattern impairs search quality

# Single-Term Feedback II: Top-3 Facets

| Facet Ranking | Candidate List Extraction | Parsed Docs | nDCG @10 | nDCG @20 |
|---|---|---|---|---|
| No facets | | | 0.0672 | 0.0759 |
| QF-I | HTML | 20 | 0.0824 | 0.0919 |
| QF-I | HTML | 50 | 0.0737 | 0.0915 |
| QF-I | HTML + Meta | 20 | 0.0919 | 0.0954 |
| QF-I | HTML + Meta | 50 | 0.0780 | 0.0911 |
| NAV | HTML | 20 | 0.0808 | 0.0929 |
| NAV | HTML | 50 | 0.0857 | 0.0932 |
| NAV | HTML + Meta | 20 | 0.0800 | 0.0915 |
| NAV | HTML + Meta | 50 | 0.0911 | 0.0960 |

Figure: Single term feedback performance using top-3 facets

- ▶ NAV and QF-I achieve similar quality
- ▶ NAV requires top-50 documents to be effective
- ▶ Meta Pattern is beneficial

# Single-Term Feedback III: Mean Number Facet Terms

| Facet Ranking | Candidate List Extraction | Parsed Docs | # Terms per Facet |
|---------------|---------------------------|-------------|-------------------|
| QF-I | HTML | 20 | 6.29 |
| QF-I | HTML | 50 | 7.69 |
| QF-I | HTML + Meta | 20 | 6.63 |
| QF-I | HTML + Meta | 50 | 8.26 |
| NAV | HTML | 20 | 7.51 |
| NAV | HTML | 50 | 6.94 |
| NAV | HTML + Meta | 20 | 7.62 |
| NAV | HTML + Meta | 50 | 7.27 |

Figure: Mean number of facet terms of the top-3 facets

▶ Increasing number of documents is required to assess the NAV features correctly
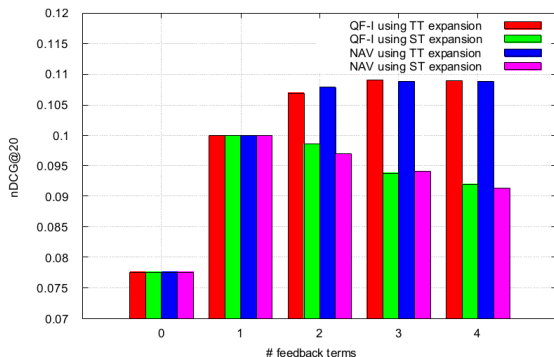
# Multi Term Feedback



Figure: Results for multi term feedback using top-5 facets

▶ ST is not capable of utilizing more than one feedback term

# Conclusion

- Facets generated by NAV, compared to QF-I facets, provide at least the same extrinsic utility
- Each baseline retrieval model might require its specific soft ranking expansion model
- Meta pattern HTML extraction algorithm yields lists that improve facet extraction significantly

# Bibliography

Wisam Dakka and Panagiotis G Ipeirotis.
Automatic extraction of useful facet hierarchies from text databases.
In *Proc. ICDE*. IEEE, 2008.

Zhicheng Dou and et al.
Finding dimensions for queries.
In *Proc. CIKM*. ACM, 2011.

Weize Kong and James Allan.
Extracting query facets from search results.
In *Proc. SIGIR*. ACM, 2013.

Weize Kong and James Allan.
Extending faceted search to the general web.
In *Proc. CIKM*. ACM, 2014.

Chengkai Li and et al.
Facetedpedia: dynamic generation of query-dependent faceted interfaces for wikipedia.
In *Proc. WWW*. ACM, 2010.