



User profiling

Group 3: Anton Lu, Björn Lindqvist, Chris Lin, Sergio Liberman

Background and related methods

Modern commercial search engines like Google compute search results to a given query not just based on measures like tf-idf and PageRank, but also on the user's past behavior. In particular, what links users have selected in preceding searches will heavily influence the ranking of the results of new searches.

In this project we have explored and implemented methods to manipulate search rankings based on past search history of individual people.

Methods proposed by others include:

Personalised PageRank

The method adjusts the random surfer model so the surfer is biased toward some pages, such as those found in the users history and bookmarks. In this way, each user has their own set of unique PageRank scores. [1]

Another method is to assign each page a categorical PageRank score, and use the score from a specific category for ranking, depending on the topics related to the user's query. [2]

Session-based user profiling

Shen, Tan and Zhai [3] proposed categorizing user preferences as either ephemeral or long-lasting. Focusing on ephemeral preferences, they developed a Bayesian model in which each query result would be affected by the previous ones in the same session. In this way, given a predecessor query *cgi programming*, the query *java* is more likely to be about programming

User embeddings

Representing user preferences as vectors and having their similarities with document vectors influence search rankings have been rigorously explored.

User representations have previously been explored with both standard tf-idf embeddings [4,5], as well as representing users and documents with topic vectors, trained using deep learning methods. [6]

Problem statement

In this project we assumed that we were able to record the queries individuals input into a search engine, as well as which pages the user clicks from the search results. We also presumed that the user had a social profile, representing the user's personal preferences.

Using this information, the goal was construct an user profile, and to over time influence the user's search results, so that results more aligned with the user's immediate profile, as well as previous searches were more highly ranked than others.

Methods

I. User and document representation

Each user profile is represented with three document profiles in tf-idf-space:

\mathbf{u}_t : Terms of document titles clicked on

\mathbf{u}_c : Terms of document categories clicked on

\mathbf{u}_x : Terms of document content clicked on

Final user and document vectors \mathbf{u} and \mathbf{d}_i are calculated as

$$\mathbf{u} = \frac{1}{Z_u} [w_t \mathbf{u}_t + w_c \mathbf{u}_c + w_x \mathbf{u}_x]$$
$$\mathbf{d}_i = \frac{1}{Z_d} [w_t \mathbf{t}_i + w_c \mathbf{c}_i + w_x \mathbf{x}_i]$$

Where the Z are normalization factors.

II. Query expansion

Using the idea of the Rocchio algorithm we enhance queries input to the system, \mathbf{q} , using the modified user profile vector $\tilde{\mathbf{u}}$ as

$$\mathbf{q}_{new} = \mathbf{q} + \gamma \tilde{\mathbf{u}}$$

Where the user profile vector \mathbf{u} has been pruned so that terms with small weighting are removed to improve system performance. Terms that already exist in \mathbf{q} were also removed from \mathbf{u} to prevent overemphasizing certain terms. In this case, $\gamma \in [0,1]$.

III. Re-ranking strategy

Re-ranking of search results is done by fetching documents from the search engine using ranked search for the active query, and then re-ranking them using the user profile. For each document \mathbf{d}_i and original score s_i , its final score is calculated as

$$r_i = \alpha s_i + (1 - \alpha) \mathbf{u}^T \mathbf{d}_i$$

Where $\alpha \in [0,1]$ controls the influence of the user profile on the final ranking.

IV. Adaptive profile updating

To emulate concentration shift and forgetting, we employed an exponential decaying scheme. The user profile vectors are then functions of time, defined as

$$\mathbf{u}_k(t) = \mathbf{u}_k(t') e^{-\lambda(t-t')}$$

Where k in this cases represents title, category, content, as presented in (I). The exponential decay constant λ is set to be different for different kinds of profile vectors. For instance, a real-world user is more likely to remember the title of an article rather than the text content, so the decay factor of the title vector is smaller than that of content. When the user clicks on a new page on the results page, the old profile vector is decayed, and then added to the new document vector of the page clicked.

V. Static profile vector

To complement the dynamically updated user profile vector \mathbf{u}_d as seen in (I), one can also introduce a static profile vector \mathbf{u}_s representing long-lasting user preferences, to simulate social profiles. The user profile vector used in (II), (III) is then instead defined as

$$\mathbf{u} = \beta \mathbf{u}_d + (1 - \beta) \mathbf{u}_s$$

Which means that the static preferences do not decay over time as in (IV).

Implementation and system setup

For this task we used a 4 vCPU Google Cloud VM with Elasticsearch 6.3.1, and indexed the whole Swedish Wikipedia dump dated 2019-04-29 as search corpus. The dump contained 11,548,339 documents.

To present search results to users, a website frontend was built using Materialize and Vue.js with a Python Flask backend.

Individual user social profiles were simulated by letting the user input age, sex, location, and selected keywords into the web interface, which were stored in the persistent user database as a static user profile.

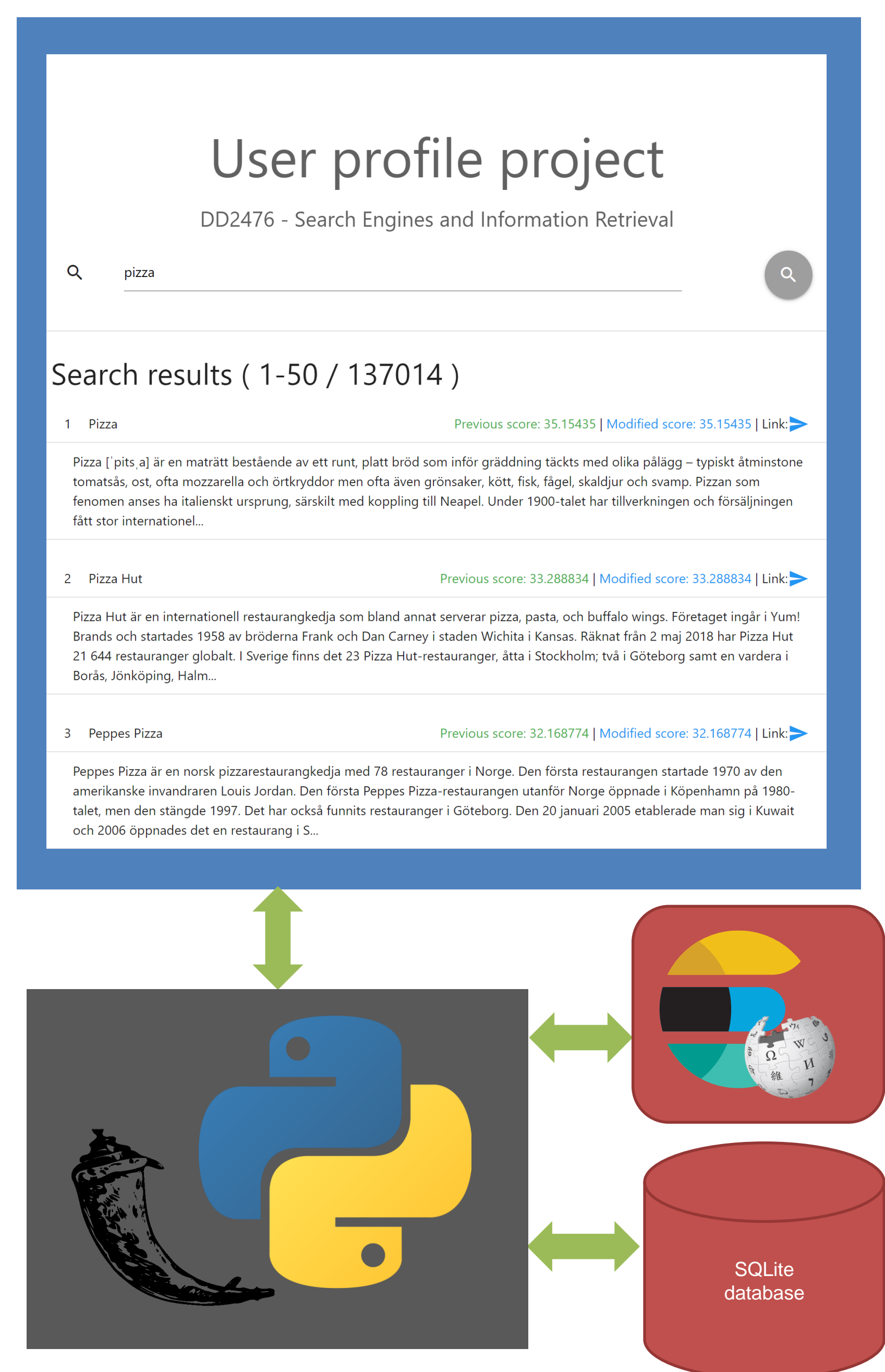


Figure 2: overview of the structure of our implementation.

When an user clicks a search result, the dynamic user profile is updated according to (IV) using the document vectors of the page and (I), and in the case when a static profile exists, it is also used to produce the final user vector.

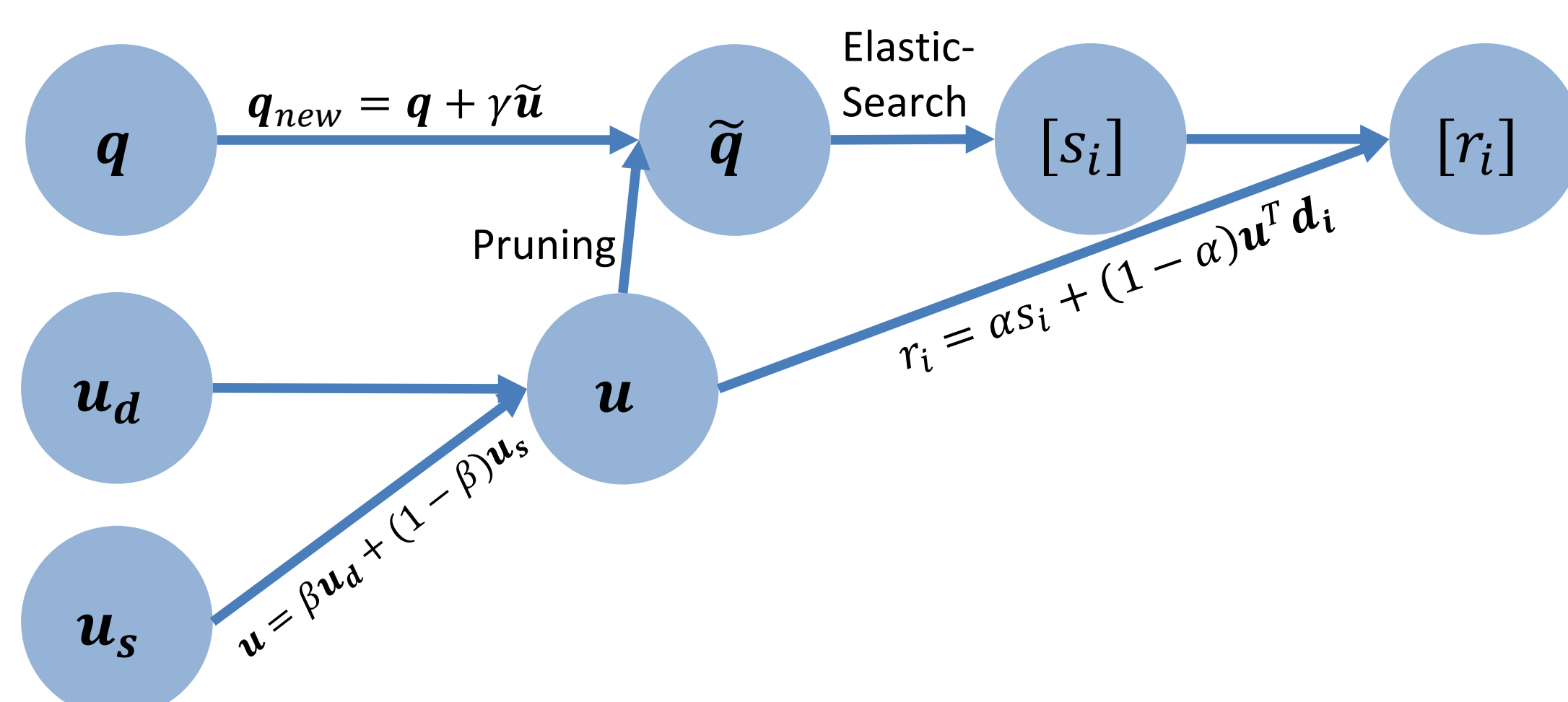


Figure 1: Simplified computational graph showing the process from query input q to final ranking scores r_i .

Evaluation

Queries were issued to the search engine, and after n number of searches, the quality of the profile-based search result were compared to that of a blank profile search.

The metric used for evaluation was the Normalized Discounted Cumulative Gain (nDCG) at position p, defined as

$$nDCG_p = \frac{DCG_p}{IDCG_p}$$

Where

$$DCG_p = \sum_{i=1}^{\rho} \frac{2^{rel_i} - 1}{\log_2(i + 1)}$$

$$IDCG_p = \sum_{i=1}^{|REL|} \frac{2^{rel_i} - 1}{\log_2(i + 1)}$$

and $|REL|$ represents a list of the p top documents, and rel_i is a relevancy score from 0 to 3.

We also used precision as metric for search result goodness.

Experiment 1

In the first evaluation, we ran test queries for 5 blank user profiles, and computed their $nDCG_{20}$. Then the profiles were being updated by clicking 10 relevant and 10 irrelevant documents, and observing what effect they had on the search results. The performance after each click was measured using $nDCG_{20}$. The documents were re-ranked with $\alpha = 0.7$, and the top 100 terms in the user profile were used to expand the search query. The static profile was disabled in this experiment, and experiment 2.

Experiment 2

To test sensitivity to query expansion, another experiment was performed by setting $\alpha = 0$ (no query expansion), and using the top 500 terms in the user profile to re-rank search results.

Experiment 3

A second experiment was run to investigate the impact for weightings between the static and dynamic user profile. We assigned an user a static profiles with somewhat unrelated keywords, and ran several queries related to “Japan”, and evaluated the performance of the search engine to help the user for example find information about sightseeing points in Japan.

Experiment 1 results

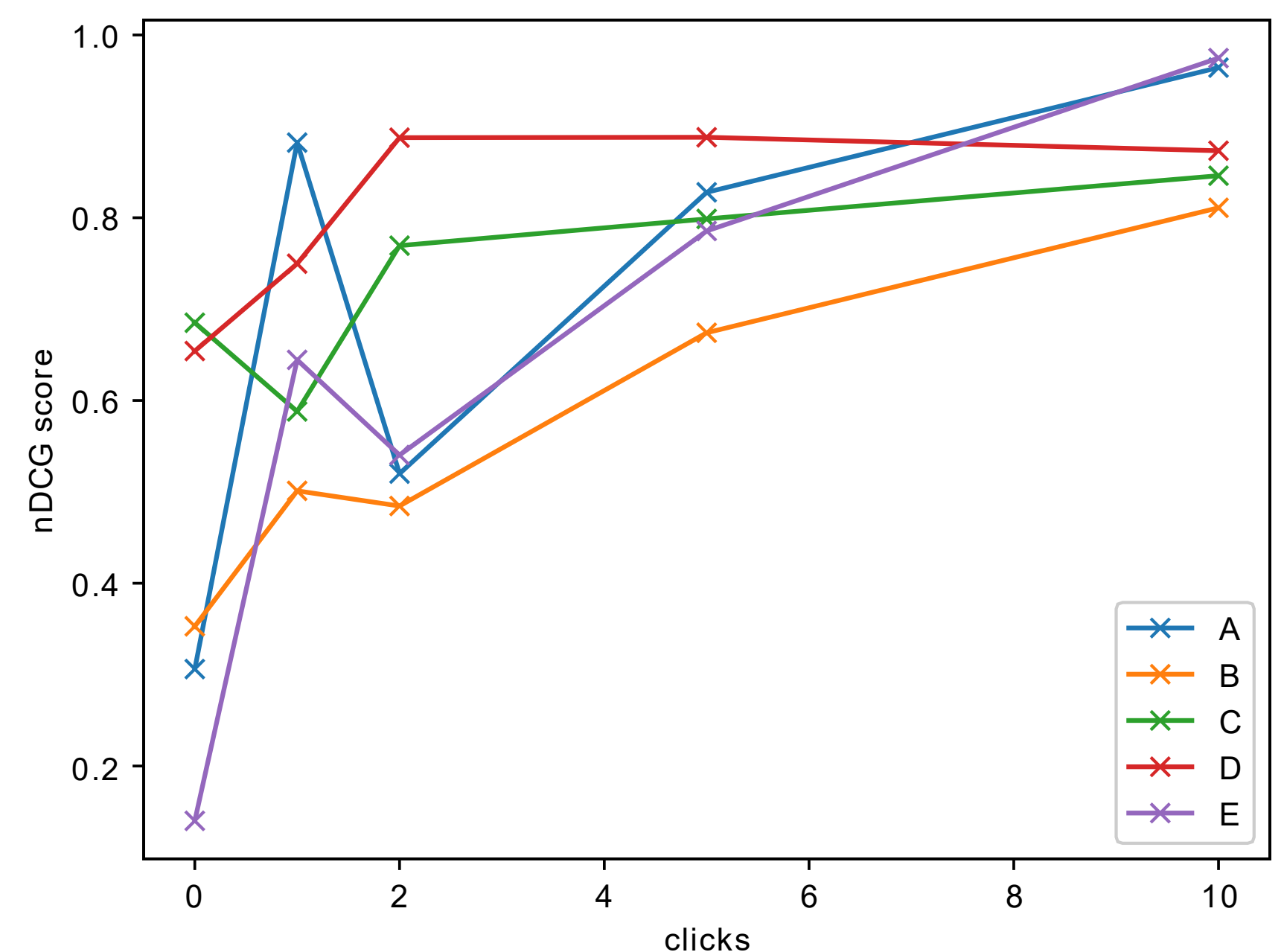


Figure 3: Experiment 1, $nDCG$ scores for 5 different users as function of number of clicks (implicit feedback). New queries were posed between clicks.

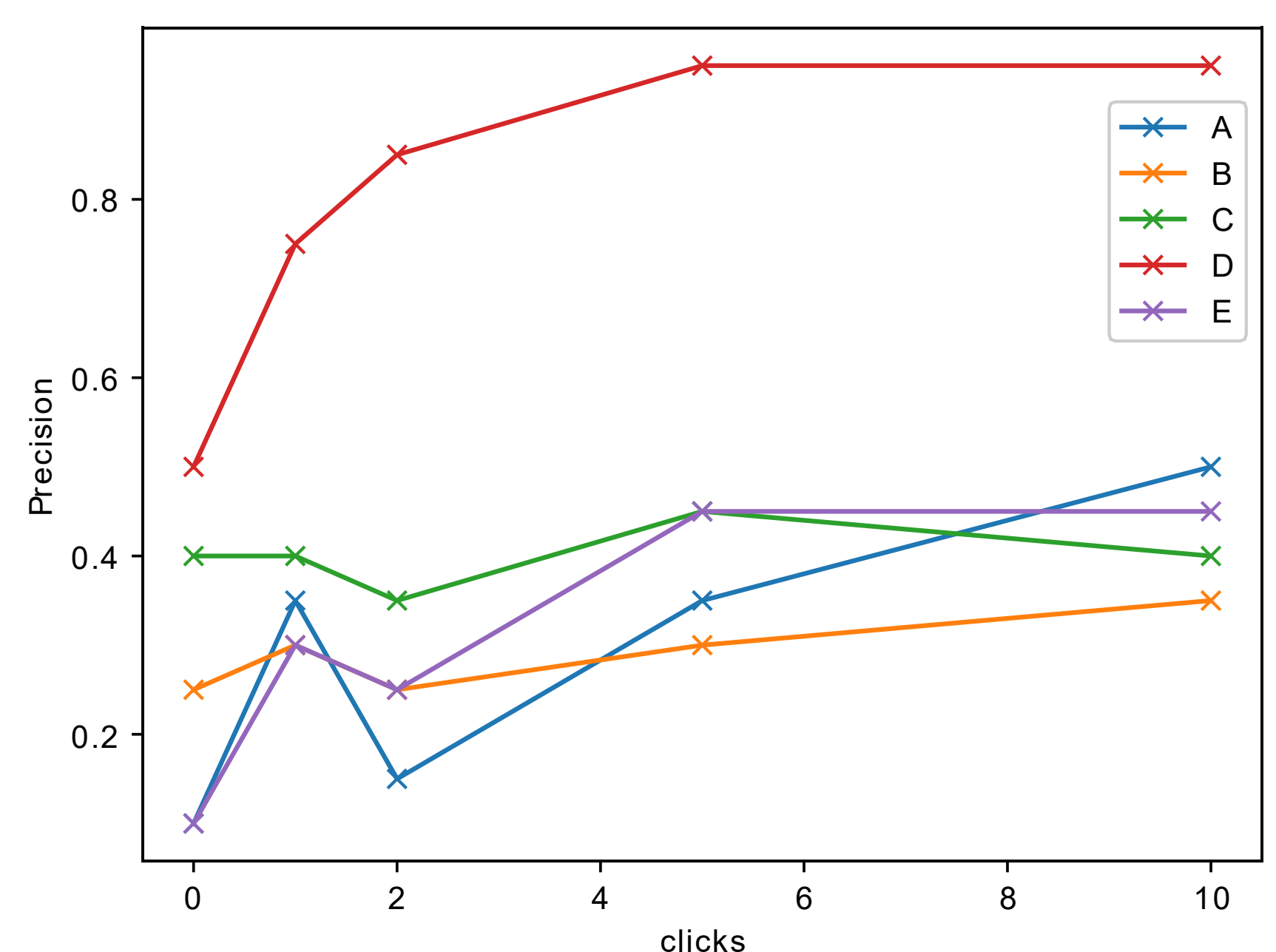


Figure 4: Experiment 1, Precision scores for 5 different users as function of number of clicks (implicit feedback). New queries were posed between clicks.



Figure 5: Experiment 1, Word cloud displaying the user profile vector for user A after 10 clicks.

A line graph showing the nDCG score (Y-axis, ranging from 0.2 to 0.9) versus the number of clicks (X-axis, ranging from 0 to 10). Five algorithms are compared: A (blue line with 'x' markers), B (orange line with 'x' markers), C (green line with 'x' markers), D (red line with 'x' markers), and E (purple line with 'x' markers). Algorithm D consistently achieves the highest nDCG score, starting at approximately 0.75 at 0 clicks and peaking at about 0.92 at 4 clicks. Algorithm C follows, starting at approximately 0.33 at 0 clicks and peaking at about 0.67 at 2 clicks. Algorithm B starts at approximately 0.50 at 0 clicks and fluctuates between 0.35 and 0.55. Algorithm E starts at approximately 0.15 at 0 clicks and fluctuates between 0.25 and 0.40. Algorithm A consistently achieves the lowest nDCG score, starting at approximately 0.17 at 0 clicks and fluctuating between 0.15 and 0.21.

clicks	A	B	C	D	E
0	0.17	0.50	0.33	0.75	0.15
1	0.21	0.35	0.58	0.87	0.28
2	0.15	0.39	0.67	0.89	0.25
3	0.17	0.52	0.54	0.92	0.39
4	0.21	0.53	0.54	0.92	0.35
5	0.21	0.52	0.61	0.89	0.39
6	0.21	0.55	0.62	0.89	0.34
7	0.21	0.52	0.61	0.88	0.36
8	0.21	0.51	0.67	0.88	0.34
9	0.21	0.51	0.67	0.88	0.34
10	0.21	0.51	0.66	0.89	0.36

For most personas in both variants experiments 1 | 2 the relevancy of the search results increased a lot. Even when not using query expansion, search performance still improved, albeit not as much as in experiment 1. Relying only on re-ranking can still be useful as it does not impact the ranking of non-related queries as much.

A bar chart comparing the nDCG score for three datasets (flygplats, slott, tempel) across four dynamic profile ratios (̢): 0.9, 0.7, 0.3, and 0.1. The y-axis represents the nDCG score from 0.0 to 1.0. The x-axis represents the dynamic profile ratio (̢). The legend indicates that flygplats is represented by blue bars, slott by orange bars, and tempel by green bars. The nDCG score generally decreases as the dynamic profile ratio decreases for all three datasets.

Dynamic profile ratio (β)	flygplats	slott	tempel
0.9	0.92	0.93	0.88
0.7	0.91	0.86	0.77
0.3	0.67	0.49	0.30
0.1	0.40	0.32	0.23

The result shows that the dynamic profile portion (large β) is excellent at capturing the user temporary profile. The more portion of the dynamic profile, the better the capture of ephemeral preferences.

Here we illustrate how an user profile evolves as an user uses the search engine

amerika
kina
programming
anime
matematik

- Japanska
- Osaka
- Himeji

DD2476 - Search Engines and Information Retrieval

DD2476 - Search Engines and Information Retrieval

1

Från slott till slott

Previous score: 24.12583 | Modified score: 16.8788091 | [Link](#)

Från slott till slott (franska: D'un château l'autre) är en roman från 1957 av den franske författaren Louis-Ferdinand Céline. Den skildrar Célines två i exil i Svingerimen i östade Tyskland tillämnade med Vöhringen i slutskedet av andra världskriget. Från slott till slott är den första delen i en triologi om dessa erfarenheter, den följes av Nord från 1960 och Ryggraden, som gavs ut postumt 1990...

2

Slott

Previous score: 23.601292 | Modified score: 16.500904 | [Link](#)

Ett slott är en byggnad som ursprungligen byggdes av kungar och furstar när borgen hade spelat ut sin roll som försvarsvärk. Slotten blev makt- och regeringscentrum efter medeltiden. Beträffningen slott kom också att användas när hopadiga låt bygga mer ståndsmässiga byggnader på sina stasgårdar. Med tiden kom också andra högre instanser, till exempel brukspatron, att bygga slott. Slottets funk...

3

Karl Fredrik Slotte

Previous score: 23.54678 | Modified score: 16.49792657 | [Link](#)

Karl Fredrik Slotte, född den 27 oktober 1848 i Neerudalen, Osterodten, död den 19 juli 1914 i Helsingfors, var en finländsk fysiker, som till Karl Johan Slotte, bror till Alexander Slotte. Slotte avlade filosofiske kandidats examen, tjänstgjorde i Åbo som lärare i matematik och fysik, blev 1882 filosofie doktor och utnämndes samma år till äldre lärare i allmän och tillämpad fysik vid Polytekniska in...

4

Slotte

Previous score: 23.46604 | Modified score: 16.420480 | [Link](#)

Slotte är ett finlandssvenskt efternamn som burits av bland andra: Alexander Slotte (1861–1927) Birger Slotte (1900–1983) Karl Johan Slotte (1827–1903) Karl-Johan Slotte (1933–2017) Karl-Fredrik Slotte (1848–1914) Per-Håkan Slotte (född 1942) Peter Slotte (född 1941)...

DD2476 - Search Engines and Information Retrieval

DD2476 - Search Engines and Information Retrieval

sløt

Search results (1 -50 / 3258559)

- 1 Himelji slott [Previous score:](#) [24.819134](#) | [Modified score:](#) [17.4048923059780](#) | [Link](#)

- Himelji slott (洞窟城、Himeji-jō) är ett japanskt slott beläget i Himeji i prefekturen Hyogo. Det är ett av de äldsta ännu existerande byggnadsvrken från det medeltida Japan och har blivit upptaget som världs kulturarv tillsammans med en av den japanska nationens kulturstäder. Tillnamnet på slottet Matsushiro slott och kommandot slott är det ett av Japan's "De berömda slotten". Himeji syns ofta på japanen.

- 2 Fåren slött till slott [Previous score:](#) [24.23623](#) | [Modified score:](#) [16.96721206500413](#) | [Link](#)

- Fårn slött till slott (thanka: D'un chateau l'autre) är en roman från 1957 av den franske författaren Louis-Ferdinand Céline. Den skildrar Célines tid i exil i Sigmaringen (tidigare tyska staden med Volkermünster i närheten av andra världskriget). Från slott till slott är den första delen i en tre-dels serie om dessa erfarelser, den följes av Nord från 1960 och Rogdon, som gavs ut postumt 1969.

- 3 Nagoya slott [Previous score:](#) [24.078291](#) | [Modified score:](#) [16.877051575006612](#) | [Link](#)

- Nagoya slott är ett japanskt slott i Nagoya i centrala Japan. Slottet byggdes under Edo-perioden av Tokugawa Ieyasu på platsen där det tidigare slottet låg. Det var 1868 som för en av uppgifterna rör tillgängligheten och shogunatet förordade mången Japen. Efter att halvöknen som kungligt slott överfördes till Meiji slöt nästan Nagoya slott. Det slet fanns på platsen redan på 1520-talet id.

- 4 Slott [Previous score:](#) [23.78797](#) | [Modified score:](#) [16.653642540431763](#) | [Link](#)

- Ett slott är en byggnad som ursprungligen byggdes av kungaer och furstar när borgarna hade svårt att sin roll som försvarsskär. Slottets bänkar blev små- och regeringsområden efter medietiden. Beredningen slott som också är användas när kongresserna till bygga mer sådana småskaliga byggnader på sina satsningar. Med tiden kom även religiösa byggnadsverk, till exempel brukstakerna eller borgar slott, att byggas upp i samma sättelager.

1. Larry Page et al. The PageRank Citation Ranking: Bringing Order to the Web. 1998
2. Kazunari Sugiyama, Kenji Hatano, and Masatoshi Yoshikawa. "Adaptive Web Search Based on User Profile Constructed Without Any Effort from Users". In: Proceedings of the 13th International Conference on World Wide Web. WWW '04. New York, NY, USA: ACM, 2004, pp. 675–684.
3. Xuehua Shen, Bin Tan, and ChengXiang Zhai. "Context-sensitive Information Retrieval Using Implicit Feedback". In: Proceedings of the 28th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval. SIGIR '05. Salvador, Brazil: ACM, 2005, pp. 43–50.
4. Nicolaas Mattheijs and Filip Radlinski. "Personalizing Web Search Using Long Term Browsing History". In: Proceedings of the Fourth ACM International Conference on Web Search and Data Mining. WSDM '11. HongKong, China: ACM, 2011, pp. 25–34.
5. Shengliang Xu et al. "Exploring Folksonomy for Personalized Search". In: Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval. SIGIR '08. Singapore, Singapore: ACM, 2008, pp. 155–162.
6. Thanh Vu et al. "Search Personalization with Embeddings". In: CoRRabs/1612.03597(2016). arXiv:1612.03597.