

Structured Collaborative Filtering



Alejandro Bellogín, Pablo Castells
Universidad Autónoma de Madrid, Spain
{alejandro.bellogin, pablo.castells}@uam.es

Jun Wang
University College London, UK
j.wang@cs.ucl.ac.uk



Vector Space Model

Query
 $Q = (w_1, qf_1; \dots; w_n, qf_n)$

Document
 $D^j = (w_1, t_{j1}; \dots; w_n, t_{jn})$

Similarity score
 $\text{sim}(Q, D_j) \propto Q \cdot D_j = \sum_k qf_k \cdot t_{jk}$

Collaborative Filtering

User
 $Q^u = (i_1, r_{i_1}^u; \dots; i_n, r_{i_n}^u)$

Item
 $I^j = (s_1^j, \dots; s_n^j)$

Predicted rating
 $\text{rat}(u, i) = \text{sim}(Q^u, I^j) \propto Q^u \cdot I^j = \sum_k r_{i_k}^u \cdot s_k^j$

with proper normalization
Item-based CF $\text{rat}(u, i) = \frac{\sum_k r_{i_k}^u \cdot s_k^j}{\sum_k s_k^j}$

Extended Boolean Model

Two connectives (and, or) are included in the query [Salton et al, 1983]:

$$Q := \underset{\text{or}}{\text{and}}(p_1) \left[\underset{\text{or}}{\text{and}}(p_2) [Q]^+, cw \right]$$

$$Q := (qf_1, \dots, qf_n)$$

These connectives are weighted by factors p_1 and p_2
The value cw is the clause weight

The document representation is the same:

$$D^j := (t_{j1}, \dots, t_{jn})$$

The similarity score depends on the connective:

$$\text{sim}(Q_{\text{and}(p)}, D^j) = 1 - \left[\frac{(qf_1(1-t_{j1}))^p + \dots + (qf_n(1-t_{jn}))^p}{(qf_1)^p + \dots + (qf_n)^p} \right]^{1/p}$$

$p \in [1, +\infty)$

$$\text{sim}(Q_{\text{or}(p)}, D^j) = \left[\frac{(qf_1 t_{j1})^p + \dots + (qf_n t_{jn})^p}{(qf_1)^p + \dots + (qf_n)^p} \right]^{1/p}$$

When $p = 1$
 $\text{sim}(Q_{\text{or}}, D^j) = \text{sim}(Q_{\text{and}}, D^j)$

Extended Vector-Space Representation for CF

Example

(a) using Boolean retrieval.

	Query items		user interest representation	
	a	b	a OR b	a AND b
item 1	1	1	1	1
item 2	1	0	1	0
item 3	0	1	1	0
item 4	0	0	0	0

(b) using extended Boolean retrieval ($p = 2$).

	Query items		user interest representation	
	a	b	a OR b	a AND b
item 1	1	1	1	1
item 2	1	0	$1/\sqrt{2}$	$1 - 1/\sqrt{2}$
item 3	0	1	$1/\sqrt{2}$	$1 - 1/\sqrt{2}$
item 4	0	0	0	0

A user shows interest for two items:

The Boolean model is too rigid: too loose or tight

Extended Boolean model: more discriminative

Representation

Only the user profile needs to change its representation

$$Q^u := \underset{\text{or}}{\text{and}}(p_1) \left[\underset{\text{or}}{\text{and}}(p_2) [Q^u]^+, cw \right]$$

$$Q^u := (r_{i_1}^u, \dots, r_{i_n}^u)$$

Now, the predicted rating takes into account the structure of the user profile:

$$\text{rat}(u, i) = \text{sim}(Q^u, I^j) \begin{cases} \text{sim}(Q_{\text{and}(p)}, I^j) \\ \text{sim}(Q_{\text{or}(p)}, I^j) \end{cases}$$

Now, item-based CF is simply represented as

Semantics of the p-values

- $p = +\infty$ and AND connective: a strict phrase has to be matched
- $p = +\infty$ and OR connective: a strict thesaurus feature is used (all the terms are substitutable)
- low p and AND connective: the presence of every term is worth more (but not compulsory) than the presence of only some of them
- low p and OR connective: the presence of several terms is more important than the presence of one of them
- $p = 1$: both connectives are equivalent

Different retrieval models obtained:
 $p=1$: Vector Space Model
 $p = +\infty$: Boolean model

User Profile Expansion

Motivation

Items that tend to occur together: movie series or by the same director

E.g., Lord of the Rings, Star Wars, etc.

These movies could be considered synonyms

Experiments

We expand every user profile with S synonym movies (i.e., the S most similar items).

Expanded terms (synonyms) are included using an inner Boolean OR (i.e., infinite p -value)

Method	P@1	P@3	P@5	P@10	NDCG@3	NDCG@5	NDCG@50	MRR
Baseline	0.002	0.004	0.005	0.007	0.002	0.003	0.009	0.027
$S = 1, \text{and}(1), \text{or}(\infty)$	0.129	0.121	0.120	0.114	0.099	0.099	0.126	0.243
$S = 2, \text{and}(1), \text{or}(\infty)$	0.149	0.138	0.130	0.119	0.113	0.110	0.128	0.260
$S = 5, \text{and}(1), \text{or}(\infty)$	0.190	0.166	0.155	0.141	0.140	0.134	0.146	0.304
$S = 10, \text{and}(1), \text{or}(\infty)$	0.197	0.171	0.161	0.147	0.146	0.140	0.151	0.313
$S = 5, \text{or}(\infty), \text{and}(1)$	0.004	0.005	0.005	0.005	0.002	0.003	0.010	0.028

Not every user needs to be expanded: dynamic user profile expansion

Method	P@1	P@3	P@5	P@10	NDCG@3	NDCG@5	NDCG@50	MRR
Threshold found by median	0.183	0.167	0.157	0.145	0.137	0.132	0.158	0.302
Threshold found by average	0.186	0.165	0.158	0.146	0.137	0.133	0.137	0.303

Discussion

A significant improvement is found, compared to plain profiles ($S=0$, baseline).

The order of the connectives is important

Profiles built using OR + AND make little sense

Dynamic expansion outperforms static one with larger cutoffs

Future Work

Additional user profile expansion methods to set dynamically the threshold

Combination of this method with recently proposed normalization methods in [Bellogín et al, 2011]

Inferring User Profile Structure

Motivation

Subprofiles in recommendation

E.g., users A and B have similar tastes in movies but different in music

User profiles could thus be decomposed into phrases

Experiments

Each profile is defined as soft OR's of cohesive subprofiles

Each subprofile is also composed of soft OR's of item ratings

Subprofiles are found by performing clustering (using the Weka library) on genre and similarity values among items:

K-means ($K=50$)

X-means (worse results than K-means)

Uniform p -value for all the subprofiles

Method	P@1	P@3	P@5	P@10	NDCG@3	NDCG@5	NDCG@50	MRR
Baseline	0.002	0.004	0.005	0.007	0.002	0.003	0.009	0.027
K-means genre or(1), or(2)	0.013	0.016	0.017	0.018	0.010	0.011	0.032	0.061
K-means genre or(2), or(2)	0.005	0.006	0.006	0.006	0.003	0.004	0.011	0.030
K-means sim or(1), or(2)	0.008	0.011	0.014	0.016	0.007	0.009	0.023	0.047
K-means sim or(2), or(2)	0.006	0.006	0.006	0.008	0.004	0.004	0.012	0.031

Discussion

Structured profiles outperform plain profiles (baseline)

No significant differences for different inner p -values

Better results with subprofiles induced based on genre information

Future Work

Different p -values for each subprofile, e.g., depending on the intracluster similarity

References

- [Bellogín et al, 2011] A. Bellogín, J. Wang, P. Castells. Text retrieval methods for item ranking in collaborative filtering. In ECIR 2011.
[Salton et al, 1983] G. Salton, E.A. Fox, H. Wu. Extended boolean information retrieval. Communications of the ACM 1983.