

# Impact of the Query Set on the Evaluation of Expert Finding Systems

Robin Brochier<sup>1,2</sup> Adrien Guille<sup>1</sup>  
Benjamin Rothan<sup>2</sup> Julien Velcin<sup>1</sup>

<sup>1</sup> Université de Lyon, Lyon 2, ERIC EA 3083, France  
{*robin.brochier,julien.velcin, adrien.guille*}@univ-lyon2.fr  
<https://eric.ish-lyon.cnrs.fr/>

<sup>2</sup> Digital Scientific Research Technology, Lyon, France {*robin, benjamin*}@peer.us  
<https://peer.us/>

07/12/2018

# Table of Contents

BIRNDL  
2018

Robin  
Brochier,  
Adrien Guille  
, Benjamin  
Rothan,  
Julien Velcin

Introduction

The Usual  
Evaluation  
(topic-query)

Another  
Evaluation  
(document-  
query)

Results and  
Comparisons

Ongoing  
Works

Références

- 1 Introduction
- 2 The Usual Evaluation (topic-query)
- 3 Another Evaluation (document-query)
- 4 Results and Comparisons
- 5 Ongoing Works

# Introduction

BIRNDL  
2018

Robin  
Brochier,  
Adrien Guille  
, Benjamin  
Rothan,  
Julien Velcin

Introduction

The Usual  
Evaluation  
(topic-query)

Another  
Evaluation  
(document-  
query)

Results and  
Comparisons

Ongoing  
Works

Références

- PhD in collaboration with Digital Scientific Research Technology<sup>1</sup>.
- Peerus Review : a tool for publishers to find potential reviewers<sup>2</sup>.

peerusreview

*Review of Impact of the Query Set on the Evaluation of Expert Finding Systems*

**Journals:** Not documented

**Authors:** Not documented

**Abstract:**

Expertise is a loosely defined concept that is hard to formalize. Much research has focused on designing efficient algorithms for expert finding in large databases in various application domains. The evaluation of such recommender systems lies most of the time on human-annotated sets of experts associated with topics. The protocol of evaluation consists in using the namings or short descriptions of these topics as raw queries in order to rank the available set of candidates. Several measures taken from the field of information retrieval are then applied to rate the rankings of candidates against the ground truth set of experts. In this paper, we apply this topic-query evaluation methodology with the AMiner data and explore a new document-query methodology to evaluate experts retrieval from a set of queries sampled directly from the experts documents. Specifically, we describe two datasets extracted from AMiner, three baseline algorithms from the literature based on several document representations and provide experiment results to show that using a wide range of more realistic queries provides different evaluation results to the usual topic-queries.

**Reviewers:**

We are currently searching for reviewers, this operation may take few minutes.

1. <https://peer.us/>
2. <https://review.peer.us/>

# Motivations

BIRNDL  
2018

Robin  
Brochier,  
Adrien Guille  
, Benjamin  
Rothan,  
Julien Velcin

Introduction

The Usual  
Evaluation  
(topic-query)

Another  
Evaluation  
(document-  
query)

Results and  
Comparisons

Ongoing  
Works

Références



**Topic-Query** : "Natural Language Processing"

**Document-Query** : "Learning Subjective Language : Subjectivity in natural language refers to aspects of language used to express opinions, evaluations, and speculations. There are numerous natural language processing applications for which subjectivity analysis is relevant, including information extraction and text categorization..."



**Topic-Query** : "Backend Engineer"

**Document-Query** : "You write well designed, testable and reliable code as you ship features that help support amazing causes world-wide ; maintain a highly responsive system by ensuring code quality, stability and scalability ; support fellow developers by providing advice, encouraging best practices and performing design and code reviews..."



**Topic-Query** : "Python Tensorflow Keras"

**Document-Query** : "My input is a merge of datasets from few different experiments. My custom loss function needs firstly to split that data and perform different calculations in each subset, then sum results. Unfortunately, I can't split it inside custom loss. I have prepared a piece of code that works with `tf.boolean_mask()`..."

**Is the usual topic-query evaluation statistically significant ?**

**Does an algorithm perform the same way for both types of queries ?**

# An Example of a Dataset

BIRNDL  
2018

Robin  
Brochier,  
Adrien Guille,  
Benjamin  
Rothan,  
Julien Velcin

Introduction

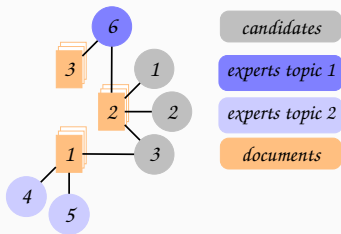
The Usual  
Evaluation  
(topic-query)

Another  
Evaluation  
(document-  
query)

Results and  
Comparisons

Ongoing  
Works

Références



(a) Bipartite graph linking candidates and documents.

$$G = \begin{pmatrix} 0 & 0 & 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}_{|D| \times |C|}$$

$$X = \begin{pmatrix} 0.4 & \dots & 0.8 \\ 0.9 & \dots & 0.1 \\ 0.3 & \dots & 0.3 \end{pmatrix}_{|D| \times |X|}$$

(b) Adjacency matrix  $G$  of the bipartite graph and the features  $X$  matrix of the documents.

Figure – Hypothetical example of a dataset for expert finding.

BIRNDL  
2018

Robin  
Brochier,  
Adrien Guille  
, Benjamin  
Rothan,  
Julien Velcin

Introduction

The Usual  
Evaluation  
(topic-query)

Another  
Evaluation  
(document-  
query)

Results and  
Comparisons

Ongoing  
Works

Références

AMiner provides a dataset of scientific authors and their publications (titles+abstracts).

Citation V1<sup>3</sup> + Expert Finding (human-annotated)<sup>4</sup> :

- 532,968 candidates
- 480,630 documents
- 1,100,698 candidates-documents links
- 210 experts
- 7 topics

---

3. <https://aminer.org/citation>

4. [https://aminer.org/lab-datasets/expertfinding/  
#expert-list](https://aminer.org/lab-datasets/expertfinding/#expert-list)

# Usual Expert Finding Evaluation.

BIRNDL  
2018

Robin  
Brochier,  
Adrien Guille  
, Benjamin  
Rothan,  
Julien Velcin

Introduction

The Usual  
Evaluation  
(topic-query)

Another  
Evaluation  
(document-  
query)

Results and  
Comparisons

Ongoing  
Works

Références

---

## Algorithm 1 The usual way of evaluating.

---

```
Require: Ranking_Algorithm
scores ← [ ]
for all topics do
  candidates_ranking = Ranking_Algorithm(current_topic_textual_expression)
  current_score ← Evaluate(candidates_ranking, ground_truth_topic_experts_set)
  scores.append(current_score)
end for
return Mean(scores), STD(scores)
```

---

Topic-queries in AMiner dataset :

- “intelligent agents”
- “planning”
- “machine learning”
- “natural language processing”
- “information extraction”
- “semantic web”
- “support vector machine”

Metrics used :

- precision at k
- average precision
- reciprocal rank
- area under the roc curve

# Limitations ?

BIRNDL  
2018

Robin  
Brochier,  
Adrien Guille  
, Benjamin  
Rothan,  
Julien Velcin

Introduction

The Usual  
Evaluation  
(topic-query)

Another  
Evaluation  
(document-  
query)

Results and  
Comparisons

Ongoing  
Works

Références

- Only 7 queries : might be statistically not significant.
- Queries are keywords : bias due to their short length ?
- Real case applications : users have different behaviors.



# Another Evaluation (document-query)

BIRNDL  
2018

Robin  
Brochier,  
Adrien Guille  
, Benjamin  
Rothan,  
Julien Velcin

Introduction

The Usual  
Evaluation  
(topic-query)

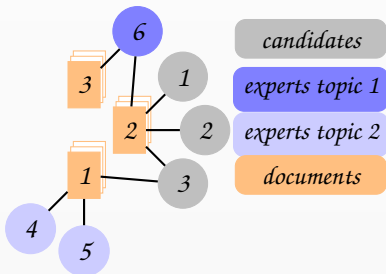
Another  
Evaluation  
(document-  
query)

Results and  
Comparisons

Ongoing  
Works

Références

## Sampling documents to use them as queries :



for topic 1 :

for expert 6 :

we sample document 3.

we sample document 2.

for topic 2 :

for expert 4 :

we add document 1.

for expert 5 :

we add document 1.

# Sampling Documents-Queries from Ground Truth Experts

BIRNDL  
2018

Robin  
Brochier,  
Adrien Guille  
, Benjamin  
Rothan,  
Julien Velcin

Table – Document-Queries Sampled

Topic	Exps	Docs
intelligent agents	29	685
planning	33	510
machine learning	38	628
natural language processing	38	453
information extraction	18	290
semantic web	41	521
support vector machine	29	323
<b>TOTAL</b>	<b>210</b>	<b>3410</b>

Introduction

The Usual  
Evaluation  
(topic-query)

Another  
Evaluation  
(document-  
query)

Results and  
Comparisons

Ongoing  
Works

Références

# Baseline Algorithms

BIRNDL  
2018

Robin  
Brochier,  
Adrien Guille  
, Benjamin  
Rothan,  
Julien Velcin

Introduction

The Usual  
Evaluation  
(topic-query)

Another  
Evaluation  
(document-  
query)

Results and  
Comparisons

Ongoing  
Works

Références

- **P@noptic** : Craswell et al., “P@ noptic expert : Searching for experts not just for documents”  
=> A simple algorithm creating meta-documents for each candidate by aggregating their documents (**only text**).
- **Vote** : Macdonald et Ounis, “Voting for candidates : adapting data fusion techniques for an expert search task”  
=> A voting system where each document votes for its candidates weighted by their similarity to the query (**text and first-order proximity**).
- **Propagation** : Serdyukov, Rode et Hiemstra, “Modeling multi-step relevance propagation for expert finding”  
=> PageRank-like algorithm initialized with query-documents similarities (**text and N-order proximity**).

# Topic-query Evaluation Results

BIRNDL  
2018

Robin  
Brochier,  
Adrien Guille  
, Benjamin  
Rothan,  
Julien Velcin

Introduction

The Usual  
Evaluation  
(topic-query)

Another  
Evaluation  
(document-  
query)

Results and  
Comparisons

Ongoing  
Works

Références

		TF	TF-IDF	LSI
P@noptic	AUC	0.777±0.099	0.778±0.102	0.832±0.083
	P@10	0.662±0.262	<b>0.685±0.260</b>	<b>0.615±0.301</b>
	AP	0.398±0.193	0.415±0.204	<b>0.395±0.233</b>
	RR	1.615±0.923	1.538±0.746	1.769±0.890
Vote	AUC	0.714±0.137	0.714±0.138	0.800±0.107
	P@10	0.608±0.312	0.608±0.287	0.538±0.325
	AP	0.373±0.212	0.381±0.209	0.390±0.247
	RR	2.308±2.398	1.538±0.929	2.769±3.765
Prop ( $\eta = 0.1$ )	AUC	0.834±0.093	0.834±0.096	0.824±0.085
	P@10	0.669±0.270	0.677±0.264	<b>0.615±0.298</b>
	AP	<b>0.458±0.239</b>	<b>0.473±0.243</b>	0.389±0.230
	RR	1.462±0.929	1.538±1.082	<b>1.692±1.136</b>
Propagation ( $\eta = 0.5$ )	AUC	<b>0.842±0.086</b>	<b>0.842±0.088</b>	<b>0.833±0.083</b>
	P@10	<b>0.677±0.255</b>	<b>0.685±0.274</b>	<b>0.615±0.296</b>
	AP	0.457±0.232	0.472±0.242	<b>0.395±0.231</b>
	RR	<b>1.308±0.606</b>	<b>1.385±0.738</b>	<b>1.692±1.136</b>

Table – Results of the Topic-query evaluations.

# Document-query Evaluation Results

BIRNDL  
2018

Robin  
Brochier,  
Adrien Guille,  
Benjamin  
Rothan,  
Julien Velcin

Introduction

The Usual  
Evaluation  
(topic-query)

Another  
Evaluation  
(document-  
query)

Results and  
Comparisons

Ongoing  
Works

Références

(a) Scores and STD.

		TF	TF-IDF	LSI
P@noptic	AUC	0.593±0.131	0.618±0.134	0.620±0.139
	P@10	0.255±0.266	0.335±0.294	<b>0.302</b> ±0.298
	AP	0.150±0.117	0.181±0.133	0.174±0.133
	RR	9.153±16.028	6.210±13.323	8.663±16.262
Vote	AUC	<b>0.606</b> ±0.131	<b>0.630</b> ±0.137	<b>0.637</b> ±0.142
	P@10	<b>0.284</b> ±0.263	0.322±0.280	0.275±0.276
	AP	<b>0.169</b> ±0.131	<b>0.193</b> ±0.145	<b>0.187</b> ±0.152
	RR	<b>6.819</b> ±15.421	<b>5.783</b> ±13.983	<b>8.438</b> ±16.987
Propagation ( $\eta = 0.1$ )	AUC	0.591±0.140	0.617±0.143	0.612±0.145
	P@10	0.256±0.253	<b>0.336</b> ±0.292	0.300±0.288
	AP	0.152±0.114	0.183±0.132	0.173±0.127
	RR	9.035±16.718	6.773±14.641	8.948±17.422
Propagation ( $\eta = 0.5$ )	AUC	0.598±0.141	0.623±0.142	0.621±0.147
	P@10	0.255±0.245	0.333±0.283	0.298±0.282
	AP	0.155±0.115	0.185±0.133	0.177±0.130
	RR	8.994±17.224	6.419±14.088	9.089±18.358

(b) Inter-topic STD.

TF	TF-IDF	LSI
0.114	0.112	0.119
0.174	0.191	0.195
0.097	0.102	0.104
6.063	3.676	6.088
<b>0.109</b>	<b>0.109</b>	<b>0.111</b>
0.166	0.181	<b>0.174</b>
0.102	0.110	0.115
<b>3.400</b>	<b>2.789</b>	<b>4.168</b>
0.128	0.126	0.131
0.156	0.187	0.184
<b>0.093</b>	0.102	<b>0.101</b>
9.175	6.139	9.245
0.125	0.122	0.128
<b>0.147</b>	<b>0.177</b>	0.175
<b>0.093</b>	<b>0.101</b>	0.102
8.488	5.059	8.944

Table – Results of the Document-query evaluations.

# Lessons Learned

BIRNDL  
2018

Robin  
Brochier,  
Adrien Guille  
, Benjamin  
Rothan,  
Julien Velcin

Introduction

The Usual  
Evaluation  
(topic-query)

Another  
Evaluation  
(document-  
query)

Results and  
Comparisons

Ongoing  
Works

Références

- Not the same algorithms that perform best for both evaluations.
- Document-query inter-topic standard deviations are lower than the overall query standard deviations.
- The (short) namings of the topics and their low number have a too strong influence on the evaluation.
- Worse results for the document-query evaluation highlights the need of annotations for documents.

# Ongoing Works

BIRNDL  
2018

Robin  
Brochier,  
Adrien Guille  
, Benjamin  
Rothan,  
Julien Velcin

Introduction

The Usual  
Evaluation  
(topic-query)

Another  
Evaluation  
(document-  
query)

Results and  
Comparisons

Ongoing  
Works

Références

- Human annotations of AMiner document-queries.
- Expert finding evaluation python library based on AMiner and StackExchange<sup>5</sup> data.
- Online evaluations based on volunteers of `peer.us`.
- Algorithm to embed experts and documents in the same vector space with the ability to embed previously unseen documents.

---

5. <https://stackexchange.com/>

# Thank You

BIRNDL  
2018

**Robin  
Brochier,  
Adrien Guille  
, Benjamin  
Rothan,  
Julien Velcin**

Introduction

The Usual  
Evaluation  
(topic-query)

Another  
Evaluation  
(document-  
query)

Results and  
Comparisons

**Ongoing  
Works**

Références

Questions ?

*GitHub : [github.com/brochier/impact\\_query\\_expert\\_finding](https://github.com/brochier/impact_query_expert_finding)*



# References I

BIRNDL  
2018

Robin  
Brochier,  
Adrien Guille  
, Benjamin  
Rothan,  
Julien Velcin



Nick Craswell et al. “P@ noptic expert : Searching for experts not just for documents”. In : *Ausweb Poster Proceedings, Queensland, Australia*. T. 15. 2001, p. 17.

Introduction

The Usual  
Evaluation  
(topic-query)

Another  
Evaluation  
(document-  
query)

Results and  
Comparisons

Ongoing  
Works

Références



Craig Macdonald et Iadh Ounis. “Voting for candidates : adapting data fusion techniques for an expert search task”. In : *Proceedings of the 15th ACM international conference on Information and knowledge management*. ACM. 2006, p. 387–396.



Pavel Serdyukov, Henning Rode et Djoerd Hiemstra. “Modeling multi-step relevance propagation for expert finding”. In : *Proceedings of the 17th ACM conference on Information and knowledge management*. ACM. 2008, p. 1133–1142.