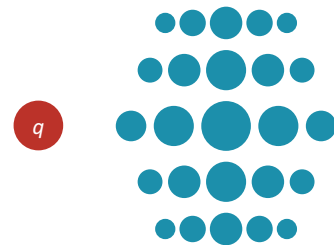


Ranking Models

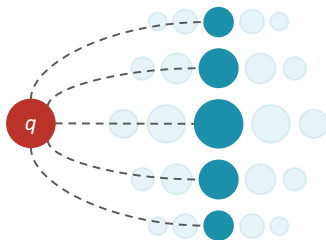
Feedback Models

Rodrygo L. T. Santos
rodrygo@dcc.ufmg.br

The ranking problem



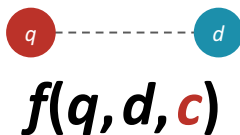
The ranking problem



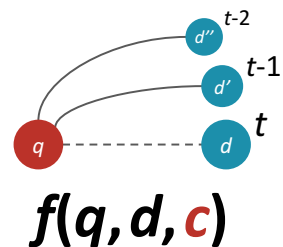
The ranking problem



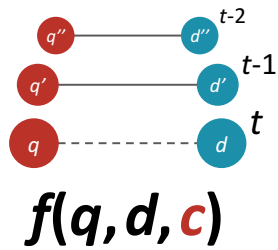
The ranking problem



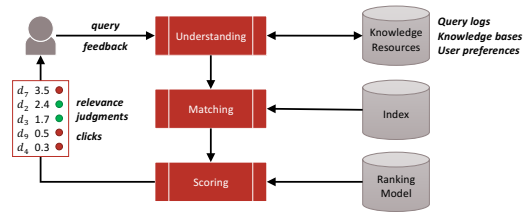
Exploiting interactions



Exploiting interactions



Eliciting feedback



Eliciting feedback

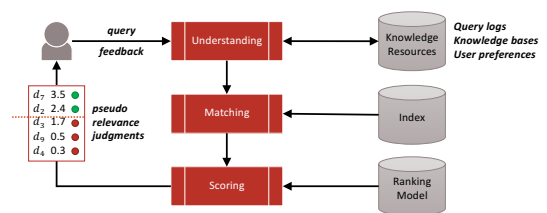
Explicit feedback

- Explicit relevance judgments
- Reliable, but costly

Implicit feedback

- Positive-only "judgments" (e.g., clicks, dwell time)
- Noisy and biased, but cheap and abundant

Simulating feedback



Simulating feedback

Pseudo-relevance feedback

- Top- k results are assumed to be relevant
- Very sensitive to ranking quality, but automatic

Exploiting feedback

Machine-learned ranking

- User feedback can be treated as supervision for learning effective ranking models
- See classes on **Learning to Rank**

Exploiting feedback

Query expansion

- Feedback documents can help enhance the user's query by providing related expansion terms

Example: [information retrieval]

- Relevant or pseudo-relevant documents may provide related terms like "search engine", "ranking"

Feedback in vector space models

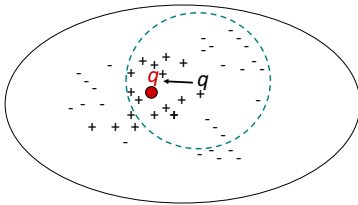
General idea: query modification

- Adding new (weighted) terms
- Adjusting weights of old terms

Rocchio (1971): most well-known approach

- Also effective and robust in practice

Rocchio method



Rocchio method

Standard operation in vector space

$$\vec{q}_m = \alpha \vec{q} + \frac{\beta}{|G|} \sum_{\forall \vec{d}_i \in G} \vec{d}_i - \frac{\gamma}{|\bar{G}|} \sum_{\forall \vec{d}_j \in \bar{G}} \vec{d}_j$$

Labels in the diagram:
 - \vec{q}_m : Modified query
 - \vec{q} : Original query
 - G : Rel docs (relevant documents)
 - \bar{G} : Non-rel docs (non-relevant documents)
 - α, β, γ : Parameters

Rocchio example

$V = \{\text{news, about, presidential, campaign, food, text}\}$

$\vec{q} = \{1, 1, 1, 1, 0, 0\}$

		news	about	pres.	campaign	food	text
-	d_1	{ 1.5	0.1	0.0	0.0	0.0	0.0 }
-	d_2	{ 1.5	0.1	0.0	2.0	2.0	0.0 }
+	d_3	{ 1.5	0.0	3.0	2.0	0.0	0.0 }
+	d_4	{ 1.5	0.0	4.0	2.0	0.0	0.0 }
-	d_5	{ 1.5	0.0	0.0	6.0	2.0	0.0 }

Rocchio example

		news	about	pres.	campaign	food	text
-	d_1	{ 1.5	0.1	0.0	0.0	0.0	0.0 }
-	d_2	{ 1.5	0.1	0.0	2.0	2.0	0.0 }
+	d_3	{ 1.5	0.0	3.0	2.0	0.0	0.0 }
+	d_4	{ 1.5	0.0	4.0	2.0	0.0	0.0 }
-	d_5	{ 1.5	0.0	0.0	6.0	2.0	0.0 }
+	C_r	{ $\frac{1.5+1.5}{2}$	0.0	$\frac{3.0+4.0}{2}$	$\frac{2.0+2.0}{2}$	0.0	0.0 }
-	C_n	{ $\frac{1.5+1.5+1.5}{3}$	$\frac{0.1+0.1+0.0}{3}$	0.0	$\frac{0.0+2.0+6.0}{3}$	$\frac{0.0+2.0+2.0}{3}$	0.0 }

Rocchio example

	{	news	about	pres.	campaign	food	text	}
+	C_r	{	$\frac{1.5+1.5}{2}$	0.0	$\frac{3.0+4.0}{2}$	$\frac{2.0+2.0}{2}$	0.0	0.0 }
-	C_n	{	$\frac{1.5+1.5+1.5}{3}$	$\frac{0.1+0.1+0.0}{3}$	0.0	$\frac{0.0+2.0+6.0}{3}$	$\frac{0.0+2.0+2.0}{3}$	0.0 }

$$\vec{q} = \{1, 1, 1, 1, 0, 0\}$$

$$\vec{q}_m = \alpha \cdot \vec{q} + \beta \cdot C_r - \gamma \cdot C_n$$

$$= \{\alpha + 1.5\beta - 1.5\gamma, \alpha - 0.067\gamma, \alpha + 3.5\beta, \alpha + 2\beta - 2.67\gamma, -1.33\gamma, 0\}$$

Rocchio in practice

Negative examples are not very important

- Non-relevant documents lack coherence

High-dimensional centroid estimation

- Restrict the vector onto a lower dimension

Training set is small and noisy and may be biased

- Keep relatively high weight on the original query (α)

Feedback in probabilistic models

Binary independence assumption

- Documents and queries as incidence vectors
 - $q = (x_1, x_2, \dots, x_{|V|})$ with $x_i \in \{0,1\}$
 - $d = (y_1, y_2, \dots, y_{|V|})$ with $y_i \in \{0,1\}$
- No association between terms

Binary Independence Model (BIM)

$$f(q, d) = O(G|d, q)$$

$$\approx \sum_{t \in d} \log \frac{p_t(1 - u_t)}{u_t(1 - p_t)}$$

where $p_t = P(y_t|G)$
 $u_t = P(y_t|\vec{G})$

Probability estimates without feedback

$$p_t = P(y_t|G)$$

- Cannot estimate without actual relevance data
→ assume constant (e.g., $p_t = 0.5$)

$$u_t = P(y_t|\vec{G})$$

- Non-relevant documents are the vast majority
→ assume collection statistics ($u_t = n_t/n$)

Probability estimates without feedback

$$\begin{aligned} f(q, d) &\approx \sum_{t \in d} \log \frac{p_t(1 - u_t)}{u_t(1 - p_t)} \\ &= \sum_{t \in d} \log \frac{0.5(1 - n_t/n)}{n_t/n(1 - 0.5)} \\ &= \sum_{t \in d} \log \frac{n - n_t}{n_t} \approx \sum_{t \in d} \log \frac{n}{n_t} \end{aligned}$$

Contingency table

	G	\bar{G}	Total
y_t	r_t	$n_t - r_t$	n_t
\bar{y}_t	$r - r_t$	$n - n_t - r + r_t$	$n - n_t$
Total	r	$n - r$	n

With relevance data

- $p_t = (r_t + 0.5)/(r + 1)$
- $u_t = (n_t - r_t + 0.5)/(n - r + 1)$

Probability estimates with feedback

$$f(q, d) \approx \sum_{t \in d} \log \frac{p_t(1 - u_t)}{u_t(1 - p_t)}$$

$$= \sum_{t \in d} \log \frac{(r_t + 0.5)/(r - r_t + 0.5)}{(n_t - r_t + 0.5)/(n - n_t - r + r_t + 0.5)}$$

BM25 adds tf saturation and query tf

Feedback in language models

Query likelihood model

$$\begin{aligned} \circ f(q, d) &= P(q|\theta_d) \\ &\propto \log P(q|\theta_d) \\ &= \sum_{t \in q} \text{tf}_{t,q} \log P(t|\theta_d) \end{aligned}$$

Difficulty

- Models documents, not query, i.e., $p(t|\theta_d)$ is fixed

Relevance models

Language model representing information need

- Query and feedback documents are samples
- $P(d|\theta_G)$: probability of generating the text in a document given a relevance model
- Document likelihood model
- Less effective than query likelihood due to difficulties comparing across documents of different lengths

Divergence-based ranking

Estimate relevance model from query and feedback

- Rank documents by similarity to relevance model

Kullback-Leibler divergence (KL-divergence)

$$\circ f(q, d) = -D_{\text{KL}}(\theta_G || \theta_d)$$

Divergence-based ranking

$$\begin{aligned} f(q, d) &= -D_{\text{KL}}(\theta_G || \theta_d) \\ &= - \sum_t P(t|\theta_G) \log \frac{P(t|\theta_G)}{P(t|\theta_d)} \\ &= \sum_t P(t|\theta_G) \log P(t|\theta_d) - \sum_t P(t|\theta_G) \log P(t|\theta_G) \end{aligned}$$

document independent

Divergence-based ranking

$$f(q, d) \propto \sum_t P(t|\theta_G) \log P(t|\theta_d)$$

- Without feedback, under MLE: $P(t|\theta_G) \propto \text{tf}_{t,q}$
- Relevance model degenerates to query likelihood

Estimating relevance models

Probability of pulling a word t out of the “bucket” representing the relevance model depends on the query terms we have just pulled out

$$\begin{aligned} P(t|\theta_G) &\approx P(t|q) \\ &= \frac{P(t, q)}{P(q)} \end{aligned}$$

Estimating relevance models

$$\begin{aligned} P(t, q) &= \sum_{d \in F} p(d) P(t, q|d) \\ &= \sum_{d \in F} p(d) P(t|q, d) P(q|d) \\ &\approx \sum_{d \in F} p(d) P(t|d) \prod_{t_i \in q} P(t_i|d) \end{aligned}$$

Estimating relevance models

$$P(t, q) \approx \sum_{d \in F} p(d) P(t|d) \prod_{t_i \in q} P(t_i|d)$$

Assuming uniform $P(d)$

- $P(t, q)$ is an average of query likelihood scores across feedback documents, weighted by $P(t|d)$

Example from top 10 docs

<i>president lincoln</i>	<i>abraham lincoln</i>	<i>fishing</i>	<i>tropical fish</i>
lincoln	lincoln	fish	fish
president	america	farm	tropic
room	president	salmon	japan
bedroom	faith	new	aquarium
house	guest	water	water
white	abraham	wild	species
america	new	caught	aquatic
guest	room	tag	fair
serve	christian	china	coral
bed	history	time	source
washington	public	eat	tank
old	bedroom	raise	reef
office	war	city	animal
war	politics	people	tarpon
long	old	fishermen	fishery
abraham	national	boat	

Example from top 50 docs

<i>president lincoln</i>	<i>abraham lincoln</i>	<i>fishing</i>	<i>tropical fish</i>
lincoln	lincoln	fish	fish
president	president	water	tropic
america	america	catch	water
new	abraham	reef	storm
national	war	fishermen	species
great	man	river	boat
white	civil	new	sea
war	new	year	river
washington	history	time	country
clinton	two	bass	tuna
house	room	boat	world
history	booth	world	million
time	time	farm	state
center	politics	angle	time
kennedy	public	fly	japan
room	guest	trout	mile

Summary

Acquiring feedback

- Explicit, implicit, simulated (pseudo) feedback

Exploiting feedback via query expansion

- Rocchio for VSM
- Observed relevance for BMI and BM25
- Feedback language models for LM

Challenges

Long queries are inefficient for typical search engines

- Only reweight certain prominent terms

Users are often reluctant to provide explicit feedback

- Effective pseudo-relevance feedback is challenging
- Implicit feedback is abundant, yet often biased

Feedback is also useful as a learning signal

References

[Text Data Management: A Practical Introduction to Information Retrieval and Text Mining](#), Ch. 6

Zhai and Massung, 2016

[Introduction to Information Retrieval](#), Ch. 9

Manning et al., 2008

[Search Engines: Information Retrieval in Practice](#), Ch. 7

Croft et al., 2009

References

[Relevance feedback in information retrieval](#)

Rocchio, 1971

[Relevance based language models](#)

Lavrenko and Croft, SIGIR 2001

[Model-based feedback in the language modeling approach to information retrieval](#)

Zhai and Lafferty, CIKM 2001

References

[A survey of automatic query expansion in information retrieval](#)

Carpineto and Romano, ACM Comp. Surveys 2012

UF  G UNIVERSIDADE FEDERAL DE MINAS GERAIS

Coming next...

Structural Models

Rodrygo L. T. Santos
rodrygo@dcc.ufmg.br