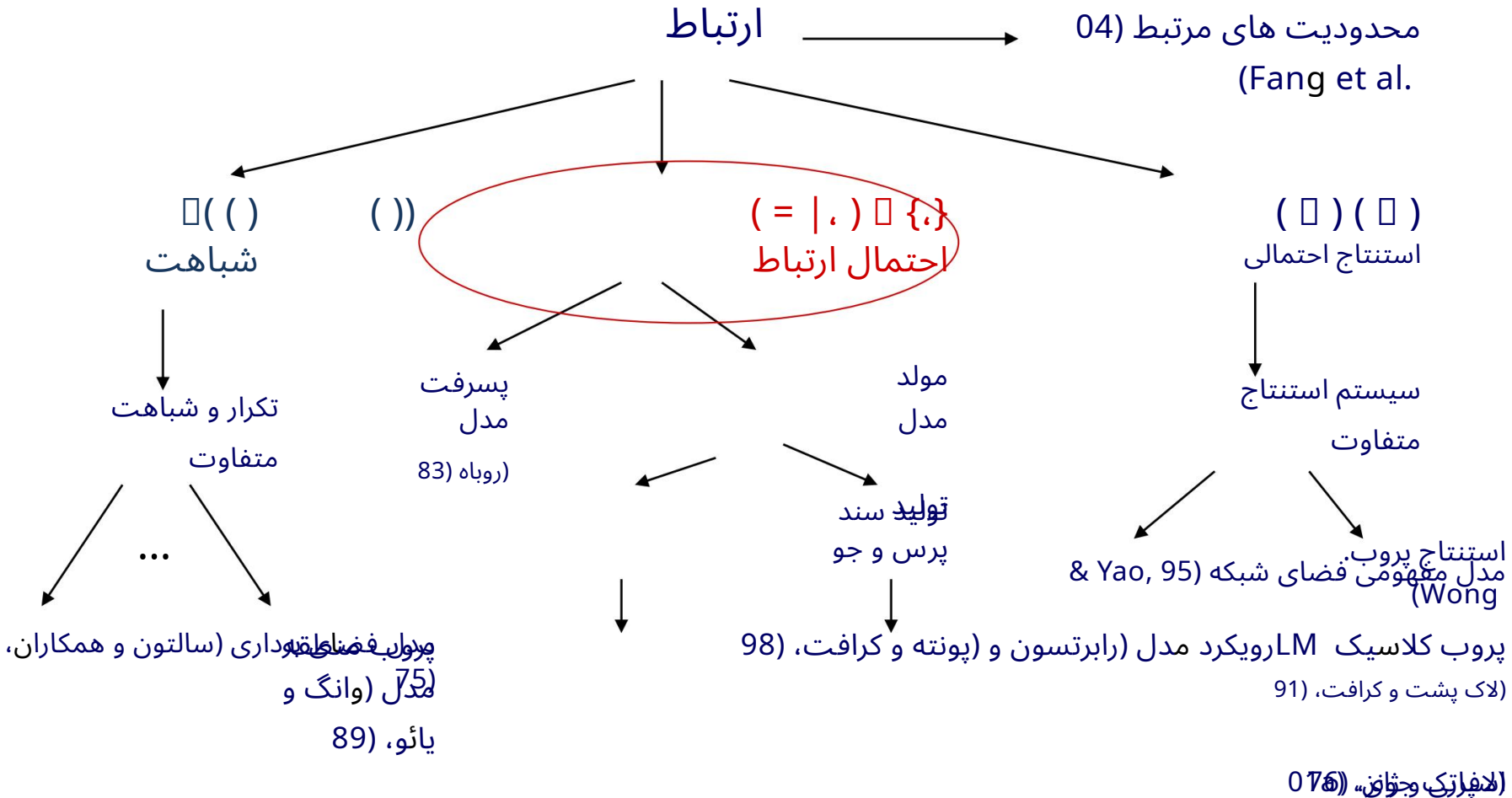


مدل های بازیابی:

احتمالی

بازیابی هوشمند اطلاعات

مفهوم ارتباط



اصل رتبه بندی احتمال

[رابرتسون 77]

• برگرداندن یک لیست رتبه بندی شده از اسناد به ترتیب احتمال نزولی که یک سند با پرس و جو مرتبط است، استراتژی بهینه تحت دو فرض زیر است:

1. سودمندی یک سند برای کاربر است مستقل از کاربرد هر سند دیگری.

2. یک کاربر نتایج را به ترتیب مرور می کند

طبق، PRP، تنها چیزی که نیاز داریم این است

"یک تابع اندازه گیری ارتباط"

که راضی می کند

$$(1, 1) > (2, 1) \quad | \quad (1, 2) > (2, 2)$$

برای

سوال اساسی

احتمال اینکه THIS سند با THIS query مرتبط باشد
چقدر است؟

به طور رسمی...

3 متغیر تصادفی: پرس و جو، ارتباط سند $\{0,1\}$

، با توجه به یک پرس و جو خاص، یک سند خاص $? = 1 = = =$
(| ،)

انحراف ...

بررسی مختصر احتمال

مفاهیم اساسی در احتمال

• آزمایش تصادفی: آزمایشی با نتیجه نامشخص
(به عنوان مثال، پرتاب یک سکه، انتخاب یک کلمه از متن)

• فضای نمونه: همه نتایج ممکن، به عنوان مثال، -پرتاب 2 سکه منصفانه، $\{ , \}$

اگر نتیجه در $\{ \}$ = (همه سرها) $\{ \}$ = (همان صورت) اتفاق می افتد

• رویداد: \square

، به عنوان مثال،

،

• احتمال رویداد: $P(\square) = 1$ (نتیجه همیشه در S)

$P(\square) = P(\square) + P(\square)$ اگر $\square = \square$ (مثلاً = صورت یکسان، = چهره متفاوت)

مفاهیم اساسی Prob. (ادامه)

• احتمال شرطی: $(A|B) = (A \cap B) / (B)$

$$- (A \cap B) = (B)(A|B) = (A)(B|A)$$

- بنابراین، $(A|B) = (A \cap B) / (B)$ (قاعده بیز)

- $(A|B) = (A \cap B) / (B)$ (قاعده بیز)

• احتمال سگین اگر پارتیشنی از $(A) = (A_1) + (A_2) + \dots + (A_n)$ تشکیل دهید

$$- (A) = (A_1) + (A_2) + \dots + (A_n)$$

- بنابراین، $(A|B) = (A \cap B) / (B)$

$$(A|B) = (A \cap B) / [(A_1 \cap B) + (A_2 \cap B) + \dots + (A_n \cap B)]$$

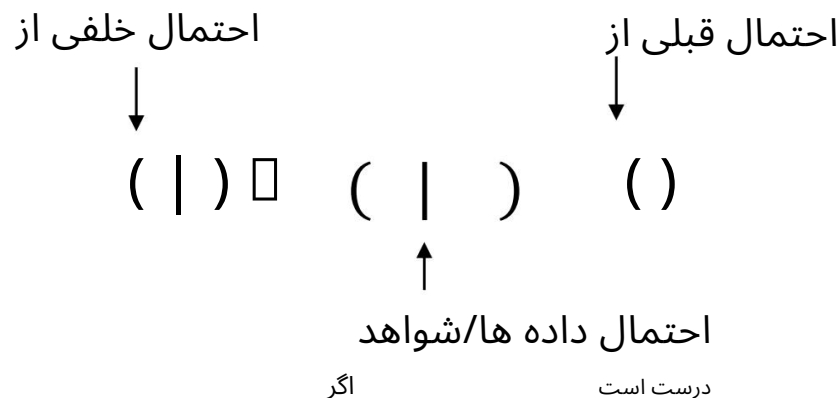
- این به ما امکان می دهد $(A|B)$ را بر اساس $(A_i|B)$ محاسبه کنیم.

تفسیر قانون بیز

{شواهد و مدارک: فضای فرضیه: $1, \dots$ } =

$$P(H_i | D) = \frac{P(D | H_i) P(H_i)}{P(D)}$$

می توانیم رها کنیم () ، اگر بخواهیم محتمل ترین فرضیه را انتخاب کنیم



انتهای مسیر انحرافی ...

حالا بیایید به بازیابی برگردیم!

مدلهای بازیابی احتمالی: شهود

فرض کنید تعداد زیادی قضاوت مرتبط داریم (مثلاً تعداد کلیک: «1» کلیک؛ «0» رد شده)

پرس و جو (Q)			ما می توانیم اسناد را بر اساس (R) و (D) رد کنیم		
Q1		1	$(\cdot, \cdot) =$		$(\cdot, \cdot, \cdot = 1)$
Q1	D2	1	$(\cdot = 1 \cdot, \cdot)$		(\cdot, \cdot)
Q1	D3	0	$(\cdot \cdot)$		$= 1 \cdot 1, 1 \neq 2$
Q1	D4	0			$(= 1 1, 2) \neq 2$
Q1	D5	1			$(= 1 1, 3) \neq 2$
...					
Q1	D1	0			
Q1	D2	1			
Q1	D3	0			
Q2	D3	1			
Q3	D1	1			
Q4	D2	1			
Q4	D3	0			
...					

چه می شود اگر ما به اندازه کافی نداشته باشیم
ورود به سیستم جستجو؟

اسناد/پرس و جوهای دیده نشده؟

ما می توانیم تقریبی کنیم $(= 1 | \cdot, \cdot)$

احتمال ارتباط

• سه متغیر تصادفی

- پرس و جو

- سند

- ارتباط $\{0,1\}$

• هدف: رتبه بر اساس $(= 1 | ,)$

- ارزیابی $(= 1 | ,)$

- در واقع، فقط باید $(= 1 | , 1)$ را با $(= 1 | , 2)$ مقایسه کنید،
یعنی اسناد را رتبه بندی کنید

• چندین روش مختلف برای پالایش $(= 1 | ,)$

پالایوشن | (F) مدل های شرطی

• ایده اصلی: ارتباط بستگی به این دارد که یک پرس و جو چقدر با یک سند مطابقت دارد

به عنوان مثال، #مطابقات، بالاترین -ویژگی ها را در \times تعریف کنید
IDF یک ترم منطبق، doclen یا حتی
(،) با توجه به هر تابع بازیابی دیگری، ...

$$- (= 1 | ,) =$$

$$(1 (,), 2 (,), \dots , (,))$$

-استفاده از داده های آموزشی (قضاوت های مرتبط شناخته شده) برای
پارامتر برآورد

-از مدل برای رتبه بندی اسناد جدید استفاده کنید

پالایش (2:)

مدل های تولیدی

• ایده پایه

-تعریف کردن ، ()
(, | 1 =) -را با استفاده از قانون بیز محاسبه کنید

$$(= 1 | ,) = \frac{(= 1 | ,)}{(= 0 | ,)} = \frac{(, | = 1)}{(, | = 0)} \times \frac{(= 1)}{(= 0)}$$

برای رتبه D نادیده گرفته شد

• موارد خاص

(- (, |))
- Query "generation": (, |)

تولید سند

$$\left(\begin{array}{c} = 1 \\ = 0 \end{array} \middle| \begin{array}{c} , \\ , \end{array} \right) = \frac{\left(\begin{array}{c} , \\ , \end{array} \middle| \begin{array}{c} , \\ , \end{array} \right) = 1 \left(\begin{array}{c} , \\ , \end{array} \middle| \begin{array}{c} , \\ , \end{array} \right) = 0 \left(\begin{array}{c} , \\ , \end{array} \middle| \begin{array}{c} , \\ , \end{array} \right)}{\left(\begin{array}{c} , \\ , \end{array} \middle| \begin{array}{c} , \\ , \end{array} \right) = 0 \left(\begin{array}{c} , \\ , \end{array} \middle| \begin{array}{c} , \\ , \end{array} \right)}$$

$$\square \frac{\left(\begin{array}{c} , \\ , \end{array} \middle| \begin{array}{c} , \\ , \end{array} \right) = 1 \left(\begin{array}{c} , \\ , \end{array} \middle| \begin{array}{c} , \\ , \end{array} \right) = 0 \left(\begin{array}{c} , \\ , \end{array} \middle| \begin{array}{c} , \\ , \end{array} \right)}{\left(\begin{array}{c} , \\ , \end{array} \middle| \begin{array}{c} , \\ , \end{array} \right) = 0 \left(\begin{array}{c} , \\ , \end{array} \middle| \begin{array}{c} , \\ , \end{array} \right)}$$

$$\square \frac{\left(\begin{array}{c} , \\ , \end{array} \middle| \begin{array}{c} , \\ , \end{array} \right) = 1 \left(\begin{array}{c} , \\ , \end{array} \middle| \begin{array}{c} , \\ , \end{array} \right) = 0 \left(\begin{array}{c} , \\ , \end{array} \middle| \begin{array}{c} , \\ , \end{array} \right)}{\left(\begin{array}{c} , \\ , \end{array} \middle| \begin{array}{c} , \\ , \end{array} \right) = 0 \left(\begin{array}{c} , \\ , \end{array} \middle| \begin{array}{c} , \\ , \end{array} \right)}$$

مدل اسناد مربوطه برای Q

مدل اسناد غیر مرتبط برای Q

تولید سند (ادامه)

- ویژگی های مستقل را فرض کنید ... که در آن $\{ , \}$ مقدار ویژگی است
- اجازه دهید $=$ (به طور مشابه $= \dots$) ، ...

$$\dots \square \frac{, = 1)}{, = 0)} \quad (| (|$$

$$\begin{array}{c} | (= | , = 1) \\ \hline (= , = 0) \\ =1 \end{array}$$

$$\begin{array}{c} (= 1 | , (= 1 |) \\ \hline = 0) \\ =1, =1 \end{array} \quad \square \frac{(= 0 | , (= 0 |)}{= 0)}{=1, =0}$$

تولید سند (ادامه)

$$\dots \quad \underset{=1, =1}{=} \frac{(\quad = 1 \mid , (\quad = \bar{1} \mid))}{= 0)} \quad \square_{=1, =0} \frac{(\quad = 0 \mid , (\quad = \bar{0} \mid))}{= 0)}$$

$$\times \square_{=1, =1} \frac{(\quad = 0 \mid , \quad = 1)}{(\quad = 0 \mid , \quad = 0)} \quad \times \square_{=1, =1} \frac{(\quad = 0 \mid , (\quad = \bar{0} \mid))}{= 1)}$$

$$\underset{=1, =1}{=} \frac{(\quad = 1 \mid , (\quad = \bar{1} \mid))}{= 0)} \cdot \frac{(\quad = 0 \mid , \quad = 0)}{(\quad = 0 \mid , \quad = 1)}$$

$$\underset{=1}{\boxed{\frac{(\quad = 0 \mid , (\quad = \bar{0} \mid))}{= 0)}}}$$

برای رتبه بندی نادیده بگیرید

0 فرض کن اکتید عبارت 1 در پرسی و جو(ظاهر نمی شود) $(= 0)$

$$\dots \quad \begin{array}{c} \text{=} \\ \text{=1, =1} \end{array} \quad \frac{(\text{=1} \mid, (\text{=1} \mid \text{=1}))}{(\text{=0} \mid, \text{=0})} \cdot \frac{(\text{=0} \mid, \text{=0})}{(\text{=0} \mid, \text{=1})} \quad \Bigg| \quad \frac{(\text{=0} \mid, (\text{=0} \mid \text{=1}))}{(\text{=0} \mid, \text{=1})}$$

برای رتبه بندی نادیده بگیرید

[illegible]

مدل رابرتسون-اسپارک جونز

(رابرتسون و اسپارک جونز (76)

$$\text{مدل (RSJ)} \quad \frac{(1-)}{(1-)} \quad \text{ورود به سیستم } \square \square \text{ (} = 1 | , \text{)}$$

دو پارامتر برای هر عبارت: $(= | = (= |$

، $\text{prob.} = \text{این اصطلاح در یک سند}$
 ، $\text{غیر مرتبط وجود دارد}$

چگونه پارامترها را تخمین بزنیم؟
 فرض کنید قضاوت های مرتبطی داریم،

$$\square = \frac{\#(\cdot \quad h \#) + \#(50.5)}{\#(\cdot \quad) + 1}$$

مدل RSJ بدون اطلاعات مرتبط

(کرافت و هارپر 79)

$$\text{مدل (RSJ)} = \frac{(1 - \alpha) (1 - \beta)}{\alpha + \beta} \text{ ورود به سیستم } (= 1 | ,)$$

$$= 1, = 1$$

چگونه پارامترها را تخمین بزنیم؟
فرض کنید قضاوت مربوطه نداریم، - ما ثابت فرض می کنیم - با فرض غیر مرتبط بودن همه اسناد، تخمین بزنید.

$$\frac{\alpha + 0.5}{\alpha + 0.5} \text{ ورود به سیستم } (= 1 | ,)$$

$$= 1, = 1$$

اسناد در مجموعه : # اسنادی که عبارت در آنها وجود دارد

مدل: RSI خلاصه

• مهم ترین پروب کلاسیک. مدل • IR فقط از عبارت حضور/غیبت استفاده کنید، بنابراین به عنوان مدل استقلال باینری نیز شناخته می شود

• اساسا بیز ساده برای رتبه بندی اسناد • طبیعی ترین بازخورد مرتبط/شبه • وقتی بدون قضاوت مرتبط بودن، پارامترهای مدل باید به صورت موقتی تخمین زده شوند.

• عملکرد به خوبی مدل VS تنظیم شده نیست

بهبود: IRSJ اضافه کردن TF

$$\frac{\text{PRQD} \quad \text{سند پایه مدل نسل:}}{\text{PDQR} \quad \text{PDQR}} \quad \frac{\text{PDQR} \quad 1)}{\text{PDQR}}$$

تعداد دفعات ترم کجاست ... ، بگذار =

$$\frac{\text{PRQD} \quad \text{PA d QR} \quad 1)}{\text{PRQD} \quad \text{PA d QR} \quad 0)} \quad \frac{\text{PA d QR} \quad 1)}{\text{PA d QR} \quad 0)} \quad \frac{\text{PAQR} \quad 1)}{\text{PAQR} \quad 0)} \quad \frac{\text{PAQR} \quad 1)}{\text{PAQR} \quad 0)} \quad \frac{\text{PAQR} \quad 1)}{\text{PAQR} \quad 0)} \quad \frac{\text{PAQR} \quad 1)}{\text{PAQR} \quad 0)}$$

$$\frac{\text{2-مدل مخلوط بواسون}}{p(A f | Q, R) p(E | Q, R) p(A f | E) p(E | Q, R) p(A f | E)}$$

بسیاری از پارامترهای بیشتر برای تخمین زدن!

تقریبی BM25/Okapi

(رابرتسون و همکاران 94)

Idea: تقریبی (, | 1 =) با یک تابع ساده تر که دارای ویژگی های مشابه است

• مشاهدات:

(, | 1 =) - log مجموع وزن های عبارت است ، = 0 - اگر = 0

-یکنواخت با افزایش می یابد - دارای حد جانبی است

$$\frac{(1+1)}{1+} \cdot \frac{(1-)}{(1-)} \text{ ورود به سیستم}$$

• تابع ساده = است

افزودن Doc. طول و پرس و جو TF

- گنجاندن طول سند

- "با دقت" سند طولانی را جریمه کنید

- ترکیب پرس و جو TF

- یک تبدیل TF مشابه

- فرمول نهایی BM25 نامیده می شود که به بالا می رسد
عملکرد TREC

فرمول BM25

$$\sum_{T \in Q} w^{(1)} \frac{(k_1 + 1)tf}{K + tf} \frac{(k_3 + 1)qtf}{k_3 + qtf} \quad (1)$$

where

Okapi TF/BM25 TF

Q is a query, containing terms T

$w^{(1)}$ is the Robertson/Sparck Jones weight [5] of T in Q

$$\log \frac{(r + 0.5)/(R - r + 0.5)}{(n - r + 0.5)/(N - n - R + r + 0.5)} \quad (2)$$

N is the number of items (documents) in the collection

n is the number of documents containing the term

R is the number of documents known to be relevant to a specific topic

r is the number of relevant documents containing the term

K is $k_1((1 - b) + b \cdot dl / avdl)$

k_1 , b and k_3 are parameters which depend on the on the nature of the queries and possibly on the database; k_1

and b default to 1.2 and 0.75 respectively, but smaller values of b are sometimes advantageous; in long queries k_3

is often set to 7 or 1000 (effectively infinite)

tf is the frequency of occurrence of the term within a specific document

qtf is the frequency of the term within the topic from which Q was derived

dl and $avdl$ are respectively the document length and average document length measured in some suitable unit.

پسوندهای

مدل های "نسل سند".

• وابستگی اصطلاحی را ضبط کنید (Rijsbergen & Harper 78)

• روش های جایگزین برای ترکیب TF (Croft 83، Kalt96)

• انتخاب ویژگی/مدت برای بازخورد (TREC Okapi)

گزارش ها)

• برآورد مدل ربط بر اساس

بازخورد شبه [Lavrenko & Croft 01]

سوالات؟