

به نام خدا



دانشگاه تهران

دانشکدگان فنی

دانشکده مهندسی برق و
کامپیوتر



درس بازیابی هوشمند اطلاعات

پاسخ بخش تئوری تمرين ۱

نام و نام خانوادگی: سهیل حاجیان منش

شماره دانشجویی: 810100119

مهر ماه ۱۴۰۳

فهرست

اظهارنامه	3
پاسخ سوال اول	3
پاسخ بخش اول	3
فرمولbm25:	3
فرمولPL2:	3
پاسخ بخش دوم	3
پاسخ سوال دوم	5
پاسخ بخش اول	5
پاسخ بخش دوم	5
پاسخ بخش سوم	5
پاسخ سوال سوم	5
پاسخ بخش اول	5
پاسخ بخش دوم	6
پاسخ بخش سوم	6
پاسخ بخش چهارم	6
پاسخ سوال چهارم	6
پاسخ بخش اول	6
پاسخ بخش دوم	6
پاسخ بخش سوم	7
پاسخ بخش چهارم	7
پاسخ بخش پنجم	8
پاسخ سوال پنجم	8
پاسخ بخش اول	8
پاسخ بخش دوم	9
پاسخ بخش سوم	9
پاسخ بخش چهارم	9

اظهار نامه

تأیید میکنم که از LLM ها مطابق با دستورالعملهای بارگذاری شده در سامانه Elearn درس به طور مسئولانه استفاده کرده ام. تمام اجزای کار خود را درک میکنم و آماده بحث شفاهی درباره آنها هستم

پاسخ سوال اول

پاسخ بخش اول

: فرمول BM25

$$TF \text{ weighting} = \frac{(k_1 + 1) \cdot c(w, d) \cdot (k_3 + 1) \cdot c(w, d)}{(k_1((1-b)+b \cdot \frac{|d|}{avdl}) + c(w, d))(k_3 + c(w, q))}$$

$$IDF \text{ weighting} = \text{LN}(\frac{N - df(w) + 0.5}{df(w) + 0.5})$$

$$Length \text{ Normalization} = \frac{1}{k_1 \left((1 - b) + b \cdot \frac{|d|}{avdl} \right)}$$

: فرمول PL2

$$Length \text{ Normalization} \Rightarrow tfn_w^d$$

داخل این فرمول تاثیر منفی (Penalize) افزایش طول سند اتفاق میفتند. در کل فرمول از این تابع استفاده شده تا این تاثیر منفی بالا بودن طول سند یا تاثیر مثبت پایین بودن طول آن کنترل شود.

$$IDF \text{ Weighting} \Rightarrow \lambda_w^d$$

چون λ از فرکانس term در کل collection میاد، عملأ نقش inverse doc frequency رو ایفا میکنه. هر چه term کمیاب تر باشد مقدار λ بزرگتر میشود و بلعکس.

TF weighting در کل قسمت های فرمول نقش موثر را ایفا می کند.

پاسخ بخش دوم

همانطور که در قسمت قبل گفته شود بارامتر C موجود در L_2 فرمول Length Normalization بر روی $\|f\|_2$ تاثیر می‌گذارد و کنترل می‌کند تا چه میزان طول اسناد نرمال سازی شود.

هر چقدر C بزرگتر شود، نرمال سازی بیشتری انجام می‌شود و تفاوت‌های بین طول سند با میانگین طول سند با ضرب C شدن بیشتر می‌شود. در نتیجه اسناد با طول بیشتر از میانگین، شدیدتر $penalize$ می‌شوند و اسناد با طول کمتر بیشتر بهشون امتیاز داده می‌شوند.

اما هرچقدر C کوچکتر باشد، نرمال سازی کمتری انجام می‌شود و رفتار به TF خام نزدیک تر می‌شود و سند‌های بزرگتر ممکن است امتیاز بیشتری از سند‌های کوچکتر بگیرند با این اینکه از لحاظ محتوا فرقی ندارند. در واقع خود TF اثر پر رنگ تری اینجا دارد. هنگامی که اسنادمان تفاوت چندانی در طول با یکدیگر ندارند، میتوان مقدار C را کمتر گذاشت.

پاسخ سوال دوم

پاسخ بخش اول

نیاز به Relevance judgement هایی داریم که مشخص شود که هر سند به پرس و جو مرتبط است یا خیر تا سپس بتوانیم به ازای هر A_i موجود در یک سند ببینیم احتمال اینکه A_i در سند های مرتبط آمده باشد و احتمال اینکه در سند های نامرتبط آمده باشد چقدر است.

اگر موجود نباشد، فرض میکنیم:

۱. p_i یک مقدار ثابت است.
۲. q_i را تخمین میزنیم با فرض اینکه تمامی اسناد موجود نامرتبط هستند. یعنی q_i برابر است با اینکه هر A_i در تمامی اسناد ظاهر شده است تقسیم بر کل اسناد.

پاسخ بخش دوم

مدل بصورت زیر ساده می شود:

$$\begin{aligned} \sum_{i=1, di=q_i=1}^k \log \frac{p_i(1-q_i)}{q_i(1-p_i)} &\approx \sum_{i=1, di=q_i=1}^k \log \frac{1-q_i}{q_i} = \sum_{i=1, di=q_i=1}^k \log \frac{1-\frac{ni}{N}}{\frac{ni}{N}} = \sum_{i=1, di=q_i=1}^k \log \frac{N-ni}{ni} \\ &\approx \sum_{i=1, di=q_i=1}^k \log \frac{N-ni+0.5}{ni+0.5} \end{aligned}$$

پاسخ بخش سوم

هر چه n_i بزرگتر شود صورت لگاریتم کوچکتر و مخرج آن بزرگتر می شود که در نتیجه احتمال مرتبط بودن و وزن کلمه کاهش می یابد. این رفتار مشابه محاسبه IDF در الگوریتم BM25 است که کلمات هر چه در داکیومنت های بیشتری آمده باشند، مقدار IDF آنها کمتر می شود و وزن کمتری دارند.

پاسخ سوال سوم

پاسخ بخش اول

کافیست $b = 0$ باشد. چون طبق فرمول در مخرج داریم :

$$k_1 \left((1 - b) + b \cdot \frac{|d|}{avdl} \right)$$

که اگر $b = 0$ باشد این مقدار همیشه برای k_1 می شود.

پاسخ بخش دوم

خیر مانند بخش اول همواره k_1 جواب عبارت نوشته شده می شود و فرمول به صورت زیر ساده می گردد :

$$\text{LN}\left(\frac{N - df(w) + 0.5}{df(w) + 0.5}\right) \cdot \frac{(k_1 + 1) \cdot c(w, d)}{k_1 + c(w, d)} \cdot \frac{(k_3 + 1) \cdot c(w, q)}{k_3 + c(w, q)}$$

پاسخ بخش سوم

چون b در مخرج است و محور x نمودار است پس ضریب آن اگر مثبت باشد، شیب نمودار منفی و ضریب آن اگر منفی باشد، شیب نمودار مثبت می شود. اگر $|d|$ بزرگتر باشد، ضریب b منفی و شیب خط مثبت و اگر $|d|$ کوچکتر باشد، ضریب b مثبت و شیب خط منفی میشود. در صورت برابر بودن نیز شیب خط صفر است.

پاسخ بخش چهارم

با افزایش k_1 از 1.0 تا 3 تاثیر تغییرات مقدار b بر امتیاز نهایی شدیدتر می شود و در نقطه $k_1 = 3$ بیشترین تاثیر تغییر b را داریم. دلیل آن هم این است که k_1 در ضریب پشت متغیر b ضرب می شود و تاثیر مستقیم بر شیب دارد. هر چه k_1 بیشتر باشد، ضریب پشت b بزرگتر می شود.

پاسخ سوال چهارم

پاسخ بخش اول

بله به وضوح سند اول تعریف مختصری از ماشین لرنینگ ارائه می دهد و سند دوم یکی از زیر مجموعه های ماشین لرنینگ را توضیح می دهد. پس سند اول با پرس و جو مرتبط تر است و به درستی تشخیص داده شده است.

متوسط طول اسناد 23 و طول سند اول 26 و طول سند دوم 20 است.

با افزایش مقدار پارامتر b میتوان ترتیب رتبه بندی اسناد را تغییر داد. چون پارامتر b میزان جریمه کردن امتیاز اسناد بزرگ و افزایش امتیاز اسناد کوچکتر از میانگین را تنظیم می کند و هرچقدر بیشتر باشد، اسناد بزرگ تر از میانگین بیشتر جریمه می شوند. پس با رسیدن مقدار b از عدد 0.9 به بعد، سند دوم (که کوچکتر از میانگین است) در جایگاه 1 قرار می گیرد و سند اول (که بزرگتر از میانگین است) در جایگاه دوم قرار می گیرد. در این حالت امتیاز سند دوم 1.4373 و سند اول 1.4367 می باشد.

پاسخ بخش دوم

بله کماکان رتبه بندی درست است و سند اول ارتباط معنایی بیشتری با پرس و جو دارد و به درستی در جایگاه اول قرار گرفته است.

این می توان با افزایش پارامتر k_1 ترتیب را بهم زد. پارامتر k_1 کنترل می کند تا چه اندازه تکرار کلمه در سند تاثیر گذار باشد. چون در سند دوم هم کلمه *nachi ne* و هم کلمه *algorithm* بیشتر از سند اول تکرار شده و تنها کلمه *al* در سند اول بیشتر از سند دوم تکرار شده است، پس می توان با افزایش k_1 از مقدار 2.1 به بعد از آن سند دوم را در جایگاه اول قرار داد.

در این حالت ($k=2.1$) امتیاز سند دوم 1.9691 و امتیاز سند اول 1.9657 می باشد.

چون هر دو سند طول برابر دارند پارامتر b در اینجا بی ثمر است.

پاسخ بخش سوم

Sample 3 : سند دوم طول بیشتر از میانگین و سند اول طول کمتر از میانگین دارد در باقی ویژگی های هر کلمه هر دو سند رفتار یکسانی از خود نشان می دهد. پس هر چقدر هم که b کوچکتر باشد، باز هم امتیاز سند کوچکتر بهتر می شود یا در بهترین حالت ($b=0$) امتیاز یکسانی دارند. این موضوع به این دلیل است که BM25 به جایگاه و ترتیب کلمات در کنار هم کاری ندارد.

Sample 4 : هیچ کاری به معنی سند ندارد و تنها به میزان تکرار کلمات در سند توجه می کند به همین دلیل می توان با تکرار بی دلیل یک کلمه به میزان قابل توجه امتیاز بالایی در جست و جو کسب کرد. سند اول امتیازش به مراتب از سند دوم پایینتر است چون کلمه آیفون دو بار در آن تکرار شده و کلمه جدید هم در آن تکرار نشده است.

Sample 5 : باز هم تا حدودی مشابه قبل است، نمیتوان تشخیص داد که منظور از آیفون جدید، موبایل آیفون جدید است پس صرفا به کمک تعداد تکرار کلمات در اسناد رتبه بندی را مشخص می کند و چون $b=0$ است پس سند دوم رتبه بهتری دارد چون تعداد آیفون بیشتر در آن تکرار شده است.

Sample 6 : این الگوریتم تکرار های مترادف های یک کلمه (ماشین = خودرو، برقی = الکتریکی) را در نظر نمی گیرد به همین خاطر سند اول که کلمات مشابه بیشتری دارد رتبه بهتری گرفته است. (مشکل vocabulary mismatch)

پاسخ بخش چهارم

معمولا Title شامل اطلاعات مهم تری است و باید وزن دهی بیشتری نسبت به بدن سند داشته باشد. به عنوان مثال $TF=1$ در title به طور معمول مهم تر از $TF=1$ در بدن سند. بنابراین نمیتوان امتیاز Title و بدن سند را با یک فرمول حساب کرد. پس ایده تجمیع امتیاز ها بصورت خام کار نمی کند همانطور که در مثال هم مشخص است (سند 1 مرتبط تر تشخیص داده شده است با وجود اینکه سند 2 در مورد آیفون 15 صحبت می کند).

یک ایده می تواند این باشد که دو فرمول جداگانه برای محاسبه امتیاز title و body داشت و در نهایت امتیاز سند برابر است با تجمیع این دو امتیاز. در فرمول title length normalization اخیلی کمتری نسبت به فرمول سابق داشته باشیم چون انتظار می رود title معمولا در یک محدوده طولی باشد. از طرفی مقدار $k1$ را بزرگتر بگیریم تا تکرار هر کلمه تاثیر بیشتری روی امتیاز title بگذارد.

مشکلات BM25 در اسناد چندبخشی

پایه ای فرض می کند سند یک متن بدون ساختار است.

ولی در وب / مقالات علمی / ... اسناد چند بخش دارند.

وقتی بخواهیم BM25 را روی این اسناد اجرا کنیم چند مشکل بزرگ ایجاد می شود:

1. اگر فیلدها را جدا score بدھیم و جمع کنیم → غیر خطی بودن TF از بین می رود → دوباره term independence بر می گردد (که غلط است).
2. IDF برای هر فیلد ثابت نیست، چون title کوتاه است و body بلند. در نتیجه آمار جهانی term در هر field متفاوت و unstable می شود.
3. اگر وزن همه field ها = 1 بگیریم، باز همان حالت unstructured را بر نمی گرداند چون رفقار TF غیرخطی است.
4. مشخص نیست length normalization را باید روی تک تک فیلدها بدھیم یا مجموع سند؟
5. پارامتر tuning هم مشکل می شود (برای هر field باید b $k1$, b $optimal$ شود)

BM25F

هر field یک weight دارد مثلاً $title=6$, $body=2$

TF field های با وزن ضرب می‌شوند و یک pseudo-frequency ساخته می‌شود

سپس BM25 معمولی فقط روى همین pseudo-frequency اعمال می‌شود

اگر همه weight ها = 1 باشد، رفتار دقیقاً مثل یک سند unstructured واحد می‌شود.

توضیحات خلاصه شده BM25F از مقاله [A Tutorial on the BM25F Model](#) گرفته شده است.

پاسخ بخش پنجم

رتبه بندی درست نیست اما چون کلمه پرواز IDF یکسانی با کلمات مثل تهران، مالزی و vip دارد ولی تعداد تکرار آن 2 است پس امتیاز بیشتری گرفته و همان باعث شده که رتبه سند دو بالاتر شود.

اهمیت توکن مالزی بیشتر است چون از لحاظ مفهومی فرد لزوماً دنبال پرواز نیست و حتماً دنبال پرواز به سمت مالزی است. پس مالزی مهم تر است ولی در معیار BM25 هر دو IDF یکسانی گرفته اند. پس لزوماً وجود IDF نمیتواند تعیین کننده اهمیت یک term در جمله باشد و مفهوم کلمات در کنار هم نیز اهمیت دارد.

پاسخ سوال پنجم

پاسخ بخش اول

Ranking 1

Item B	↑	↓	Rel: 0
Item D	↑	↓	Rel: 3
Item F	↑	↓	Rel: 1
Item E	↑	↓	Rel: 2
Item C	↑	↓	Rel: 0
Item A	↑	↓	Rel: 0

Evaluation Metrics

NDCG@5: 0.683

DCG@5 = $0/\log_2(2) + 3/\log_2(3) + 1/\log_2(4) + 2/\log_2(5) + 0/\log_2(6) = 3.254$, IDCG@5 = $3/\log_2(2) + 2/\log_2(3) + 1/\log_2(4) + 0/\log_2(5) + 0/\log_2(6) = 4.762$

P@1: 0.000

P@1 = 0

P@5: 0.600

P@5 = 3/5

R@5: 1.000

R@5 = 3/3

MRR: 0.500

MRR = $1/2 = 0.500$

AP: 0.639

P@2 = 0.500, P@3 = 0.667, P@4 = 0.750

R-Precision: 0.667

RP = 2/3

Ranking 2

Item F	↑	↓	Rel: 1
Item B	↑	↓	Rel: 0
Item A	↑	↓	Rel: 0
Item C	↑	↓	Rel: 0
Item E	↑	↓	Rel: 2
Item D	↑	↓	Rel: 3

Evaluation Metrics

NDCG@5: 0.372

DCG@5 = $1/\log_2(2) + 0/\log_2(3) + 0/\log_2(4) + 0/\log_2(5) + 2/\log_2(6) = 1.774$, IDCG@5 = $3/\log_2(2) + 2/\log_2(3) + 1/\log_2(4) + 0/\log_2(5) + 0/\log_2(6) = 4.762$

P@1: 1.000

P@1 = 1

P@5: 0.400

P@5 = 2/5

R@5: 0.667

R@5 = 2/3

MRR: 1.000

MRR = $1/1 = 1.000$

AP: 0.633

P@2 = 1.000, P@5 = 0.400, P@6 = 0.500

R-Precision: 0.333

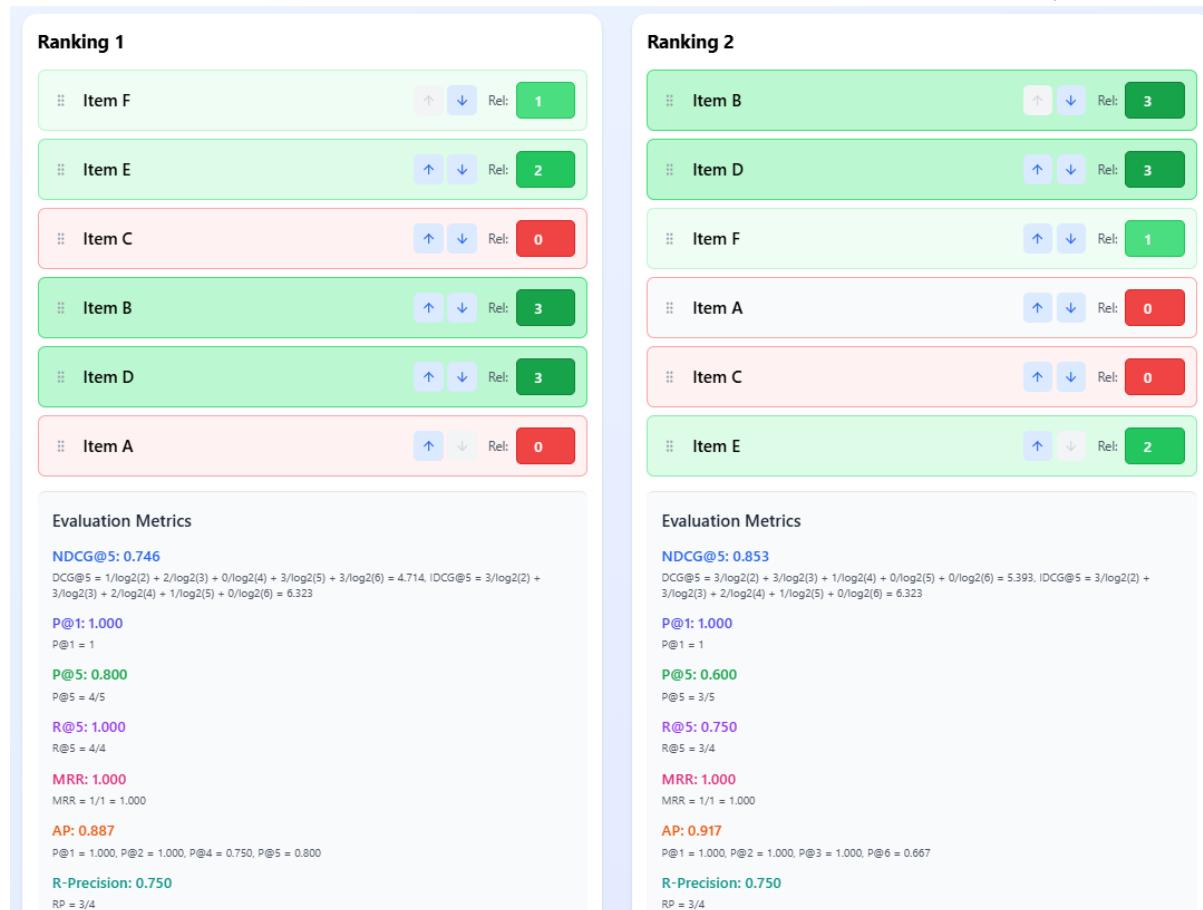
RP = 1/3

پاسخ بخش دوم

MRR زمانی مفید است که هدف پیدا کردن سریع اولین پاسخ درست باشد، مثل سیستم‌های Question Answering یا چت‌بات‌ها که مهم است جواب صحیح در رتبه‌های اول ظاهر شود.

پاسخ بخش سوم

بله مطالق دو مدل پیشنهاد شده عکس زیر :



پاسخ بخش چهارم

معیار NDCG محاسبه می‌کند که رتبه بندی انجام شده تا چه اندازه به رتبه ایده آل اسناد برای آن کوئری مشابهت دارد و یک امتیازی برای هر رتبه می‌دهد.

فرمولش به صورت رو به رو است :

$$DCG = \sum_{i=1}^k \frac{2^{rel_i} - 1}{\log i + 1}$$

مقدار IDCG هم با همین فرمول ولی با اسناد بصورت ایده آل مرتب شده حساب می‌شود.

از تقسیم DCG بر IDCG مقدار NDCG رتبه بند محاسبه می‌گردد.

پاسخ بخش پنجم

دو مدل پیشنهادی :

Ranking 1

⋮ Item F	↑ ↓	Rel: 1
⋮ Item E	↑ ↓	Rel: 2
⋮ Item C	↑ ↓	Rel: 0
⋮ Item B	↑ ↓	Rel: 0
⋮ Item A	↑ ↓	Rel: 0
⋮ Item D	↑ ↓	Rel: 0

Evaluation Metrics

NDCG@5: 0.860

DCG@5 = $1/\log_2(2) + 2/\log_2(3) + 0/\log_2(4) + 0/\log_2(5) + 0/\log_2(6) = 2.262$, IDCG@5 = $2/\log_2(2) + 1/\log_2(3) + 0/\log_2(4) + 0/\log_2(5) + 0/\log_2(6) = 2.631$

P@1: 1.000

P@1 = 1

P@5: 0.400

P@5 = 2/5

R@5: 1.000

R@5 = 2/2

MRR: 1.000

MRR = 1/1 = 1.000

AP: 1.000

P@1 = 1.000, P@2 = 1.000

R-Precision: 1.000

RP = 2/2

Ranking 2

⋮ Item B	↑ ↓	Rel: 3
⋮ Item D	↑ ↓	Rel: 3
⋮ Item F	↑ ↓	Rel: 0
⋮ Item A	↑ ↓	Rel: 0
⋮ Item C	↑ ↓	Rel: 0
⋮ Item E	↑ ↓	Rel: 2

Evaluation Metrics

NDCG@5: 0.830

DCG@5 = $3/\log_2(2) + 3/\log_2(3) + 0/\log_2(4) + 0/\log_2(5) + 0/\log_2(6) = 4.893$, IDCG@5 = $3/\log_2(2) + 3/\log_2(3) + 2/\log_2(4) + 0/\log_2(5) + 0/\log_2(6) = 5.893$

P@1: 1.000

P@1 = 1

P@5: 0.400

P@5 = 2/5

R@5: 0.667

R@5 = 2/3

MRR: 1.000

MRR = 1/1 = 1.000

AP: 0.833

P@1 = 1.000, P@2 = 1.000, P@6 = 0.500

R-Precision: 0.667

RP = 2/3

: AP فرمول

$$AP = \frac{1}{R} \sum_{k=1}^n P(k).rel(k)$$

به recall حساس است چون تعداد کل استناد مرتبط در محاسبه میانگین دقت وارد می شود و اگر تعداد استناد مرتبط بیشتری در کل موجود باشد، هر سند مرتبط رتبه پایین تر تأثیر کمتری روی ap خواهد داشت.