

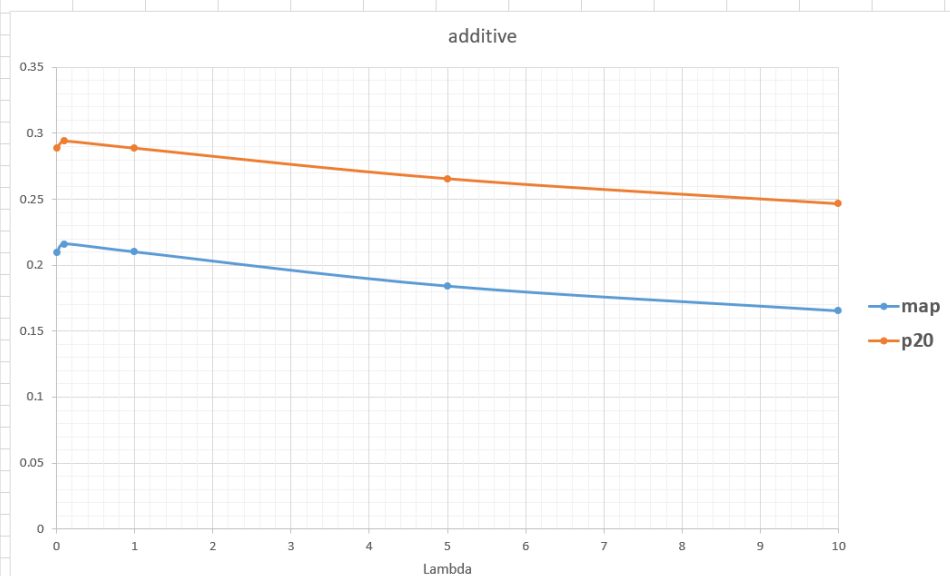
بهترین مقدار برای حدود 0.275 یا 0.35 است.

در این روش به کمک پارامتر لاندا ترکیبی از تاثیر دو عامل maximum likelihood و مدل مرجع را داریم. در این مجموعه لاندا کمی کمتر از 0.5 خوب جواب داده که انگار کلمات زیادی هستند که در سند نیستند و باید ارزش گذاری شوند و اتفاقاً می توانند مهم نیز باشند! مشاهده می شود که با لاندا 1 (حذف تاثیر مدل مرجع) نتایج بدتر از زمانی که حذف تقریبی maximum likelihood را داریم شده است.



:additive smoothing

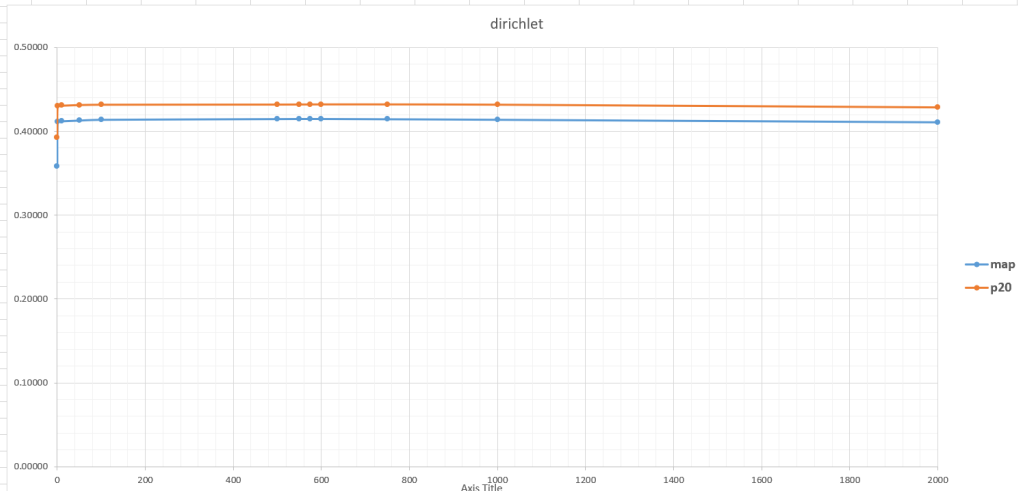
additive	map	p20
0.01	0.20942	0.28841
0.1	0.21598	0.29412
1	0.20989	0.28876
5	0.18425	0.26554
10	0.16568	0.24652



در این مدل با 0.1 بهترین نتیجه را تقریباً داریم. در این روش به تعداد کلمات مقداری افزوده می شود. گویا بیشتر کلمات به اندازه کافی در متن هستند و اضافه کردن حتی یک کلمه باعث می شود نسبت تعداد کلمات به طول سند بیش از حد زیاد شود که مطلوب نیست.

:dirichlet

dirichlet	map	p20
0	0.35832	0.39252
1	0.41165	0.43055
10	0.41192	0.43070
50	0.41281	0.43146
100	0.41370	0.43186
500	0.41456	0.43206
550	0.41470	0.43216
575	0.41473	0.43221
600	0.41469	0.43226
750	0.41436	0.43226
1000	0.41376	0.43191
2000	0.41075	0.42870

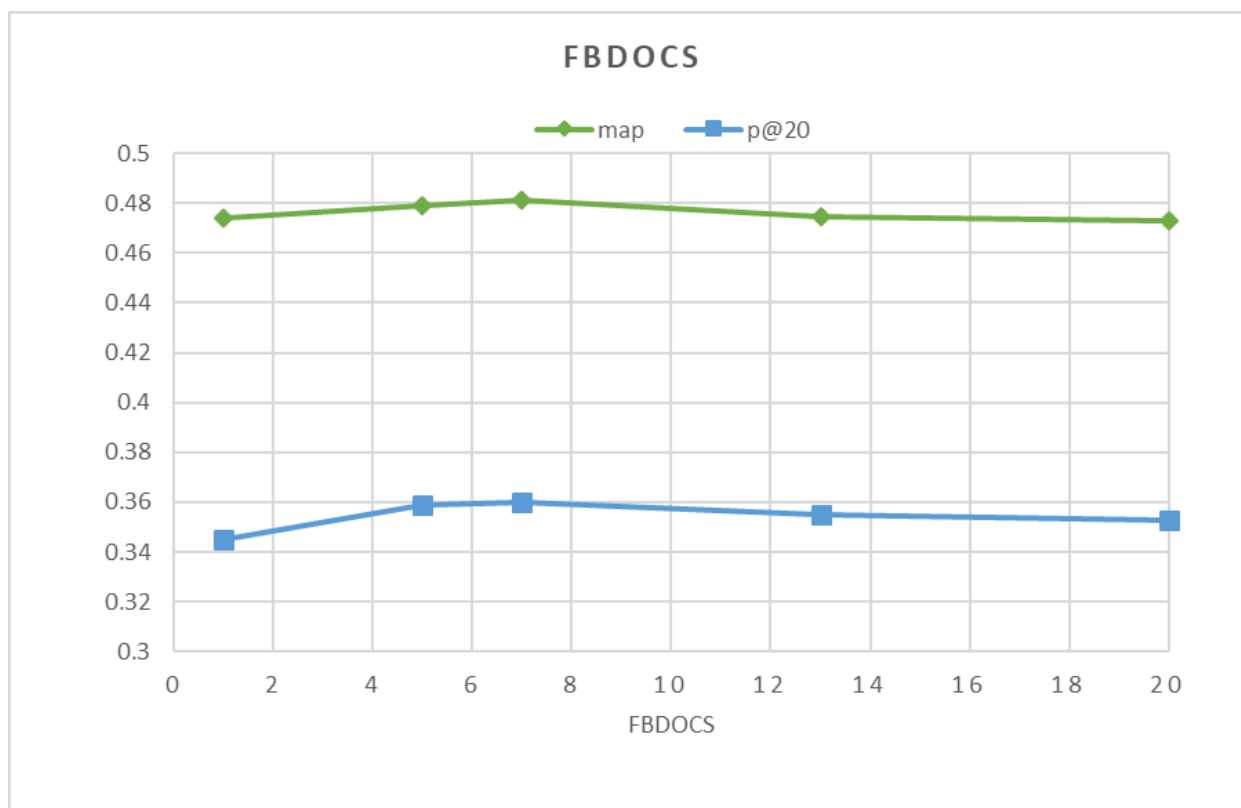


همانطور که دیده می شود در تغییر از 0 به 1 برای مو پیشرفت داریم اما در مقادیر بیشتر پیش چشمی دیده نمی شود. این روش smoothing زمانی که تعداد کلمات دیده نشده یا کم دیده شده در سند زیاد باشد به درد میخورد و هرچه طول سند کمتر باشد تاثیر آن بیشتر است. در مورد ما احتمالاً تعداد خوبی از کلمات پرس و جو در سند ها دیده می شوند و یا با تاثیر مقدار کمی مو نیز قابل تمییز از یکدیگر می شوند و smoothing از یه جا به بعد خیلی تاثیر نمی گذارد.

best = 575 - 600

(سوال 2)

fbDocs (الف)

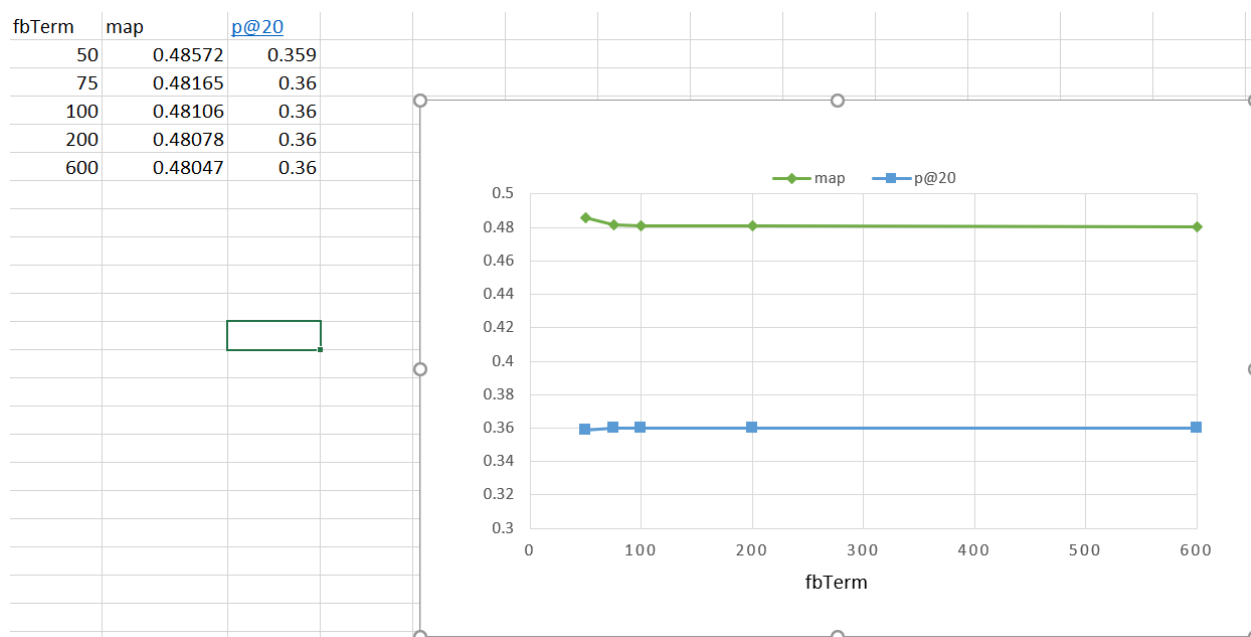


برای $fbDocs = 7$ مقادیر بهینه اند.



اگر تعداد اسناد بازخورد خیلی کم باشد عملاً از فیدبک در آپدیت مدل کوئری استفاده خاصی ن
شاید تاثیر مدل مرجع بیشتر دیده می شود. اگر هم خیلی زیاد باشد ممکن است تعداد خوبی از کلمات موجود
در فیدبک ها که خیلی هم اهمیت ندارند وزن زیادی بگیرند و در مدل $teta f$ وزن بسیاری بگیرند.

fbTerms (ب)



برای 50 در موارد آزمایشی ما نتیجه کمی بهتر است. نباید خیلی به همه کلمات اسناد بازخورد اکتفا کرد.
لزوما مدل کوئری را بهتر نمی کنند. اینکه نتیجه از مقداری ب بعد تغییر نکرده شاید ب این دلیل است که
تعداد کلمات اسناد بازخورد تمام شده و هرچه بیشتر کنیم تاثیری ندارد. (در تست fbTerm از مقدار بهینه
قسمت قبل یعنی $fbDoc = 7$ استفاده کردیم)



- B

(الف

$$q = \text{"بازیابی اطلاعات اسناد"} \quad |d_1| = 8 \quad |d_2| = 8$$

در نگاه اول ، به نظر می رسد سندی که کلیات کلیدی "بازیابی" ، "اطلاعات" و "اسناد" و شاید حتی "هوشمند" را بیشتر دارد ، برای کاربر مورد انتظار و مورد پسند است :

کلمه "بازیابی" در سند دوم 2 بار بیشتر از سند اول ، کلمه "هوشمند" در سند دوم یک بار تکرار شده اما در سند اول نیامده ، کلمه "اطلاعات" در سند دوم 2 بار و در سند اول 3 بار آمده ، کلمه "اسناد" در هر دو سند به تعداد برابر آمده .

انتظار می رود سند دوم به دلیل تاکید بیشتر بر روی واژه "بازیابی" و تعداد تکرار تقریباً یکسان در دیگر کلمات ، نسبت به سند اول بهتر باشد .

$$1) d_2$$

$$2) d_1$$

(ب

$$p(w|d) = \frac{c(w,d)}{|d|}$$

$$p(w|d_1) = p(w = \text{"بازیابی"} | d_1) \times p(w = \text{"اطلاعات"} | d_1) \times p(w = \text{"اسناد"} | d_1) =$$

$$\frac{c(\text{بازیابی}, d_1)}{|d_1|} \times \frac{c(\text{اطلاعات}, d_1)}{|d_1|} \times \frac{c(\text{اسناد}, d_1)}{|d_1|} =$$

$$\frac{2}{9} \times \frac{3}{9} \times \frac{1}{9} \cong 0.0082304$$

$$p(w|d_2) = p(w = \text{"بازیابی"} | d_2) \times p(w = \text{"اطلاعات"} | d_2) \times p(w = \text{"اسناد"} | d_2) =$$

$$\frac{c(\text{بازیابی}, d_2)}{|d_2|} \times \frac{c(\text{اطلاعات}, d_2)}{|d_2|} \times \frac{c(\text{اسناد}, d_2)}{|d_2|} =$$



$$\frac{4}{8} \times \frac{2}{8} \times \frac{1}{8} \cong 0.0156 \rightarrow d_2 \gg d_1$$

(ج)

$$d = 0.5 \quad p(w|d) = \frac{c(w,d)+0.5}{|d|+0.5 \times |v|} \quad v: \text{vocab size}$$

$$p(w|d_1) = p(w = \text{"بازیابی"} | d_1) \times p(w = \text{"اطلاعات"} | d_1) \times p(w = \text{"اسناد"} | d_1) =$$

$$\frac{2+0.5}{9+0.5 \times 10} \times \frac{3+0.5}{9+0.5 \times 10} \times \frac{1+0.5}{9+0.5 \times 10} = \frac{2.5}{14} \times \frac{3.5}{14} \times \frac{1.5}{14} \cong 0.0047$$

$$p(w|d_2) = \frac{4+0.5}{8+0.5 \times 10} \times \frac{2+0.5}{8+0.5 \times 10} \times \frac{1+0.5}{8+0.5 \times 10} = \frac{4.5}{13} \times \frac{2.5}{13} \times \frac{1.5}{13}$$

$$\cong 0.0076 \rightarrow d_2 \gg d_1$$

(د)

$$p(w|d) = \frac{(c(w,d) - s, 0) + d |d|_u p(w)c}{|d|}$$

$|d|_u$: number of unique words in d

$$d = 0.5$$

$$p(w|d_1) = p(w = \text{"بازیابی"} | d_1) \times p(w = \text{"اطلاعات"} | d_1) \times p(w = \text{"اسناد"} | d_1) =$$

$$\frac{1.5+0.5 \times 5 \times 0.6}{9} \times \frac{2.5+0.5 \times 5 \times 0.1}{9} \times \frac{0.5+0.5 \times 5 \times 0.1}{9} =$$

$$\frac{3}{9} \times \frac{2.75}{9} \times \frac{0.75}{9} = \frac{6.1875}{729} \cong 0.0084$$



$$p(w|d_2) = p(w = \text{"بازیابی"} | d_2) \times p(w = \text{"اطلاعات"} | d_2) \times p(w = \text{"اسناد"} | d_2) =$$

$$\frac{3.5+0.5 \times 4 \times 0.6}{8} \times \frac{1.5+0.5 \times 4 \times 0.1}{8} \times \frac{0.5+0.5 \times 4 \times 0.1}{8}$$

$$\frac{4.7}{8} \times \frac{1.7}{8} \times \frac{0.7}{9} = \frac{5.593}{512} \cong 0.0109 \rightarrow d_2 \gg d_1$$

ه (بله ممکن است. ممکن است کلمه ای را داشته باشیم که اصلاً در سند نیست و $c(w, d)$ آن برابر 0 است اما $p(w, c)$ آن بزرگ است و عام است.

در این حالت

$$c(w, d) - S + s|d|_u p(w, c)$$

کلمه دیگری که در سند وجود دارد بیشتر شود.

$$w_1 = \text{زیبا} \quad c(w_1, d) = 0 \quad p(w_1|c) = 0.9$$

$$w_2 = \text{علی} \quad c(w_2, d) = 1 \quad p(w_2|c) = 0.1 \quad s(\text{lambda}) = 0.5 \quad |d| = 10 \quad |d|_u = 5$$

$$p(w_1|d) = \frac{0.5 \times 5 \times 9}{10} > \frac{0.5 + 0.5 \times 5 \times 0.1}{10} = p(w_2|d)$$

-A

۱_ غلط است.

در این مدل به ازای یک μ ثابت، هرچه طول سند بیشتر باشد smoothing کمتر انجام می شود.

اما μ برآوردی از میانگین طول سند نیست.

۲_ درست است.

این الگوریتم در هر مرحله likelihood را افزایش می دهد تا به یک max local optima برسد. پس برای نقاط مختلف از یک مقدار شروع کرده و به یک max محلی می رساند و minimum optimisation انجام نمی دهد.



۳_ درست است.

Smoothing دو هدف دارد. یکی نسبت دادن وزن برای کلماتی که در سند نیامده اند. و دیگری نزدیک کردن وزن کلمات عام مانند stopword ها در دو داکيومنت.

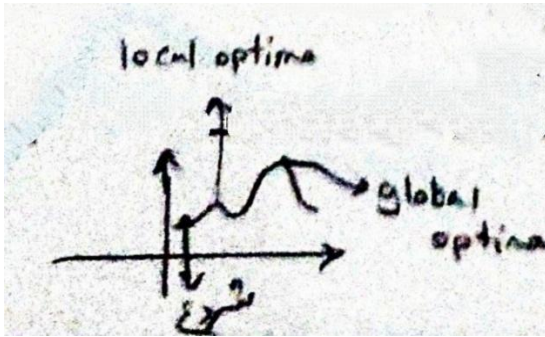
۴_ غلط است.

این الگوریتم در هر مرحله مقدار هدف خود را بیشتر می کند و likelihood را در هر مرحله افزایش می دهد تا اهداف به local optima خود برسد.

۵_ غلط است.

این الگوریتم برای وزن هر کلمه به طور اتفاقی از یک جا شروع می شود آنقدر iteration ها را ادامه می دهد تا به یک local optima برسد و پس از آن دیگر ادامه نمی دهد.

ممکن است آن max محلی، بیشترین مقدار ممکن نباشد!



سوال 4)

طبق توضیحی که در باره maxRank داریم .

$\text{maxRank}(n) - \text{maxRank}(n+1)$ = تعداد کلماتی که n بار تکرار میشوند

طبق قانون زیپ داریم :

$$r * \text{freq}() = A * N$$

maxRank کلمات با n بار تکرار زیرا در فرمول فوق جایگزاری می کنیم :



n

$$\text{Rank}(n) \times n = A * N \rightarrow \max \text{Rank} =$$

$\max \text{Rank}(n) - \max \text{Rank}(n+1) =$ بار تکرار میشوند n تعداد کلماتی که

$$\frac{A*N}{n} - \frac{A*N}{n+1}$$

$$= A * N \left(\frac{1}{n} - \frac{1}{n+1} \right) = A * N \left(\frac{1}{n(n+1)} \right) = \frac{A*N}{n(n+1)}$$

حال باید $A * N$ را حساب کنیم.

فرض میکنیم کمترین تعداد تکرار یک کلمه 1 بار می باشد (کلمات unique)

کلمه با کمترین تکرار (1 بار) در رتبه آخر رتبه بندی کلمات می باشد.

اگر فرض کنیم اندازه کلمات مجزا 7 باشد.

طبق قانون زیپ داریم :

$$r * \text{freq}(r) = A * N \rightarrow V \times 1 = A * N \}$$

فرمول تعداد کلمات با n بار تکرار صفحه قبل را با مقدار یافت شده برای

$A * N$ باز نویسی می کنیم :

$$\text{تعداد کلمات با } n \text{ بار تکرار} = \frac{A*N}{n(n+1)} - \frac{V}{n(n+1)}$$

الف) چه نسبتی از عبارات فقط یک بار ظاهر می شود؟

$$\frac{\text{تعداد کلمات، با 1 بار تکرار}}{\text{کلمات تعداد مجزا}} = \frac{V \times \left(\frac{1}{1(1+1)} \right)}{V} = \frac{1}{2}$$

ب) چه نسبتی 10 بار ظاهر می شود ؟



$$\frac{\text{تعداد کلمات، با}}{\text{کلمات، مجزا تعداد}} = \frac{V \left(\frac{1}{10(11)} \right)}{V} = \frac{1}{110}$$

(ج)

$$\frac{\text{تعداد کلمات، با } k \text{ بار تکرار}}{\text{تعداد کلمات، مجزا}} = \frac{7 \left(\frac{1}{k(k+1)} \right)}{7} = \frac{1}{k(k+1)}$$

