

به نام خدا



دانشگاه تهران
دانشکدگان فنی
دانشکده مهندسی برق
و کامپیوتر



درس بازیابی هوشمند اطلاعات

پاسخ بخش تئوری تمرین ۲

نام و نام خانوادگی: روزین پنا

شماره دانشجویی: ۲۲۰۷۰۱۰۴۶

مهر ماه ۱۴۰۳

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.

فهرست

۳	پاسخ سوال اول
۳	پاسخ بخش الف
۵	پاسخ بخش ب
۷	پاسخ بخش ج
۱۱	پاسخ سوال دوم
۱۱	پاسخ بخش الف
۱۵	پاسخ بخش ب
۱۶	پاسخ سوال سوم
۱۶	پاسخ بخش الف
۱۷	پاسخ بخش ب
۱۸	پاسخ بخش ج
۱۹	پاسخ بخش د

پاسخ سوال اول

پاسخ بخش الف

"الگوریتم های پیشرفته یا دگیری ماشین و کاربردهای آن در سیستم های توصیه‌گر و بازیابی اطلاعات" = D_3

: Jelinek-Mercer smoothing . ۱

$$p(w_i | d) = (1 - \lambda) \cdot \frac{c(w_i, d)}{|d|} + \lambda \cdot p(w_i | C)$$

$$p("يادگیری" | D_3) = (1 - \lambda) \cdot \frac{c("يادگیری", D_3)}{|D_3|} + \lambda \cdot p("يادگیری" | C)$$

$$p("يادگیری" | D_3) = 0.4 \cdot \frac{1}{13} + 0.6 \cdot 0.003$$

$$p("يادگیری" | D_3) = 0.031 + 0.0018$$

$$p("يادگیری" | D_3) \approx 0.0328$$

: Dirichlet smoothing . ۲

$$p(w_i | d) = \frac{c(w_i, d) + \mu \cdot p(w_i | C)}{|d| + \mu}$$

$$p("يادگیری" | D_3) = \frac{c("يادگیری", D_3) + \mu \cdot p("يادگیری" | C)}{|D_3| + \mu}$$

$$p("يادگیری" | D_3) = \frac{1 + 2000 \cdot 0.003}{13 + 2000}$$

$$p("يادگیری" | D_3) = \frac{1 + 6}{2013}$$

$$p("يادگیری" | D_3) \approx 0.0035$$

: Absolute Discounting . ۳

$$p(w_i | d) = \frac{\max(c(w_i, d) - \delta, 0) + \delta \cdot |d|_u \cdot p(w_i | C)}{|d|}$$

$$p("يادگیری" | D_3) = \frac{\max(c("يادگیری", D_3) - \delta, 0) + \delta \cdot |D_3|_u \cdot p("يادگیری" | C)}{|D_3|}$$

در "الگوریتم های پیشرفته یادگیری ماشین و کاربردهای آن در سیستم های توصیه‌گر و بازیابی اطلاعات" ← تنها واژه "و" تکرار شده:

$$|D_3|_u = 12 \quad \text{تعداد واژه های کتا}$$

$$p("يادگیری" | D_3) = \frac{\max(1 - 0.5, 0) + 0.5 \cdot 12 \cdot 0.003}{13}$$

$$p("يادگیری" | D_3) = \frac{0.5 + 0.5 \cdot 12 \cdot 0.003}{13}$$

$$p("يادگیری" | D_3) = \frac{0.5 + 0.018}{13}$$

$$p("يادگیری" | D_3) \approx 0.45$$

: Additive (Laplace) . ۴

$$p(w_i | d) = \frac{c(w_i, d) + \alpha}{|d| + \alpha \cdot |V|}$$

$$p("يادگیری" | D_3) = \frac{c("يادگیری", D_3) + \alpha}{|D_3| + \alpha \cdot |V|}$$

$$p("يادگیری" | D_3) = \frac{1 + 1}{13 + 1.10000}$$

$$p("يادگیری" | D_3) = \frac{2}{10013}$$

$$p("يادگیری" | D_3) \approx 0.0002$$

JM	DIR	AD	Additive	
0.328	0.0035	0.45	0.0002	$p("يادگیری" D_3)$

روش Absolute Discounting بیشترین احتمال را دارد، چون هم تعداد وقوع کلمه «یادگیری» را نگه می‌دارد ($c(w_i, d) - \delta$) پس اگر کلمه در سند آمده باشد، سهم مستقیمش از بین نمی‌رود، هم یک مقدار اضافه از مدل مجموعه به آن می‌دهد این به خاطر تعداد واژه‌های یکتای سند بزرگتر می‌شود، پس عدد کلی AD بزرگتر می‌شود.

در JM وزن زیادی به داده های کلی داده شده ($\lambda = 0.6$) و سهم سند کم می شود، برای همین مقدارش کمتر است.

Dirichlet بیشتر به عده های مدل مجموعه تکیه می کند، وقتی m خیلی بزرگ است، یکبار وقوع کلمه "یادگیری" در سند در برابر m تقریباً ناچیز می شود و نتیجه نزدیک $p(w | C)$ می ماند؛ بنابراین مقدار نهایی پایین است.

Laplace هم به همه کلمه ها یک مقدار ثابت اضافه می کند، اما مشکلش این است که **مخرج خیلی بزرگ** می شود $|d| + \alpha |V|$. با واژگان بزرگم ($|V|=10000$)، حتی اگر شمارش کلمه "یادگیری" یک باشد، نسبت به مخرج عدد خیلی کوچک می شود و برای بازیابی اطلاعات مناسب نیست؛ چون واژه های دیده شده و ندیده را تقریباً یکسان رقیق می کند.

پاسخ بخش ب

$q = "یادگیری بازیابی"$

$D_1 = "بازیابی اطلاعات و یادگیری ماشین در علوم داده"$

"الگوریتم های پیشرفته یادگیری ماشین و کاربردهای آن در سیستم های توصیه گر و بازیابی اطلاعات" $= D_3$

$$\log p(q | d) = \sum_{i=1}^m \log p(w_i | d) = \sum_{w \in V, c(w|d) > 0} c(w, q) \log p(w | d)$$

باید به جای $p(w | d)$ از نسخه هموارشده آن استفاده کنیم با روش Dirichlet Prior

$$p(w | d) = \frac{c(w, d) + \mu \cdot p(w | C)}{|d| + \mu}$$

$$p(w | d) = \frac{|d|}{|d| + \mu} \cdot \frac{c(w, d)}{|d|} + \frac{\mu}{|d| + \mu} \cdot p(w | C)$$

محاسبه $: p(q | D_1)$

$$p("بازیابی" | D_1) = \frac{c("بازیابی", D_1) + \mu \cdot p("بازیابی" | C)}{|D_1| + \mu}$$

$$p("بازیابی" | D_1) = \frac{1 + 2000 \cdot 0.002}{8 + 2000} = 0.0024900$$

$$p("یادگیری" | D_1) = \frac{c("یادگیری", D_1) + \mu \cdot p("یادگیری" | C)}{|D_1| + \mu}$$

$$p("یادگیری" | D_1) = \frac{1 + 2000 \cdot 0.003}{8 + 2000} = 0.0034860$$

محاسبه $p(q | D_3)$:

$$(p("بازیابی" | D_3) = \frac{c("بازیابی", D_3) + \mu \cdot p("بازیابی" | C)}{|D_3| + \mu})$$

$$p("بازیابی" | D_3) = \frac{1 + 2000 \cdot 0.002}{13 + 2000} = 0.0024839$$

$$(p("یادگیری" | D_3) = \frac{c("یادگیری", D_3) + \mu \cdot p("یادگیری" | C)}{|D_3| + \mu})$$

$$p("یادگیری" | D_3) = \frac{1 + 2000 \cdot 0.003}{13 + 2000} = 0.0034774$$

جایگذاری:

$$\log p(q | D_1) = \sum_{w \in V, c(w|d) > 0} c(w, q) \log p(w | D_1)$$

$$\log p(q | D_1) = c("یادگیری", q) \log p("بازیابی" | D_1) + c("بازیابی", q) \log p("یادگیری" | D_1)$$

$$\log p(q | D_1) = 1 \cdot \log p("بازیابی" | D_1) + 1 \cdot \log p("یادگیری" | D_1)$$

$$\log p(q | D_1) = \log 0.0024900 + \log 0.0034860$$

$$p(q | D_1) = 0.0024900 \times 0.0034860 = 0.0000086801$$

$$\log p(q | D_3) = c("یادگیری", q) \log p("بازیابی" | D_3) + c("بازیابی", q) \log p("یادگیری" | D_3)$$

$$\log p(q | D_3) = 1 \cdot \log p("بازیابی" | D_3) + 1 \cdot \log p("یادگیری" | D_3)$$

$$\log p(q | D_3) = \log 0.0024839 + \log 0.0034774$$

$$p(q | D_3) = 0.0024839 + 0.0034774 = 0.0000086375$$

$$p(q | D_1) > p(q | D_3)$$

چرا $p(q | D_1) > p(q | D_3)$

در Dirichlet پارامتر ($\mu = 2000$) خیلی بزرگ است، بنابراین سهم «شمارش واقعی در سند» نسبت به «احتمال زمینه (مدل زبانی مجموعه)» کوچک می‌شود.

$$p(w | d) = \frac{1}{|d| + \mu} \cdot \underbrace{\frac{|d| \cdot c(w, d)}{|d|}}_{\substack{\text{یکسان در } D_1 \text{ و } D_3 \\ \text{یکسان در } D_1 \text{ و } D_3}} + \frac{\mu}{|d| + \mu} \cdot \underbrace{\frac{p(w | C)}{D_1 \cup D_3}}_{\substack{\text{یکسان در } D_1 \text{ و } D_3}}$$

وقتی طول سند کوتاه‌تر باشد ($|D_1|=8$)، اثر مدل مجموعه و تعداد وقوع واژه در سند نسبت به مخرج $d | + \mu$ کمی بزرگ‌تر است؛ برای D_3 که طولش بیشتر است ($|D_3|=13$) مدل مجموعه و تعداد وقوع واژه کم اثرتر می‌شود.

$$0.996016 = \frac{2000}{2008} = \frac{\mu}{|D_1| + \mu} > \frac{\mu}{|D_3| + \mu} = \frac{2000}{2013} = 0.993541$$

$$0.0004980 = \frac{1}{2008} = \frac{1}{|D_1| + \mu} > \frac{1}{|D_3| + \mu} = \frac{1}{2013} = 0.0004968$$

یعنی تفاوت تنها در مخرج است. مخرج بزرگتر در D_3 باعث کاهش اندکی در $p(w|c)$ و $c(w,d)$ نسبت به D_1 شده.

پاسخ بخش ج

مدل ترکیبی:

$$P(w|d) = \alpha P_{JM}(w|d) + \beta P_{Dir}(w|d) + \gamma P_{AD}(w|d) + \theta P_{Add}(w|d)$$

به طوری که

$$\alpha + \beta + \gamma + \theta = 1 \quad \alpha, \beta, \gamma, \theta \geq 0$$

۱. چه زمانی هر یک از مولفه‌ها نقش مهم‌تری دارن؟

:Dirichlet •

$$P_{Dir} = \frac{|d|}{|d| + \mu} \cdot \frac{c}{|d|} + \frac{\mu}{|d| + \mu} \cdot P_c$$

$$P_{Dir} \rightarrow \frac{c}{|d|} \quad \text{در نتیجه } \frac{|d|}{|d| + \mu} \rightarrow 1 \quad \circ \quad \text{اگر } |d| \rightarrow \infty \text{ آنگاه:}$$

$$P_{Dir} \approx P_c \quad \frac{|d|}{|d| + \mu} \approx 0 \quad \circ \quad \text{اگر } |d| \ll \mu \text{ آنگاه:}$$

هر چه $|d|$ نسبت به μ بزرگتر شود، Dirichlet بیشتر به مشاهدات و مدل سند تکیه می‌کند؛ در غیر این صورت به احتمال زمینه و مدل مجموعه تکیه می‌کند.

:Jelinek-Mercer •

$$P_{JM} = (1 - \lambda) \frac{c}{|d|} + \lambda P_c$$

$$P_{JM} \rightarrow \frac{c}{|d|} \quad \circ \quad \text{اگر } \lambda \rightarrow 0 \text{ آنگاه:}$$

$$P_{JM} \rightarrow P_c \quad \circ \quad \text{اگر } \lambda \rightarrow 1 \text{ آنگاه:}$$

در سند‌های طولانی که $\frac{c}{|d|}$ قابل اطمینان است، انتخاب λ کوچک باعث می‌شود P_{JM} نقش مهم‌تری داشته باشد؛ تأثیر $|d|$ غیرمستقیم و از طریق نسبت $\frac{c}{|d|}$ است، ولی کنترل اصلی را پارامتر λ می‌کند.

: Absolute Discounting •

$$P_{AD} = \frac{\max(c - \delta, 0) + \delta |d|_u P_C}{|d|}$$

$$P_{AD} = \frac{c - \delta}{|d|} + \frac{\delta |d|_u}{|d|} P_C \quad \text{اگر } c \geq \delta \quad \circ$$

دو ترم داریم:

$$P_{Seen}(w|d) = \frac{c - \delta}{|d|} \quad \text{: discounted ML estimate . I}$$

$$\alpha_d P(w|C) = \frac{\delta |d|_u}{|d|} P_C \quad \text{: Collection Model . II}$$

$$\frac{\partial P_{AD}}{\partial |d|_u} = \frac{\delta}{|d|} P_C \quad \text{اگر } |d|_u \text{ مشتق نسبت به } |d| = \text{constant} \quad \checkmark$$

$$\frac{\partial P_{AD}}{\partial |d|_u} > 0 \quad \text{چون } \delta > 0, |d| > 0, P_C > 0 \quad \text{داریم}$$

(با افزایش $|d|_u$ و ثابت ماندن $|d|$) آنگاه: $P_{AD} \leftarrow$ افزایش

P_{AD} به صورت خطی و صعودی افزایش می‌یابد و AD به نفع و اژه‌هایی که در سند کم‌تکرارند شانس بیشتری می‌دهد.

$$\frac{\delta |d|_u}{|d|} P_C = \delta \left(\frac{|d|_u}{|d|} \right) P_C \quad \text{اگر اگر } |d| \neq \text{constant} \quad \text{افزایش:} \\ \text{با افزایش } |d|_u \text{ یا کاهش } |d| \text{ آنگاه: } P_{AD} \leftarrow$$

$$P_{AD} = \frac{\delta |d|_u}{|d|} P_C \quad \text{اگر } c < \delta \quad \text{آنگاه:} \\ \text{باز هم استدلال بالا میتوانیم داشته باشیم.}$$

$$P_{AD} \rightarrow \frac{c}{|d|} \quad \text{اگر } \delta \rightarrow 0 \quad \text{آنگاه:}$$

$$P_{AD} \rightarrow \frac{\delta |d|_u}{|d|} P_C \quad \text{اگر } \delta \rightarrow 1 \quad \text{آنگاه:}$$

P_{AD} زمانی مهم‌تر است که سند تعداد یکتا‌های نسبتاً زیادی داشته باشد یا وقتی می‌خواهیم جرم احتمال discounted ML estimate را به شکل متناسب $|d|_u$ با توزیع کنیم.

:Additive (Laplace) •

$$P_{Add} = \frac{c + \alpha}{|d| + \alpha |V|}$$

$P_{Add} \rightarrow \frac{1}{ V }$	اگر $ V \rightarrow \infty$ آنگاه :
$\lim_{ V \rightarrow 0} P_{Add}(w d) = \frac{c+\alpha}{ d }$	اگر $ V \rightarrow 0$ آنگاه :
$\lim_{\alpha \rightarrow \infty} P_{Add}(w d) = \lim_{\alpha \rightarrow \infty} \frac{\alpha(1+c/\alpha)}{\alpha(V + d /\alpha)} = \frac{1}{ V }$	اگر $ V > 0$ و $\alpha \rightarrow \infty$ آنگاه :
$\lim_{\alpha \rightarrow \infty} \frac{c+\alpha}{ d +\alpha V } = \lim_{\alpha \rightarrow \infty} \frac{\alpha(1+o(1))}{\alpha(V +o(1))} = \frac{1}{\lim_{\alpha \rightarrow \infty} V }$	اگر $\alpha V \rightarrow \infty$ و $\alpha \rightarrow \infty$ آنگاه :
$\lim_{\alpha V \rightarrow 0} P_{Add}(w d) = \frac{c+\alpha}{ d }$	اگر $\alpha V \rightarrow 0$ آنگاه :

پس اگر فضای واژگان کوچک باشد، Additive میتواند خوب باشد؛ ولی در \mathbb{IR} عملی با $|V|$ بزرگ این مؤلفه معمولاً رقیق‌کننده و کم‌تأثیر است.

۲. چیکار کنیم $\alpha, \beta, \gamma, \theta$ خودکار تعریف بشن؟

باید وزن‌های ترکیبی چهار مؤلفه مدل زبان را طوری یاد بگیریم که ترکیب آن‌ها بیشترین احتمال را روی مجموعه نگهداری (validation) بدهد. وزن‌ها باید مثبت باشند و جمع‌شان برابر ۱ شود؛ برای این کار پارامترهای آزاد z را یاد می‌گیریم و آن‌ها را با softmax به $\alpha, \beta, \gamma, \theta$ نگاشت می‌زنیم. در ادامه روش را با مسیرهای forward و backward حل می‌کنیم.

تعریف مدل و نگاشت وزن‌ها

• مدل ترکیبی :

$$P(w | d) = \alpha P_{JM}(w | d) + \beta P_{Dir}(w | d) + \gamma P_{AD}(w | d) + \theta P_{Add}(w | d)$$

• نگاشت وزن‌ها (پارامترهای آزاد z_1, z_2, z_3, z_4) :

$$\omega_i = \frac{e^{z_i}}{\sum_{j=1}^4 e^{z_j}} \quad \text{که } \omega \in \{\alpha, \beta, \gamma, \theta\}.$$

مسیر - Forward

برای هر جفت (q, d) در مجموعه اعتبار و برای هر واژه w در q ، چهار مؤلفه $P_{JM}, P_{Dir}, P_{AD}, P_{Add}$ را محاسبه می‌کنیم. (پارامترهای داخلی هر مؤلفه مثل $\lambda, \mu, \delta, \alpha$ قبلاً تعیین شده‌اند).

۱. با z فعلی وزن‌ها را با softmax می‌سازم و برای هر w مقدار ترکیب $P(w | d)$ را محاسبه می‌کنم:

$$P(w | d; z) = \sum_{i=1}^4 \omega_i(z) P_i(w | d)$$

۹. امتیاز پرسوچو در سند را با جمع لگاریتم احتمال کلمات حساب میکنیم:

$$\log P(q | d) = \sum_{w \in q} c(w, q) \log P(w | d; z)$$

۱۰. تابع هدف کل روی مجموعه اعتبار برابر است با جمع این امتیازها:

$$\mathcal{L}(z) = \sum_{(q,d) \in D_{val}} \sum_{w \in q} (w, q) \log P(w | d; z)$$

۱۱. برای پایداری یک جمله منظم‌ساز اضافه می‌کنم:

$$\mathcal{L}_{reg}(z) = \mathcal{L}(z) - \frac{\lambda}{2} \| z \|^2$$

عدد $\mathcal{L}_{reg}(z)$ نشان می‌دهد وزن‌های فعلی چقدر خوبند.

مسیر - Backward

اگر برای یک واژه-سند، $P_j(w | d)$ از میانگین ترکیب $P(w | d)$ بهتر باشد، باید وزن مؤلفه z را بیشتر کنیم؛ اگر بدتر باشد، باید کمتر کنیم. نسبت مؤلفه به ترکیب این شهود را مشخص می‌کند:

$$R_j(w, d) = \frac{P_j(w | d)}{P(w | d)}.$$

۱۲. مشتق لگاریتم برای هر z_j (نتیجه عملی محاسبات گرادیان) چنین است:

$$\frac{\partial \mathcal{L}}{\partial z_j} = \sum_{(q,d)} \sum_{w \in q} c(w, q) \omega_j \left(\frac{P_j(w | d)}{P(w | d)} - 1 \right)$$

۱۳. اضافه کردن منظم‌سازی به گرادیان:

$$\nabla_{z_j} \mathcal{L}_{reg} = \frac{\partial \mathcal{L}}{\partial z_j} - \lambda z_j$$

۱۴. به روزرسانی پaramترهای آزاد z با یک گام گرادیان یا با Adam/L-BFGS :

$$z_j \leftarrow z_j + \eta \nabla_{z_j} \mathcal{L}_{reg}$$

توضیح ساده: اگر عبارت داخل جمع مثبت باشد، وزن مؤلفه z را بیشتر میکنیم؛ اگر منفی باشد، کمتر میکنیم.

حلقه یادگیری کامل:

$$z = (0,0,0,0) \rightarrow \omega_i = 0.25$$

۱. مقداردهی اولیه:

۲. تکرار تا همگرایی:

$\mathcal{L}_{reg}(z)$: محاسبه $P(w | d; z)$ و $P_i(w | d)$ ، سپس محاسبه $\mathcal{L}_{reg}(z)$: Forward.

محاسبه $\nabla_z \mathcal{L}_{reg}$: Backward.

به روزرسانی z با نرخ یادگیری η (یا استفاده از Update برای خودتنظیم نرخ Adam).

۳. شرط توقف: تغییرات \mathcal{L}_{reg} یا z بسیار کوچک شود یا تعداد تکرارها به حد مشخص برسد.

۴. خروجی: وزن‌ها $\alpha, \beta, \gamma, \theta = \text{softmax}(z)$ که به صورت خودکار یاد گرفته میشوند.

پاسخ سوال دوم

پاسخ بخش الف

جدول اسناد و امتیاز اولیه:

امتیاز اولیه	محتوا	سند
12.5	الگوریتم‌های یادگیری عمیق و شبکه‌های عصبی	D ₁
11.8	یادگیری عمیق در بینایی ماشین	D ₂
8.3	الگوریتم‌های ژنتیک و بهینه‌سازی	D ₃
10.2	یادگیری تقویتی عمیق	D ₄
9.7	معماری‌های عمیق در پردازش زبان	D ₅

$$Z_i \in \{1(\text{background}), 0(\text{topic})\}$$

احتمال این که واژه از منبع عمومی باشد :

$$P^{(n)}(Z_i = 1 | w_i) = \frac{\lambda \cdot P(w_i | C)}{\lambda \cdot P(w_i | C) + (1 - \lambda) \cdot P^{(n)}(w_i | \theta_F)}$$

$$\lambda = 0.3 : \text{LAMBDA}$$

{ عميق، يا دگيري، الگوريتم } : WORD

$P^{(0)}(w | \theta_F)$: INITIAL FEEDBACK MODEL يکنواخت

$$P^{(0)}(w | \theta_F) = \frac{1}{3} \approx 0.333333$$

$$: n = 0$$

الگوريتم : W =

$$P^{(0)}(Z_i = 1 | \text{الگوريتم}) = \frac{0.3 \cdot P(\text{الگوريتم} | C)}{0.3 \cdot P(\text{الگوريتم} | C) + 0.7 \cdot P^{(0)}(\text{الگوريتم} | \theta_F)}$$

$$P^{(0)}(Z_i = 1 | \text{الگوريتم}) = \frac{0.3 \cdot 0.001}{0.3 \cdot 0.001 + 0.7 \cdot 0.334}$$

$$P^{(0)}(Z = 1 | \text{الگوريتم}) \approx \frac{0.0003}{0.2341} \approx 0.00128$$

يا دگيري : W =

$$P^{(0)}(Z_i = 1 | \text{يا دگيري}) = \frac{0.3 \cdot P(\text{يا دگيري} | C)}{0.3 \cdot P(\text{يا دگيري} | C) + 0.7 \cdot P^{(0)}(\text{يا دگيري} | \theta_F)}$$

$$P^{(0)}(Z_i = 1 | \text{يا دگيري}) = \frac{0.3 \cdot 0.002}{0.3 \cdot 0.002 + 0.7 \cdot 0.334}$$

$$P^{(0)}(Z = 1 | \text{يا دگيري}) \approx \frac{0.0006}{0.2344} \approx 0.00256$$

عميق : W =

$$P^{(0)}(Z_i = 1 | \text{عميق}) = \frac{0.3 \cdot P(\text{عميق} | C)}{0.3 \cdot P(\text{عميق} | C) + 0.7 \cdot P^{(0)}(\text{عميق} | \theta_F)}$$

$$P^{(0)}(Z_i = 1 | \text{عميق}) = \frac{0.3 \cdot 0.0015}{0.3 \cdot 0.002 + 0.7 \cdot 0.334}$$

$$P^{(0)}(Z = 1 | \text{عميق}) \approx \frac{0.00045}{0.23425} \approx 0.00192$$

: M-step

با استفاده از تخمین‌های مرحله‌ی قبل، مدل موضوعی رو به روزرسانی می‌کنیم:

$$P^{(n+1)}(w_i | \theta_F) = \frac{c(w_i, F) \cdot (1 - P^{(n)}(z_i = 1 | w_i))}{\sum_{w_j} c(w_j, F) \cdot (1 - P^{(n)}(z_j = 1 | w_j))}$$

: که

$$\begin{aligned} \sum_{w_j} c(w_j, F) \cdot (1 - P^{(n)}(z_j = 1 | w_j)) &= c(\text{الگوریتم}, F) \cdot (1 - 0.00128) + \\ &\quad c(\text{یادگیری}, F) \cdot (1 - 0.00256) + \\ &\quad c(\text{عمیق}, F) \cdot (1 - 0.00192) \end{aligned}$$

طبق جدول اسناد داریم:

$$\sum_{w_j} c(w_j, F) \cdot (1 - P^{(n)}(z_j = 1 | w_j)) = 2 \cdot (0.99872) + 3 \cdot (0.99744) + 4 \cdot (0.99808)$$

$$\sum_{w_j} c(w_j, F) \cdot (1 - P^{(n)}(z_j = 1 | w_j)) = 1.99744 + 2.99232 + 3.99232$$

$$\sum_{w_j} c(w_j, F) \cdot (1 - P^{(n)}(z_j = 1 | w_j)) = 8.98208$$

$$P^{(1)}(\text{الگوریتم} | \theta_F) = \frac{c(\text{الگوریتم}, F) \cdot (1 - P^{(0)}(z_i = 1 | w_i))}{\sum_{w_j} c(w_j, F) \cdot (1 - P^{(0)}(z_j = 1 | w_j))}$$

$$P^{(1)}(\text{الگوریتم} | \theta_F) = \frac{1.99744}{8.98208} \approx 0.22238$$

$$P^{(1)}(\text{یادگیری} | \theta_F) = \frac{c(\text{یادگیری}, F) \cdot (1 - P^{(0)}(z_i = 1 | w_i))}{\sum_{w_j} c(w_j, F) \cdot (1 - P^{(0)}(z_j = 1 | w_j))}$$

$$P^{(1)}(\text{یادگیری} | \theta_F) = \frac{2.99232}{8.98208} \approx 0.33314$$

$$P^{(1)}(\text{عمیق} | \theta_F) = \frac{c(\text{عمیق}, F) \cdot (1 - P^{(0)}(z_i = 1 | w_i))}{\sum_{w_j} c(w_j, F) \cdot (1 - P^{(0)}(z_j = 1 | w_j))}$$

$$P^{(1)}(\text{عمیق} | \theta_F) = \frac{3.99232}{8.98208} \approx 0.44448$$

برای مجموعه بازخورد F به صورت زیر تعریف می‌شود:

$$\log p(F^n | \theta) = \sum_i \sum_w c(w; d_i) \log [(1 - \lambda) P^{(n)}(w | \theta) + \lambda p(w | C)]$$

: $P(w) = \lambda P(w | C) + (1 - \lambda) P(w | \theta_F)$ **MIXTURE MODEL FOR FEEDBACK** طبق فرمول

$$\log p(F^n | \theta) = \sum_i \sum_w c(w; d_i) \log P^{(n)}(w)$$

$$P^{(1)}(\text{الگوریتم} | \theta_F) = \lambda P(C) + (1 - \lambda) P^{(1)}$$

$$= 0.3 \cdot 0.001 + 0.7 \cdot 0.22238 = 0.155966 \approx 0.16$$

$$P^{(1)}(\text{يادگیری} | \theta_F) = \lambda P(C) + (1 - \lambda) P^{(1)}$$

$$= 0.3 \cdot 0.002 + 0.7 \cdot 0.33314 = 0.233798 \approx 0.23$$

$$P^{(1)}(\text{عمیق} | \theta_F) = \lambda P(C) + (1 - \lambda) P^{(1)}$$

$$= 0.3 \cdot 0.0015 + 0.7 \cdot 0.44448 = 0.311586 \approx 0.31$$

$$\begin{aligned} \log p(F^1 | \theta) &= \sum_w c(w; d_1) \log P^{(1)}(w) + \sum_w c(w; d_2) \log P^{(1)}(w) + \sum_w c(w; d_3) \log P^{(1)}(w) \\ &\quad + \sum_w c(w; d_4) \log P^{(1)}(w) + \sum_w c(w; d_5) \log P^{(1)}(w) \end{aligned}$$

$$\begin{aligned} \log p(F^1 | \theta) &= \sum_i c(d_i) \log P^{(1)}(\text{الگوریتم}; d_i) \\ &\quad + \sum_i c(d_i) \log P^{(1)}(\text{يادگیری}; d_i) \\ &\quad + \sum_i c(d_i) \log P^{(1)}(\text{عمیق}; d_i) \end{aligned}$$

$$\begin{aligned} \log p(F^1 | \theta) &= c(F) \log P^{(1)}(\text{الگوریتم}, F) \\ &\quad + c(d_{\text{يادگیری}}, F) \log P^{(1)}(\text{يادگیری}, F) \\ &\quad + c(d_{\text{عمیق}}, F) \log P^{(1)}(\text{عمیق}, F) \end{aligned}$$

$$\sum_w c(w; d_1) \log P^{(1)}(w) = \text{الگوریتم}(d_1) \log P^{(1)}(\text{الگوریتم}(d_1)) + c(\text{عمیق}(d_1); \text{یادگیری}) \log P^{(1)}(\text{عمیق}(d_1))$$

$$\sum_w c(w; d_1) \log P^{(1)}(w) = \log 0.16 + \log 0.23 + \log 0.31 = -1.943$$

$$\sum_w c(w; d_2) \log P^{(1)}(w) = 0 + \log 0.23 + \log 0.31 = -1.147$$

$$\sum_w c(w; d_3) \log P^{(1)}(w) = \log 0.16 + 0 + 0 = -0.796$$

$$\sum_w c(w; d_4) \log P^{(1)}(w) = 0 + \log 0.23 + \log 0.31 = -1.147$$

$$\sum_w c(w; d_5) \log P^{(1)}(w) = 0 + 0 + \log 0.31 = -0.509$$

$$\log p(F^1 | \theta) = (-1.943) + (-1.147) + (-0.796) + (-1.147) + (-0.509)$$

$$\log p(F^1 | \theta) = -5.542$$

$$\log p(F^0 | \theta) = 9 \log 0.334 = -4.286$$

نتیجه نهایی:

word	#	$P(w C)$	Iteration 0		Iteration 1	
			$p(w \theta)$	$P(z=1)$	$p(w \theta)$	$P(z=1)$
الگوریتم	2	0.001	0.334	0.00128	0.22238	
یادگیری	3	0.002	0.334	0.00256	0.33314	
عمیق	4	0.0015	0.334	0.00192	0.44448	
Log-Likelihood			-4.286		-5.542	

پاسخ بخش ب

کدام کلمات وزن بیشتری می‌گیرند و چرا؟

در میان واژه‌های بررسی شده، کلمه «عمیق» بیشترین وزن را به دست آورده است. دلیل این موضوع این است که این واژه بیشترین تعداد تکرار را در اسناد بازخورد داشته و در محاسبات مرحله‌ی E نیز مقدار $P(Z=1 | w)$ برای آن بزرگتر بوده است، بنابراین سهم بیشتری در مدل موضوعی پیدا کرده است. پس از آن واژه «یادگیری» قرار می‌گیرد که با سه بار تکرار و مقدار نسبتاً بزرگ $P(Z=1 | w)$ وزن متوسطی گرفته است. در نهایت واژه «الگوریتم» کمترین وزن را دارد زیرا تنها دو بار در

اسناد دیده شده و احتمال تعلق آن به مدل زمینه بیشتر بوده است، در نتیجه سهم کمتری در مدل موضوعی دارد.

اگر λ افزایش یابد، تأثیر آن چیست؟

اگر مقدار λ افزایش یابد، وزن مدل زمینه بیشتر می‌شود و در نتیجه احتمال تعلق واژه‌ها به مدل زمینه افزایش پیدا می‌کند. این تغییر باعث می‌شود که مقدار

$1 - P(Z = 1 | w)$ کاوش می‌یابد و در مرحله‌ی M وزن واژه‌ها در مدل بازخورد کمتر می‌شود. پس، هر چه λ بزرگتر باشد، مدل بازخورد محافظه‌کارتر عمل می‌کند و بیشتر به توزیع مجموعه‌ی زمینه نزدیک می‌شود. این یعنی واژه‌های عمومی‌تر تقویت می‌شوند و واژه‌های خاصتر که در اسناد بازخورد اهمیت دارند، ممکن است نادیده گرفته شوند.

آیا D_3 باعث نویز می‌شود؟ چگونه می‌توان آن را شناسایی کرد؟

سند D_3 فقط جمله‌ی «الگوریتم‌های ژنتیک و بهینه‌سازی» است، در عمل فقط کلمه‌ی «الگوریتم» رو تقویت می‌کنه و هیچ کمکی به واژه‌های اصلی موضوع یعنی «یادگیری» و «عمیق» نمی‌کنه. چون این دو واژه در اون سند وجود ندارن، سند از نظر محتوا ای با موضوع اصلی یعنی «یادگیری عمیق» هم راستا نیست و به همین دلیل می‌توانه نویز باشد.

برای اینکه بفهمیم نویزه یا نه، می‌توانیم چند چیز رو نگاه کنیم: اول اینکه ببینیم چقدر با واژه‌های کلیدی موضوع همپوشانی داره؛ اگر خیلی کم باشه، احتمال نویز بودن زیاده. دوم اینکه بررسی کنیم سهم این سند در محاسبات کلی چقدر؛ اگر تأثیرش روی احتمال‌ها و لاغرایکلیهود خیلی کم یا حتی منفی باشه، نشونه‌ی نویزه. سوم اینکه امتیاز اولیه‌ی سند رو با بقیه مقایسه کنیم؛ اگر پایین‌تر باشه، باز هم نشونه‌ی کم ربط بودنشه.

پاسخ سوال سوم

پاسخ بخش الف

QUERY LIKELIHOOD MODEL (مدل احتمال پرس‌وجو) :

این مدل فرض می‌کنه که هر سند یک مدل زبانی داره، و بررسی می‌کنه که احتمال تولید پرس‌وجو توسط مدل زبانی سند چقدر هست. یعنی:

$$P(Q | D) = \prod_{w \in Q} P(w | D)$$

که در اون $P(w | D)$ با استفاده از فراوانی واژه در سند و SMOOTH (JELLINEK-MERCER) محاسبه می‌شه.

KL-DIVERGENCE MODEL (مدل واگرایی کولبک-لایبلر) :

این مدل فاصله‌ی آماری بین توزیع واژه‌های پرس‌وجو و سند رو اندازه‌گیری می‌کنه. فرمول اصلی:

$$D_{KL}(P(Q) || P(D)) = \sum_{w \in V} P(w | Q) \cdot \text{LOG} \frac{P(w | Q)}{P(w | D)}$$

هرچی این فاصله کمتر باشه، سند به پرسوچو نزدیکتر و رتبه‌ی بالاتری می‌گیره.

در مدل‌های شباهت توزیعی مثل **QUERY LIKELIHOOD** و **KL-DIVERGENCE** هدف اینه که شباهت آماری بین توزیع واژه‌ها در پرسوچو و سند بررسی بشه نه صرفًا تطابق معنایی.

این شباهت باعث می‌شه سند‌هایی که از نظر آماری به پرسوچو نزدیکترن، رتبه‌ی بالاتری بگیرن.

اما چون این مدل‌ها فقط به توزیع واژه‌ها نگاه می‌کنن، ممکنه اسناد بلندتر یا کلیتر که واژه‌های بیشتری دارن، امتیاز بیشتری بگیرن حتی اگر سند کوتاه‌تری وجود داشته باشه که دقیقاً به مفهوم پرسوچو اشاره کرده باشه.

- در **QUERY LIKELIHOOD**، سند بلند احتمال بیشتری داره که واژه‌ها پرسوچو رو شامل بشه، حتی اگر اون واژه‌ها در زمینه‌ی اصلی سند نباشن.
- در **KL-DIVERGENCE**، سند بلندتر می‌تونه توزیع واژه‌ی متنوعتری داشته باشه که به توزیع پرسوچو نزدیکتر بشه، حتی اگر معنای دقیق نداشته باشه.
- کردن **SMOOTH** باعث می‌شه واژه‌های عمومی‌تر (از مدل زمینه‌ای) وارد محاسبه بشن و به نفع سند‌های کلیتر تهموم بشه.

پاسخ بخش ب

در مدل‌های آماری مثل **KL-Divergence** و **Query Likelihood**، هدف مقایسه‌ی توزیع واژه‌ها بین پرسوچو و سند هست. اما یه مشکل رایج اینه که بعضی واژه‌های پرسوچو ممکنه اصلًا در سند ظاهر نشده باشن، و این باعث می‌شه احتمال $P(w|D)$ صفر بشه. برای حل این مشکل، از **هموارسازی (Smoothing)** استفاده می‌کنیم.

شرايطی که هموارسازی باعث افت عملکرد می‌شود:

١. وقتی سند کوتاه باشه:
 - تعداد واژه‌ها کمتره، پس $P_{ML}(w|D)$ برای واژه‌های پرسوچو ممکنه صفر یا خیلی کم باشه.
 - هموارسازی باعث می‌شه $P(w|C)$ وارد محاسبه بشه، که معمولاً واژه‌های عمومی‌تر رو تقویت می‌کنه.
 - نتیجه: سند کوتاه که واژه‌های خاص داره، ممکنه امتیاز پایین‌تری بگیره چون مدل به واژه‌های عمومی‌تر وزن می‌ده.

٢. وقتی سند خام باشه:

- اگر واژه‌های سند خیلی تخصصی باشن و در مجموعه کمتر دیده شده باشن، $P(w|C)$ برای اون‌ها کم می‌شه.

- هموارسازی باعث می‌شود این واژه‌های خاص ضعیف بشن و سند از نظر مدل کمتر مرتبط به نظر بیاد حتی اگر معنایی دقیق باشد.

۳. تنظیم نادرست پارامتر هموارسازی:

- اگر مقدار λ یا μ خیلی زیاد باشد، مدل بیشتر به مجموعه کلی تکیه می‌کند و سند‌های خاص را نادیده می‌گیرد.
- اگر خیلی کم باشد، مدل حساسیت زیادی به داده‌های سند پیدا می‌کند و ممکن است ناپایدار بشد.

پاسخ بخش ج

Query Drift به پدیده‌ای گفته می‌شود که در اون توزیع زبانی پرس‌وجو به‌طور ناخواسته تغییر می‌کند، طوری که تمرکز مدل بازیابی از موضوع اصلی پرس‌وجو منحرف می‌شود. این اتفاق معمولاً در فرآیند بازخورد مرتبط (Relevance Feedback) یا گسترش پرس‌وجو (Query Expansion) (رخ می‌دهد، وقتی واژه‌هایی به پرس‌وجو اضافه می‌شون که از نظر آماری رایج ولی از نظر معنایی بی‌ربط هستند.

در مدل‌های شباهت توزیعی، مثل **Query Likelihood**، شباهت بین پرس‌وجو و سند بر اساس احتمال تولید پرس‌وجو توسط مدل زبانی سند محاسبه می‌شود.

اگر واژه‌های عمومی مثل "information" ، "system" ، "data" به پرس‌وجو اضافه بشن، چون این واژه‌ها در اکثر اسناد وجود دارند، مقدار $P(w | D)$ برای اون‌ها بالاست. در نتیجه، سند‌هایی که این واژه‌های عمومی رو دارند حتی اگر به موضوع اصلی پرس‌وجو ربطی نداشته باشند، امتیاز بیشتری می‌گیرند. این باعث می‌شود اسناد دقیق‌تر و تخصصی‌تر رتبه‌ی پایین‌تری بگیرند.

در مدل **KL-Divergence** هم همین اتفاق می‌افتد. شباهت بین توزیع واژه‌های پرس‌وجو و سند با فرمول زیر سنجیده می‌شود:

$$D_{KL}(P(Q) \parallel P(D)) = \sum_{w \in V} P(w | Q) \cdot \log \frac{P(w | Q)}{P(w | D)}$$

اگر توزیع $P(Q)$ به سمت واژه‌های عمومی متمايل بشود، فاصله KL با سند‌هایی که اون واژه‌ها رو دارند کمتر می‌شود، حتی اگر اون سند‌ها به موضوع اصلی پرس‌وجو بی‌ربط باشند. این یعنی رتبه‌بندی به نفع اسناد عمومی‌تر تغییر می‌کند.

۳ روش برای جلوگیری:

۱. مدل ترکیبی با کنترل نویز (Generative Mixture Model)

در این مدل، واژه‌ها از دو منبع تولید می‌شون: مدل زمینه‌ای $P(w | C)$ و مدل موضوعی $P(w | \theta)$. با تنظیم مناسب پارامتر λ ، می‌توانیم تأثیر واژه‌های عمومی رو کاوش بدیم:

$$P(w) = (1 - \lambda) \cdot P(w | \theta) + \lambda \cdot P(w | C)$$

هرچی λ کمتر باشد، تمرکز مدل روی واژه‌های موضوعی بیشتر می‌شود و Drift اتفاق می‌افتد.

۲. فیلتر واژه‌های عمومی با وزن‌دهی معکوس (IDF Filtering)

قبل از اعمال مدل، میتونیم واژه‌هایی که در مجموعه زیاد دیده می‌شون (دارای IDF پایین) را حذف یا وزن‌شون رو کم کنیم. این باعث می‌شود واژه‌های خاص و مهم‌تر تأثیر بیشتری در رتبه‌بندی داشته باشند و Drift کنترل بشوند.

۳. Query Term Weighting (وزن‌دهی پویا به واژه‌های پرس‌وجو)

به جای اینکه همه واژه‌های پرس‌وجو وزن برابر داشته باشند، میتوانیم وزن واژه‌ها را بر اساس اهمیت آماری یا معنایی تنظیم کنیم. واژه‌های عمومی وزن کمتر می‌گیرند و واژه‌های خاص وزن بیشتر. این کار باعث می‌شود Drift کنترل بشوند چون واژه‌های عمومی نمیتوانن توزیع پرس‌وجو را منحرف کنند.

پاسخ بخش د

Distribution Shift سیستم بازیابی با هاش مواجه می‌شود، با توزیعی که مدل زبانی بر اساس اون آموخته دیده متفاوت باشد. این تغییر میتوانه ناشی از موارد زیر باشد:

- ورود واژه‌ها یا اصطلاحات جدید (مثل "LLM" یا "ChatGPT")
- تغییر سبک نگارش (مثل از رسمی به محاوره‌ای)
- تغییر موضوعات رایج (مثل از پزشکی به هوش مصنوعی)
- تغییر زبان یا ترکیب زبان‌ها (مثل پرس‌وجوی انگلیسی در مجموعه فارسی)

در بازیابی اطلاعات، این یعنی پرس‌وجوها یا اسناد جدید ممکنه شامل واژه‌ها، ساختارها یا مفاهیمی باشند که مدل قبلاً ندیده یا به‌ندرت دیده. در نتیجه، مدل نمیتوانه به درستی احتمال واژه‌ها رو تخمین بزنند یا شباهت معنایی رو تشخیص بدهد.

مدل‌های زبانی مثل Query Likelihood یا $P(w | D)$ یا KL-Divergence به داده‌های قبلی وابسته‌اند. وقتی واژه‌ای جدید وارد می‌شود:

- در مدل‌های کلاسیک: احتمال واژه صفر می‌شود یا باید با smoothing تخمین زده بشوند، که معمولاً مقدار خیلی کمی می‌دهند.
- در مدل‌های یادگیری عمیق embedding: واژه ممکنه وجود نداشته باشد یا ضعیف باشد، و مدل نمیتوانه ارتباط معنایی رو درک کند.

نتیجه: سندهایی که واژه‌های جدید دارند، حتی اگر مرتبط باشند، ممکنه رتبه‌ی پایینی بگیرند چون مدل نمیتوانه شباهت آماری یا معنایی رو درست تشخیص بدهد.

دو راهکار مهندسی یا طراحی برای افزایش پایداری در برابر Distribution Shift

۱. استفاده از مدل‌های تطبیقی (Adaptive Language Models)

مدل‌های تطبیقی به‌طور مداوم خودشون رو با داده‌های جدید هماهنگ می‌کنند. دو روش مهم:

میش مدل زبانی به صورت منظم با داده های جدید (مثل Fine-tuning اسناد تازه وارد یا پرسو جوهای اخیر) آموزش داده می شه. این باعث می شه مدل واژه ها و ساختارهای جدید رو یاد بگیره.

Continual Learning: مدل به جای آموزش مجدد کامل، به صورت افزایشی یاد می گیره. یعنی دانش قبلی حفظ می شه و داده های جدید بهش اضافه می شن. این روش برای محیط هایی که داده ها دائماً تغییر می کنن خیلی مؤثره.

مثال: اگر در یک موتور جستجو، کاربران شروع به استفاده از واژه هی جدیدی مثل "LLM" کنن، مدل تطبیقی می تونه با استفاده از پرسو جوهای اخیر و اسناد مرتبط، معنای این واژه رو یاد بگیره و در رتبه بندی لحاظ کنه.

۲. افزایش پوشه واژگانی با استفاده از مدل های زمینه ای (Contextual Embedding)

مدل هایی مثل BERT، RoBERTa، T5 به جای اینکه برای هر واژه یک embedding ثابت بسازن، بردار معنایی واژه رو بر اساس جمله یا زمینه تولید می کنن. این باعث می شه حتی واژه های جدید یا نادر، در زمینه های مناسب معنا پیدا کنن.

مثال: واژه هی "prompt" ممکنه در زمینه هی "writing prompt" معنای آموزشی داشته باشه، ولی در جمله هی "prompt engineering in LLMs" معنای فنی پیدا کنه. مدل های زمینه ای می تونن این تفاوت رو تشخیص بدن حتی اگر واژه قبلاً دیده نشده باشه.