

شبکه های اجتماعی و اقتصادی

دانشکده مهندسی کامپیوتر

مریم رمضانی

بهار ۱۴۰۴



تمرین چهارم

سیستم های توصیه گر، گراف، شبکه

تاریخ انتشار: ۱۰ خرداد ۱۴۰۴

۱. سوالات خود در مورد این تمرین را در کوئرا مطرح کنید.

۲. لطفا پاسخ خود را با توضیحات کافی ارائه دهید.

تاریخ تحویل: ۲۳ خرداد ۱۴۰۴

سوالات تئوری (۸۰ نمره)

پرسش ۱ (۲۰ نمره)

سوال: تحلیل گراف های ناهمگن و کاربرد Meta-Path در خوشه بندی مدل سازی گراف های ناهمگن:

گرافی ناهمگن تعریف کنید که شامل نودهای متفاوت مانند مقاله ها، ژورنال ها و نویسندگان باشد. همچنین گراف دیگری شامل فیلم ها، کارگردان ها و بازیگران ایجاد نمایید. در هر دو گراف، نوع نودها و نوع یال ها (روابط میان نودها) را مشخص کنید.

تعریف Meta-Path:

مفهوم Meta-Path را در گراف های ناهمگن توضیح دهید. چرا این مفهوم برای تحلیل رابطه های غیر مستقیم بین نودها مهم است؟

کاربرد Meta-Path در مدل سازی گراف:

نحوه ی استفاده از Meta-Path ها برای مدل سازی روابط میان نودهای هم نوع (مثلاً نویسنده با نویسنده یا فیلم با فیلم) را توضیح دهید. چگونه می توان از این مدل سازی برای اعمال الگوریتم های Community Detection استفاده کرد؟

طراحی Meta-Path برای گراف فیلم ها: حداقل سه Meta-Path مختلف پیشنهاد دهید که برای خوشه بندی (Clustering) فیلم ها مفید باشند. برای هر Meta-Path توضیح دهید چه نوع ارتباطی میان فیلم ها برقرار می شود و چرا برای کشف جامعه های معنایی مناسب است.

تعریف Meta-Path برای گراف مقالات جهت خوشه بندی نویسندگان: چگونه می توان از Meta-Path ها برای تعریف شباهت میان نویسندگان استفاده کرد؟ چند نمونه Meta-Path ارائه دهید که خوشه بندی نویسندگان بر اساس آن ها معنادار باشد (مثلاً همکاری مشترک، انتشار در ژورنال مشابه و ...).

پاسخ

۱ مدل سازی گراف های ناهمگن

۱.۱ گراف مقالات

نودها:

- مقاله
- نویسنده
- ژورنال

یال ها:

- نویسنده → نوشته است ← مقاله
- مقاله → منتشر شده در ← ژورنال

۲.۱ گراف فیلم ها

نودها:

- فیلم
- کارگردان
- بازیگر

یال ها:

- کارگردان → کارگردانی کرده ← فیلم
- بازیگر → بازی کرده در ← فیلم

۲ تعریف Meta-Path

- **Meta-Path** دنباله‌ای از انواع نودها و انواع یال‌ها است که مسیرهای معنادار در گراف ناهمگن را تعریف می‌کند.
- **Meta-Path** امکان تحلیل روابط غیر مستقیم را فراهم می‌کند. در بسیاری موارد، نودها ارتباط مستقیم ندارند ولی از طریق زنجیره‌ی معنایی می‌توان مشابهت یا نزدیکی آن‌ها را سنجید.

۳ کاربرد Meta-Path در مدل‌سازی گراف

با تعریف **Meta-Path**، می‌توان شباهت میان نودهای هم‌نوع را مدل‌سازی کرد. این شباهت‌ها می‌توانند مبنای ایجاد ماتریس شباهت (**Similarity Matrix**) باشند. سپس الگوریتم‌های **Community Detection** روی این ماتریس اعمال می‌شوند و خوشه‌های معنایی استخراج می‌شوند.

۴ طراحی Meta-Path برای گراف فیلم‌ها

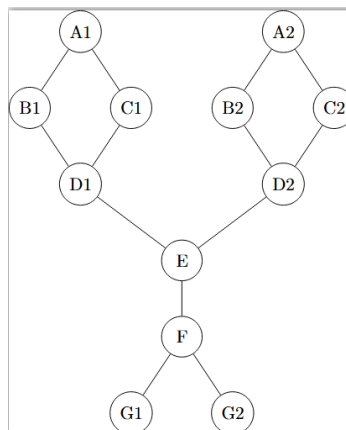
- فیلم-بازیگر-فیلم: یک بازیگر در هر دو فیلم بازی کرده است. هر بازیگر سبک خاصی دارد و مناسب بازی در برخی فیلم‌ها است.
- فیلم-کارگردان-فیلم: هر دو فیلم توسط یک کارگردان کارگردانی شده‌اند. کارگردان تأثیر زیادی در سبک فیلم دارد و یکی بودن کارگردان نشان‌دهنده شباهت زیاد بین دو فیلم است.
- فیلم-ژانر-فیلم: دو فیلم از یک ژانر هستند. در این حالت لازم است ژانر نیز به نودهای گراف اضافه شود.

۵ تعریف Meta-Path برای گراف مقالات

در گراف مقاله-ژورنال-نویسنده، می‌توان **Meta-Path** های زیر را تعریف کرد:

- نویسنده-مقاله-نویسنده: اگر دو نویسنده یک کار مشترک داشته باشند به هم مرتبط می‌شوند.
 - نویسنده-مقاله-ژورنال-مقاله-نویسنده: اگر دو نویسنده در یک ژورنال مشترک مقاله منتشر کرده باشند به هم مرتبط می‌شوند.
- این روابط می‌توانند نشان دهند دو نویسنده حوزه‌های کاری نزدیک دارند و زیاد بودن این مسیرها احتمال هم‌حوزه بودن آن‌ها را تقویت می‌کند.

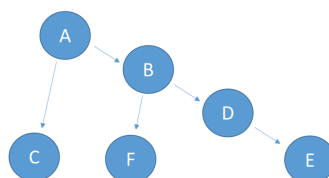
پرسش ۲ (۲۰ نمره)



اجرای الگوریتم‌های **Node2Vec**، **GraphSAGE** و **GAT** را برای گراف داده شده به صورت گام به گام تا سه مرحله توضیح دهید. بر اساس پیچیدگی روابط بین رئوس، اندازه گراف و زمانی بودن توضیح دهید در هر حالت کدام الگوریتم مناسب‌تر است و چرا.

پاسخ

پرسش ۳ (۲۰ نمره)



احتمال	یال
۰,۹	$A \rightarrow B$
۰,۴	$A \rightarrow C$
۰,۶	$B \rightarrow D$
۰,۳	$B \rightarrow F$
۰,۵	$D \rightarrow E$

مرکزیت در فضای تعبیه شده (Embedding Space) در برابر مرکزیت ساختاری

هر گره در یک فضای دوبعدی نهفته تعبیه شده که از طریق راهپیمایی‌های بایاس دار (Node2Vec) به دست آمده است. مختصات گره‌ها به صورت زیر است:

گره	(x, y)
A	(9.0, 1.0)
B	(8.0, 2.0)
C	(7.0, 0.0)
D	(6.0, 3.0)
E	(5.0, 4.0)
F	(7.5, 2.5)

(آ) فاصله اقلیدسی گره F تا سایر گره‌ها را محاسبه کرده و میانگین فاصله‌ها را گزارش دهید.

(ب) آیا مرکزیت در فضای تعبیه شده لزوماً با پتانسیل تأثیرگذاری در گراف ساختاری همبستگی دارد؟ با ذکر مثال مشخص از این گراف توضیح دهید.

تحلیل انتشار تأثیر (مدل IC)

(آ) با استفاده از مدل انتشار مستقل (IC) تعداد مورد انتظار گره‌های فعال شده با شروع از گره A را محاسبه کنید. فرض کنید فرایند انتشار در حداکثر سه مرحله انجام می‌شود. فرمول $\text{Expected}(X \rightarrow Y) = P(X \text{ activates } Y) \times P(X \text{ is active})$ را نیز محاسبه و تحلیل کنید:

• احتمال فعال شدن گره E

• تفاوت در انتظار تأثیر زمانی که گره B به جای A به عنوان seed انتخاب می‌شود

(ب) توضیح دهید که چرا از الگوریتم‌های حریصانه در بیشینه‌سازی تأثیر استفاده می‌شود و به صورت خلاصه نشان دهید که تابع تأثیر زیرمجموعه‌گرا (submodular) است.

دینامیک - GRL هم‌رخدادی متنی در Node2Vec

فرض کنید تعبیه‌های فوق با استفاده از الگوریتم Node2Vec به دست آمده‌اند با تنظیمات زیر: طول راهپیمایی: ۳، دو راهپیمایی برای هر گره، پارامتر بازگشت $p = 1$ ، پارامتر in-out برابر $q = 0.25$ (تمایل به اکتشاف)

(آ) دو راهپیمایی نمونه از گره A را شبیه‌سازی کنید. وزن انتخاب‌ها و مسیرها را مشخص کنید. مثال: $A \rightarrow B \rightarrow F \rightarrow B$

(ب) بر اساس راهپیمایی‌ها، ماتریس هم‌رخدادی بین گره‌ها را بسازید. این ماتریس چه تأثیری بر شباهت نقطه‌ای بین A و سایر گره‌ها در آموزش skip-gram دارد؟

(ج) توضیح دهید که چگونه تغییر پارامتر q در گراف‌های پراکنده و متراکم تعبیه‌ها را تحت تأثیر قرار می‌دهد. این مسئله چه ارتباطی با سیستم‌های توصیه‌گر حساس به تأثیر دارد؟

ملاحظات راهبردی و نمایش چندوجهی

(آ) محدودیت‌های استفاده از تعبیه‌های ایستا (مانند Node2Vec یا DeepWalk) در فرآیندهای پویا مثل انتشار تأثیر را توضیح دهید. یک روش جایگزین پیشنهاد دهید (مانند GRL زمانی یا مدل‌های مبتنی بر توجه).

(ب) در گراف‌های چندوجهی که تأثیر تنها ساختاری نیست (مثلاً وابسته به محتوا یا زمان یا اعتماد اجتماعی است)، چگونه می‌توان معماری‌های GRL را گسترش داد؟ یک اسکچ از نحوه انتقال پیام در GNN را که شامل احتمال تأثیر است، ارائه دهید.

به تمامی پرسش‌ها با استدلال دقیق پاسخ دهید. در صورت نیاز به مفاهیم نظری مانند زیرمجموعه‌گرایی، توزیع پایدار راهپیمایی تصادفی یا گرادیان‌های تابع هزینه skip-gram ارجاع دهید.

پاسخ

پرسش ۳ - مرکزیت در تعبیه، انتشار، IC و Node2Vec

(۱) مرکزیت در فضای تعبیه

مختصات تعبیه دوبعدی فرضی:

$$A(9, 1), \quad B(8, 2), \quad C(7, 0), \quad D(6, 3), \quad E(5, 4), \quad F(7.5, 2.5).$$

برای هر $X \neq F$ داریم $d(F, X) = \sqrt{(x_F - x_X)^2 + (y_F - y_X)^2}$:

$$\begin{aligned} d(F, A) &= \sqrt{(7/5 - 9)^2 + (2/5 - 1)^2} \approx 2/1213, \\ d(F, B) &= \sqrt{(7/5 - 8)^2 + (2/5 - 2)^2} \approx 0/7071, \\ d(F, C) &= \sqrt{(7/5 - 7)^2 + (2/5 - 0)^2} \approx 2/5495, \\ d(F, D) &= \sqrt{(7/5 - 6)^2 + (2/5 - 3)^2} \approx 1/5811, \\ d(F, E) &= \sqrt{(7/5 - 5)^2 + (2/5 - 4)^2} \approx 2/9155. \end{aligned}$$

میانگین فاصله F از بقیه تقریباً $1/9749$ است.

نکته تحلیلی. مرکزی بودن در تعبیه لزوماً به معنی قدرت اثرگذاری در شبکه انتشار نیست؛ چون انتشار به جهت یال‌ها و احتمال‌های فعال‌سازی وابسته است، نه صرفاً نزدیکی در تعبیه.

(۲) مدل انتشار مستقل (IC) تا ۳ گام با بذر A

یال‌ها و احتمال‌ها:

$$A \rightarrow B : 0/9, \quad A \rightarrow C : 0/4, \quad B \rightarrow D : 0/6, \quad B \rightarrow F : 0/3, \quad D \rightarrow E : 0/5.$$

احتمال فعال بودن گره‌ها:

$$\begin{aligned} P(A) &= 1, \quad P(B) = 0/9, \quad P(C) = 0/4, \\ P(D) &= 0/9 \times 0/6 = 0/54, \quad P(F) = 0/9 \times 0/3 = 0/27, \quad P(E) = 0/54 \times 0/5 = 0/27. \end{aligned}$$

تعداد مورد انتظار فعال‌ها:

$$1 + 0/9 + 0/4 + 0/54 + 0/27 + 0/27 = 3/38.$$

$$P(E | A) = 0/9 \times 0/6 \times 0/5 = 0/27 \quad \text{همچنین}$$

اگر بذر B باشد:

$$P(B) = 1, \quad P(D) = 0/6, \quad P(F) = 0/3, \quad P(E) = 0/6 \times 0/5 = 0/3 \Rightarrow 1 + 0/6 + 0/3 + 0/3 = 2/20.$$

$$\text{بنابراین تفاوت انتظاری: } 3/38 - 2/20 = 1/18.$$

چرایی روش حریصانه. تابع گسترش مورد انتظار در IC زیرمدولار است (بازده کاهنده)، لذا انتخاب حریصانه k بذر تضمین $1 - 1/e$ دارد.

(۳) Vec2Node با $p = 1, q = 0/25$

در گذار مرتبه دوم، وزن‌ها:

$$\alpha_{pq}(t, x) = \begin{cases} 1/p & x = t, \\ 1 & (t, x) \in E, \\ 1/q & (t, x) \notin E. \end{cases} \Rightarrow 1/q = 4.$$

پس از گام $A \rightarrow B$ ، همسایه‌های B برابر $\{A, D, F\}$ است و

$$w(A) = 1, \quad w(D) = 4, \quad w(F) = 4 \Rightarrow P(A) = \frac{1}{9}, \quad P(D) = \frac{4}{9}, \quad P(F) = \frac{4}{9}.$$

(الف) دو راه‌روی نمونه طول ۳ از هر گره

$$\begin{aligned} A: & A \rightarrow B \rightarrow F \rightarrow B, \quad A \rightarrow B \rightarrow D \rightarrow E \\ B: & B \rightarrow D \rightarrow E \rightarrow D, \quad B \rightarrow F \rightarrow B \rightarrow D \\ C: & C \rightarrow A \rightarrow B \rightarrow D, \quad C \rightarrow A \rightarrow B \rightarrow F \\ D: & D \rightarrow E \rightarrow D \rightarrow B, \quad D \rightarrow B \rightarrow F \rightarrow B \\ E: & E \rightarrow D \rightarrow B \rightarrow D, \quad E \rightarrow D \rightarrow B \rightarrow F \\ F: & F \rightarrow B \rightarrow D \rightarrow E, \quad F \rightarrow B \rightarrow A \rightarrow B \end{aligned}$$

(ب) ماتریس هم‌وقوعی با پنجره $w = 2$

جدول زیر با شمردن هم‌وقوعی گره‌ها در فاصله موضعی $2 \leq$ داخل ۱۲ راهرو بالا به دست آمده است.

	F	E	D	C	B	A	
A	۳	۰	۲	۲	۶	۰	
B	۱۰	۶	۱۰	۲	۰	۶	
C	۰	۰	۰	۰	۲	۲	
D	۴	۸	۰	۰	۱۰	۲	
E	۰	۰	۸	۰	۶	۰	
F	۰	۰	۴	۰	۱۰	۳	

تفسیر. هم‌وقوعی‌های بزرگ بین (B, D) ، (B, F) و (A, B) دیده می‌شود؛ بنابراین در تعبیه، نمایش A به این گره‌ها نزدیک‌تر خواهد شد. مقدار $q < 1$ رفتار مشابه DFS و اکتشاف دورتر را تقویت می‌کند؛ درحالی‌که $q > 1$ به رفتار محلی‌تر شبیه BFS منجر می‌شود.

پرسش ۹ (۲۰ نمره) سیستم‌های پیشنهاددهنده

یک گراف دو بخشی user-item را در نظر بگیرید که در آن هر یال بین کاربر U و آیت I نشان‌دهنده‌ی این است که کاربر U آیت I را پسندیده است. همچنین ماتریس امتیازدهی را برای این مجموعه از کاربران و آیت‌ها با R نمایش می‌دهیم، که در آن هر سطر از R مربوط به یک کاربر و هر ستون مربوط به یک آیت است. اگر کاربر i آیت j را پسندیده باشد، آنگاه $R_{i,j} = 1$ ، در غیر این صورت $R_{i,j} = 0$. همچنین فرض می‌کنیم m کاربر و n آیت داریم، بنابراین ماتریس R ابعادی برابر با $m \times n$ خواهد داشت.

ماتریس P با ابعاد $m \times m$ را تعریف می‌کنیم به عنوان یک ماتریس قطری که در آن عنصر قطری i ام برابر است با درجه‌ی رأس کاربر i ، یعنی تعداد آیت‌هایی که کاربر i آن‌ها را پسندیده است. به صورت مشابه، ماتریس Q با ابعاد $n \times n$ یک ماتریس قطری است که در آن عنصر قطری j ام برابر است با درجه‌ی رأس آیت j ، یا به عبارتی تعداد کاربرانی که آیت j را پسندیده‌اند. برای مثال به شکل زیر توجه کنید.

(الف)

ماتریس شباهت کاربر به صورت غیر نرمال شده را به صورت $T = RR^T$ تعریف می‌کنیم (یعنی حاصل ضرب ماتریس R و ترانپوز آن). مفهوم عناصر $T_{i,j}$ (برای $i \neq j$) را از منظر ساختارهای گراف دو بخشی (مانند درجه‌ی رأس‌ها، مسیر میان رأس‌ها و ...) توضیح دهید (به شکل ۱ مراجعه کنید).

شباهت کسینوسی:

یادآوری می‌کنیم که شباهت کسینوسی بین دو بردار \mathbf{u} و \mathbf{v} به صورت زیر تعریف می‌شود:

$$\cos - \text{sim}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}$$

(ب)

ماتریس شباهت آیت را به صورت S_I با ابعاد $n \times n$ تعریف می‌کنیم، به طوری‌که عنصر سطر i و ستون j برابر است با شباهت کسینوسی بین آیت i و آیت j که متناظر با ستون‌های i و j از ماتریس R هستند. نشان دهید که:

$$S_I = Q^{-1/2} R^T R Q^{-1/2}$$

که در آن $Q^{-1/2} = 1/\sqrt{Q_{rc}}$ به صورت $Q^{-1/2}$ برای تمام درایه‌های غیر صفر تعریف شده است و برای سایر درایه‌ها برابر صفر است.

همین سؤال را برای ماتریس شباهت کاربران، یعنی S_U تکرار کنید؛ به طوری‌که عنصر سطر i و ستون j برابر است با شباهت کسینوسی بین کاربر i و کاربر j که متناظر با سطرهای i و j از ماتریس R هستند. عبارت S_U را نیز به صورت یک عملیات ماتریسی بر حسب R ، P و Q بیان کنید؛ توجه داشته باشید که نباید به صورت درایه به درایه S_U را تعریف نمایید.

$$\mathbf{R} = \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\mathbf{P} = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\mathbf{Q} = \begin{bmatrix} 3 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 3 \end{bmatrix}$$

* شکل ۱: گراف دو بخشی کاربر-آیت.

پاسخ شما باید نحوه‌ی استخراج این رابطه‌ها را نشان دهد. (نکته: برای جذر درایه به درایه‌ی یک ماتریس، می‌توان آن را به توان $1/2$ نوشت.) (ج)

روش پیشنهادی بر اساس فیلترینگ مشارکتی کاربر-کاربر برای کاربر u به این صورت است: برای تمام آیتم‌ها s محاسبه کنید

$$r_{u,s} = \sum_{x \in \text{کاربران}} \cos(\text{sim}(x, u)) \cdot R_{xs}$$

و سپس k آیتی را پیشنهاد دهید که مقدار $r_{u,s}$ در آن‌ها بیشینه است.

به صورت مشابه، روش پیشنهادی بر اساس فیلترینگ مشارکتی آیت-آیت برای کاربر u به این صورت است: برای تمام آیتم‌ها s محاسبه کنید

$$r_{u,s} = \sum_{x \in \text{آیتم‌ها}} R_{ux} \cdot \cos(\text{sim}(x, s))$$

و سپس k آیتی را پیشنهاد دهید که مقدار $r_{u,s}$ در آن‌ها بیشینه است.

ماتریس پیشنهادی Γ را تعریف می‌کنیم که ابعادی برابر با $m \times n$ دارد، به طوری که $\Gamma(i, j) = r_{i,j}$. ماتریس Γ را برای هر دو روش فیلترینگ کاربر-کاربر و فیلترینگ آیت-آیت بر حسب P, R, Q بیابید. پاسخ نهایی شما باید شامل عملیات در سطح ماتریسی باشد و نه تعریف درایه به درایه. راهنمایی: در حالت آیت-آیت داریم:

$$\Gamma = RQ^{-1/2}R^T.$$

پاسخ شما باید روند رسیدن به این عبارت‌ها را (حتی در مورد آیت-آیت که فرمول نهایی داده شده) توضیح دهد. (د)

جدول امتیازدهی زیر را بین پنج کاربر و شش آیت در نظر بگیرید:

شناسه آیت	۱	۲	۳	۴	۵	۶
۱	۵	۶	۷	۴	۳	؟
۲	۴	؟	۳	؟	۵	۴
۳	؟	۳	۴	۱	۱	؟
۴	۷	۴	۳	۶	؟	۴
۵	۱	؟	۳	۲	۲	۵

جدول ۱: امتیازدهی پنج کاربر (ردیف‌ها) به شش آیت (ستون‌ها).

- (د-۱) مقادیر نامشخص امتیازدهی کاربر ۲ را با استفاده از الگوریتم فیلترینگ مشارکتی مبتنی بر کاربر پیش‌بینی کنید. از ضریب همبستگی پیرسون با میانگین‌گیری استفاده کنید. فرض کنید اندازه گروه هم‌تایان حداکثر ۲ است و همبستگی‌های منفی را حذف نمایید.
- (د-۲) مقادیر نامشخص امتیازدهی کاربر ۲ را با استفاده از الگوریتم فیلترینگ مشارکتی مبتنی بر آیت پیش‌بینی کنید. از شباهت کسینوسی تعدیل‌شده استفاده کنید. فرض کنید اندازه گروه هم‌تایان حداکثر ۲ است و همبستگی‌های منفی را حذف نمایید.

اکنون، یک سامانه پیشنهاددهنده فیلم را در نظر بگیرید که در آن فیلم‌ها با ژانرها مرتبط هستند و امتیازدهی یک کاربر خاص مشخص است:

شناسه فیلم	کمدی	درام	عاشقانه	هیجانی	اکشن	ترسناک	دوست داشتن یا نداشتن
۱	۱	۰	۱	۰	۰	۰	دوست ندارد
۲	۱	۱	۱	۰	۱	۰	دوست ندارد
۳	۱	۱	۰	۰	۰	۰	دوست ندارد
۴	۰	۰	۰	۱	۱	۰	دوست دارد
۵	۰	۱	۰	۱	۱	۱	دوست دارد
۶	۰	۰	۰	۰	۱	۱	دوست دارد
آزمون-۱	۰	۰	۰	۱	۰	۱	؟
آزمون-۲	۰	۱	۱	۰	۰	۰	؟

جدول ۲: رابطه ژانر-فیلم و بازخورد کاربر در قالب دوست داشتن یا نداشتن.

تمام قوانین انجمنی را با حداقل پشتیبانی ۳۳٪ و اطمینان ۷۵٪ استخراج کنید. بر اساس این قوانین، آیا آیت آزمون-۱ یا آزمون-۲ را به کاربر پیشنهاد می‌کنید؟

موارد مورد نیاز برای تحویل:

- (i) تفسیر T_{ij} و T_{ii} .
(ii) بیان S_I و S_U برحسب ماتریس‌های R ، P و Q به همراه توضیح.
(iii) بیان Γ برحسب R ، P و Q به همراه توضیح.
(iv) محاسبات مبتنی بر تکنیک‌های فیلترینگ مشارکتی.
(v) پاسخ به سوالات.

پاسخ

پاسخ سؤال ۹

(الف)

می‌دانیم $R_{ij} \in \{0, 1\}$ و $R_{ij}^\top = R_{ji}$. بنابراین می‌توانیم T_{ii} و T_{ij} را به صورت زیر محاسبه کنیم:

$$T_{ii} = \sum_{k=1}^n R_{ik} R_{ki}^\top = \sum_{k=1}^n R_{ik}^\top.$$

از آن‌جا که اگر $R_{ik} = 1$ آنگاه $R_{ik}^\top = 1$ و در غیر این صورت ۰ است، خواهیم داشت

$$T_{ii} = \sum_{k=1}^n R_{ik}.$$

پس T_{ii} برابر است با تعداد آیت‌هایی که کاربر i آن‌ها را دوست دارد (درجه گره کاربر i).
برای $j \neq i$:

$$T_{ij} = \sum_{k=1}^n R_{ik} R_{kj}^\top = \sum_{k=1}^n R_{ik} R_{jk}.$$

عبارت $R_{ik} R_{jk} = 1$ تنها وقتی برقرار است که هر دو کاربر i و j آیت k را دوست داشته باشند؛ بنابراین T_{ij} تعداد آیت‌های مشترکی را می‌شمارد که هر دو کاربر i و j می‌پسندند (و معادل تعداد مسیرها بین i و j در گراف کاربر-آیت است).

(ب)

بگذارید R_i^\top بردار سطری i ام در R^\top (امتیازهای آیت i) باشد. سپس $R^{*\top}$ را طوری تعریف می‌کنیم که

$$R_{ij}^{*\top} = \frac{R_{ij}^\top}{\|R_i^\top\|}, \quad \|R_i^\top\| = \sqrt{\sum_{j=1}^m (R_{ij}^\top)^2}.$$

از آن‌جا که $\|R_i^\top\|$ ریشه تعداد کاربرانی است که آیت i را دوست دارند، و این تعداد همان درجه گره آیت i است، می‌توان نوشت

$$\|R_i^\top\| = \sqrt{Q_{ii}}.$$

ماتریس $Q^{-1/2}$ را چنین تعریف می‌کنیم:

$$(Q^{-1/2})_{rc} = \begin{cases} \frac{1}{\sqrt{Q_{rc}}} & \text{برای درایه‌های نامساوی با صفر،} \\ 0 & \text{در غیر این صورت.} \end{cases}$$

در نتیجه

$$R^{*\top} = Q^{-1/2} R^\top.$$

شباهت کسینوسی آیت‌های i و j برابر است با

$$\frac{\sum_{k=1}^m R_{ik}^\top R_{jk}^\top}{\|R_i^\top\| \|R_j^\top\|} = \frac{\sum_{k=1}^m R_{ik}^\top R_{jk}^\top}{\sqrt{Q_{ii}} \sqrt{Q_{jj}}} = R_i^{*\top} \cdot (R_j^{*\top})^\top.$$

پس ماتریس شباهت آیت‌ها داریم

$$S_I = R^{*\top} (R^{*\top})^\top.$$

از آن‌جا که $Q^{-1/2}$ متقارن است، $(R^{*\top})^\top = (Q^{-1/2} R^\top)^\top = R Q^{-1/2}$ ، در نتیجه

$$S_I = Q^{-1/2} R^\top R Q^{-1/2}.$$

به طور مشابه، R_i را بردار سطری i ام در R (امتیازهای کاربر i) بگیرید و R^* را طوری تعریف کنید که

$$R_{ij}^* = \frac{R_{ij}}{\|R_i\|}, \quad \|R_i\| = \sqrt{\sum_{j=1}^n R_{ij}^2} = \sqrt{P_{ii}},$$

که در آن P_{ii} درجه کاربر i است. بگذارید $P^{-1/2}$ همانند قبل تعریف شود؛ آنگاه $R^* = P^{-1/2} R$ و شباهت کسینوسی کاربران i و j برابر است با

$$\frac{\sum_{k=1}^n R_{ik} R_{jk}}{\|R_i\| \|R_j\|} = \frac{\sum_{k=1}^n R_{ik} R_{jk}}{\sqrt{P_{ii}} \sqrt{P_{jj}}} = R_i^* \cdot (R_j^*)^\top.$$

پس ماتریس شباهت کاربران

$$S_U = R^* (R^*)^\top = P^{-1/2} R R^\top P^{-1/2}.$$

(ج)

(۱) فیلترگذاری مشارکتی کاربر-کاربر. برای پیش‌بینی $r_{u,s}$ داریم

$$r_{u,s} = \sum_{x \in \text{users}} \text{cos-sim}(x, u) R_{x,s} = \sum_{x \in \text{users}} (P^{-1/2} R R^\top P^{-1/2})_{u,x} R_{x,s} = (P^{-1/2} R R^\top P^{-1/2})_u \cdot R_s.$$

در نتیجه، ماتریس توصیه

$$\Gamma = P^{-1/2} R R^\top P^{-1/2} R.$$

(۲) فیلترگذاری مشارکتی آیتم-آیتم. برای پیش‌بینی $r_{u,s}$ داریم

$$r_{u,s} = \sum_{x \in \text{items}} R_{u,x} \text{cos-sim}(x, s) = \sum_{x \in \text{items}} R_{u,x} (Q^{-1/2} R^\top R Q^{-1/2})_{x,s} = R_u \cdot (Q^{-1/2} R^\top R Q^{-1/2})_s.$$

پس ماتریس توصیه

$$\Gamma = R Q^{-1/2} R^\top R Q^{-1/2}.$$

توضیح: در این فرمول‌ها، P و Q ماتریس‌های قطری درجات کاربران و آیتم‌ها هستند.

(د) (ترجمه پاسخ سؤال ۲: پیش‌بینی امتیازهای کاربر ۲ با جداول کامل)

فرض کنید جدول امتیازها بین پنج کاربر و شش آیتم به صورت زیر است (علامت «؟» یعنی مقدار نامشخص):

۶	۵	۴	۳	۲	۱	⇒ Item-Id
						⇓ User-Id
؟	۳	۴	۷	۶	۵	۱
۴	۵	؟	۳	؟	۴	۲
؟	۱	۱	۴	۳	؟	۳
۴	؟	۶	۳	۴	۷	۴
۵	۲	۲	۳	؟	۱	۵

گام اول: جدول میان‌مرکز شده (برای پیرسون و adjusted cosine). برای هر کاربر، میانگین سطر را از امتیازهای همان سطر کم می‌کنیم:

Mean	۶	۵	۴	۳	۲	۱	⇒ Item-Id
۵	؟	-۲	-۱	۲	۱	۰	۱
۴	۰	۱	؟	-۱	؟	۰	۲
۲.۲۵	؟	-۱.۲۵	-۱.۲۵	۱.۷۵	۰.۷۵	؟	۳
۴.۸	-۰.۸	؟	۱.۲	-۱.۸	-۰.۸	۲.۲	۴
۲.۶	۲.۴	-۰.۶	-۰.۶	۰.۴	؟	-۱.۶	۵

(الف) روش کاربر-کاربر (Pearson) + (mean-centering) ضریب‌های پی‌رسون بین کاربر ۲ و سایر کاربران (بر اساس آیت‌های مشترک ارزیابی شده) به صورت زیر محاسبه می‌شوند:

$$\text{Pearson}(1, 2) = \frac{(2)(-1) + (-2)(1)}{\sqrt{2^2 + (-2)^2} \sqrt{(-1)^2 + 1^2}} = -1,$$

$$\text{Pearson}(3, 2) \approx \frac{(-1/75)(-1) + (-1/25)(1)}{\sqrt{1/75^2 + (-1/25)^2} \sqrt{(-1)^2 + 1^2}} = \frac{-3}{3/4} \approx -0/99,$$

$$\text{Pearson}(4, 2) \approx \frac{1/8}{2/95} \approx 0/61, \quad \text{Pearson}(5, 2) \approx \frac{-0/4 - 0/6}{\sqrt{2} \sqrt{1/6^2 + 0/4^2 + 0/6^2 + 2/4^2}} \approx -0/23.$$

با توجه به شرط «همسایگی حداکثر ۲ و حذف ضرایب منفی»، نزدیک‌ترین همسایه معتبر فقط کاربر ۴ است. بنابراین از سطر میان‌مرکز شده کاربر ۴ برای آیت‌های نامشخص کاربر ۲ استفاده می‌کنیم:

$$\hat{r}_{2,2}^{(mc)} = -0/8, \quad \hat{r}_{2,4}^{(mc)} = 1/2.$$

با افزودن میانگین کاربر ۲ (یعنی ۴):

$$\hat{r}_{2,2} = 4 + (-0/8) = 3/2, \quad \hat{r}_{2,4} = 4 + 1/2 = 5/2.$$

(ب) روش آیت-آیت (Adjusted Cosine). برای آیت ۲، نزدیک‌ترین آیت‌ها ۳ و ۶ با شباهت‌های ۰/۹۹۷۷ و ۰/۷۰۷۱ هستند؛ با وزن‌دهی امتیازهای کاربر ۲ روی آن‌ها (۳ و ۴):

$$\hat{r}_{2,2} = \frac{0/9977 \times 3 + 0/7071 \times 4}{0/9977 + 0/7071} \approx 3/41.$$

برای آیت ۴، نزدیک‌ترین آیت‌ها ۱ و ۵ با شباهت‌های ۰/۷۹۱ و ۰/۹۴ هستند و امتیازهای کاربر ۲ روی آن‌ها ۴ و ۵ است:

$$\hat{r}_{2,4} = \frac{0/791 \times 4 + 0/94 \times 5}{0/791 + 0/94} \approx 4/54.$$

جمع‌بندی. در هر دو رویکرد، ترتیب ترجیح یکسان است ($\hat{r}_{2,4} > \hat{r}_{2,2}$)، ولی مقادیر مطلق متفاوت‌اند. برای تکمیل گزارش، می‌توانید در جدول اول (خام) مقادیر پیش‌بینی شده $\hat{r}_{2,2}$ و $\hat{r}_{2,4}$ را جای «۴»‌ها وارد کنید.

(ه) قوانین انجمنی و ترجیح کاربر فعال

قوانین استخراج شده.

- کمدی، درام \Rightarrow نپسندیدن
- کمدی، عاشقانه \Rightarrow نپسندیدن
- کمدی \Rightarrow نپسندیدن
- عاشقانه \Rightarrow نپسندیدن
- دلهره/هیجان‌انگیز \Rightarrow پسندیدن
- اکشن \Rightarrow پسندیدن
- ترسناک \Rightarrow پسندیدن
- دلهره/هیجان‌انگیز، اکشن \Rightarrow پسندیدن
- ترسناک، اکشن \Rightarrow پسندیدن

نتیجه‌گیری. به سادگی دیده می‌شود که Test-1 تمام قوانینی را که گزاره نتیجه آن‌ها «پسندیدن» است فعال می‌کند، در حالی که Test-2 همه قوانینی را که نتیجه‌شان «نپسندیدن» است شلیک می‌کند. بنابراین کاربر فعال، Test-1 را به Test-2 ترجیح خواهد داد.

توضیح تکمیلی.

- وجود قوانین عام مانند «کمدی \Rightarrow نپسندیدن» باعث می‌شود قوانین خاص‌تر هم‌راستا (مثل «کمدی، عاشقانه \Rightarrow نپسندیدن») نیز بدون تناقض عمل کنند.
- اگر چند قانون هم‌زمان فعال شوند و همه آن‌ها یک نتیجه داشته باشند (همه «پسندیدن» یا همه «نپسندیدن»)، تصمیم نهایی پایدار است. در مسائل واقعی، معمولاً برای حل تعارض یا اولویت‌بندی از وزن‌دهی بر اساس *confidence* و *support* هر قانون یا از رای‌گیری وزن‌دار استفاده می‌شود.

تاریخ تحویل:

سوالات عملی (۴۰ نمره)

پرسش ۱ (۲۰ نمره) به فایل جویپتر Q1.ipynb مراجعه کنید.

پاسخ

پرسش ۲ (۲۰ نمره) به فایل جویپتر Q2.ipynb مراجعه کنید.

پاسخ