



باسم‌هه تعالی

تاریخ: ۱۴۰۴/۱۰/۱۷

پروژه‌ی سوم درس مبانی علم داده

## PageRank with Restarts on WikiVote: Dense vs Sparse and Deadend Handling

در این پروژه شما الگوریتم PageRank با restart را روی شبکه‌ی واقعی WikiVote پیاده‌سازی و تحلیل می‌کنید. ابتدا روی دو زیرگراف کوچک (برای اجرای dense و فهم رفتار الگوریتم) کار می‌کنید، سپس نسخه‌ی sparse را برای اجرای مقیاس‌پذیر روی گراف کامل پیاده‌سازی می‌کنید. بخش کلیدی پروژه، تشخیص و اصلاح مشکل deadend/dummy page و تفسیر نتایج در متن شبکه‌ی رأی دهن است.

### قوانين پروژه (کتابخانه‌های مجاز/غیرمجاز)

مجاز: numpy, pandas, scipy, matplotlib, gzip, time, os, collections  
غیرمجاز: networkx.pagerank و هر تابع آماده‌ی PageRank/PPR (مثل networkx.pagerank/PPR یا هر پیاده‌سازی آماده مشابه)

### معرفی دیتاست‌ها

#### دیتاست اصلی: WikiVote

یک گراف جهت‌دار از رأی‌دادن کاربران ویکی‌پدیا در فرآیند adminship. هر یال  $A \rightarrow B$  یعنی: کاربر A به کاربر B رأی داده است.

#### زیرگراف A (نمونه‌گیری BFS/Frontier)

• فایل‌ها:

wikivote\_subgraph\_A.edges.tsv –  
(old\_id ↔ new\_id) نگاشت wikivote\_subgraph\_A.mapping.tsv –

#### زیرگراف B (نمونه‌گیری Top-k induced)

• فایل‌ها:

wikivote\_subgraph\_B.edges.tsv –  
wikivote\_subgraph\_B.mapping.tsv –

نکته: فایل‌های edges.tsv با شناسه‌های reindexed (از ۰ تا  $n - 1$ ) هستند تا ساخت ماتریس آسان باشد. فایل mapping برای ترجمه‌ی شناسه‌ها به ID اصلی WikiVote استفاده می‌شود.

### تعاریف کلیدی

#### ۱. (Compressed Sparse Row) csr\_matrix

یک روش ذخیره‌سازی ماتریس تنک است که فقط عناصر غیرصفر را نگه می‌دارد. برای گراف‌ها یعنی: بهجای نگهداشتن کل

ماتریس  $n \times n$ , فقط یال‌های موجود ذخیره می‌شوند.

**مزیت کلیدی:** ضرب ماتریس-بردار روی csr معمولاً حدوداً متناسب با  $\text{nnz}$  (تعداد عناصر غیرصفر (تقریباً برابر با تعداد یال‌ها)) انجام می‌شود، نه  $n^2$ .

## ۲. adjacency matrix (ماتریس مجاورت)

اگر یال  $v \rightarrow u$  وجود داشته باشد، آنگاه  $A[u, v] = 1$  در غیر این صورت ۰.

## ۳. in-degree / out-degree

$(u)$ : تعداد یال‌های خروجی از  $u$

$(v)$ : تعداد یال‌های ورودی به  $v$

## ۴. deadend / dangling

نودی که  $0 = \text{out-degree}$  دارد

## ۵. row-normalization (نرمال‌سازی سطری)

ساختن ماتریس انتقال  $P$  از روی  $A$ : هر سطر تقسیم بر  $\text{out-degree}$  همان نود (تا مجموع سطر ۱ شود؛ به جز danglers).

## ۶. (r) restart probability / teleport

همان پارامتر کتاب (مثالاً ۰.۱۵)

## ۷. residual (معیار توقف)

$\|p^{(t+1)} - p^{(t)}\|_1$ : residual

اگر residual از  $\epsilon$  کوچک‌تر شد، می‌گوییم همگرا شده است.

## ۸. Personalized PageRank (PPR)

مثل PageRank است، فقط توزیع restart یکنواخت نیست؛ روی یک seed یا مجموعه‌ای از seeds متمرکز می‌شود.

# گام ۰) آماده‌سازی و شناخت داده

کارهایی که باید انجام دهید

۱. دیتاست را بخوانید و تعداد node یکتا و تعداد edge را گزارش کنید.

۲. آمار پایه‌ی درجه‌ها را بدهید:

• in-degree و out-degree: میانگین، میانه، max، درصد صفرها

• تعداد node با  $0 = \text{out-degree}$  (یعنی deadend/dangling)

۳. بررسی کنید آیا self-loop دارید یا نه. اگر دارید، تصمیم بگیرید حذف شود یا نه (و یک دلیل کوتاه بنویسید).

# خروجی‌های این گام

• یک جدول خلاصه‌ی آمار

• Top-20 بر اساس in-degree و out-degree

• نمودار Rank-Degree(log-log): محور x رتبه‌ی گره‌ها بر اساس درجه ( $1 = \text{بیشترین درجه}$ ), محور y مقدار درجه. دو منحنی جدا برای in-degree و out-degree.

# سؤال‌های تحلیلی

۱. شبکه‌ی WikiVote بیشتر شبیه شبکه‌ی وب است یا شبکه‌ی ایمیل؟ چرا؟

۲. اگر تعداد زیادی out-degree صفر وجود داشته باشد، چه مشکلی برای تفسیر «random walk with restart» ایجاد می‌شود؟ ( فقط پیش‌بینی کنید).

## پارامترها

0.15 (restart/teleport):  $r$   
 $10^{-9}$  (پیشنهادی):  $\varepsilon$   
500 :**max\_iters**

## گام ۱) PageRank (Dense نسخه) روی دو زیرگراف کوچک از WikiVote

### داده‌های ورودی

- زیرگراف A (روش نمونه‌گیری BFS/Frontier)
- زیرگراف B (روش نمونه‌گیری Top-k induced)

کارهایی که باید انجام دهید (برای A و B جداگانه)

- ۱) بارگذاری و ساخت ماتریس مجاورت
  - فایل edges.tsv را بخوانید.
  - $n$  را از فایل mapping یا  $\text{max}(\text{id}) + 1$  به دست آورید.
  - یک ماتریس  $A \in \mathbb{R}^{n \times n}$  بسازید که اگر  $v \rightarrow u$  وجود داشت، آنگاه  $A[u, v] = 1$  شود.

خروجی مورد انتظار:

- چاپ  $n$  و تعداد یال‌ها
- چاپ تعداد self-loop

### ۲) ساخت ماتریس انتقال (Row-normalized)

- out-degree هر نود را حساب کنید.
- ماتریس انتقال را بسازید.
- نودهای با  $\text{out-degree} == 0$  را فقط بشمارید (فعلاً درستش نکنید).

خروجی مورد انتظار:

- تعداد deadend‌ها ( $\text{out\_degree} == 0$ )

### ۳) اجرای Basic — PageRank (بدون رفع deadend نسخه)

- $r = 0.15$  و بردار شروع را یکنواخت بگیرید (یعنی به همه نودها مقدار برابر دهید).
- در هر iteration مقدار residual (نرم  $L_1$ ) را حساب و ذخیره کنید.
- با شرط  $\varepsilon < \text{residual}$  توقف کنید.

خروجی مورد انتظار:

- تعداد iteration تا توقف

residual\_last •

$\sum_i p_i$  •

(new\_id نود (بر حسب Top-10 •

محور  $y$  لگاریتمی residual بر حسب iteration نمودار •

#### ۴) اجرای PageRank — نسخه‌ی اصلاح شده (رفع نشستی deadend)

- همان PageRank را اجرا کنید، اما طوری اصلاح کنید که در هر iteration مجموع احتمال‌ها باید تقریباً برابر ۱ بماند:  $\sum_i p_i \approx 1$ . (نشستی را جبران کنید).
- خروجی مورد انتظار:
  - تعداد iteration
  - $\sum_i p_i \approx 1$
- Top-10 نودها و مقایسه با نسخه‌ی Basic (آیا تغییر کرد یا خیر؟)
- نمودار residual مثل قبل

#### سوال‌های تحلیلی

۱. چرا در زیرگراف A در نسخه‌ی Basic مقدار  $\sum_i p_i$  کمتر از ۱ شد؟ (ارتباط با deadend/dummy/restart)
۲. چرا در زیرگراف B نسخه‌ی Basic و اصلاح شده تقریباً یکسان شد؟ (به تعداد deadend‌های B اشاره کنید.)
۳. اگر  $\epsilon$  را از  $10^{-9}$  به  $10^{-8}$  تغییر دهید، Top-10 چقدر عوض می‌شود و چرا؟

## گام (۲) PageRank نسخه‌ی Sparse و مقایسه‌ی زمان/حافظه

کارهایی که باید انجام دهید

#### ساخت ماتریس Sparse

برای هر زیرگراف (A و B) و سپس برای گراف کامل:

- از روی لیست یال‌ها یک ماتریس مجاورت  $\text{csr\_matrix}$  بسازید.

#### نرمال‌سازی ردیفی

- ماتریس انتقال  $P$  را با نرمال‌سازی ردیفی (row-normalized) بسازید.
- تعداد نودهای با  $= 0$  out-degree را گزارش کنید.
- نکته: در این گام فقط گزارش کنید؛ رفع مشکل dangling مربوط به گام ۳ است.

#### اجرای (Power Iteration) Sparse با ضرب PageRank

- PageRank را با تکرارهای متوالی (Power Iteration) اجرا کنید.
- شرط توقف پروژه را اعمال کنید

#### مقایسه‌ی زمان اجرا

برای زیرگراف‌های A و B زمان این بخش‌ها را جداگانه اندازه‌گیری کنید:

- زمان ساخت  $\text{csr}$
- زمان نرمال‌سازی
- زمان اجرای iteration
- زمان کل

## اجرای Sparse روی گراف کامل

همین نسخه Sparse را روی FullGraph اجرا کنید و گزارش دهید:

- تعداد iteration تا همگرایی

● زمان کل iteration (یا زمان متوسط هر iteration)

$$\sum_i p_i$$

برای زیرگراف‌های A و B، زمان اجرای PageRank را در دو حالت Sparse و Dense مقایسه کنید. (اختیاری)

## خروجی‌های مورد انتظار

جدول زمان برای A و B با ستون‌های زیر:

graph, build\_csr\_s, normalize\_s, iterate\_s, total\_s, iters, sum\_p, nnz

گزارش عددی برای FullGraph

n, nnz, iters, time\_iterate\_s, sum(p)

## سؤال‌های تحلیلی

۱. با استفاده از عده‌های واقعی دیتاست WikiVote (تعداد نودها و تعداد یال‌ها)، توضیح دهید چرا ساخت و نگهداری ماتریس Dense برای گراف کامل عملی نیست، اما نمایش Sparse (CSR) عملی است.؟ جواب را هم با O-Big و هم با اعداد واقعی WikiVote توضیح دهید.

۲. اگر  $\sum_i p_i$  در FullGraph کمتر از ۱ شد، این نشانه‌ی چیست و چه چیزی در گراف باعث آن می‌شود؟

## گام ۳) رفع مشکل Personalized PageRank + deadend/dummy

### (FullGraph) روی deadend (A)

ورودی: همان  $P_{full}$  و  $P_{full}^T$  که در گام ۲ ساخته‌اید.  $P$  ماتریس انتقال (نمایش سطری) و  $P^T$  (ترانهاده‌ی آن).

کارهایی که باید انجام دهید

- دوباره PageRank را اجرا کنید ولی این بار رفع deadend را اضافه کنید.

- دو نسخه را اجرا کنید و مقایسه کنید:

– q-fix: نشت را به همان بردار teleport (یعنی  $q$ ) برگردانید.

– u-fix: نشت را یکنواخت بین همه نودها پخش کنید.

## خروجی‌های مورد انتظار

● جدول residual\_last, iters, sum\_p, method

● Top-20 برای هر روش

● میزان همپوشانی Top-20 بین روش‌ها (overlap)

## سؤال تحلیلی

- چرا با رفع  $\sum_i p_i$ , deadend ۱ شود؟

## (FullGraph) کاربردی، روی Personalized PageRank (B)

کارهایی که باید انجام دهید

Seed = old\_id = 4037 •

.(seed restart با روی PPR را اجرا کنید •

خروجی‌های مورد انتظار

$\sum_i p_i$  •

Top-20 •

### سوالات تحلیلی

۱. PPR در WikiVote چه نوع «اطرافیان مرتبط» را برجسته می‌کند؟

۲. چرا ممکن است بعضی نودها عملاً به seed ربطی نگیرند (حس disconnected)؟

### سؤالات مفهومی

۱. در گراف WikiVote یا B → A یعنی «A به B رأی داده». پس PageRank بالا برای B دقیقاً چه تعییر اجتماعی دارد؟

۲. «چطور ممکن است کسی رأی‌های کمتری گرفته باشد ولی PageRank بالاتری بگیرد؟»

۳. توضیح دهید چرا اگر دو کاربر به هم رأی بدھند (دوطرفه)، می‌تواند PageRank را بالا ببرد.

۴. چرا restart با  $r = 0.15$  جلوی «بی‌نهایت بالا رفتن» را می‌گیرد؟

۵. اگر فقط اجازه داشته باشید یک معیار برای «مهم بودن» تعریف کنید، کدام را انتخاب می‌کنید و چرا؟ گزینه‌ها: in-degree, out-degree, PageRank, PPR.

پاسخ باید شامل ۲ مزیت و ۱ محدودیت باشد.