



به نام خدا  
سیستم‌های توزیع شده  
۲-۱۳۹۷

تمرین چهارم  
مدرس: صابر صالح

### نکات مهم

لطفا ابتدا به نکات زیر توجه کنید:

– برای پیاده‌سازی این تمرین از زبان Python استفاده نمایید.

– پاسخ تمرین را فقط در Quera آپلود کنید.

– به شیوه‌ی ورودی و خروجی دقت کنید.

– مهلت ارسال تمرین تا پایان روز ۱۳۹۸/۳/۲۸ می‌باشد.

موفق باشید.

## ۱ سوال ۱

در این تمرین، کد spark ای بنویسید که به وسیله‌ی آن بتوان در یک شبکه اجتماعی به کاربران پیشنهادهایی برای ایجاد دوستی<sup>۱</sup> بدهد. برای اینکار الگوریتم باید دو نفری را که تعداد زیادی دوستان مشترک دارند به هم پیشنهاد دهد.

## ۱.۱ داده‌ها

◀ داده‌ها در فایل q1\_data.txt قرار دارند.

◀ هر خط از این فایل به شکل <friends> <TAB> <user> است. در ابتدای هر خط نام کاربر (<user>) قرار دارد. سپس با یک فاصله، اسامی دوستان (<friends>) این شخص پشت سر هم قرار دارد که با یک ویرگول جدا شده است. دقت کنید که در این داده‌ها دوستی‌ها دو طرفه است یعنی یال‌های گراف کاربران جهت دار نیست.

## ۲.۱ توضیحات بیشتر

◀ برای پیاده‌سازی این الگوریتم برای هر کاربر ۱۰ شخص را به عنوان پیشنهاد دوستی مشخص کنید که در حال حاضر با هم دوست نیستند و بیشترین تعداد دوستان مشترک را با هم دارند.

◀ اگر شخصی کمتر از ۱۰ دوست درجه ۲ داشت (یعنی مجموعه دوستان دوستانش به ۱۰ نفر نمی‌رسید) همه‌ی آن‌ها را نمایش دهید (به صورت نزولی از تعداد دوست مشترک بیشتر به کمتر). اگر شخصی هیچ دوستی نداشت لیست خالی قرار دهید.

◀ به صورت کوتاه توضیح دهید که چگونه الگوریتم خود را پیاده کردید.

<sup>1</sup>Friendship recommendation

◀ برای تست کد خود، جواب خروجی برای کاربر شماره ۱۱ به شکل زیر است.

ID: 11\_ 27552, 27667, 27620, 33192, 27600, 32072, 27573, 7785, 27589, 27617

دقت کنید که خروجی برای هر کاربر به شکل نمونه آورده شده باشد. یعنی به شکل

ID:<user>\_ <recommendation>}

◀ دقت کنید که کدها دارای خطوط خالی و یا جدا شدن بدون <TAB> هستند، در هنگام پردازش داده‌ها شروطی را قرار دهید که این داده‌ها پردازش نشوند.

## ۲ PageRank

PageRank (PR) الگوریتمی می‌باشد که توسط گوگل سرچ استفاده می‌شود تا در search engine خود به webpage ها rank بدهد. این الگوریتم یک روش برای اندازه گیری اهمیت website ها می‌باشد. اساس کار PageRank بوسیله شمردن تعداد و نیز کیفیت لینک‌ها به یک page می‌باشد که بر این اساس تشخیص می‌دهد که یک website به چه اندازه مهم است. در این روش فرض می‌شود که website های مهم احتمالاً لینک‌های بیشتری از بقیه سایت‌ها دریافت می‌کنند.

### ۱.۲ الگوریتم

شما الگوریتم خود را بدینگونه شروع می‌کنید که در ابتدا برای rank هر page یک مقدار اولیه در نظر می‌گیرید. این مقدار اولیه برابر  $\frac{1}{n}$  می‌باشد که در این رابطه  $n$  بیان کننده تعداد node ها در گراف می‌باشد. سپس برای هر وب سایت  $x$  شما باید PageRank آن را به صورت iterative محاسبه نمایید که از فرمول زیر بدست می‌آید:

$$PR(x) = \frac{1-d}{N} + d \sum_{y \rightarrow x} \frac{PR(y)}{out(y)} \quad (1)$$

در فرمول بالا  $d$  یک عدد بین 0 و 1 می‌باشد (که بصورت معمول مقدار آن را 0.85 در نظر می‌گیرند).

همچنین  $out(y)$  بیان کننده تعداد یال‌هایی می‌باشد که از راس  $y$  خارج شده‌اند.

حال شما باید این کار را برای همه node ها در گراف انجام بدهید و این کار را باید به تعداد  $iteration$

ها که یکی از ورودی‌های الگوریتم است تکرار کنید.

## Pseudo-code ۲.۲

### Algorithm 1 PageRank Algorithm

<pre> 1: <b>procedure</b> PAGERANK(<math>G, iteration</math>) 2:   <math>d \leftarrow 0.85</math> 3:   <math>oh \leftarrow G</math> 4:   <math>ih \leftarrow G</math> 5:   <math>N \leftarrow G</math> 6:   <b>for</b> all <math>p</math> in the graph <b>do</b> 7:     <math>opg[p] \leftarrow \frac{1}{N}</math> 8:   <b>while</b> <math>iteration &gt; 0</math> <b>do</b> 9:     <math>dp \leftarrow 0</math> 10:    <b>for</b> all <math>p</math> that has no out-links <b>do</b> 11:      <math>dp \leftarrow dp + d \times \frac{opg[p]}{N}</math> 12:    <b>for</b> all <math>p</math> in the graph <b>do</b> 13:      <math>npg[p] \leftarrow dp + \frac{1-d}{N}</math> 14:      <b>for</b> all <math>ip</math> in <math>ih[p]</math> <b>do</b> 15:        <math>npg[p] \leftarrow npg[p] + \frac{d \times opg[ip]}{oh[ip]}</math> 16:      <math>opg \leftarrow npg</math> 17:      <math>iteration \leftarrow iteration - 1</math> </pre>	<pre> ▷ <math>G</math>: inlink file, <math>iteration</math>: # of iteration ▷ damping factor: 0.85 ▷ get outlink count hash from <math>G</math> ▷ get inlink hash from <math>G</math> ▷ get # of pages from <math>G</math> ▷ initialize PageRank ▷ get PageRank from pages without outlink ▷ get PageRank from random jump ▷ get PageRank from inlink ▷ update PageRank </pre>
---	--

### ۳ سوال ۳ (امتیازی)

#### ۱.۳ مقدمه

در این سوال قصد داریم الگوریتم MST (Minimum Spanning Tree) را پیاده کنیم. هدف این است که برای یک گراف بدون جهت، درختی را پیدا کنیم که تمام گره‌های گراف را شامل شود و مجموع وزن یال‌های این درخت کمترین مقدار بین تمام درخت‌های ممکن با این دو ویژگی باشد.

#### ۲.۳ MST

در اینجا فرض اولیه این است که هر گره وزن یال‌های مجاور خود را می‌داند. همچنین هر گره دارای یک `<uid>` است و هیچ دو وزن یالی مانند هم نیست. در نتیجه درخت خروجی برای این الگوریتم درختی منحصر به فرد است. پایان الگوریتم زمانیست که هر گره تمام یال‌های متصل به خود را در دو گروه "یال‌های جزو درخت" و "یال‌های خارج درخت" قرار داده باشد. در ادامه توصیه می‌شود که بخش 15.5 از کتاب مرجع را به دقت مطالعه کنید.

#### ۳.۳ گروه بندی یال‌ها

در هر گره، ابتدا تمام یال‌های مجاورش را به صورت basic تعریف می‌کنید. زمانی که مشخص شود این یال جزو MST نیست، این یال را به گروه rejected که همان گروه "یال جزو درخت نیست" اضافه می‌کنیم. همچنین این یال را از گروه basic حذف می‌کنیم. زمانی که مشخص شد یالی جزو درخت MST است، آن را به گروه branch اضافه می‌کنیم و از گروه basic حذف می‌کنیم. در پایان اجرای الگوریتم، هر گره باید تمام یال‌های موجود در گروه basic را به یکی از گروه‌های rejected یا branch تخصیص داده باشد.

## ۴.۳ پیام‌ها

پیام‌های زیر را تعریف می‌کنیم. علاوه بر پیام‌های زیر می‌توانید در صورت نیاز یک یا دو پیام دیگر با کاربردی که خود تعیین می‌کنید تعریف کنید.

## initiate ۱.۴.۳

این پیام توسط رهبر هر کامپوننت فرستاده می‌شود. این پیام به گره‌ها می‌گویند که عمل پیدا کردن MWOE<sup>۲</sup> خود را شروع کنند. همچنین مقادیر core و level نیز توسط این پیام به تمام اعضای کامپوننت فرستاده می‌شود.

```
{"uid": <node_id>, "type": "initiate", "value": <(core, level)>}
```

## report ۲.۴.۳

بعد از اینکه هر گره MWOE خود و تمام descendant‌های خود در کامپوننت را فهمید، مقدار کمینه این مقادیر به همراهی که این پیام طی می‌کند را به گره parent خود (گره‌ای که پیام initiate را داده است) می‌فرستد. به طور مثال اگر گره ۵ بعد از دریافت پیام report از تمام descendant‌هایش متوجه شد که گره ۶ کمترین مقدار MWOE را داراست، در صورتی که پیام initiate از گره ۴ به گره ۵ رسیده باشد، گره ۵ پیام report را با مقادیر "6 5" path = و  $mwoe = MOWE_6$  به گره ۴ می‌فرستد.

```
{"uid": <node_id>, "type": "report", "value": <(path, mwoe)>}
```

ساختار متغیر path را به دلخواه تعیین کنید. اما دقت کنید که مسیری که از رهبر به گره دارای

<sup>2</sup>Minimum Weight Outgoing Edge

MWOE است باید مشخص باشد و پیام `change_root` نباید به صورت `broadcast` فرستاده شود.

test ۳.۴.۳

زمانی که گره  $i$  دنبال mwoe خود می‌گردد، به تمام یال‌های موجود در گروه basic خود پیام test را می‌فرستد. جواب این پیام مشخص می‌کند که هر گره  $j$  در طرف دیگر یال با گره  $i$  در یک کامپوننت است یا خیر.

```
{"uid": <node_id>, "type": "test", "value": <(core, mwoe)>}
```

accept ۴.۴.۳

در صورتی که گره  $i$  با لول بزرگتر یا مساوی گره  $j$ ، پیام test را از گره  $j$  دریافت کند، در صورتی که در یک کامپوننت نبودند (مقدار level و core متفاوتی داشتند) این پیام را به گره  $j$  ارسال می‌کند.

```
{"uid": <node_id>, "type": "accept", "value": <(core, level)>}
```

reject ۵.۴.۳

زمانی که گره  $i$  پیام تست را از گره  $j$  که در یک کامپوننت با گره  $i$  است دریافت می‌کند این پیام را به گره  $j$  می‌فرستد.

```
{"uid": <node_id>, "type": "reject", "value": None}
```

change\_root ۶.۴.۳

زمانی که گره رهبر مقدار MWOE کل کامپوننت را تعیین کرد، در صورتی که این مقدار یال مجاور خودش نباشد، این پیام را با استفاده از مقدار path دریافتی توسط پیام report به گره دارای MWOE

می‌فرستد. این پیام به گره دارای MWOE می‌گوید که به کامپوننت مجاور MWOE وصل شود.

```
{"uid": <node_id>, "type": "change_root", "value": path}
```

connect ۷.۴.۳

این پیام توسط گره متصل به یال MWOE به گره طرف دیگر این یال فرستاده می‌شود. در صورتی که پیام connect از هر دو سر این یال بگذرد عمل merge اتفاق می‌افتد. در غیر این صورت عمل absorb اتفاق می‌افتد.

```
{"uid": <node_id>, "type": "change_root", "value": <(core,level)>}
```

end ۸.۴.۳

زمانی که گره رهبر مقدار MWOE را بی‌نهایت تشخیص داد، این پیام را از طریق درخت MST به تمام گره‌های دیگر Broadcast می‌کند. هر گره بعد از Broadcast کردن پیام end توسط خودش، پایان برنامه خودش را اعلام می‌کند و می‌ایستد. بعد از دریافت پیام end توسط تمام گره‌ها، الگوریتم پایان می‌یابد.

```
{"uid": <node_id>, "type": "end", "value": None}
```

۵.۳ حالت‌های هر گره

برای راحتی کار برای هر گره ۳ حالت را تعریف می‌کنیم: sleep، search و found. نحوه‌ی تغییر بین این حالت‌ها را به دلخواه خود تعریف کنید. اما به طور کلی زمانی که یک گره در حال یافتن MWOE است در حالت search و زمانی که پیام report را به گره parent خود می‌فرستد، به حالت found می‌رود. همچنین در ابتدا همه‌ی گره‌ها در حالت sleep قرار دارند.

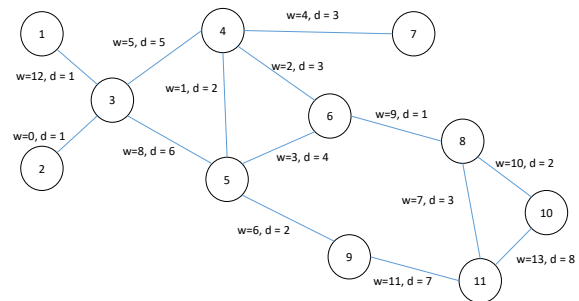


## ۶.۳ ورودی و خروجی

بخشی از فایل وردی که در input.txt قرار دارد، در اینجا آورده شده است. گرافی که در این فایل قرار دارد همان گرافی است که در اسلاید های آموزشی شما موجود است. هر خط از ورودی به شکل زیر است. دقت شود که گراف بدون جهت است و هر دو گره دو سر هر یال باید وزن و تاخیر یال را بدانند.

```
{<head_node> <tail_node> <weight> <delay>}
```

```
1 3 12 1
2 3 0 1
3 4 5 5
3 5 8 6
```



## ۷.۳ نکات نهایی

سعی کنید که قبل از پیاده سازی به طور کامل الگوریتم را بفهمید، حالت های مختلف را بررسی کنید. کد را به صورت کلاس بندی شده بنویسید و سعی کنید در زمان نوشتن برای هر تابع توضیحات مناسبی بنویسید. همراه با کد گزارشی را نیز بنویسید و تمام حالت های خاص که ممکن است اتفاق بیفتد را در آن ذکر کنید. همچنین مشکلاتی که در حین پیاده سازی با آن مواجه شده‌اید را نیز در آن ذکر کنید. پیشنهاد می‌شود که برای خوانایی کد ماژول pylint را نصب و نحوه‌ی استفاده از آن را مشاهده کنید. این ماژول مشکلات کدی که زدید مانند نبود کامنت یا متغیرها با نامگذاری اشتباه را به شما می‌گوید. همچنین پیشنهاد می‌شود این لینک را نیز مشاهده کنید. علت تاکید بر خوانایی کد به این دلیل است که حجم کد ممکن است زیاد شود و برای رفع باگ به مشکل بر بخورید.