

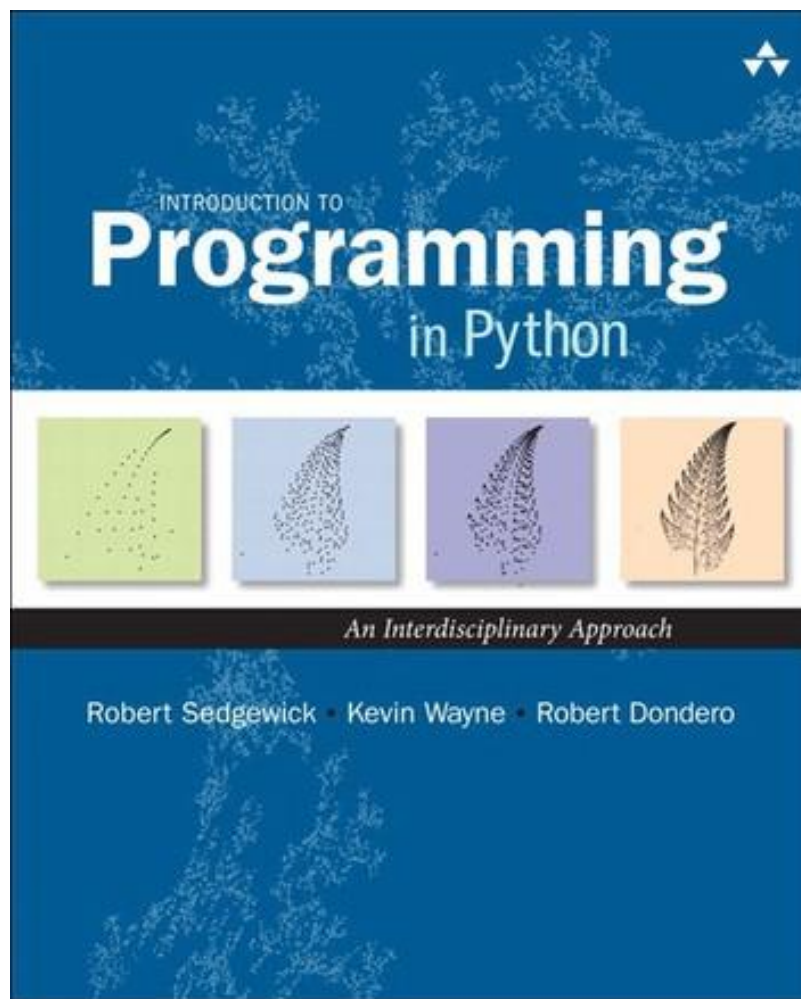
عناصر برنامه نویسی: کاوشگر تصادفی وب

سید ناصر رضوی www.snrazavi.ir

۱۳۹۸

۱-۶ کاوشگر وب

۲



جستجوی وب

۳

- مرتبط بودن. آیا سند مفروض شبیه واژه‌های پرس و جو است؟
- اهمیت. آیا سند مفروض برای کاربرهای گوناگون مفید است؟



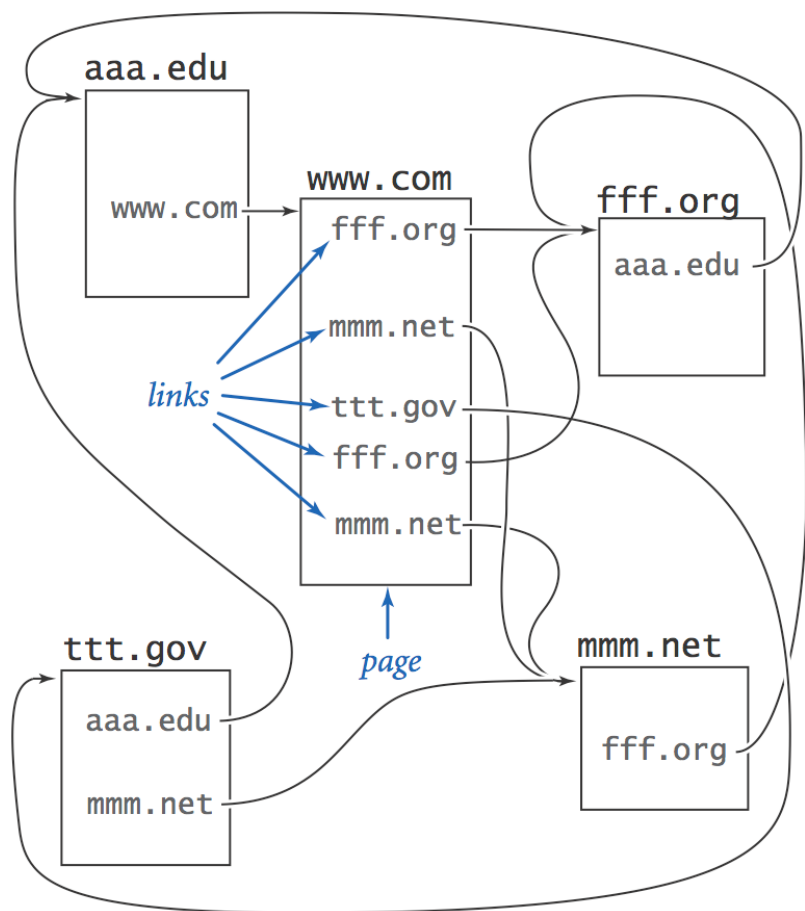
الگوریتم PageRank

۴

□ الگوریتم PageRank گوگل. [سرچی برین و لری پیج، ۱۹۹۸]

□ اندازه‌گیری محبوبیت صفحات بر اساس ساختار پیوندی وب.

□ یک انقلاب در زمینه دسترسی به اطلاعات.



- **مدل.** انتخاب صفحه بعدی به وسیله کاوشگر وب:
 - ۹۰ درصد مواقع، کاوشگر بر روی یک **پیوند تصادفی** کلیک می کند.
 - ۱۰ درصد مواقع، کاوشگر به یک **صفحه کاملاً تصادفی** می رود.
- **هشدارها.** یک مدل خام اما مفید برای کاوش صفحات وب.
 - هیچ کس پیوندها را با احتمال تصادفی انتخاب نمی کند.
 - قاعده ۹۰-۱۰ تنها یک حدس است.
 - این مدل دکمه بازگشت به صفحه قبل و همچنین صفحات نشانه گذاری شده را در نظر نمی گیرد.
 - تنها می تواند از عهده یک نمونه بسیار کوچک از وب برآید.
 - ...

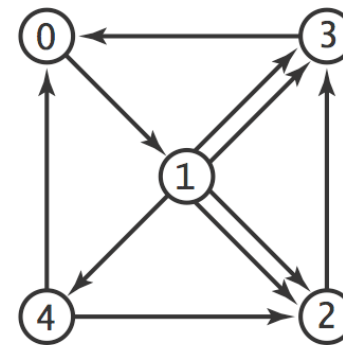
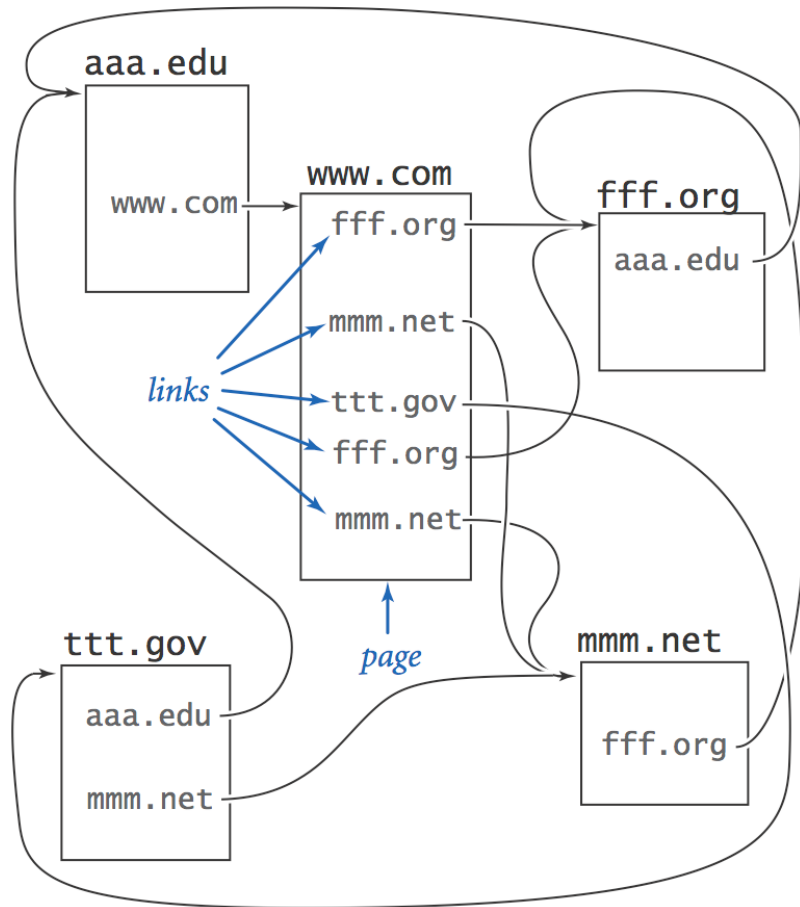
قالب ورودی برای گراف وب

۶

□ قالب ورودی.

□ N صفحه با شماره‌های 0 تا $N - 1$

□ نمایش هر پیوند به وسیله یک زوج از اعداد صحیح



% more tiny.txt

5 ← N

0 1

1 2 1 2

1 3 1 3 1 4

2 3

3 0

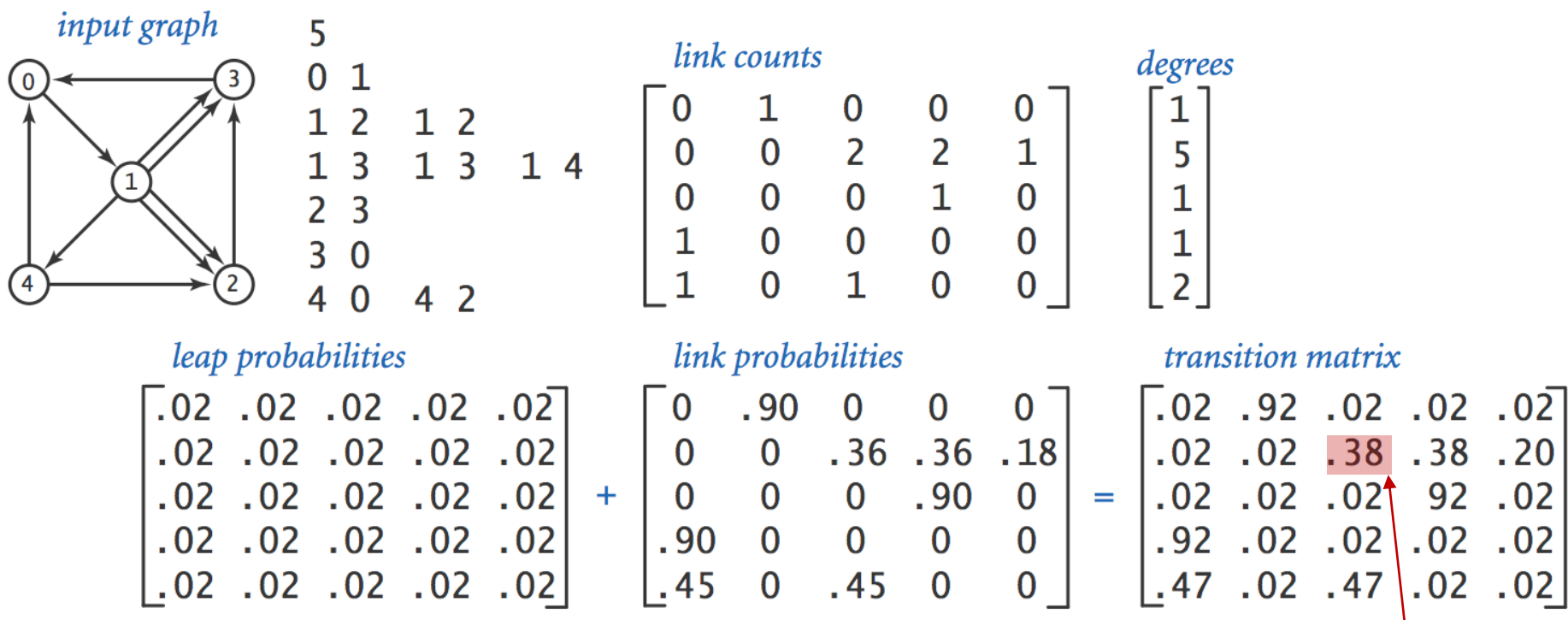
4 0 4 2

← links

ماتریس انتقال

۷

□ ماتریس انتقال. $p[i][j]$ = احتمال رفتن از صفحه i به صفحه j



احتمال رفتن از صفحه ۱ به صفحه ۲

محاسبه ماتریس انتقال از روی گراف وب

۸

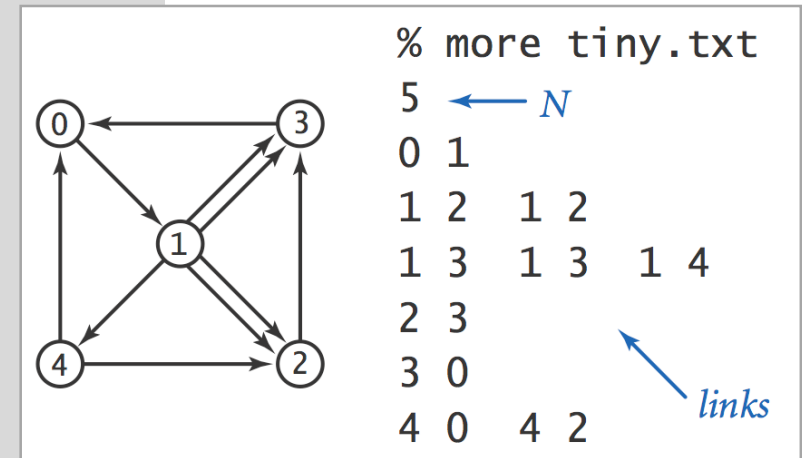
```
import stdio

n = stdio.readInt()
link_counts = [[0] * n for i in range(n)]
out_degrees = [0] * n

while not stdio.isEmpty():
    i = stdio.readInt()
    j = stdio.readInt()
    out_degrees[i] += 1
    link_counts[i][j] += 1

print('%d %d' % (n, n))

for i in range(n):
    for j in range(n):
        p = .9 * (link_counts[i][j] / out_degrees[i]) + .1 / n
        print('%8.5f' % p, end=' ')
    print()
```



```
% python transition.py < tiny.txt
5 5
0.02000 0.92000 0.02000 0.02000 0.02000
0.02000 0.02000 0.38000 0.38000 0.20000
0.02000 0.02000 0.02000 0.92000 0.02000
0.92000 0.02000 0.02000 0.02000 0.02000
0.47000 0.02000 0.47000 0.02000 0.02000
```

قاعده ۹۰-۱۰

شبیه‌سازی مونت کارلو

۹

□ شبیه‌سازی مونت کارلو.

□ کاوشگر از صفحه صفر شروع می‌کند.

□ به صورت تکراری و بر اساس ماتریس انتقال صفحه بعدی را انتخاب می‌کند.

□ سپس محاسبه می‌کند از هر صفحه به چه میزان بازدید شده است.

چگونه؟ اسلاید بعد را ببینید



.02	.92	.02	.02	.02
.02	.02	.38	.38	.20
.02	.02	.02	.92	.02
.92	.02	.02	.02	.02
.47	.02	.47	.02	.02

ماتریس انتقال

کاوشر تصادفی

۱۰

□ حرکت تصادفی. فرض کنید کاوشگر در صفحه page است؟ صفحه بعدی j چگونه انتخاب می شود؟

□ سطر page از ماتریس انتقال شامل احتمالات مربوط به این صفحه است.

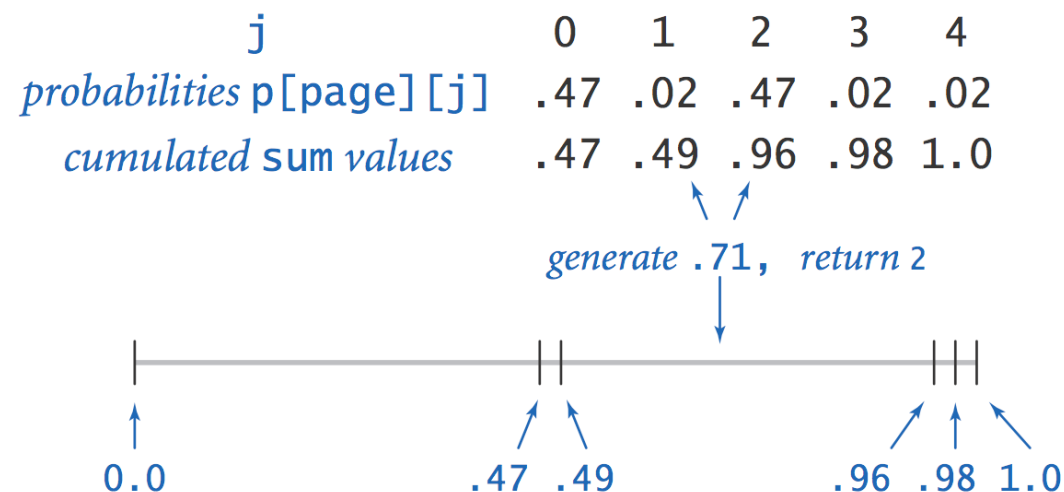
□ احتمال های **تجمعی** برای این صفحه محاسبه می شود.

□ یک عدد تصادفی مانند r در بازه 0.0 و 1.0 تولید می شود.

□ صفحه بعدی j بر اساس بازه های که r در آن قرار گرفته انتخاب می شود.

	j				
	probabilities p[page][j]				
	cumulated sum values				
page	.02	.92	.02	.02	.02
	.02	.02	.38	.38	.20
	.02	.02	.02	.92	.02
	.92	.02	.02	.02	.02
	.47	.02	.47	.02	.02

ماتریس انتقال



کاوشگر تصادفی

۱۱

□ حرکت تصادفی. فرض کنید کاوشگر در صفحه `page` است؟ صفحه بعدی `j` چگونه انتخاب می‌شود؟

□ سطر `page` از ماتریس انتقال شامل احتمالات مربوط به این صفحه است.

□ احتمال‌های **تجمعی** برای این صفحه محاسبه می‌شود.

□ یک عدد تصادفی مانند `r` در بازه `0.0` و `1.0` تولید می‌شود.

□ صفحه بعدی `j` بر اساس بازه‌ای که `r` در آن قرار گرفته انتخاب می‌شود.

```
r = random.random()
```

```
total = 0.0
```

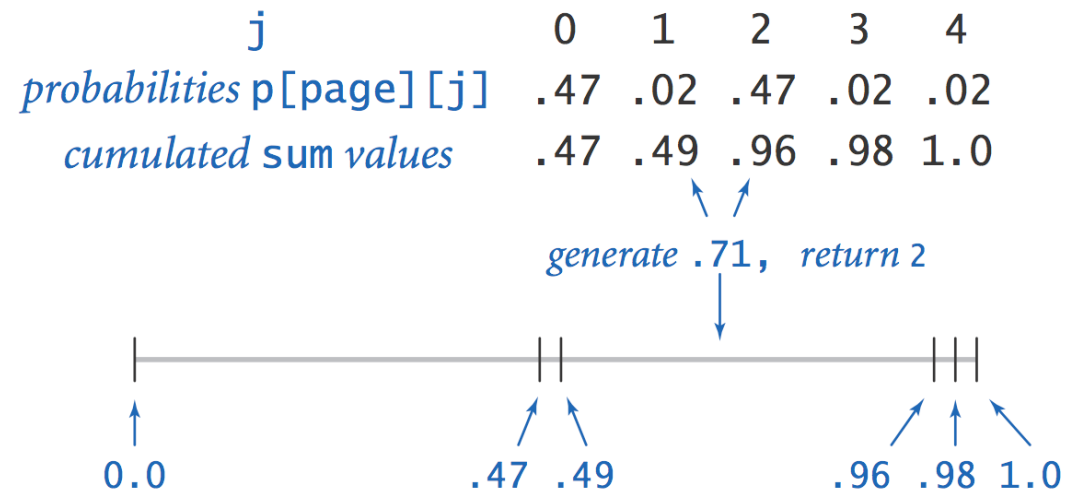
```
for j in range(n):
```

```
    total += p[page][j]
```

```
    if r < total:
```

```
        page = j
```

```
        break
```



كاوشگر وب: شبیه‌سازی مونت کارلو (۱)

۱۲

```
import sys
import random
import stdio

moves = int(sys.argv[1])

n = stdio.readInt()
stdio.readInt()

p = [[0.0] * n for i in range(n)]

for i in range(n):
    for j in range(n):
        p[i][j] = stdio.readFloat()
...

```

```
% python transition.py < tiny.txt
5 5
0.02000 0.92000 0.02000 0.02000 0.02000
0.02000 0.02000 0.38000 0.38000 0.20000
0.02000 0.02000 0.02000 0.92000 0.02000
0.92000 0.02000 0.02000 0.02000 0.02000
0.47000 0.02000 0.47000 0.02000 0.02000

% python transition.py < tiny.txt | randomsurfer.py 1000000

```

کاوشگر وب: شبیه‌سازی مونت کارلو (۲)

۱۳

```
hits = [0] * n
page = 0
for i in range(moves):
    r = random.random()
    total = 0.0
    for j in range(n):
        total += p[page][j]
        if r < total:
            page = j
            break
```

```
hits[page] += 1
```

```
for v in hits:
    print('%8.5f' % (1. * v / moves), end='')
```

```
% python transition.py < tiny.txt
```

```
5 5
```

```
0.02000 0.92000 0.02000 0.02000 0.02000
0.02000 0.02000 0.38000 0.38000 0.20000
0.02000 0.02000 0.02000 0.92000 0.02000
0.92000 0.02000 0.02000 0.02000 0.02000
0.47000 0.02000 0.47000 0.02000 0.02000
```

```
% python transition.py < tiny.txt | randomsurfer.py 1000000
```

```
0.27480 0.26180 0.14860 0.25030 0.06450
```

حرکت از صفحه `page` به صفحه `j`

همگرایی: کمی ریاضیات

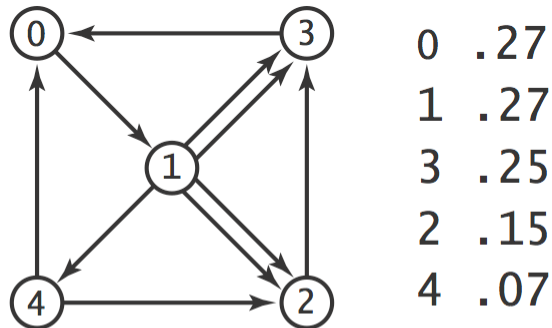
۱۴

□ همگرایی. برای مدل کاوشگر تصادفی، مدت زمانی که کاوشگر بر روی هر صفحه صرف می کند، به یک **توزیع منحصر به فرد** همگرا می شود که مستقل از صفحه شروع است.

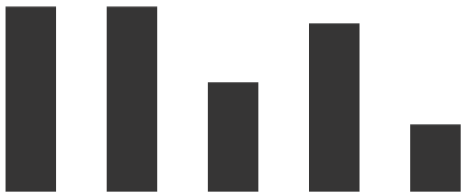
رتبه صفحه

توزیع ایستای زنجیره مارکوف

بردارهای ویژه اصلی ماتریس انتقال



$$\begin{bmatrix} \frac{428,671}{1,570,055} & \frac{417,205}{1,570,055} & \frac{229,519}{1,570,055} & \frac{388,162}{1,570,055} & \frac{106,498}{1,570,055} \end{bmatrix}$$



همگرایی: کمی ریاضیات

۱۵

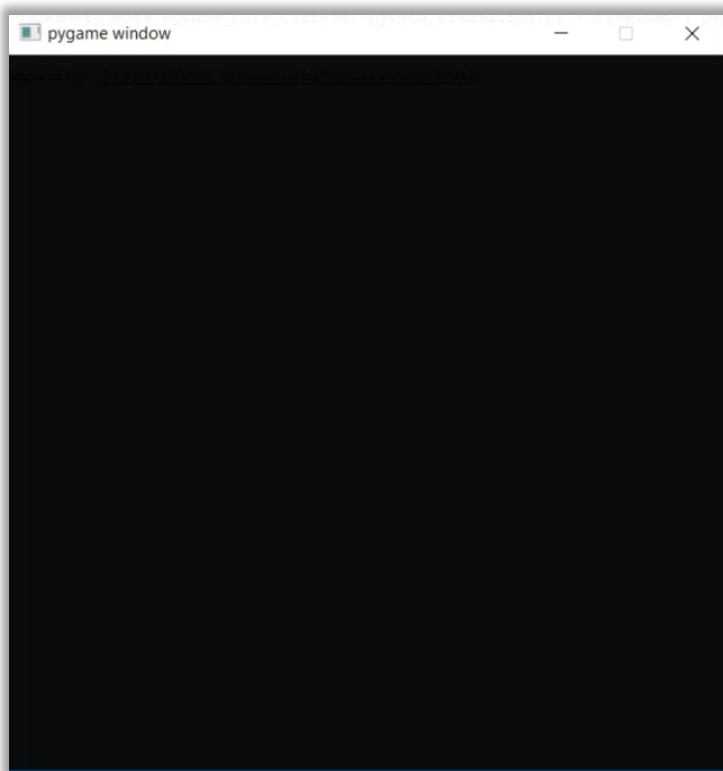
□ همگرایی. برای مدل کاوشگر تصادفی، مدت زمانی که کاوشگر بر روی هر صفحه صرف می کند، به یک توزیع منحصر به فرد همگرا می شود که مستقل از صفحه شروع است.

رتبه صفحه

توزیع ایستای زنجیره مارکوف

بردارهای ویژه اصلی ماتریس انتقال

$$\begin{bmatrix} \frac{428,671}{1,570,055} & \frac{417,205}{1,570,055} & \frac{229,519}{1,570,055} & \frac{388,162}{1,570,055} & \frac{106,498}{1,570,055} \end{bmatrix}$$



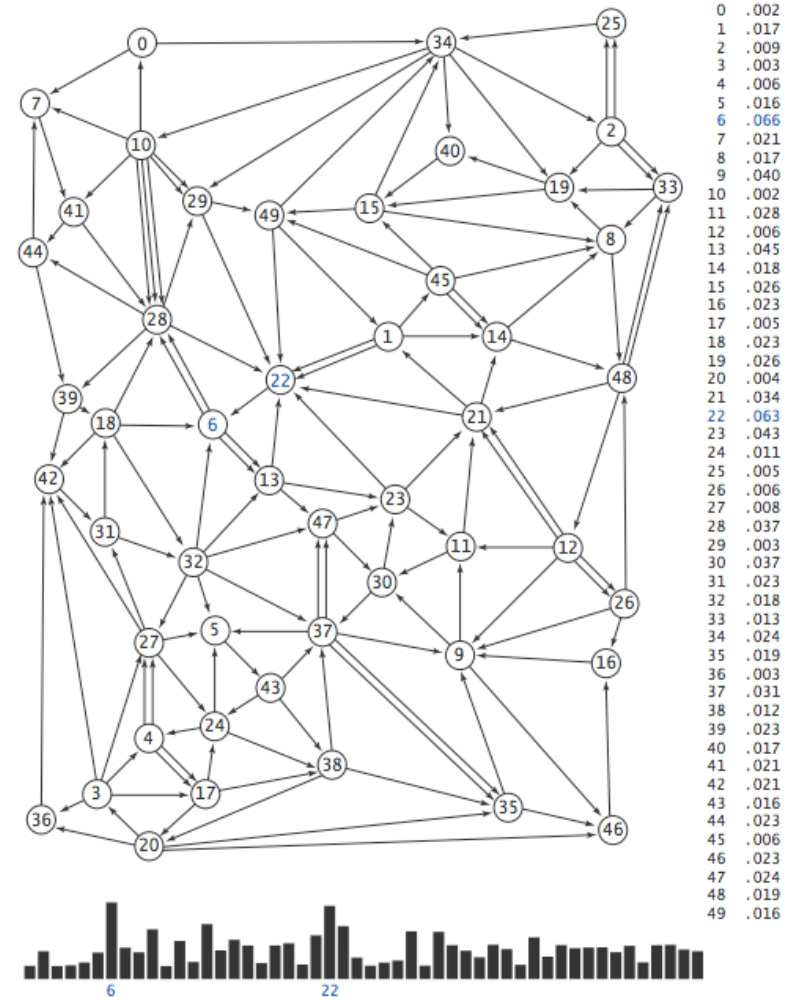
یک گراف بزرگ تر

۱۶

% more medium.txt

50

0	7	0	34				
1	14	1	22	1	22	1	45
2	19	2	25	2	33		
3	4	3	17	3	27	3	36
4	17	4	17	4	27	4	27
5	43						
6	13	6	13	6	28		
7	41						
8	19	8	48				
9	11	9	30	9	46		
...							



Page ranks with histogram for a larger example