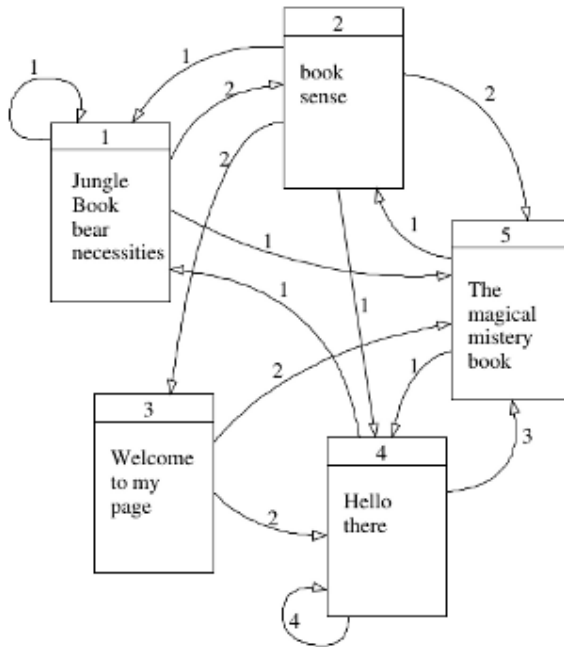


## Google Nasil Isler?

Lineer Cebir hocalari Google'a mutessekkir olmalı, cunku bu unlu arama motorunun kullandigi PageRank tekniginin ozu aslinda lineer cebirin temelini olusturan kavramlardan ozdeger / ozvektor hesabi. Ogrenciler "niye bu kavramlari ogreniyoruz hocam?" diye sorunca artik cevap kolay: "bu yontemi Google da kullaniyor!".

Simdi arama motorunun yapmasi gerekenleri dusunelim: Google'a bir kelime yazdigimizda geri gelen sonuclar nasil kararlaltirilir? İlk akla gelen yontem tabii ki Web'deki tum sayfalarin (milyarlarca sayfa) sayfalar izerindeki kelimelerin o sayfa ile iliskilendirilmesi ve arama yapilınca kelimeye gore sayfa geri getirilmesi. Mesela alttaki ornekte "book (kitap)" yazınca geriye 1., 2. ve 5. sayfalar geri gelecek. Fakat hangi sirada? Bu sayfalardan hangisi digerlerinden daha onemli?



PageRank'in temelinde daha fazla referans edilen sayfalarin daha ustte cikmasi yatar. Hatta o referans eden sayfalarin kendilerine daha fazla referans var ise bu etki ta en sondaki sayfaya kadar yansitilir, hatta bu zincir bastan sona her seviyede hesaplanabilir. Peki bu nasil gerceklestirilir?

PageRank Web sayfalarini bir Markov Zinciri olarak gorur. Markov Zincirleri seri halindeki  $X_n, n = 0, 1, 2, \dots$  rasgele degiskenini modeller ve bu degiskenler belli sayidaki konumların birinde olabilirler. Mesela konumları bir dogal sayı ile ilintilendirirsek  $X_n = i$  olabilir ki  $i = \{0, 1, \dots\}$  diye kabul edelim.

Markov Zincirlerinde (MZ)  $i$  konumundan  $j$  konumuna gecis olasiligini,  $P_{ij}$ , biliriz ve bu  $P(X_{n+1} = j | X_n = i)$  olarak acilabilir. Acilimdan gorulecegi uzere bir MZ sonraki adima gecis olasiligi icin sadece bir onceki adima bakar. Bu tur once/sonra yapısındaki iki boyutlu hal, cok rahat bir sekilde matrisine ceviri-

ilebilir / gosterilebilir. Onceki konum satirlar, sonraki konum kolonlar olarak betimlenir mesela.

### Ornek

Bir sonraki gunde yagmur yagmayacagini bir MZ olarak tasarlayalim. Bir sonraki gunde yagmur yagmayacagini sadece bugun etkiliyor olsun. Eger bugun yagmur yagiyorsa yarin yagmur yagmasi 0.7, eger bugun yagmiyor ise yarin yagmasi 0.4. MZ soyle

$$P = \begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix}$$

Gecis olasiliklarindan bahsettigimize gore ve elimizde sinirli / belli sayida konum var ise, bir MZ'nin her satirindaki olasiliklarin toplami tabii ki 1'e esit olmalidir.

MZ'lerin ilginç bir özelliği  $n$  adım sonra  $i, j$  gecisinin  $P^n$  hesabıyla yapılabilmesidir. Yani  $P$ 'yi  $n$  defa kendisiyle carpip  $i, j$  kordinatına bakarsak  $n$  adım sonrasını rahatca gorebiliriz. Bunun ispatını burada vermeyeceğiz.

Mesela üstteki örnekte, eger bugun yagmur yagiyorsa 4 gün sonra yagmur yagma olasılığı nedir?

```
import numpy.linalg as lin
P = np.array([[0.7,0.3],[0.4,0.6]])
P4 = lin.matrix_power(P,4)
print P4

[[ 0.5749  0.4251]
 [ 0.5668  0.4332]]
```

Aradığımız gecis için kordinat 0,0'a bakıyoruz ve sonuc 0.5749. Numpy `matrix_power` bir matrisi istedigimiz kadar kendisiyle carpmamızı sağlıyor.

### Duragan Dagilim (Stationary Distribution)

Eger yagmur örneğindeki matrisi carpmaya devam edersek, mesela 8 kere kendisiyle carpsak sonuc ne olurdu?

```
import numpy.linalg as lin
P = np.array([[0.7,0.3],[0.4,0.6]])
P8 = lin.matrix_power(P,8)
print P8

[[ 0.57145669  0.42854331]
 [ 0.57139108  0.42860892]]
```

Dikkat edilirse, her satir bir degere yaklasmaya basladi. Bu deger MZ'nin duragan dagilimidir, belli kosullara uyan her MZ'nin boyle bir duragan dagilimi vardır. Bu kosullar MZ'nin periyodik olmayan (aperiodic) ve tekrar eden (recurrent) olmasıdır. Bu sartlar çok “özel” sartlar degildir aslında, daha çok “normal” bir MZ'yi tarif ediyor diyebiliriz. Tüm konumları tekrar eden yapmak kolaydır, MZ

tek bagli (singly connected) hale getirilir, yani her konumdan her diger konuma bir gecis olur, ve periyodik olmamasi icin ise MZ'ye olmadigi zamanlarda bir konumdan kendisine gecis saglanir (az bir gurultu uzerinden).

Nihayetinde duraganlik su denkleme sebebiyet verir,

$$\pi = \pi P$$

Burada duragan dagilim  $\pi'$  dir. Bu denklem tanidik geliyor mu? Devrigini alarak soyle gosterelim, belki daha iyi taninir,

$$P^T \pi^T = \pi^T$$

Bir sey daha ekleyelim,

$$P^T \pi^T = 1 \cdot \pi^T$$

Bu ozdeger/vektor formuna benzemiyor mu? Evet! Bu form

$$Ax = \lambda x$$

MZ denklemi sunu soyluyor, 1 degerindeki ozdegere ait ozvektor bir MZ'nin duragan dagilimidir! Bu arada, MZ gecis matrisi P'nin en buyuk ozdegerinin her zaman 1 oldugunu biliyoruz (cunku ustteki tarif ettigimiz ozel sartlara sahip olan turden matrisler boyle ozdegerlere sahip olmalidir). Bu durumda en buyuk ozdegere ait ozvektoru hesaplamak yeterli olacaktır. Bunu yapmayi zaten *Lineer Cebir Ders 21'*de ogrenmistik, ust metot (power method) sayesinde bu hesap kolayca yapilabiliyor.

Simdi en basktaki Web sayfalarina ait gecisleri yazalim,

```
P = [[1./4, 2./4, 0, 0, 1./4],
      [1./6, 0, 2./6, 1./6, 2./6],
      [0, 0, 0, 2./4, 2./4],
      [1./8, 0, 0, 4./8, 3./8],
      [0, 1./2, 0, 1./2, 0]]
```

```
P = np.array(P)
print P
```

```
[[ 0.25      0.5      0.      0.      0.25
   0.16666667 0.      0.33333333 0.16666667 0.33333333]
 [ 0.      0.      0.      0.5      0.5
   0.125      0.      0.      0.5      0.375
   0.      0.5      0.      0.5      0.] ]]
```

Simdi ust metodu kullanarak duragan dagilimi hesaplayalim. Herhangi bir baslangic vektorunu P ile 20 kere carpmak yeterli olur.

```
import numpy.linalg as lin
x=np.array([.5, .3, .1, .1, 0]) # herhangi bir vektor
for i in range(20):
    x = np.dot(x,P)
print 'pi = ', x

pi = [ 0.10526316  0.18421053  0.06140351  0.38596491  0.2631579 ]
```

Not: Aslında cebirsel olarak  $P$ 'yi 20 kere kendisiyle carpmak ve sonucu baslangic vektoru ile bir kere carpmak ta dusunulebilirdi. Fakat 20 kere vektor / matris carpimlari yapmak, 20 kere matris / matris carpimi yapmaktan daha verimli olacaktir. Buyuk Veri ortamı için de bu soylenebilir.

Neyse, eger ozvektor hesabini kendimiz elle yapmak yerine direk kutuphane cagrisi kullansaydik,

```
import numpy.linalg as lin
evals, evec = lin.eig(P.T)
pi = evec[:,0] / np.sum(evec[:,0])
print np.abs(pi)

[ 0.10526316  0.18421053  0.06140351  0.38596491  0.26315789]
```

Ayni sonuca ulastik. Bu sonuc gosteriyor ki “book” yazdigizda Google bize 5. sayfayi en basta olacak sekilde sonuc dondurmeli, cunku onun duragan dagilimi 1,2,5 sayfalarinin arasinda en yuksek.

### Duragan Dagilima Bakis

MZ ve duragan dagilimin PageRank'le alakasini bir daha dusunelim. MZ ile  $n$  adim sonrasini hesaplayabiliyoruz, duragan dagilim ise sonsuz adim sonrasini ifade ediyor. Ve bu dagilim, bir anlamda, sonsuz yapilan adimlar sirasinda *en fazla hangi konumlarda* zaman gecirilecegini gosteriyor. Konum yerine sayfa dersek duragan dagilimin niye en onemli sayfalari belirlemek için onemli oldugunu anlariz.

Kullanici herhangi bir sayfada iken hangi diger sayfalara gidecegi o sayfa uzerinde baglantilar uzerinden anlasilir, PageRank bu baglantinin mevcudiyetine bakar sadece, o mevcudiyet uzerinden bir gecis olasiligi hesaplar, ve bu olasiliga gore (raslantisal sekilde) baglantinin takip edilecegi dusunulur. Bu arada cogunlukla sayfalar arasindaki baglantilarin agirligi 1 olacaktir, fakat ornek amacli 2,3 gibi sayilar da kullaniliyor.

### Rasgele Sayfa Gecisi

Google veri temsili uzerinde bazi ekler yapmaktadir, mesela kullanicinin hicbir baglanti takip etmeyip tarayiciya direk URL girerek baska bir sayfaya ziplamasi (teleporting) bir sekilde temsil edilmelidir. Ayrica hic disa baglanti vermeyen sayfalar (olu noktalar) hesaba katilmalidir. Simdi  $\pi^T$  yerine  $p$ ,  $P$  yerine  $N$  kullanalim, PageRank ozyineli algoritmasi

$$p = N^T p$$

olarak gosterilebilir.

Bu her iki durum icin formül su sekilde ikiye ayirilir,

$$p = (1 - d)N^T p + dN_f^T p$$

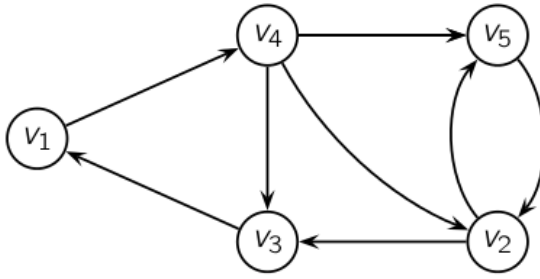
$$= ((1 - d)N^T + dN_f^T)p$$

$$= M^T p$$

ki,

$$M = (1 - d)N^T + dN_f^T$$

olacaktır.  $N_f$  bir normalize edilmiş “ziplama” matrisidir, yani her sayfadan her diger sayfaya bir baglanti “varmis gibi” yapar, mesela 5x5 boyutunda tum ogeleri 0.20 olacaktır.  $d$  bir agirlik sabitidir, Google’in bunu 0.85 olarak tanımladigi duyulmustur, ve gercek baglanti matrisi ve rasgele ziplama matrisi arasinda bir agirlik tanımlar, her ikisinde de birazcik alarak (daha cok ana  $N$ ’den tabii) niahi matrisi olusturur. Ornek olarak su grafige bakalim,



```

N = [[0, 0, 0, 1., 0],
      [0, 0, 1./2, 0, 1./2],
      [1, 0, 0, 0, 0],
      [0, 1./3, 1./3, 0, 1./3],
      [0, 1, 0, 0, 0]]

```

```

N = np.array(N)

```

```

Nf = 0.20 * np.ones((5,5))

```

```

d = 0.85

```

```

M = d*N + (1-d)*Nf

```

```

x=np.array([.5, .3, .1, .1, 0]) # herhangi bir vektor

```

```

for i in range(20):

```

```

    x = np.dot(x,M)

```

```

print 'result = ', x

```

```

result = [ 0.18959094  0.24375097  0.18775335  0.19115138  0.18775335]

```

Sonuca gore  $v_2$  en yuksek PageRank degerine sahip.

#### Mimari

Google tabii ki arama sonuclarini iyilestirmek icin yillar icinde diger ek fonksiyonlari motoruna ekledi. Duyumlarimiza gore artik PageRank gibi onlarca ek kriter kullanilmaktadir; fakat PageRank hala cok onemli ve sirketin kurulusu baglaminda Google'i Google yapan algoritmaydi, onun diger motorlara nazaran elindeki avantajı, en buyuk ilerlemesiydi.

Sistem kodlamasi acisindan PageRank'e tum Web sayfaları ve onların arasindaki iliskiler verilmelidir, bu milyarlarca sayfa ve onların arasindaki baglantılar demektir. Google bunu yapabilmek icin Web "agını" orumcek (spider) programları ile surekli geziyor, ve bu veriyi alip, PageRank'e hesap icin aktariyor.

[1] Murphy, K., CS340: Machine Learning Lecture Notes, [www.ugrad.cs.ubc.ca/~cs340](http://www.ugrad.cs.ubc.ca/~cs340)

[2] Ross, S., Introduction to Probability Models, 8th Edition

[3] Zaki, Mair, Fundamentals of Data Mining Algorithms