

Technische Universiteit Delft
Faculteit Elektrotechniek, Wiskunde en Informatica
Delft Institute of Applied Mathematics

De wiskunde achter Google
Mathematics behind Google

Verslag ten behoeve van het
Delft Institute for Applied Mathematics
als onderdeel ter verkrijging

van de graad van

BACHELOR OF SCIENCE
in
TECHNISCHE WISKUNDE

door

Nana Cao

Delft, Nederland
December 2008



BSc verslag TECHNISCHE WISKUNDE

“De wiskunde achter Google”

“Mathematics behind Google”

Nana Cao

Technische Universiteit Delft

Begeleider

Dr. M.H.A Haase

Overige commissieleden

Dr.

J.G. Spandaw

Dr.

C. Kraaikamp

December, 2008

Delft

De wiskunde achter Google

Inhoudsopgave

1	Inleiding	2
2	Zoekmachines in het algemeen	3
3	Het PageRank algoritme	4
3.1	Het basis idee	5
3.2	Non-Unique Rankings	7
3.3	Dangling node	8
3.4	Verwisselen van index	11
3.5	Een remedie voor $\dim(V_1(A)) > 1$	11
4	Berekenen van de eigenvector en de tweede eigenvector van Google	14
4.1	Machtsmethode	14
4.2	Convergentiesnelheid en de tweede eigenwaarde	15
4.3	Aangepaste matrix M	19
5	De stelling van Perron	21
6	Toekomst van het zoeken	24

1 Inleiding

Het bachelorproject van de bachelor technische wiskunde dient als afsluiting van de bachelorfase. De bedoeling van het Bachelorproject Analyse is om zelfstandig een wiskundige onderwerp te bestuderen door middel van een wiskundige onderzoek of literatuurstudie.

Hoewel men een literatuur onderzoek Analyse vaak als abstract beschouwen, kent de Analyse eigenlijk veel praktische en belangrijke toepassingen. Voor dit project kies ik voor het bestuderen van zoekmachines op het world wide web. Het Internet is immers bijna niet meer weg te denken uit ons dagelijks leven. Het web biedt een enorm schat van informatie, daarmee is het belangrijk snel de juiste informatie te kunnen vinden. Een zoekmachine op het web speelt hierbij een belangrijke rol.

Voor ons is de vraag wat wiskunde of Analyse te maken heeft met zo'n zoekmachine. Het doel van dit project is om een kijkje te nemen achter de schil van de zoekmachine Google.

We zullen eerst uitzoeken hoe Google werkt in het algemeen en daarna het wiskundige model van het zoeken beschouwen. We zullen zien dat Google gebaseerd is op een fundamenteel algoritme dat een eigenvector van een gigantische matrix uitrekent. De wiskunde hierachter zal te maken hebben met een belangrijke stelling van Perron over positieve matrices en hun eigenwaarden.

In dit verslag zullen we het wiskundige model opstellen en het PageRank algoritme van Google zo nauwkeurig mogelijk beschrijven. De wiskunde hierachter met in het bijzonder het bewijs van Perron gaan we zo precies mogelijk uitwerken. Verder zullen de numerieke methode voor het uitrekenen van de oplossing beschrijven en de convergentiesnelheid van het algoritme bestuderen wat te maken zal hebben met de tweede eigenwaarde van Google.

2 Zoekmachines in het algemeen

Zoekmachines zijn al jarenlang een middel om informatie te vinden op het internet. Ze wijzen ons de weg op internet met de ‘kennis’ die ze beschikken over het volledige web. Maar hoe werkt een zoekmachine?

Wat gebruikers vooral zien van een zoekmachine is de grafische schil, de website waar de zoekopdrachten worden ingetikt en de resultaten getoond.

De drie belangrijke onderdelen achter de schil zijn:

1. De spider

De spider (ook wel spin, verkenners of crawler genoemd) is een programma dat zoveel mogelijk pagina's op internet bezoekt, een spider 'leest' webpagina's en volgt de daarop voorkomende links om naar nieuwe pagina's te gaan. De inhoud van alle pagina's wordt 'gelezen'. De tekst, de afbeeldingen, de aangetroffen documenten enzovoorts gaan mee naar de database (het tweede deel van de zoekmachine). De hyperlinks naar andere pagina's of andere sites worden gevolgd om ook die pagina's binnen te halen, enzovoorts.

2. De database

In de database van een zoekmachine wordt de inhoud van de gespiderde webpagina's op een slimme manier opgeslagen. Naast de tekst gaan er zoveel mogelijk additionele gegevens mee. Zoals de datum van creatie, gegevens over kleuren op de pagina, soorten documenten die zijn aangetroffen, enzovoorts. In deze databank kan snel worden gezocht. De meeste wachttijd gaat verloren met het transport van de gegevens van en naar de computer van de gebruiker. Een zoektochtje bij Google naar een combinatie van twee veelvoorkomende woorden in 339 miljoen documenten duurt minder dan een halve seconde.

3. De rangschikking van resultaten (ranking)

Het rangschikken van de resultaten is het derde belangrijke onderdeel van een zoekmachine. Een goede rangschikking is heel belangrijk zodat de beste zoekresultaten op de eerste paar pagina's komen te staan. Het heeft geen zin om door tientallen pagina's met zoekresultaten te moeten bladeren.

De selectiecriteria om te bepalen welke link waardevol voor de gebruiker is waren vroeger vrij eenvoudig. Webmasters konden met behulp van zogenaamde meta-tags onder meer omschrijvingen en keywords toevoegen aan hun pagina's. Zoekmachines keken simpelweg welke meta-tags overeenkwamen met de zoekopdracht. Maar van dit systeem werd er veel misbruik van

gemaakt. Doordat informatie kan worden toegevoegd die voor de bezoeker niet zichtbaar is, kon men populaire keywords toevoegen die feitelijk niets met de website te maken hebben. Tegenwoordig beoordelen zoekmachines alleen nog de inhoud die de bezoeker daadwerkelijk te zien krijgt. Keywords worden dan voornamelijk uit de teksten en titels gehaald.

De precieze selectiecriteria van zoekmachines blijft meestal geheim, niet alleen uit concurrentie-overwegingen, maar ook om te voorkomen dat webmasters hun pagina's zo inrichten dat ze ongeacht de inhoud van die pagina altijd bovenaan komen te staan.

Uniek aan Google is het algoritme dat PageRank wordt genoemd. Daarbij wordt behalve de inhoud ook gekeken naar het aantal links van en naar een pagina. Verder is er een reeks andere elementen die de relevantie bepalen.

Zo speelt de titel van een pagina een grote rol bij het plaatsen op een ranglijst. Daarnaast zijn woorden die in het begin van een document worden gevonden veel belangrijker dan woorden aan het eind van een pagina. Andere factoren zijn bijvoorbeeld :

1. De lettergrootte van een woord, hoe groter de lettergrootte hoe belangrijker.
2. Het aantal malen dat een woord voorkomt (woordfrequentie) en de woordafstand tussen twee of meerdere gezochte woorden (woord *proximity*);
3. De woordlengte van een pagina (echter, erg korte en erg lange pagina's krijgen weer een lagere beoordeling).
4. De frequentie dat de inhoud wordt ververs: hoe vaker je je pagina ververs met nieuwe inhoud hoe hoger de score.

3 Het PageRank algoritme

Eind jaren 90 kwam het concept van link analysis op, volgens dit principe wordt een website belangrijker als er door veel andere websites naar wordt verwezen. Google werd groot door handig van dit concept gebruik te maken. In september 1998 werd Google opgericht door twee promovendi aan de Stanford-universiteit, Larry Page en Sergey Brin. Door de unieke zoektechnologie groeide het bedrijf razendsnel, en werd al binnen een paar jaar de meest gebruikte zoekmachine op het web.

Uniek aan Google was het algoritme dat PageRank wordt genoemd, een methode dat in staat is elke pagina in het web een waarde toe te kennen aan de hand van zijn relevantie.

3.1 Het basis idee

We bekijken bij het algoritme eerst het aantal links dat naar een pagina verwijst. Een link dat is gemaakt op een pagina heet een **backlink** van dat pagina. Het is als het ware een verkiezing waar een webpagina op een ander pagina kan stemmen, de hoogte van het aantal stemmen geeft de waarde van een pagina aan. Een waardetoekenning is dus altijd niet negatief. Bekijk een web met n pagina's, zij x_k de waarde van pagina k . Dan wordt de belangrijkheidscore van dat pagina berekend met de volgende formule

$$x_k = \sum_{j \in L_k} \frac{x_j}{n_j},$$

hier is $L_k \subset \{1, 2, \dots, n\}$ de verzameling van pagina's bijbehorend bij de backlinks van pagina k en n_j het aantal links dat pagina j maakt in het web. De formule is afgeleid met de volgende aannames:

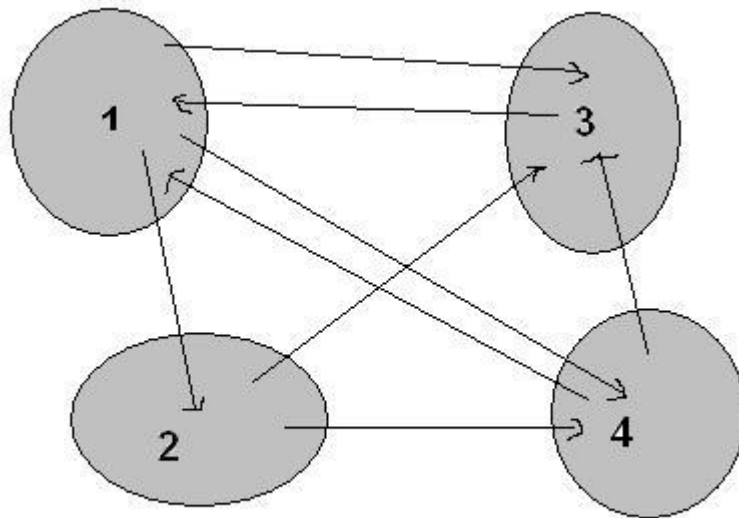
1. Een link van een pagina op zichzelf zal niet worden meegeteld.
2. Een link gemaakt door een belangrijkere pagina heeft een grotere wegingsfactor.
3. Als een webpagina alleen naar jouw pagina verwijst, is dit van meer waarde dan wanneer die ook naar (veel) andere pagina's verwijst.

We krijgen nu dus een stelsel van n lineaire vergelijkingen dat geschreven kan worden als $Ax = x$ met $x = [x_1 \ x_2 \ \dots \ x_n]^T$, en A wordt de matrix met alle wegingen van de pagina's van de volgende vorm:

$$A = \begin{pmatrix} 0 & A_{12} & A_{13} & \dots & A_{1n} \\ A_{21} & 0 & A_{23} & \dots & A_{2n} \\ A_{31} & A_{32} & 0 & \dots & \vdots \\ A_{41} & A_{42} & A_{43} & \ddots & \vdots \\ \vdots & \dots & \ddots & \ddots & \vdots \\ A_{n1} & A_{n2} & \dots & A_{n(n-1)} & 0 \end{pmatrix}$$

met A_{ij} gelijk aan de wegingsfactor van pagina j op pagina i , oftewel $\frac{1}{n_j}$ als het gelinkt is, anders 0.

We passen dit toe op een voorbeeld: Gegeven is een simpel web, dat wordt voorgesteld door de volgende graaf. Elke cirkel is een pagina in het web, en elke pijl is een link. Zo is er een link van pagina 1 naar 2, maar niet de andere kant op.



Voor pagina 1 bijvoorbeeld krijgen we: $x_1 = x_3/1 + x_4/2$ omdat pagina 3 en 4 backlinks zijn van pagina 1 en pagina 3 in totaal maar 1 link heeft gemaakt en pagina 4 twee. Door analoge afleiding krijgen $x_2 = x_1/3$, $x_3 = x_1/3 + x_2/2 + x_4/2$ en $x_4 = x_1/3 + x_2/2$. Deze lineaire vergelijkingen kunnen we schrijven als $Ax = x$ met $x = [x_1, x_2, x_3, x_4]^T$ en

$$A = \begin{pmatrix} 0 & 0 & 1 & \frac{1}{2} \\ \frac{1}{3} & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{3} & \frac{1}{2} & 0 & 0 \end{pmatrix}$$

Wij zoeken dus eigenlijk een niet-negatieve eigenvector x met eigenwaarde 1 voor de matrix A . In het algemeen geldt dat de matrix van elke web zonder **danglingnode** (web pagina zonder uitgaande links) 1 als een eigenwaarde moet hebben. Dit kunnen we laten zien met behulp van de *kolomstochastische* eigenschap van de matrix A .

We definiëren het volgende:

3.1 Definitie Een vierkante matrix A is **kolomstochastisch** als alle elementen van A niet negatief zijn en de som van de elementen van elke kolom gelijk is aan 1.

De matrix A voor een web zonder dangling node is kolomstochastisch, er geldt de volgende:

3.2 Propositie *Elke kolomstochastische matrix heeft 1 als een eigenwaarde en behoort tot een niet negatieve eigenvector (ook wel de Perron vector genoemd).*

We gaan nu laten zien dat 1 een eigenwaarde moet zijn van een kolomstochastische matrix.

Bewijs. Zij A een $n \times n$ kolomstochastische matrix en $\mathbf{e} = [1, \dots, 1]^T$ de n -dimensionale kolomvector met all zijn elementen gelijk aan 1. Omdat $\det(A - \lambda I) = \det(A - \lambda I)^T = \det(A^T - \lambda I)$ hebben A en A^T hetzelfde karakteristieke polynoom en daarom ook dezelfde eigenwaarden. Omdat A kolomstochastisch is, is het niet moeilijk in te zien dat $A^T \mathbf{e} = \mathbf{e}$, 1 is dus een eigenwaarde voor A^T , en daarmee ook voor A . \square

Zo is 1 een eigenwaarde van ons matrix A in de voorgaande voorbeeld, en elke meervoud van de vector $[12, 4, 9, 6]^T$ een eigenvector bij die eigenwaarde. In dit geval krijgen we $x_1 = 12/31$, $x_2 = 4/31$ enz. Wat opvalt is dat hoewel pagina 3 meer backlinks bezit, pagina 1 toch een hogere score krijgt, maar dat komt overeen met onze tweede aanname.

3.2 Non-Unique Rankings

Met onze formule in de vorige paragraaf voor x_k kunnen we een probleem krijgen wanneer de dimensie van de eigenruimte behorend bij de eigenwaarde 1 niet gelijk is aan 1. We kunnen dan een ranking krijgen die niet uniek is.

3.3 Definitie Een web is **sterk verbonden** als je van elke pagina uit dat web via links naar een ander willekeurige pagina van dat web kunt gaan in eindig veel stappen.

3.4 Definitie De verzameling van alle eigenvectoren corresponderend met een vaste eigenwaarde λ van een vierkante matrix A vormt samen met de nulvector (oftewel $\text{Ker}(\lambda I - A)$) de **eigenruimte** van λ , we noteren het als $V_\lambda(A)$.

We gaan nu laten zien waarom een web bestaande uit niet verbonden subwebben meerdere lineair onafhankelijke eigenvectoren bij de eigenwaarde 1 moet hebben, oftewel $\dim(V_1(A)) > 1$.

Neem aan dat een web W n pagina's en 2 onverbonden subwebben W_1, W_2 bevat. Zij n_i het aantal pagina's in W_i , $i = 1, 2$, nummer de paginas in W_1 met indices 1 tot en met n_1 , de pagina's in W_2 met $n_1 + 1$ tot en met $n_1 + n_2$. We krijgen nu een matrix A die een blokdiagonale structuur heeft

$$A = \begin{pmatrix} A_1 & 0 \\ 0 & A_2 \end{pmatrix}.$$

Hier zijn A_i de link matrix voor W_i omdat we W_i als een zelfstandige web op zichzelf mogen beschouwen. Elke $n_i \times n_i$ matrix A_i is kolomstochastisch en heeft dus een eigenvector $v^i \in \mathbb{R}^n$ behorend bij eigenwaarde 1. Construeer nu voor $i = 1, 2$ een vector $w^i \in \mathbb{R}^n$ die 0-componenten heeft voor alle elementen behorende tot de blokken anders dan blok i , dus in dit geval

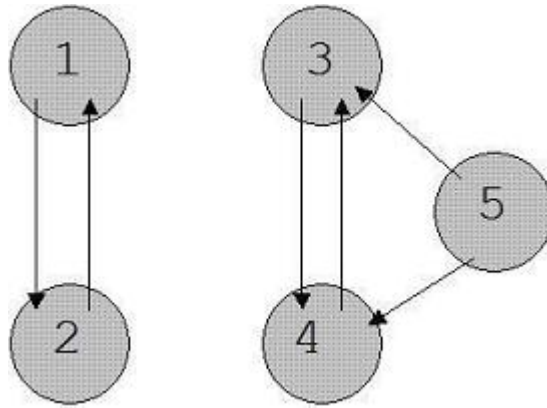
$$w^1 = \begin{pmatrix} v^1 \\ 0 \end{pmatrix}, w^2 = \begin{pmatrix} 0 \\ v^2 \end{pmatrix}$$

De vectoren w_1, w_2 zijn nu lineair onafhankelijke eigenvectoren van A met eigenwaarde 1 omdat

$$Aw^i = w^i \quad i = 1, 2.$$

Er bestaat dus geen unieke positieve eigenvector.

Dit kunnen we ook zien aan de hand van het volgende voorbeeld. Beschouw het volgende web :



De linkmatrix is dus

$$A = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & \frac{1}{2} \\ 0 & 0 & 1 & 0 & \frac{1}{2} \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}.$$

Hierbij is $x = [1/2, 1/2, 0, 0, 0]^T$ een mogelijke eigenvector bijbehorende bij eigenwaarde 1 maar ook $y = [0, 0, 1/2, 1/2, 0]$, en dus ook elke lineaire combinatie hiervan. Er bestaat dus geen unieke eigenvector.

3.3 Dangling node

Een ander probleem dat men tegen zal komen als we de linkmatrix A gebruiken is een web met dangling nodes. Een web met dangling nodes bevat pagina's

die geen uitgaande links bevatten en geeft daardoor een linkmatrix A die kolommen kunnen bevatten met alleen nullen als elementen. In dat geval is de matrix A *kolomsubstochastisch*,

3.5 Definitie Een vierkante matrix is **kolomsubstochastisch** als alle elementen van die matrix niet negatief zijn en de som van alle elementen van elke kolom kleiner of gelijk is aan 1.

Zo'n matrix A heeft alleen eigenwaarden die kleiner of gelijk zijn aan 1 in absolute waarde, maar 1 hoeft niet altijd een eigenwaarde te zijn van A . Om dit te laten zien definiëren we eerst het volgende:

3.6 Definitie Zij $\sigma(A) := \{\lambda_1, \lambda_2, \dots, \lambda_n\}$ de verzameling eigenwaarden (het **spectrum**) van een matrix $A \in \mathbb{C}^{n \times n}$, waarbij we aannemen dat $|\lambda_1| \geq |\lambda_2| \dots \geq |\lambda_n|$. De **spectrale straal** van A wordt gedefinieerd als $r(A) := \max_i(|\lambda_i|)$.

3.7 Definitie zij \mathbb{K} het veld van de reële of de complexe getallen. Beschouw de ruimte \mathbf{M}_n van alle $n \times n$ vierkante matrices met n met elementen uit \mathbb{K} . We noemen een functie $\|\cdot\| : \mathbf{M}_n \rightarrow \mathbb{R}$ een **matrixnorm** op $\mathbb{K}^{n \times n}$ als voor alle $A, B \in \mathbb{K}^{n \times n}$ geldt dat

1. $\|A\| \geq 0, \|A\| = 0$ desda $A = 0$
2. $\|cA\| = |c|\|A\|$ voor alle $c \in \mathbb{K}$
3. $\|A + B\| \leq \|A\| + \|B\|$
4. $\|AB\| \leq \|A\|\|B\|$

Omdat de ruimte die we bekijken eindig dimensionaal is, zijn alle matrixnormen op die ruimte equivalent. Dit komt door de equivalentie van vectornormen op eindig dimensionale lineaire vectorruimten.

Verder hebben we het volgende verband tussen de spectrale straal en de norm van een matrix.

3.8 Stelling Zij $\|\cdot\|$ een matrixnorm, dan $r(A) \leq \|A\|$.

Bewijs. Als λ een eigenwaarde is van A dan $|\lambda| \leq r(A)$ en bovendien is er volgens de definitie van de spectrale straal minstens een eigenwaarde van A waarvoor geldt $|\lambda| = r(A)$. Als $Ax = \lambda x, x \neq 0$ en $|\lambda| = r(A)$, neem nu een matrix X met al zijn kolommen gelijk aan de eigenvector x , dan krijgen we $AX = \lambda X$. Zij $\|\cdot\|$ een matrixnorm, dan $|\lambda|\|X\| = \|\lambda X\| = \|AX\| \leq \|A\|\|X\|$

en dus $|\lambda| = r(A) \leq \|A\|$. \square

3.9 Stelling Zij $\|\cdot\|$ een vectornorm op \mathbb{K}^n . Definieer $\|\cdot\|$ op $\mathbb{K}^{n \times n}$ met $\|A\| \equiv \max_{\|x\|=1} \|Ax\|$, dan is $\|\cdot\|$ een matrixnorm. We noemen $\|\cdot\|$ de geïnduceerde norm van de vectornorm $\|\cdot\|$.

Bewijs. We bekijken de vier axioma's van een matrixnorm

1. $\|A\| \geq 0, \|A\| = 0$ volgt direct uit het feit dat $\|A\|$ gedefinieerd is als het maximum van een verzameling niet negatieve functiewaarden. En $\|A\| = 0$ desda $A = 0$ volgt uit het feit dat $Ax = 0$ voor alle x juist wanneer $A = 0$.

2. Het tweede axioma volgt uit de volgende berekening

$$\|cA\| = \max \|cAx\| = \max |c| \|Ax\| = |c| \max \|Ax\| = |c| \|A\|.$$

Met $|c|$ de absolute waarde van c .

3. Analooft voor de derde axioma geldt nu dat

$$\|A+B\| = \max \|(A+B)x\| = \max \|Ax+Bx\| \leq \max(\|Ax\| + \|Bx\|) \leq \max \|Ax\| + \max \|Bx\| = \|A\| + \|B\|.$$

4. En het laatste axioma geldt omdat

$$\|AB\| = \max \frac{\|ABx\|}{\|x\|} = \max \frac{\|ABx\|}{\|Bx\|} \frac{\|Bx\|}{\|x\|} \leq \max \frac{\|Ay\|}{\|y\|} \max \frac{\|Bx\|}{\|x\|} = \|A\| \|B\|.$$

Waarbij we zonder verlies van algemeenheid aannemen dat het maximum genomen is over die x die niet in de nulruimte van B zitten.

\square

Neem nu de geïnduceerde matrixnorm $\|\cdot\|$ van de vector kolomsomnorm.

$$\|A\|_1 \equiv \max_{\|x\|_1=1} \|Ax\|_1.$$

Bekijk $\|Ax\|_1$,

$$\|Ax\|_1 = \sum_{i=1}^n |[Ax]_i| = \sum_{i=1}^n \left| \sum_{j=1}^n A_{ij}x_j \right| \leq \sum_{i,j} A_{ij}|x_j| = \sum_j \left(\sum_i A_{ij} \right) |x_j|$$

Omdat A kolomsubstochastisch is, is $\sum_i A_{ij} \leq 1$. Dit geeft dat

$$\|Ax\|_1 = \sum_j (\sum_i A_{ij}) |x_j| \leq \sum_j |x_j| = \|x\|_1 = 1 \rightarrow \|A\| \leq 1.$$

Uit de vorige stelling volgt dan dat elke eigenwaarde van zo'n matrix kleiner of gelijk moet zijn aan 1 terwijl 1 niet altijd een eigenwaarde hoeft te zijn van zo'n matrix.

De kans bestaat dat er dangling nodes op het web aanwezig zijn, aangezien de grootte van ons world wide web. De aangepaste matrix dat we gaan introduceren in Paragraaf 3.5 kan een oplossing zijn voor dit probleem. In dat geval bestaat er altijd een unieke rangschikking.

3.4 Verwisselen van index

Tot nu toe hebben we stilzwijgend aangenomen dat de volgorde van de nummering van pagina's geen invloed hebben op de rangschikking. Dat gaan we nu laten zien door te bekijken wat er gebeurt als we de indices van paginas i en j verwisselen:

Stap 1. Laat \tilde{A} de linkmatrix zijn van de hergenummerde web. Dan $\tilde{A} = PAP$ met P de matrix die men krijgt na verwisselen van rij i en j van de $n \times n$ identiteitsmatrix. Want we weten van de lineaire algebra dat $A \rightarrow PA$ de rijen i en j van A verwisseld en $A \rightarrow AP$ de kolommen, verder geldt er dat $P^2 = I$.

Stap 2. Laat x een eigenvector zijn van A , dus $Ax = \lambda x$ voor een zekere λ , dan is $y = Px$ een eigenvector voor \tilde{A} met eigenwaarde λ , want

$$\tilde{A}y = PAPy = PAPPx = PAx = P\lambda x = \lambda Px = \lambda y.$$

Dit laat dus zien dat de volgorde van de nummering geen invloed heeft op de rangschikking, omdat we voor de nieuwe nummering een zelfde eigenvector krijgen met i en j verwisseld, de rangschikking blijft dus invariant.

3.5 Een remedie voor $\dim(V_1(A)) > 1$

Ons world wide web bestaat uit miljarden van webpagina's, de grote hoeveelheid rekenwerk die men uit moet voeren om een eigenvector te krijgen is dan ook gigantisch. Het is dus heel belangrijk om te weten dat er een unieke ranking bestaat.

We weten nu dat een niet sterk verbonden web geen unieke eigenvector heeft. Maar het world wide web kan uit vele disjuncte componenten bestaan, het is dus zeer belangrijk om een oplossing te vinden voor dit probleem.

Het idee van de oplossing die we nu gaan geven komt van een speciaal geval van de Perron-Frobenius stelling. We bewijzen in dit deel alleen de delen die

van toepassing zijn voor ons probleem, in een aparte hoofdstuk over de stelling zal uitgebreid het algemene geval bekeken en bewezen worden.

We nemen aan dat ons web geen dangling nodes bevat maar wel disjuncte componenten kan hebben. Voor zo'n web hebben we een oplossing. Het idee is als volgt, we hebben besloten dat een pagina belangrijker is als dat pagina meer backlinks heeft, maar dat betekent eigenlijk dat men tijdens het willekeurig klikken op links op het web een grotere kans heeft om op dat pagina te komen dan een minder belangrijke pagina. In dit proces van klikken op links kunnen we nooit bij een disjuncte web komen omdat er geen link is gemaakt. Maar een gebruiker hoeft niet per se via een link naar een pagina te komen. Door willekeurig een pagina te openen kunnen we alsnog in een disjuncte web komen. Door in ons model zo'n virtuele links te maken kunnen we een mogelijke oplossing krijgen.

Laat S een $n \times n$ matrix zijn met al zijn elementen gelijk aan $1/n$. De matrix S is totaal positief, kolomstochastisch en $V_1(S)$ is 1-dimensionaal omdat S rang 1 heeft. Als oplossing vervangen we nu de oude linkmatrix A door de nieuwe matrix

$$M = (1 - m)A + mS, \quad S = \frac{1}{n} \mathbf{e} \mathbf{e}^T$$

met een zekere $m \in (0, 1)$. M is als het ware het gewogen gemiddelde van A en S , de m moeten we niet al te groot nemen om dat er anders veel informatie van onze echte linkmatrix verloren kan gaan. De oorspronkelijke waarde van m die Google gebruikt is 0.15. De keuze van deze waarde zal te maken hebben met het snelheid van convergentie tijdens het uitrekenen van zo'n eigenvector. In het volgende hoofdstuk zullen we meer aandacht hieraan besteden.

We gaan nu eerst laten zien waarom de nieuwe matrix ons probleem op kan lossen. Merk eerst op dat men de eigenvector zo kan normeren zodat $\sum_i x_i = 1$. Op die manier kunnen we als we denken aan kansen een betere beeld krijgen van onze rangschikking.

Zij $s = \frac{1}{n} \mathbf{e}$. Dan $Sx = s$ omdat $\sum_i x_i = 1$. De vergelijking $x = Mx$ kunnen we dus ook schrijven als

$$x = (1 - m)Ax + ms.$$

We gaan nu bewijzen dat $V_1(M)$ 1-dimensionaal moet zijn door een speciaal geval van de Perron-Frobenius stelling te bekijken zoals eerder genoemd. Merk eerst op dat alle elementen M_{ij} van de matrix M strikt positief zijn. We geven de volgende definitie

3.10 Definitie Een matrix M is **positief** als $M_{ij} > 0$ voor alle i en j .

3.11 Stelling *Als M positief en kolomstochastisch is, dan is $V_1(M)$ 1-dimensionaal, en er bestaat een eigenvector in $V_1(M)$ met positieve elementen.*

Het is voldoende om alleen reële eigenvectoren te bekijken. (Stel er is een eigenvector $z = x + iy$ met $Az = z$, dan zijn x en y ieder op zichzelf ook eigenvectoren van A . Dus als we kunnen laten zien dat alle reële eigenvectoren veelvoud van elkaar zijn, dan zijn we klaar.) Van Propositie 2.1.2 weten we dat $V_1(M)$ niet nul is omdat M kolomstochastisch is. Voor het eerste deel van het bewijs moeten we nog alleen laten zien dat er geen onafhankelijke eigenvectoren in $V_1(M)$ kunnen bestaan.

3.12 Propositie *Zij M positief en kolomstochastisch, dan hebben alle reële eigenvectoren in $V_1(M)$ de eigenschap dat al hun elementen van hetzelfde teken moeten zijn. Er zijn dus geen vectoren die elementen groter dan 0 en elementen kleiner dan 0 bevatten.*

Bewijs. Bekijk eerst de standaard driehoeksongelijkheid

$$\left| \sum_i y_i \right| \leq \sum_i |y_i|, \quad y \in \mathbb{R}^n.$$

Deze ongelijkheid is strikt als y_i van gemengd teken zijn (dat wil zeggen dat de eigenvectoren zowel elementen kleiner als groter dan 0 bevatten). Stel $x \in V_1(M)$ bevat elementen van verschillende teken, van $x = Mx$ hebben we $x_i = \sum_{j=1}^n M_{ij}x_j$. Nu moet $M_{ij}x_j$ dus ook van gemengd teken zijn omdat M positief is. We krijgen de volgende strikte ongelijkheid

$$|x_i| = \left| \sum_{j=1}^n M_{ij}x_j \right| < \sum_{j=1}^n M_{ij}|x_j|.$$

Sommeer nu het linker en rechterlid van i tot n en verwissel de i en j sommatie. We krijgen dan

$$\sum_{i=1}^n |x_i| < \sum_{i=1}^n \sum_{j=1}^n M_{ij}|x_j| = \sum_{j=1}^n \left(\sum_{i=1}^n M_{ij} \right) |x_j| = \sum_{j=1}^n |x_j|.$$

Dit geeft een tegenspraak, hiermee hebben we dus laten zien dat alle elementen van x van hetzelfde teken moeten zijn. \square

3.13 Propositie *Zij v en w twee lineair onafhankelijke vectoren in \mathbb{R}^m , $m \geq 2$, dan bestaan er s en t niet beide gelijk 0 zodat de vector $x = sv + tw$ elementen van verschillende teken bevatten.*

Bewijs. Lineair onafhankelijkheid impliceert dat v en w beide niet gelijk zijn aan 0. Stel $d = \sum_i v_i$. Als $d = 0$ dan moet v elementen van gemengde teken bevatten. Neem nu $s = 1, t = 0$, dan hebben we zo'n vector x . Als $d \neq 0$, laat $s = -\frac{\sum_i w_i}{d}, t = 1$ en $x = sv + tw$, Omdat v en w lineair onafhankelijk zijn is $x \neq 0$, echter, $\sum_i x_i = 0$, dus x moet elementen van gemengd teken bevatten. \square

Nu zijn we toe aan ons hoofdbewijs.

Bewijs. Stel er zijn twee lineair onafhankelijke eigenvectoren v en $w \in V_1(M)$. Voor alle s en t dat niet beide nul zijn geldt dat de vector $x = sv + tw, x \neq 0$ in $V_1(M)$ moeten liggen en dus alle elementen van hetzelfde teken moet zijn. Volgens Propositie 2.5.3 moet x van gemengde teken zijn, dit geeft een tegenspraak. Hieruit kunnen we dus concluderen dat $V_1(M)$ niet twee lineair onafhankelijke vectoren kan hebben en dus 1-dimensionaal moet zijn. Nu de eigenvector niet van gemengd teken zijn kunnen we altijd een veelvoud van de eigenvector nemen zodat er een positieve rangschikking bestaat (bij een negatieve eigenvector kunnen we altijd met -1 vermenigvuldigen).

4 Berekenen van de eigenvector en de tweede eigenvector van Google

Het berekenen van een eigenvector bij een matrix met een dergelijke omvang als de linkmatrix van ons world wide web kan alleen maar bij benadering gebeuren. Een goede numerieke methode met snelle convergentie is vereist, voor Google werkt de *machtsmethode* goed genoeg.

4.1 Machtsmethode

Een goede numerieke benadering voor de eigenvector kunnen we krijgen met behulp van de machtsmethode (power method). We beginnen met een willekeurige niet negatieve beginvector x_0 , bijvoorbeeld $x_0 = 1/n e^T$. Met de formule

$$x_k = Mx_{k-1} = M^k x_0$$

laten we de k naar oneindig gaan. Afhankelijk van de grootte van de eigenvectoren kan deze rij divergeren, naar nul gaan of convergeren naar een vaste vector. Maar dat geeft precies onze rangschikingsvector omdat we bij convergentie zal hebben dat $x = \lim x_k = \lim x_{k+1} = \lim Mx_k = M \lim x_k = Mx$. We kunnen hierbij onze matrix M zien als een overgangsmatrix en de vector x

als een toestandsvector. Het idee is als volgt, stel we hebben een vaste aantal gebruikers (we moeten denken aan miljoenen) die ieder op een willekeurige pagina begint (starttoestand), vervolgens gaan de gebruikers willkeurig op een link van dat pagina klikken om naar de volgende pagina te gaan, elke keer tellen we bij elke pagina hoeveel gebruikers daar terecht zijn gekomen, door continu het proces te volgen willen we uiteindelijk in een toestand terecht komen waar het aantal gebruikers op een pagina niet meer verandert. Door het percentage van aantal gebruikers van het totale gebruikers te bepalen, kunnen we aan elke pagina een waarde toekennen. Het is de vraag of er zo'n eindtoestand bestaat (of de iteratie wel convergeert).

4.2 Convergentiesnelheid en de tweede eigenwaarde

Een belangrijke voorwaarde voor de convergentie is dat $V_1(M)$ 1-dimensionaal moet zijn en dat de rest van de eigenwaarden in absolute waarde kleiner zijn dan 1. Bovendien geeft de waarde van de grootte van de tweede eigenwaarde een schatting van de convergentiesnelheid. Met behulp van een stelling over matrixnorm en een lemma proberen we dit te laten zien.

4.1 Stelling $\forall A \forall \varepsilon > 0, \exists$ matrix norm $\|\cdot\|$ zodat $r(A) \leq \|A\| \leq r(A) + \varepsilon$.

Bewijs. Met Schur triangularization stelling weten we dat er een unitaire matrix U en een bovendiagonaalmatrix δ bestaan zodat $A = U^T \delta U$. Neem $D_t \equiv \text{diag}(t, t^2, t^3, \dots, t^n)$ en bereken nu

$$D_t \delta D_t^{-1} = \begin{pmatrix} \lambda_1 & t^{-1}d_{12} & t^{-2}d_{13} & \dots & t^{-n+1}d_{1n} \\ 0 & \lambda_2 & t^{-1}d_{23} & \dots & t^{-n+2}d_{2n} \\ 0 & 0 & \lambda_3 & \dots & t^{-n+1}d_{3n} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & \dots & t^{-n+1}d_{n-1,n} \\ 0 & 0 & 0 & 0 & \lambda_n \end{pmatrix}.$$

Door $t > 0$ groot genoeg te nemen kunnen we de som van de absolute waarden van alle elementen die niet op diagonaal staan kleiner dan ε krijgen. In het bijzonder krijgen we $\|D_t \delta D_t^{-1}\|_1 \leq r(A) + \varepsilon$ wanneer we t groot genoeg nemen. Definieer nu de matrix norm $\|\cdot\|$ als volgt

$$\|B\| \equiv \|D_t U^T B U D_t^{-1}\|_1 = \|(U D_t^{-1})^{-1} B (U D_t^{-1})\|_1$$

for alle $n \times n$ -matrices B . Als we t groot genoeg nemen, dan hebben we een matrixnorm geconstrueerd zodat $\|A\| \leq r(A) + \varepsilon$ geldt. We weten al dat $\|A\| \geq r(A)$ voor alle matrixnormen, hiermee hebben we de stelling bewezen.

4.2 Lemma Zij M een kolomstochastische matrix, dan $\mathbb{C}^n = \text{Ker}(I - M) \oplus \text{Im}(I - M)$.

Bewijs. Van de dimensiestelling weten we dat

$$\dim(\text{Ker}(I - A)) + \dim(\text{Im}(I - A)) = n.$$

We moeten dus aantonen dat de doorsnede gelijk is aan 0, daarmee is de som ook direct. We laten nu zien dat er geen andere element dan 0 in de doorsnede van de kern en de beeldruimte bestaan.

Zij $x \in \text{Ker}(I - M)$, en stel $\exists y, x = (I - M)y$ zodat x ook een element is van de beeldruimte. Nu geldt er

$$0 = (I - M)x = (I - M)^2y = y - 2My + M^2y$$

We laten eerst met behulp van de volledig inductiestelling zien dat $M^{n+1}y = My - nx \quad \forall n \geq 0$, merk eerst op dat $M^2y = 2My - y = My + My - y = My - x$.

Stap 1. Voor $n = 0$ geldt dat $M^1y = My - 0x$, dus dit klopt.

Stap 2. We nemen nu aan dat $M^ny = My - (n - 1)x$ voor een zekere $n \geq 1$ als inductie veronderstelling.

Stap 3. Als $M^ny = My - (n - 1)x$, dan

$$M^{n+1}y = M(M^ny) = M(My - (n - 1)x) = M^2y - M(n - 1)x = (My - x) - (n - 1)Mx = My - x - nx - x = My - nx.$$

Conclusie: $M^{n+1}y = My - nx \quad \forall n \geq 0$.

Omdat de matrix M kolomstochastisch is, is $\sup\|M^ny\|_1 \leq \|y\|_1$. Uit de conclusie van ons bewijs hierboven krijgen we

$$n\|x\| = \|nx\| = \|My - M^{n+1}y\| \leq 2\|y\| \quad \forall n.$$

Dit kan alleen als $\|x\| = 0$, maar dan $x = 0$. qed

Laat

$$M = (1 - m)A + mS, \quad S = \frac{1}{n}\mathbf{e}\mathbf{e}^T.$$

Met de decompositie als hierboven gedefinieerd krijgen we

$$M \sim \begin{pmatrix} I & 0 \\ 0 & B \end{pmatrix}$$

Bij de machtsmethode nemen we de machten van M oftewel

$$M^n \sim \begin{pmatrix} I^n & 0 \\ 0 & B^n \end{pmatrix}$$

Nu zijn alle machten van I gelijk aan I en er geldt dat de spectrale straal van B gelijk aan de absolute waarde van λ_2 , de tweede eigenwaarde van M . Omdat deze waarde strikt kleiner is dan 1 (dat gaan we later laten zien) geldt er dat $r(B) = |\lambda_2| < q < 1$ voor een zekere q . Zoals we eerder liet zien bestaat er een norm $\|\cdot\|$ zodat $\|B\| < r(B) + \varepsilon$. Kies ε zodanig dat $\varepsilon + r(B) = q$, dan $\|B^n\| \leq \|B\|^n < q^n \rightarrow 0$ als $n \rightarrow \infty$. Met de machtmethode krijgen we dus onze unieke eigenvector bijbehorend bij de grootste eigenwaarde, met een snelheid ongeveer gelijk als de tweede eigenwaarde.

Omdat M kolomstochastisch is weten we dat $V_1(M)$ 1-dimensionaal is met alle andere eigenwaarden strikt kleiner dan 1 in absolute waarden. We laten nu zien dat de tweede eigenwaarde gelijk moet zijn aan m .

4.3 Stelling *Stel $M = (1 - m)A + mS$, dan is de absolute waarde van de tweede eigenwaarde λ_2 (eigenwaarde van een matrix die op de dominante eigenwaarde na de grootste absolute waarde van alle eigenwaarden van die matrix heeft) van M kleiner of gelijk dan m .*

Het bewijs hiervan geven we met behulp van enkele lemma's.

4.4 Lemma *Voor de tweede eigenvector x_2 bijbehorend bij de tweede eigenwaarde λ_2 van M geldt: $\mathbf{e}^T x_2 = 0$.*

Omdat A kolomstochastisch is, geldt $\mathbf{e}^T M = \mathbf{e}^T$. Maar $\mathbf{e}^T x_2 = \mathbf{e}^T M x_2 = \mathbf{e}^T (\lambda_2 x_2) = \lambda_2 \mathbf{e}^T x_2$, dit geeft dat $(1 - \lambda_2) \mathbf{e}^T x_2 = 0$. Omdat $\lambda_2 \neq 1 \rightarrow \mathbf{e}^T x_2 = 0$.

4.5 Lemma $Sx_2 = 0$

Per definitie geldt $S = \frac{1}{n} \mathbf{e} \mathbf{e}^T$, met het vorige lemma volgt dat $Sx_2 = \frac{1}{n} \mathbf{e} \mathbf{e}^T x_2 = \frac{1}{n} \mathbf{e} 0 = 0$

4.6 Lemma *De tweede eigenvector x_2 van M is ook een eigenvector van A , met de bijbehorende eigenwaarde gelijk aan $\gamma = \lambda_2/m$.*

Bewijs.

$$mAx_2 + (1 - m)Sx_2 = Mx_2 = \lambda_2 x_2$$

met het vorige lemma krijgen we $mAx_2 = \lambda_2 x_2$ en delen door m geeft

$$Ax_2 = \frac{\lambda_2}{m} x_2$$

4.7 Lemma $|\lambda_2| \leq m$

Van het vorige lemma weten we dat λ_2/m een eigenwaarde is van A , omdat A kolomstochastisch is, is $|\frac{\lambda_2}{m}| \leq 1$ dit geeft dan dat $|\lambda_2| \leq m$

4.8 Stelling *Wanneer een web minstens twee subwebben bevat die niet aan elkaar verbonden zijn (dit is het geval voor ons world wide web), dan is de tweede eigenwaarde van de matrix M gelijk aan m .*

Voor het geval van $m = 1$ hebben we in de voorgaande hoofdstuk al laten zien dat er meerdere eigenvectoren bij eigenwaarde 1 bestaan, dus $\lambda_2 = 1$, We bewijs nog alleen voor $m < 1$. We laten eerst zien dat er een eigenwaarde van A bestaat die gelijk is aan m , nu 1 de grootste eigenwaarde is, moet de tweede eigenwaarde dus groter dan of gelijk zijn aan 1, en van de vorige stelling volgt dan $\lambda_2 = m$.

4.9 Lemma *Elke eigenvector $y_i, i = 1, \dots, n$ van A die loodrecht op \mathbf{e} staat is een eigenvector x_i van M met $\lambda_i = m\gamma_i$.*

Bewijs. We moeten bewijzen dat als $Ay = \gamma x$ en $\mathbf{e}^T \mathbf{y} = \mathbf{0}$ dan $Mx = m\gamma y$. Gegeven $\mathbf{e}^T \mathbf{y}_i = \mathbf{0}$, dan

$$Sy_i = \frac{1}{n} \mathbf{e} \mathbf{e}^T \mathbf{y}_i = \mathbf{0}.$$

Per definitie geldt

$$Ay_i = \gamma_i y_i.$$

Daarom geldt

$$My_i = (mA + (1 - m)S)y_i = mAy_i = m\gamma_i y_i.$$

□

In hoofdstuk 2 hebben we laten zien dat na henummering van index we van A naar een matrix kan 'omschrijven' van de vorm

$$A = \begin{pmatrix} C & 0 \\ 0 & D \end{pmatrix}$$

met C en D kolomstochastisch, nu heeft de matrix twee verschillende lineaire onafhankelijke eigenvectoren x, y met $Ax = x, Ay = y$. Omdat x, y lineair onafhankelijk zijn kunnen we λ, μ niet beide gelijk aan 0 vinden zodat $e \perp \lambda x + \mu y =: w$. Nu is w niet gelijk aan 0 omdat x, y lineair onafhankelijk zijn, en er geldt $Aw = w$. Maar $Sw = \frac{1}{n} \mathbf{e} \mathbf{e}^T w = 0, \Rightarrow Mw = (mA + (1 - m)S)w = mA w = mw \Rightarrow \lambda = m$.

4.3 Aangepaste matrix M

Tot nu toe hebben we bij de matrix M gekozen voor een matrix mS met allemaal gelijke en positieve waarden. In dat geval weten we wegens de kolomstochastische eigenschap dat er een unieke dominante eigenwaarde bestaat. Maar we kunnen laten zien dat de machtmethode ook werkt als de matrix niet strikt positief is. Dit is wel interessant om te weten omdat Google op dit moment bezig is met verbeteringen van het algoritme. Daarbij is bijvoorbeeld onze virtuele link zoals gedefinieerd in de voorgaande hoofdstukken niet meer gekozen als een uniform verdeelde kans maar een variabele die afhankelijk is van ieder individuele netgebruiker. Het enige wat we willen laten zien is dat voor deze matrix M ook geldt dat 1 de unieke dominante eigenwaarde is, de rest van het bewijs gaat net zo als vanaf Lemma 4.5. We gaan dus in plaats van die mS met alle gelijke waarden een nieuwe kansverdelingsvector v gebruiken dat afhankelijk is van de achtergrond van een bezoeker. We definiëren nu dus de nieuwe Google matrix M .

$$M = mA + (1 - m)E.$$

$$v = (v_1 \ v_2 \ \dots \ v_n)^T, E = v\mathbf{e}^T$$

Hierbij is A ons linkmatrix met $0 \leq m \leq 1$, \mathbf{e} de n vector met alle elementen gelijk aan 1.

Definieer eerst

$$U := \{i \mid v_i \neq 0\}$$

$$F := \{j \mid \exists i \in U, k \in \mathbb{N}, [M^k]_{ij} > 0\}$$

$$J_F := \text{span}\{e_j \mid j \in F\}$$

4.10 Definitie Een $n \times n$ vierkante matrix $A = a_{ij}$ is *reducibel* als de indices $i, j = 1, 2, \dots, n$ kunnen worden ingedeeld in twee disjuncte niet lege verzameling i_1, i_2, \dots, i_μ en j_1, j_2, \dots, j_ν (met $\mu + \nu = n$) zodat $a_{i_\alpha j_\beta} = 0$ voor $\alpha = 1, 2, \dots, \mu, \beta = 1, 2, \dots, \nu$.

Als een matrix A reducibel is dan bestaat er een permutatie matrix P bestaat zodat $P^T A P$ is bovendagonaal. Een $n \times n$ -matrix die niet reducibel is noemen we *irreducibel*. Een matrix is irreducibel als er voor elke i en j er een k bestaat zodat $(A^k)_{ij} > 0$.

4.11 Definitie Een lineaire deelruimte $V \subset \mathbb{R}^n$ noemen we *A-invariant* als $AV \subset V$, oftewel $Av \in V$ voor alle $v \in V$.

4.12 Propositie *De deelruimte J_F is M -invariant.*

Bewijs. Neem $i \in F$ en zij $[Me_i]_l > 0$, we willen nu laten zien dat $l \in F$, $[Me_i]_l > 0$ dus

$$e_l^T Me_i > 0,$$

$i \in F$ dus voor een zekere $k \in \mathbb{N}$, $j \in U$ geldt nu

$$e_i^T M^k e_j > 0$$

maar dan

$$e_l^T Me_i e_i^T M^k e_j > 0$$

omdat $e_i e_i^T \leq I$ geldt

$$0 < e_l^T Me_i e_i^T M^k e_j \leq e_l^T MIM^k e_j = e_l^T M^{k+1} e_j$$

$$\Rightarrow l \in F$$

dus J_F is M -invariant. We weten van systeemtheorie dat wanneer V M -invariant is we M in bovendiaagonaalvorm kunnen representeren.

We weten nu dus dat M met de decompositie van $J_F \oplus J_F^\perp$ na permutaties geschreven kan worden als

$$M = \begin{pmatrix} B & C \\ 0 & D \end{pmatrix}$$

Met $B := M|_{J_F}$ irreducibel. Omdat voor elke i en j van de blok geldt dat $(M^k)_{ij} > 0$. Voor het verdere bewijs introduceren we eerst het volgende :

4.13 Definitie *Een niet negatieve irreducibele matrix noemen we **primitief** als er een simpel eigenwaarde λ bestaat met $\lambda = r(M)$ op het spectrale straal.*

Een niet negatieve irreducibele matrix M is primitief als M minstens een strikt positieve diagonaal element bevat. Voor het bewijs zie pag 674-678 van [4].

Omdat $e_i^T Me_i > 0 \quad \forall i \in U$, is $\text{spoor}(B) > 0$, maar B is ook irreducibel. We kunnen hieruit dus concluderen dat B primitief is.

En nu het bewijs dat de tweede eigenwaarde van M kleiner of gelijk is aan m . Neem aan dat $\lambda \in \sigma(M)$ met

$$M \begin{pmatrix} x \\ y \end{pmatrix} = \lambda \begin{pmatrix} x \\ y \end{pmatrix}$$

$x \in J_F, y \in J_F^\perp$, x en y niet beide gelijk aan 0.

$$M \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} Bx + Cy \\ Dy \end{pmatrix} = \begin{pmatrix} \lambda x \\ \lambda y \end{pmatrix}$$

$$\Rightarrow \lambda y = Dy = \text{pr}[My] = \text{pr}[mAy + (1-m)ve^T y] = m(\text{pr})Ay.$$

Hierbij is $\text{pr} : \mathbb{R}^d \mapsto J_F^\perp$ de orthogonale projectie, en de laatste geldt omdat $v \in J_F$.

$$|\lambda||y| = \|\lambda y\| \leq m\|(\text{pr})A\||y| \leq m\|y\|$$

De ongelijkheid geldt omdat de norm van A na projectie alleen kleiner kan worden of gelijk blijven. We krijgen dus

$$y = 0 \vee |\lambda| \leq m \leq 1$$

Maar $y = 0 \Rightarrow Bx = \lambda x \wedge x \neq 0$, B is primitief en heeft dus maar 1 eigenwaarde op zijn spectrale straal. We krijgen dus $\lambda = 1 \vee |\lambda| < 1$.

We hebben laten zien dat $|\lambda| \leq m$ of $|\lambda_2| < 1$, maar als $|\lambda_2| < 1$ dan geldt ook $|\lambda_2| \leq m$. Vanaf hier kunnen we verder met Lemma 4.5, het bewijs gaat ook voor deze matrix M .

5 De stelling van Perron

We hebben voor veel bewijzen die we gegeven hebben gebruik gemaakt van het feit dat de matrix geen negatieve elementen bevat. Er is een belangrijke en handige stelling die deze eigenschappen verklaren, wij hebben het eerder al genoemd: de stelling van Perron-Frobenius.

5.1 Lemma *Zij A een matrix met al zijn elementen $A_{ij} \geq 0$, dan heeft A de spectrale straal r als een eigenwaarde, en bevat de eigenruimte van $\lambda = r$ niet negatieve eigenvectoren.*

Bewijs. We bekijken de resolvent

$$R(\lambda) = (\lambda - A)^{-1}.$$

De Neumann reeks van de resolvent wordt gegeven door

$$R(\lambda) = \frac{1}{\lambda} + \frac{A}{\lambda^2} + \frac{A^2}{\lambda^3} + \dots$$

deze reeks convergeert als $|\lambda| > r(A)$.

Dit laten we zien met behulp van Stelling 4.1.

Zij $\|\cdot\|$ een matrixnorm zogekozen zodat $r(A) \leq \|A\| \leq r(A) + \varepsilon$ voor een zekere ε . De reeks $R(\lambda) = \sum_0^\infty \frac{A^n}{\lambda^{n+1}}$ is convergent als $|\lambda| > \|A\|$, omdat er geldt dat $r(A) \leq \|A\| \leq r(A) + \varepsilon$ kunnen we $\|A\|$ willekeurig dichtbij het spectrale straal nemen. Dus voor alle $|\lambda| > r(A)$ is de reeks convergent.

Bekijk nu op de spectrale straal, we weten dat er op de spectrale straal minstens 1 eigenwaarde moet liggen, stel $\mu \in \sigma(A)$, $|\mu| = r$, $r = r(A)$. Voor een zekere $x \in \mathbb{R}^n$, $x \neq 0$ geldt nu

$$(\lambda - A)x = (\lambda - \mu)x$$

oftwel

$$R(\lambda)x = \frac{1}{\lambda - \mu}x, \quad \forall \lambda \notin \sigma(A).$$

Definieer

$$\lambda_n := r + \frac{1}{n}, \quad \mu_n := \mu \left(\frac{r + \frac{1}{n}}{r} \right)$$

dus

$$|\mu_n| = |\lambda_n| > r.$$

Er geldt voor $|x|$ (de vector x in absolute waarde) dat

$$\frac{|x|}{|\mu_n - \mu|} = |R(\mu_n)x| = \left| \sum_{k=0}^{\infty} \frac{A^k}{\mu_n^{k+1}}x \right| \leq \sum_{k=0}^{\infty} \frac{A^k}{|\mu_n|^{k+1}}|x| = R(\lambda_n)|x|.$$

Omdat $x \neq 0 \Rightarrow |x| \neq 0$, er bestaat dus een i zodat

$$[R(\lambda_n)|x|]_i \rightarrow \infty$$

Construeer

$$X_n := \frac{R(\lambda)|x|}{\|R(\lambda_n)|x|\|}$$

Omdat $(X_n)_n$ een begrensde rij is ($\|(X_n)\| = 1$) bestaat er volgens de Bolzano-Weierstrass stelling een convergente deelrij van $(X_n)_n$.

Noem de deelrij weer $(X_n)_n$. Dan $X_n \rightarrow y \geq 0$ voor een zekere $y \geq 0$, $\|y\| = 1$. Maar er geldt ook

$$AX_n = A \left(\frac{R(\lambda)|x|}{\|R(\lambda_n)|x|\|} \right) = \frac{-|x| + \lambda_n R(\lambda_n)|x|}{\|R(\lambda_n)|x|\|} = \frac{-|x|}{\|R(\lambda_n)|x|\|} + \lambda_n X_n.$$

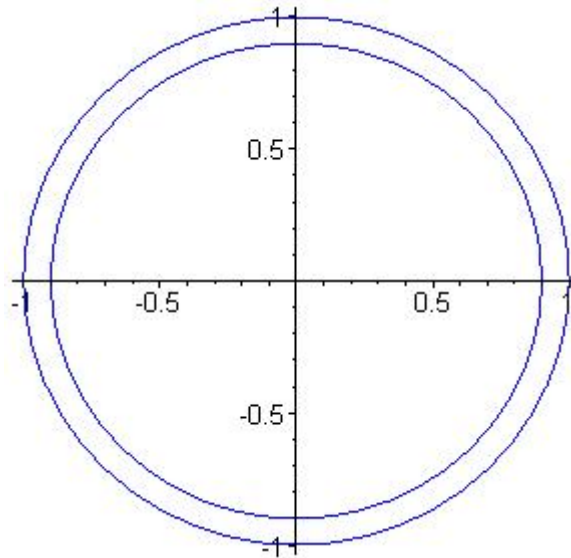
Merk op dat de bovenste geldt omdat:

$$AR(\lambda) = A(\lambda - A)^{-1} = (A - \lambda)(\lambda - A)^{-1} + \lambda(\lambda - A)^{-1} = -I + \lambda(\lambda - A)^{-1}.$$

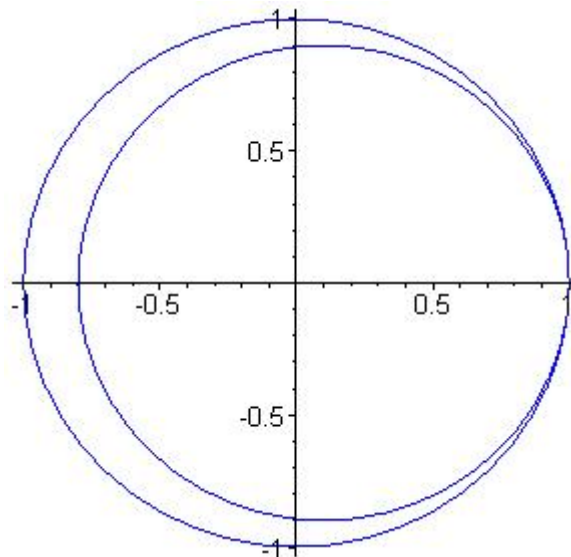
Als n naar oneindig gaat dan gaat de laatste term naar ry , maar $X_n \rightarrow y$ dus $AX_n \rightarrow Ay$, $\|y\| = 1$. De conclusie luidt dat $Ay = ry$. \square

5.2 Stelling *De eigenwaarde met de grootste absolute waarde van een positieve vierkante matrix A is simpel en positief en behoort tot een positieve eigenvector, alle andere eigenwaarden zijn kleiner in absolute waarde.*

Bewijs. Het eerste is al bewezen, we moeten nog laten zien dat er geen andere eigenwaarden op het spectrale straal liggen. Beschouw $A - \varepsilon I > 0$ voor een kleine $\varepsilon > 0$. Omdat A strikt positief is is het spectrale straal van $A - \varepsilon I > 0$ kleiner dan die van A . Zoals weergegeven in figuur hieronder



Het spectrum van deze matrix is gelijk aan $\sigma(A) - \varepsilon$. Door terugschuiven van de kleine cirkel met ε , krijgen we het spectrum $\sigma(A)$ van onze matrix A .



We zien nu dat alle andere eigenwaarden van A strikt kleiner moeten zijn dan de spectrale straal (dit geldt omdat alle digonaal elementen van A positief zijn).

6 Toekomst van het zoeken

Een belangrijke doel van Google is om alle informatie in de wereld te organiseren en universeel toegankelijk en nuttig te maken. Maar op dit moment zijn er maar 10 procent van alle in de wereld aanwezige informatie online beschikbaar volgens het hoofd van het project Google Book van het project Google Book Search. Het doel van het Google Book Search-project is om elk boek ter wereld te digitaliseren (boeken die vrij van copyright zijn). Hoe de zoekmachines de boeken inscant, is onduidelijk. Deze technologie blijft geheim.

Google probeert dus niet alleen sneller te zoeken en een betere rangschikking te bepalen, maar ook andere factoren spelen een grotere rol op de markt. Een ander project is om een rangschikking te maken die afhankelijk is van de achtergrond van een gebruiker. Dat doet Google door zoveel mogelijk informatie van een gebruiker op het web op te slaan, zo scant het bedrijf alle mails die met de e-maildienst Gmail worden verstuurd (ook om op die manier advertenties op maat op te sturen om inkomsten te maken), en slaat Google alle zoekgeschiedenis van een gebruiker.

En dat levert veel kritiek op, gebruikers willen hun privacy beschermen. Een belangrijke gegevens die Google opslaan tijdens het zoeken is ons IP-nummer. Daaruit is mogelijk uit te halen welke zoekopdrachten van welke gebruiker afkomstig is en van welke geografische positie het afkomstig is, en niet iedereen vindt dat prettig.

Op andere zoekmachines worden deze gegevens ook opgeslagen, maar meestal voor paar dagen. Maar Google sloeg ze op voor 30 jaar! Na veel protest is het nu verkort tot ongeveer 18 maanden. Maar wie daar ook problemen heeft, moet minder actief worden op internet, elke actie laat namelijk sporen achter.

Referenties

- [1] Kurt Bryan and Tanya Leise. The \$25, 000, 000, 000 eigenvector: the linear algebra behind Google. *SIAM Rev.*, 48(3):569C581 (electronic), 2006.
- [2] Taher H. Haveliwala and Sepandar D. Kamvar. The second eigenvector of the Google matrix. Stanford University Technical Report, 2003.

- [3] C. R. MacCluer. The many proofs and applications of Perrons theorem. *SIAM Rev.*, 42(3):487C498 (electronic), 2000.
- [4] Carl D. Meyer. *Matrix Analysis and Applied Linear Algebra*, SIAM, Philadelphia, 2000
- [5] Roger A. Horn and Charles R. Johnson. *Matrix Analysis*, Cambridge University Press, Cambridge, 1985