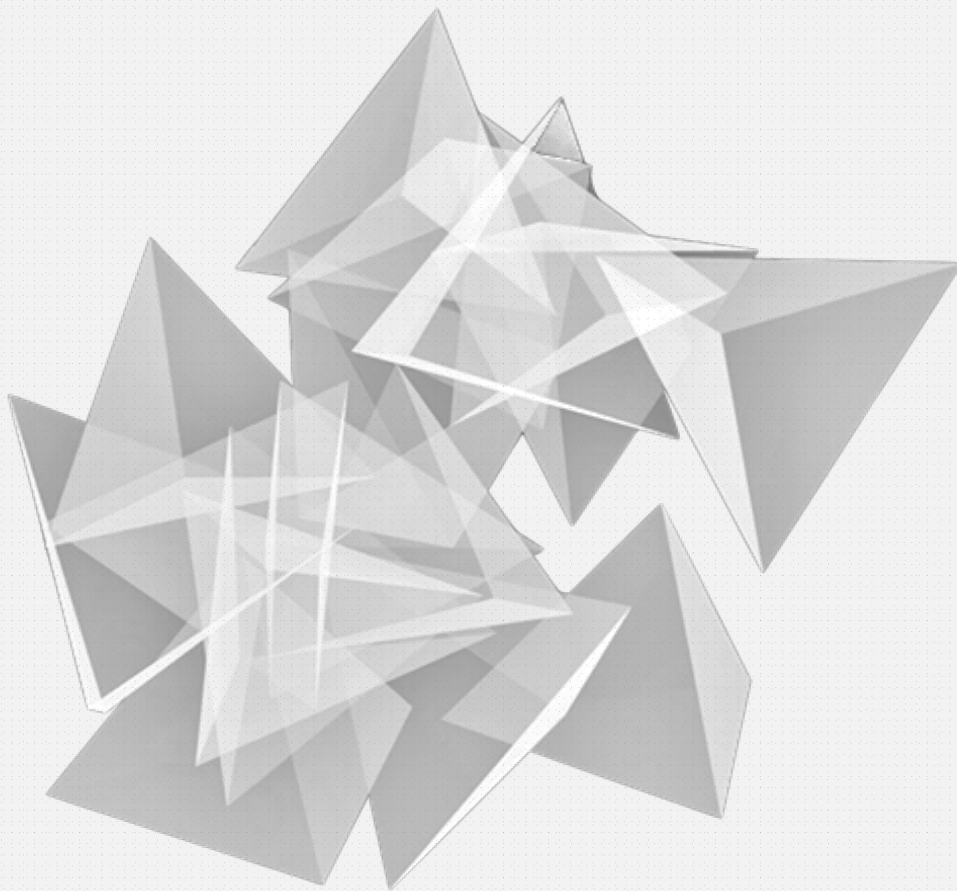


Practicum 1 - Lagerangbenaderingen

Bruno Vandekerkhove



ACADEMISCH JAAR 2019

G0Q57A: MODELLERING & SIMULATIE

Inhoudsopgave

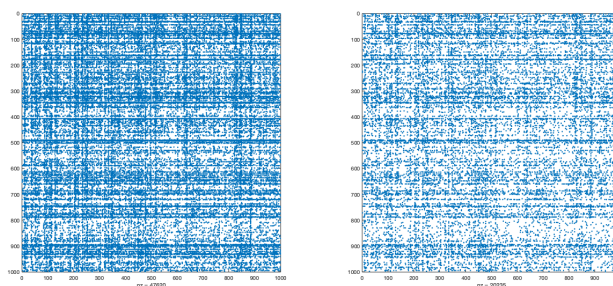
Een Aanbevelingssysteem voor Films	1
Opdracht 1	1
Opdracht 2	2
Opdracht 3	3
Opdracht 4	3
Opdracht 5	3
Opdracht 6	4
Opdracht 7	4
Opdracht 8	4
Opdracht 9	4
Opdracht 10	4
Opdracht 11	4
Opdracht 12	5
Opdracht 13	5
Opdracht 14	5
Opdracht 15	5
Opdracht 16	5
Opdracht 17	5
Opdracht 18	5
Bronnen	5
Evaluatie	6

Een Aanbevelingssysteem voor Films

De broncode bevindt zich in de **src** folder. Het algemene script (**src/s0216676_script**) is opgedeeld in secties, één per opgave. Elke opgave wordt hieronder afzonderlijk beantwoord. Aan het einde van elk antwoord wordt (indien nodig) de broncode weergegeven.

Opdracht 1

We laden de dataset in met de **load** functie. De uitvoer staat weergegeven in figuur 1.



Figuur 1: Grafische voorstelling van ijle matrices R en T .

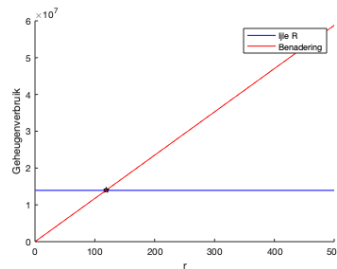
```
1 set(0, 'defaultFigurePosition', get(0, 'Screensize')); % Figuren vullen scherm
2 load('MovieLens_Subset.mat');
3 subplot(1,2,1)
4 spy(R(1:1000,1:1000))
5 subplot(1,2,2)
6 spy(T(1:1000,1:1000))
```

Opdracht 2

Stel dat gehele getallen 4 - en reële getallen 8 bytes innemen. De volle matrix `full(R)` zou dan 220388224 bytes ($\approx 210\text{MB}$) in beslag nemen ; 1 reëel getal per element. Voor de ijle matrix `sparse(R)` lijkt het op het eerste gezicht 18525728 bytes ($< 18\text{MB}$) te zijn ; 2 gehele getallen en 1 reëel getal per element dat niet nul is. MatLab verbruikt echter iets meer dan dat. Als men de juiste formule toepast voor het geheugenverbruik van ijle matrices ¹ :

$$12 \times nnz + 4 \times n$$

dan komt men uit op bijna 14 miljoen bytes. Op mijn computer was het totaal meer dan 18,6 miljoen omdat gehele getallen daar met 8 bytes worden voorgesteld.



Figuur 2: Geheugenverbruik voor ijle matrix R en lagerangsbenadering. Het snijpunt bevindt zich in $r \approx 119$.

```
1 [m,n] = size(R);
2 ratings = nnz(R);
3 int_mem = 4;
4 double_mem = 8;
5 max_r = 500;
6 %
7 fprintf('Geheugenruimte full(R) : %i\n', m * n * double_mem)
8 %
9 size_sparse_naive = ratings * (int_mem * 2 + double_mem);
10 size_sparse = 12 * ratings + 4 * n;
11 fprintf('Geheugenruimte sparse(R) : %i\n', size_sparse)
12 %
13 fullR = full(R);
14 fprintf('Matlab zelf gebruikt respectievelijk %i en %i bytes.\n', whos('fullR').bytes, ...
15         whos('R').bytes)
16 %
17 r = 1:max_r;
18 size_approx = (m + n) * double_mem * r;
19 snijpunt_r = size_sparse / ((m + n) * double_mem);
20 fprintf('Snijpunt in r = %i\n\n', snijpunt_r)
21 %
22 hold on
23 plot(r, repmat(size_sparse,1,max_r), 'b-')
24 plot(r, size_approx, 'r-')
25 plot(snijpunt_r, size_sparse, 'kp')
26 xlabel('r')
27 ylabel('Geheugenverbruik')
28 legend('Ijle R', 'Benadering', 'Location', 'northeast')
```

¹*MATLAB gebruikt het CSC formaat voor ijle matrices : “Even though MATLAB is written in C, it follows its LINPACK and Fortran predecessors and stores full matrices by columns. This organization has been carried over to sparse matrices. A sparse matrix is stored as the concatenation of the sparse vectors representing its columns. Each sparse vector consists of a floating point array of nonzero entries (or two such arrays for complex matrices), together with an integer array of row indices. A second integer array gives the locations in the other arrays of the first element in each column. Consequently, the storage requirement for an $m \times n$ reals parse matrix with nnz nonzero entries is nnz reals and $nnz + n$ integers. On typical machines with 8-byte reals and 4-byte integers, this is $12nnz + 4n$ bytes.” [1]*

Opdracht 3

Men weet dat

$$A = \sum_{i=1}^r \sigma_i u_i v_i^T \quad (1)$$

Het iteratief algoritme gaat bij elke stap σ_j , u_j en v_j bepalen zodat de Frobeniusnorm van $E_{j-1} - \sigma_j u_j v_j^T$ (de nieuwe E_j) minimaal is. Dit komt overeen met het bepalen van een afgeknotte singulierewaardenontbinding van graad 1. De σ_j dient daarbij de grootste singuliere waarde van E_{j-1} te zijn. Stel $j = 0$, dan :

$$E_1 = E_0 - \sigma_j u_j v_j^T = A - \sigma_j u_j v_j^T = \sum_{i=1}^r \sigma_i u_i v_i^T - \sigma_j u_j v_j^T = \sum_{i=2}^r \sigma_i u_i v_i^T$$

Na een aantal iteraties bekomt men uiteindelijk :

$$E_k = \sum_{i=k+1}^r \sigma_i u_i v_i^T$$

Want telkens wordt er een afgeknotte singulierewaardenontbinding bepaald van rang 1 en die komt overeen met de term horende bij de grootste singuliere waarde zodat deze term keer op keer verdwijnt na aftrekking. Totdat de laatste $r - k$ termen overblijven. Men kan opmerken dat dit niet vereist dat de u_j en v_j bekomen in stap 4 overeenstemmen met de u_i en v_i in (1). Het teken van u_j en v_j kan bijvoorbeeld verschillen.

Berekent men ten slotte $A - E_k$ dan bekomt men het gevraagde :

$$A - E_k = \sum_{i=1}^r \sigma_i u_i v_i^T - \sum_{i=k+1}^r \sigma_i u_i v_i^T = \sum_{i=1}^k \sigma_i u_i v_i^T$$

Het algoritme stopt na k stappen (mits men veronderstelt dat $k \leq r$, want anders stop the algoritme vroeger dan dat).

Opdracht 4

Gezien mijn studentnummer s0216676 is, is c gelijk aan 6. Met $\text{sizeof}(\text{double}) = 8$ zal stap 4 een totaal aan $8 * (m + n + 1)$ bytes nodig hebben. Voor stap 5 krijgt men te maken met de matrix E die maar één keer wordt gealloceerd en $8 * n * m$ bytes nodig heeft. Stap 6 verbruikt 4 bytes (voor het geheel getal). Het aantal keren dat men de stappen uitvoert maakt niet uit.

Men heeft dus te maken met $((280000 + 58000 + 1) * 8 + (280000 * 58000 * 8) + 4) / 1024^3 \approx 120$ GB wat op zich al meer is dan het optimistische $8 + 2 + (c + 1)^2 = 8 + 2 + 49 = 59$ GB. Als men met ijle voorstellingen blijft werken kan hier iets aan gedaan worden.

Opdracht 5

De drie vectoren $[i, j, x]$ die teruggegeven worden door `find(A)` verbruiken elk $\mathcal{O}(\zeta)$ geheugen. In de implementatie van de functie `s0216676_sparseModel` wordt $[x]$ zélf hergebruikt om de resultaten van de lineaire operator $P_\Omega(X)$ op te slaan. Elk element x_{ij} is gelijk aan de vermenigvuldiging van rij i van de matrix $Uk * \text{diag}(sk)$ en kolom j van Vk^T . Het volstaat dus elementsgewijs te vermenigvuldigen van rij i van Uk met sk^T en dit te vermenigvuldigen met de getransponeerde rij j van Vk .

```
1 function [P] = s0216676_sparseModel(Uk,sk,Vk,A)
2     [i,j,x] = find(A);
3     for idx = 1:nnz(A)
4         x(idx) = Uk(i(idx),:) .* sk' * Vk(j(idx),:);
5     end
6     P = sparse(i, j, x);
7 end
```

Opdracht 6

TE VRAGEN AAN ASSISTENT

```
1 B = repmat(spones(R(:)), 1, 15);
2 spy(B(1:500,:))
```

Opdracht 7

Door op voorhand $\mathcal{O}(k\zeta)$ aan geheugen vrij te maken en dan stap voor stap de kolommen van B aan te vullen met de resultaten van de `sparseModel` functie (die op zijn beurt $\mathcal{O}(\zeta)$ geheugen gebruikt) blijft het geheugengebruik binnen de opgelegde limiet.

```
1 function [s] = s0216676_optimalCoefficients(Uk,Vk,A)
2     [m,k] = size(Uk);
3     [n,~] = size(Vk);
4     B = sparse([], [], [], m * n, k, k * nnz(A)); % Allocate space
5     for j = 1:k
6         B(:,j) = reshape(s0216676_sparseModel(Uk(:,j), 1, Vk(:,j), A), [], 1); %#ok
7     end
8     [s,~] = lsqr(B, A(:));
9 end
```

Opdracht 8

Om de gemiddelde beoordeling per gebruiker te bekomen wordt de som van al diens beoordelingen genomen en gedeeld door het aantal beoordelingen dat niet nul is. Vervolgens wordt de bekomen ijle vector omgevormd tot een kolomvector.

```
1 function [mu] = s0216676_userMeans(A)
2     mu = (sum(A, 1) ./ sum(A ~= 0))'
3 end
```

Opdracht 9

Er zijn `length(mu(mu == 5))` (dwz. 13) gebruikers die een gemiddelde beoordeling hebben van 5. De laagste drie gemiddelden bedragen 0.5100, 0.5727 en 0.8983 (`sort(mu); ans(1:3)`).

Opdracht 10

Men kan de RMSE als volgt berekenen :

```
1 function [err] = s0216676_RMSE(A,B)
2     err = norm(A-B, 'fro') / sqrt(nnz(A));
3 end
```

Opdracht 11

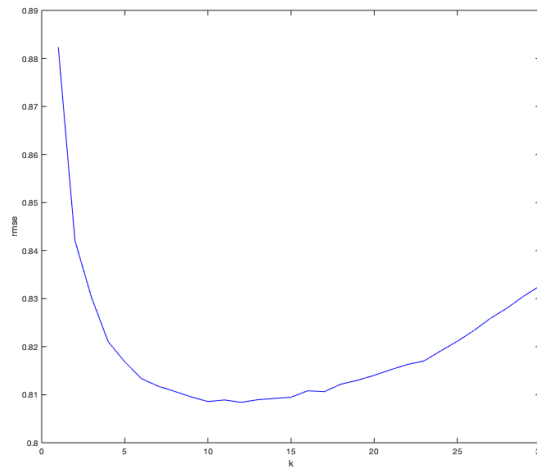
De RMSE bedraagt iets meer dan 26 :

```
1 [i,j] = find(T);
2 mu = s0216676_userMeans(R);
3 err = s0216676_RMSE(T, repmat(mu', m, 1))
4 % Prints "err = 26.6238"
```

Opdracht 12

VRAGEN ASSISTENT : DIVERGENTIE OK?

De RMSE is te zien in figuur 3. De volle berekening duurde ongeveer 250 seconden.



Figuur 3: Resultaat van het algoritme voor $k = 30$, toegepast op de MovieLens dataset. Totale duur : 256 seconden.

Opdracht 13

De vector s_{20} ziet er als volgt uit :

[1.0 806.2 386.4 383.2 284.9 248.9 192.7 199.2 179.9 171.3 133.1 152.1 136.0 143.8 122.0 120.5 128.8 116.9 120.2 124.7]

Opdracht 14

De functie kan in één lijn geschreven worden :

```
1 function [movieIDs,score] = s0216676_actualBestMovies(R)
2     [score,movieIDs] = sort(sum(R,2) ./ sum(R~=0,2), 'descend');
3 end
```

Opdracht 15

Opdracht 16

Opdracht 17

Opdracht 18

Referenties

- [1] John R. Gilbert, Cleve Moler, and Robert Schreiber. Sparse matrices in MATLAB: Design and implementation. *SIAM Journal on Matrix Analysis and Applications*, 13(1):333–356, January 1992.

Evaluatie