

نمایه گذاری معنایی پنهان *latent semantic indexing (LSI)*

نادیه آرمین

پاییز ۹۶

WT
Laboratory



آزمایشگاه فناوری وب
Web Technology Lab

Web Technology Lab
دانشگاه گیلان

فهرست مطالب

- مقدمه
 - نمایه گذاری
 - فضای برداری
- تعریف نمایه گذاری معنایی پنهان
- آشنایی با SVD
- کاربردهای SVD
- مثال
- پیاده سازی
- محاسبه شباهت ها در LSI

مقدمه_ نمایه گذاری

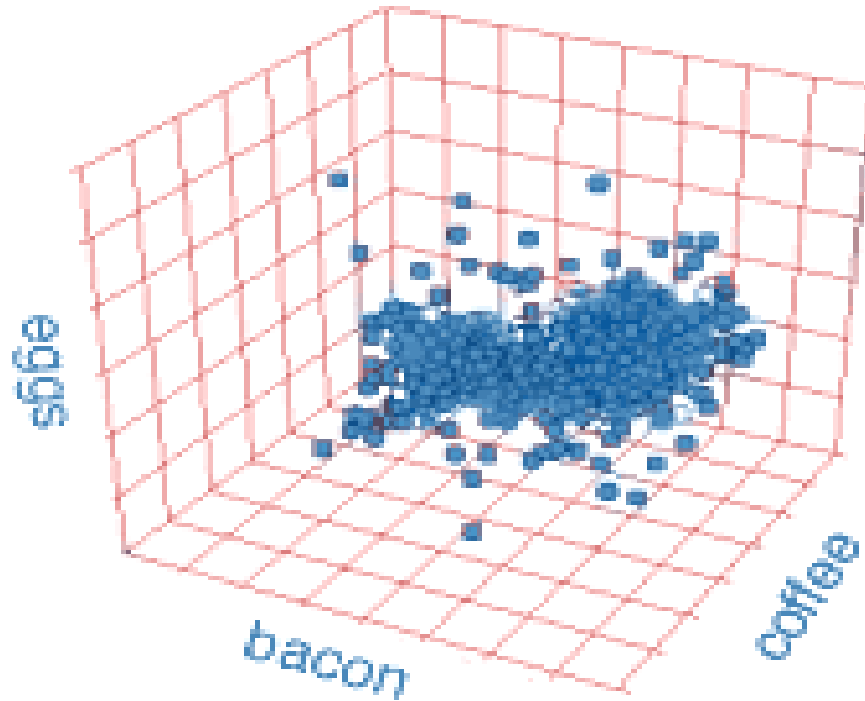
• ماتریس واژه - سند

C	d_1	d_2	d_3	d_4	d_5	d_6
ship	1	0	1	0	0	0
boat	0	1	0	0	0	0
ocean	1	1	0	0	0	0
wood	1	0	0	1	1	0
tree	0	0	0	1	0	1

- تعداد واژه

- TF_IDF

مقدمه_ فضای برداری

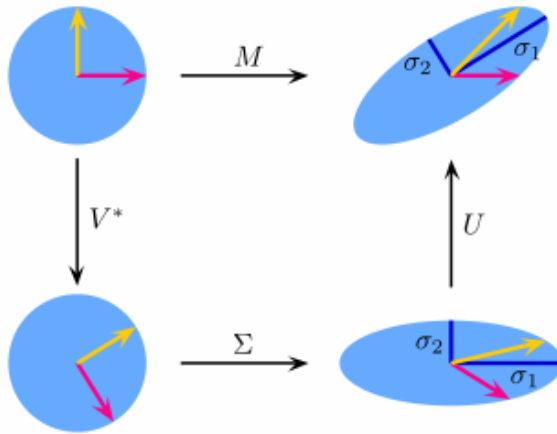


تعریف نمایه گذاری معنایی پنهان

- نگاشت ماتریس کلمه _ سند به ماتریس معنا_ سند
- مثال :
 - جستجو عبارت «درگیری های خاورمیانه»
 - احتمالاً سندی با عنوان «جنگ عراق» مفید است
 - جنگ هم معنی درگیری و عراق در خاورمیانه قرار دارد.
- مزایا:
 - عدم نیاز به استفاده از تعاریف منابع خارجی
 - پیاده سازی به وسیله روش های ریاضی

آشنایی با SVD

• تجزیه مقادیر منفرد (Singular Value Decomposition)



$$M = U \cdot \Sigma \cdot V^*$$

SVD: $M = U \Sigma V^T$ ■

■ M ماتریس سند_ واژه

■ محاسبه M'

■ M' در مقایسه با M مقایسه بهتری بین اسناد انجام میدهد

کاربردهای SVD

- کاهش ویژگی
- پردازش تصویر
- فشرده سازی
- بازیابی اطلاعات
- و ...

- SVD: $C = U\Sigma V^T$

(where C = term-document matrix)

- use the SVD to compute:

a **new, improved term-document matrix C'** .

- **better similarity** values out of C' (compared to C).

Example of $C = U\Sigma V^T$: The matrix C

C	d_1	d_2	d_3	d_4	d_5	d_6
ship	1	0	1	0	0	0
boat	0	1	0	0	0	0
ocean	1	1	0	0	0	0
wood	1	0	0	1	1	0
tree	0	0	0	1	0	1

This is a standard term-document matrix. Actually, we use a non-weighted matrix here to simplify the example.

C	d_1	d_2	d_3	d_4	d_5	d_6
ship	1	0	1	0	0	0
boat	0	1	0	0	0	0
ocean	1	1	0	0	0	0
wood	1	0	0	1	1	0
tree	0	0	0	1	0	1

=

U	1	2	3	4	5
ship	-0.44	-0.30	0.57	0.58	0.25
boat	-0.13	-0.33	-0.59	0.00	0.73
ocean	-0.48	-0.51	-0.37	0.00	-0.61
wood	-0.70	0.35	0.15	-0.58	0.16
tree	-0.26	0.65	-0.41	0.58	-0.09

×

Σ	1	2	3	4	5
1	2.16	0.00	0.00	0.00	0.00
2	0.00	1.59	0.00	0.00	0.00
3	0.00	0.00	1.28	0.00	0.00
4	0.00	0.00	0.00	1.00	0.00
5	0.00	0.00	0.00	0.00	0.39

×

V^T	d_1	d_2	d_3	d_4	d_5	d_6
1	-0.75	-0.28	-0.20	-0.45	-0.33	-0.12
2	-0.29	-0.53	-0.19	0.63	0.22	0.41
3	0.28	-0.75	0.45	-0.20	0.12	-0.33
4	0.00	0.00	0.58	0.00	-0.58	0.58
5	-0.53	0.29	0.63	0.19	0.41	-0.22

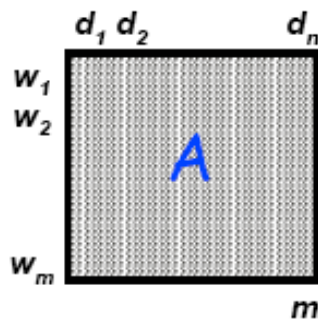
$$C = U\Sigma V^T$$

- Compute SVD of term-document matrix
- Reduce the space and compute reduced document representations
- Map the query into the reduced space
- This follows from:
- Compute similarity of q_2 with all reduced documents in $V_2 \vec{q}_2^T = \Sigma_2^{-1} U_2^T \vec{q}^T$.
- Output ranked list of docun $C_2 = U \Sigma_2 V^T \Rightarrow \Sigma_2^{-1} U^T C = V_2^T$
- Exercise: What is the fundamental problem with this approach?

محاسبه شباهت ها در LSI

- Fundamental comparisons based on SVD

- The original word-document matrix (A)



- compare two terms \rightarrow dot product of two rows of A
 – or an entry in AA^T
- compare two docs \rightarrow dot product of two columns of A
 – or an entry in $A^T A$
- compare a term and a doc \rightarrow each individual entry of A

- The new word-document matrix (A')

$U' = U_{m \times k}$
 $\Sigma' = \Sigma_k$
 $V' = V_{n \times k}$

- compare two terms $A'A'^T = (U' \Sigma' V'^T) (U' \Sigma' V'^T)^T = U' \Sigma' V'^T V' \Sigma'^T U'^T = (U' \Sigma') (U' \Sigma')^T$
 \rightarrow dot product of two rows of $U' \Sigma'$
- compare two docs $A^T A' = (U' \Sigma' V'^T)^T (U' \Sigma' V'^T) = V' \Sigma'^T U'^T U' \Sigma' V'^T = (V' \Sigma') (V' \Sigma')^T$
 \rightarrow dot product of two rows of $V' \Sigma'$
- compare a query and a doc \rightarrow each individual entry of A'

