

CS-C3240 – Machine Learning D

Feature Engineering

Stephan Sigg

Department of Communications and Networking
Aalto University, School of Electrical Engineering
stephan.sigg@aalto.fi

Version 1.0, February 14, 2022

Outline

Latent Semantic Indexing

Latent Semantic Indexing

Motivation

In information retrieval, a common task is to obtain from many documents that subset which best matches a query

Latent Semantic Indexing

Motivation

In information retrieval, a common task is to obtain from many documents that subset which best matches a query

→ Typical feature representations of documents are then term-document matrices:

Latent Semantic Indexing

Motivation

In information retrieval, a common task is to obtain from many documents that subset which best matches a query

→ Typical feature representations of documents are then term-document matrices:

Terms	Documents													
	M1	M2	M3	M4	M5	M6	M7	M8	M9	M10	M11	M12	M13	M14
abnormalities	0	0	0	0	0	0	0	1	0	1	0	0	0	0
age	1	0	0	0	0	0	0	0	0	0	0	1	0	0
behavior	0	0	0	0	1	1	0	0	0	0	0	0	0	0
blood	0	0	0	0	0	0	0	1	0	0	1	0	0	0
close	0	0	0	0	0	0	1	0	0	0	1	0	0	0
culture	1	1	0	0	0	0	0	1	1	0	0	0	0	0
depressed	1	0	1	1	1	0	0	0	0	0	0	0	0	0
discharge	1	1	0	0	0	1	0	0	0	0	0	0	0	0
disease	0	0	0	0	0	0	0	0	1	0	1	0	0	0
fast	0	0	0	0	0	0	0	0	0	1	0	1	1	1
generation	0	0	0	0	0	0	0	0	1	0	0	0	1	0
oestrogen	0	0	1	1	0	0	0	0	0	0	0	0	0	0
patients	1	1	0	1	0	0	0	1	0	0	0	0	0	0
pressure	0	0	0	0	0	0	0	0	0	0	1	0	0	1
rats	0	0	0	0	0	0	0	0	0	0	0	0	1	1
respect	0	0	0	0	0	0	0	1	0	0	0	1	0	0
rise	0	0	0	1	0	0	0	0	0	0	0	0	0	1
study	1	0	1	0	0	0	0	0	1	0	0	0	0	0

Latent Semantic Indexing

Motivation

In information retrieval, a common task is to obtain from many documents that subset which best matches a query

- Typical feature representations of documents are then term-document matrices:
- These matrices are typically huge but sparse.

Terms	Documents													
	M1	M2	M3	M4	M5	M6	M7	M8	M9	M10	M11	M12	M13	M14
abnormalities	0	0	0	0	0	0	0	1	0	1	0	0	0	0
age	1	0	0	0	0	0	0	0	0	0	0	1	0	0
behavior	0	0	0	0	1	1	0	0	0	0	0	0	0	0
blood	0	0	0	0	0	0	0	1	0	0	1	0	0	0
close	0	0	0	0	0	0	1	0	0	0	1	0	0	0
culture	1	1	0	0	0	0	0	1	1	0	0	0	0	0
depressed	1	0	1	1	1	0	0	0	0	0	0	0	0	0
discharge	1	1	0	0	0	1	0	0	0	0	0	0	0	0
disease	0	0	0	0	0	0	0	0	1	0	1	0	0	0
fast	0	0	0	0	0	0	0	0	0	1	0	1	1	1
generation	0	0	0	0	0	0	0	0	1	0	0	0	1	0
oestrogen	0	0	1	1	0	0	0	0	0	0	0	0	0	0
patients	1	1	0	1	0	0	0	1	0	0	0	0	0	0
pressure	0	0	0	0	0	0	0	0	0	0	1	0	0	1
rats	0	0	0	0	0	0	0	0	0	0	0	0	1	1
respect	0	0	0	0	0	0	0	1	0	0	0	1	0	0
rise	0	0	0	1	0	0	0	0	0	0	0	0	0	1
study	1	0	1	0	0	0	0	0	1	0	0	0	0	0

Latent Semantic Indexing

Motivation

In information retrieval, a common task is to obtain from many documents that subset which best matches a query

→ Typical feature representations of documents are then term-document matrices:

→ These matrices are typically huge but sparse.

How to identify those feature dimensions (or combinations thereof) which are most meaningful?

Terms	Documents													
	M1	M2	M3	M4	M5	M6	M7	M8	M9	M10	M11	M12	M13	M14
abnormalities	0	0	0	0	0	0	0	1	0	1	0	0	0	0
age	1	0	0	0	0	0	0	0	0	0	0	1	0	0
behavior	0	0	0	0	1	1	0	0	0	0	0	0	0	0
blood	0	0	0	0	0	0	0	1	0	0	1	0	0	0
close	0	0	0	0	0	0	1	0	0	0	1	0	0	0
culture	1	1	0	0	0	0	0	1	1	0	0	0	0	0
depressed	1	0	1	1	1	0	0	0	0	0	0	0	0	0
discharge	1	1	0	0	0	1	0	0	0	0	0	0	0	0
disease	0	0	0	0	0	0	0	0	1	0	1	0	0	0
fast	0	0	0	0	0	0	0	0	0	1	0	1	1	1
generation	0	0	0	0	0	0	0	0	1	0	0	0	1	0
oestrogen	0	0	1	1	0	0	0	0	0	0	0	0	0	0
patients	1	1	0	1	0	0	0	1	0	0	0	0	0	0
pressure	0	0	0	0	0	0	0	0	0	0	1	0	0	1
rats	0	0	0	0	0	0	0	0	0	0	0	0	1	1
respect	0	0	0	0	0	0	0	1	0	0	0	1	0	0
rise	0	0	0	1	0	0	0	0	0	0	0	0	0	1
study	1	0	1	0	0	0	0	0	1	0	0	0	0	0

Latent Semantic Indexing

Singular Value Decomposition

Any $m \times n$ matrix C can be represented as a singular value decomposition in the form $C = U\Sigma V^T$ where

U $m \times m$ matrix; columns are the orthogonal eigenvectors of CC^T

V $n \times n$ matrix; columns are the orthogonal eigenvectors of $C^T C$

Σ Diagonal Matrix with $\Sigma_{ii} = \sqrt{\lambda_i}$; $\Sigma_{ij} = 0, i \neq j$

Latent Semantic Indexing

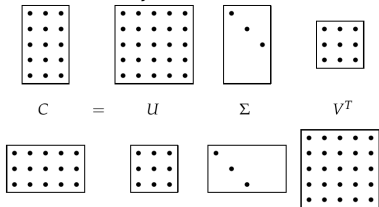
Singular Value Decomposition

Any $m \times n$ matrix C can be represented as a singular value decomposition in the form $C = U\Sigma V^T$ where

U $m \times m$ matrix; columns are the orthogonal eigenvectors of CC^T

V $n \times n$ matrix; columns are the orthogonal eigenvectors of $C^T C$

Σ Diagonal Matrix with $\Sigma_{ii} = \sqrt{\lambda_i}$; $\Sigma_{ij} = 0, i \neq j$



Latent Semantic Indexing

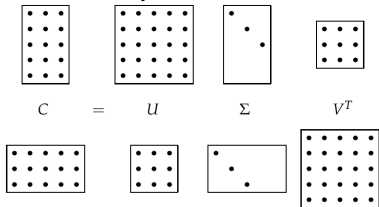
Singular Value Decomposition

Any $m \times n$ matrix C can be represented as a singular value decomposition in the form $C = U\Sigma V^T$ where

U $m \times m$ matrix; columns are the orthogonal eigenvectors of CC^T

V $n \times n$ matrix; columns are the orthogonal eigenvectors of $C^T C$

Σ Diagonal Matrix with $\Sigma_{ii} = \sqrt{\lambda_i}$; $\Sigma_{ij} = 0, i \neq j$



- First k eigenvectors map document vectors to lower dimensional representation
It can be shown that this mapping results in the k -dim. space with smallest distance to the original space

Example

	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
voyage	1	0	0	1	1	0
trip	0	0	0	1	0	1

U :

Σ :

V^T :

Example

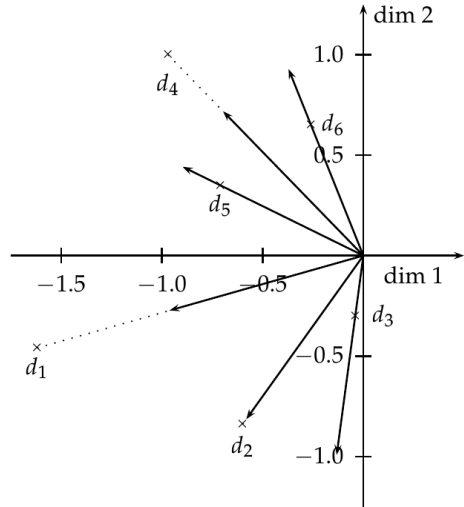
		1	2	3	4	5
U:	ship	-0.44	-0.30	0.57	0.58	0.25
	boat	-0.13	-0.33	-0.59	0.00	0.73
	ocean	-0.48	-0.51	-0.37	0.00	-0.61
	voyage	-0.70	0.35	0.15	-0.58	0.16
	trip	-0.26	0.65	-0.41	0.58	-0.09
Σ :	2.16	0.00	0.00	0.00	0.00	
	0.00	1.59	0.00	0.00	0.00	
	0.00	0.00	1.28	0.00	0.00	
	0.00	0.00	0.00	1.00	0.00	
	0.00	0.00	0.00	0.00	0.39	
V^T :		d_1	d_2	d_3	d_4	d_5
	1	-0.75	-0.28	-0.20	-0.45	-0.33
	2	-0.29	-0.53	-0.19	0.63	0.22
	3	0.28	-0.75	0.45	-0.20	0.12
	4	0.00	0.00	0.58	0.00	-0.58
	5	-0.53	0.29	0.63	0.19	0.41

Example

$\Sigma:$	2.16	0.00	0.00	0.00	0.00		
	0.00	1.59	0.00	0.00	0.00		
	0.00	0.00	0.00	0.00	0.00		
	0.00	0.00	0.00	0.00	0.00		
	0.00	0.00	0.00	0.00	0.00		
$C_2:$		d_1	d_2	d_3	d_4	d_5	d_6
	1	-1.62	-0.60	-0.44	-0.97	-0.70	-0.26
	2	-0.46	-0.84	-0.30	1.00	0.35	0.65
	3	0.00	0.00	0.00	0.00	0.00	0.00
	4	0.00	0.00	0.00	0.00	0.00	0.00
	5	0.00	0.00	0.00	0.00	0.00	0.00

Cosine-similarity

	d_1	d_2	d_3	d_4	d_5	d_6
1	-1.62	-0.60	-0.44	-0.97	-0.70	-0.26
2	-0.46	-0.84	-0.30	1.00	0.35	0.65



Questions?

Stephan Sigg

`stephan.sigg@aalto.fi`

Si Zuo

`si.zuo@aalto.fi`

Literature

- C.M. Bishop: Pattern recognition and machine learning, Springer, 2007.
- R.O. Duda, P.E. Hart, D.G. Stork: Pattern Classification, Wiley, 2001.

