

Binary Factorizations in Data Mining

Kick-off meeting
27 October 2016



Today's Agenda

- Short intro to binary factorizations
- Check attendance
- Goals of the seminar
- Organization of the seminar
- Grading & guidelines

First Things First

- A block seminar
 - Preliminary work + one (or two) day(s) of presentations
- 7 ECTS credits
- Meeting all DLs and attending all talks is **mandatory** for passing the seminar
- Attending this kick-off meeting is mandatory

Short intro

Matrix multiplication

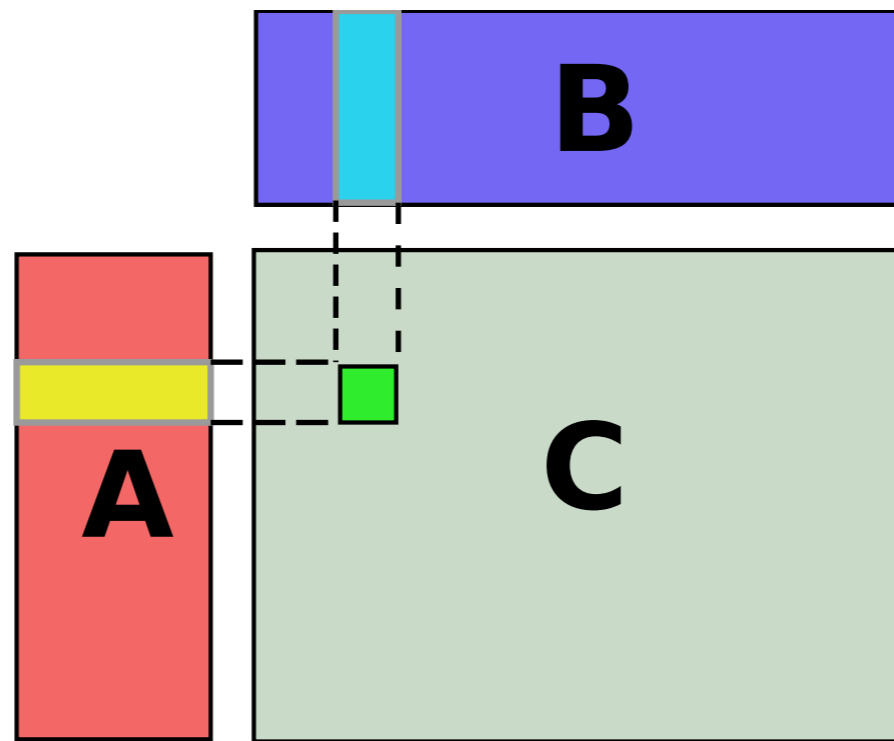
- The product of two matrices, **A** and **B**, is defined element-wise as

$$(\mathbf{AB})_{ij} = \sum_{\ell=1}^k a_{i\ell} b_{\ell j}$$

- The number of columns in **A** and number of rows in **B** must agree
 - inner dimension

Intuition for Matrix Multiplication

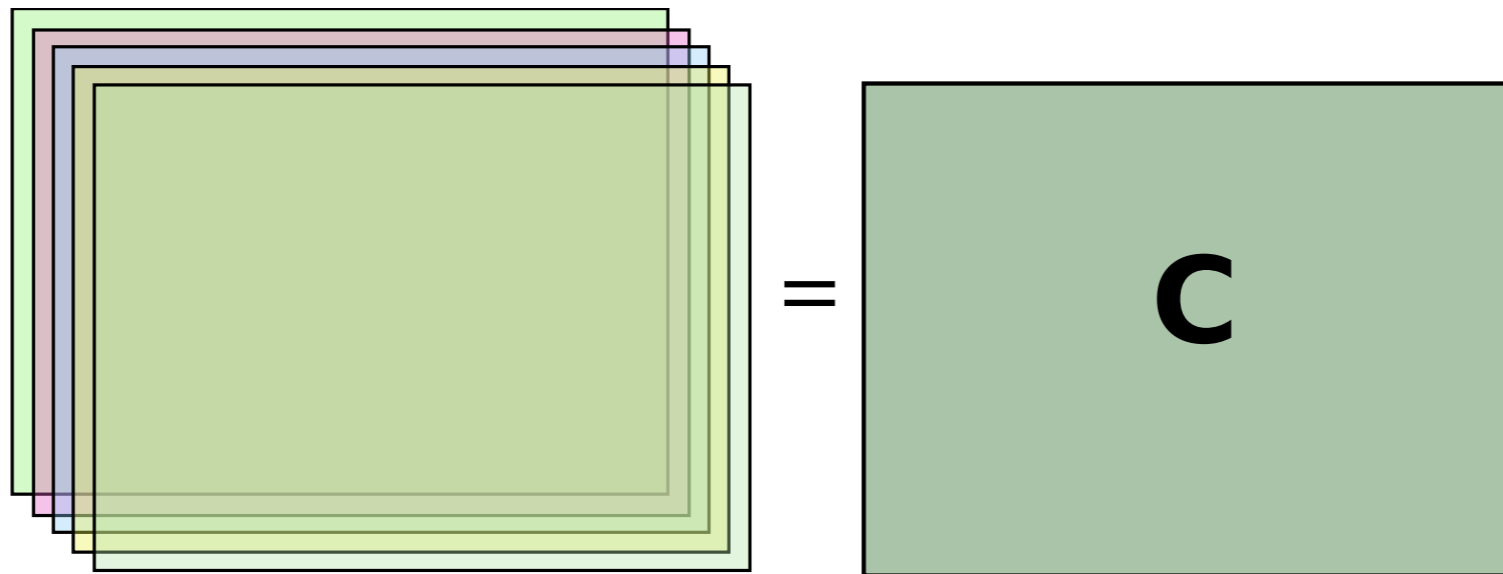
- Element $(\mathbf{AB})_{ij}$ is the inner product of row i of \mathbf{A} and column j of \mathbf{B}



$$c_{ij} = \sum_{\ell=1}^k a_{i\ell} b_{\ell j}$$

Intuition for Matrix Multiplication

- Matrix \mathbf{AB} is a sum of k matrices $\mathbf{a}_l \mathbf{b}_l^T$ obtained by multiplying the l -th column of \mathbf{A} with the l -th row of \mathbf{B}



$$\mathbf{C} = \sum_{\ell=1}^k \mathbf{a}_{\ell} \mathbf{b}_{\ell}^T$$

Matrix decompositions

- A **decomposition** of matrix **A** expresses it as a product of two (or more) **factor matrices**
 - **$A = BC$**
- Every matrix has decomposition **$A = AI$** (or **$A = IA$** if $n < m$)
- The size of the decomposition is the inner dimension of the product

Matrices in data mining

	Bread	Butter	Beer
Anna	1	1	0
Bob	1	1	1
Charlie	0	1	1

Customer transactions

	Data	Matrix	Mining
Book 1	5	0	3
Book 2	0	0	7
Book 3	4	6	5

Document-term matrix

	Avatar	The Matrix	Up
Alice		4	2
Bob	3	2	
Charlie	5		3

Incomplete rating matrix

	Jan	Jun	Sep
Saarbrücken	1	11	10
Helsinki	6.5	10.9	8.7
Cape Town	15.7	7.8	8.7

Cities and monthly temperatures

Matrix decompositions in data mining

- A common goal in data mining is to find regularities (or patterns) in the data
 - Often, to summarize the data
- A *matrix decomposition* presents the data as a sum of “simple” elements, i.e. patterns
 - but there’s also other uses... *stay tuned!*

Text mining and pLSA

- Consider a document-term matrix \mathbf{A}
 - a_{ij} is the number of times term j appears in document i

Can we find these topics automatically?

	Environmet				
	air	water	pollution	democrat	republican
doc 1	3	2	8	0	0
doc 2	1	4	12	0	0
doc 3	0	0	0	10	11
doc 4	0	0	0	8	5
doc 5	1	1	1	1	1

Politics

pLSA example

air wat pol dem rep

0.04	0.03	0.12	0	0
0.01	0.06	0.17	0	0
0	0	0	0.14	0.16
0	0	0	0.12	0.07
0.01	0.01	0.01	0.01	0.01

A

Here, A is normalized

0.39	0
0.52	0
0	0.58
0	0.36
0.09	0.06

W

How strong
the topic is
in the document?

0.48	0
0	0.52

Σ

Overall
frequency

air wat pol dem rep

0.15	0.21	0.64	0	0
0	0	0	0.53	0.47

H

How strong the
word is in the
topic?

Binary factorizations

- Often we deal with **binary** matrices
 - presence/absence data, graph adjacency matrices, market basket data, ...
- Such matrices are commonly better decomposed into binary factor matrices
 - Or one matrix binary, other numerical
 - Interpretability, sparsity, speed, ...
- Also the algebra might change to Boolean

Binary factorizations



long-haired

well-known

male

✓ 1	✓ 1	✗ 0
✓ 1	✓ 1	✓ 1
✗ 0	✓ 1	✓ 1

Example factorization

$$\begin{pmatrix} 1 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \end{pmatrix} \otimes \begin{pmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} \vee \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 1 & 1 \end{pmatrix}$$

The diagram illustrates the factorization of a 3x3 matrix into two 3x3 matrices. The original matrix is shown with red and blue boxes highlighting its structure. It is then expressed as the Kronecker product of two 2x2 matrices, also with red and blue boxes. Finally, the result is shown as the element-wise OR of two 3x3 matrices, with columns labeled A, B, and C. The first matrix has columns A and B boxed in red, and the second matrix has columns B and C boxed in blue.

Papers in this seminar...

- Factorize binary (or ternary) matrices
- Have at least one factor matrix discrete/
binary
- Either present new binary factorisation
methods, or apply the binary factorisations to
other problems
 - Or both

Head Count

Goals

- To learn how to read and understand recent research literature
- To learn how to write a concise report of a research article
- To learn how to present research
- To boldly read what no one (at this seminar) has read before
- To keep young people out of streets

Workflow

1. You read the paper (+ other papers)
2. You write a draft report and send it to me
3. I comment your report
4. You improve your report and prepare your presentation, which you also send to me
5. I comment your presentation
6. You improve your presentation and send me the final report
7. I distribute the reports to everybody
8. You read others' reports
9. You present your work and follow and discuss others' presentations

Schedule

Day	Topic
27 October	Kick-off
27 November	Report draft DL
9 December	Slides draft DL
Early January	Report DL
January/February	Seminar

Selecting the dates

- The seminar takes two full days (approx. 9:00–16:00)
 - The days have to be consecutive
- I've created a doodle where you can indicate which days would work with you
 - <http://doodle.com/poll/bgst3hwms83vdwqn>
 - Fill in the doodle **by Wednesday, 2 November**
 - Use *no* if you have hard constraints and *if-need-be* if you have soft constraints
 - The more *no* answers you have, the more likely I'm to violate them

Odds and ends

- Remember to register to HISPOS
 - You can de-register for three weeks from now, but not after that
- The papers are now behind a username and password
 - see whiteboard or contact the lecturer

Grading Overview

- Report (3–5 pages):
 - Correctness, connections, criticism, style
- Slides & presentation (20 min):
 - Delivery, clearness, presentation skills
- Discussion (5–10 min):
 - Participation, correctness, connections

Report

- 3–5 pages
- **In your own words**
 - No verbatim copy
- Explain the main ideas of the paper
 - Research questions
 - Proposed solutions and their evaluation
- Write to somebody who hasn't read the paper
- Provide also extra connections & criticism towards the approach

More on report

- Length is not a hard constraint
 - As long as needed, but no longer
 - Figures, math, choice of margins, etc. effect the final length
- You must cite your sources using the standard academic practices
 - Everything that's not your own must be cited
- I only accept reports in PDF format
- More information is send later in the semester

Presentation

- 20 minutes (+ 5–10 minutes of discussion)
- Explain the main ideas of the paper
 - What, why, and how
- Build connections and provide criticism where appropriate
- Target to audience that knows CS and basics of DM and ML
 - Some audience knows more on your topic, others less; cater for both groups

Discussion

- After you've followed a presentation, discuss
 - You know the topic a bit: you've read the report
- Ask if something was left unclear
- Tell if you know more on something
- After technical part, give constructive feedback on the presentation itself
- Actually discussing is **mandatory**

Lecture on giving presentations?

- If there's demand, I might give one lecture on giving presentations
 - No soft-skills seminar, just few tidbits on how to give a presentation and how to prepare slides in this seminar
- Schedule: November