



—

BC COMS 2710: Computational Text Analysis

—

Lecture 19 Dimensionality Reduction



- Readings 05:
 - No more this semester – congrats!

- HW02
 - Mostly done grading
 - You all did very well!!!

- HW04
 - Due Thursday



- P/D/F deadline is today
- Course Evaluations
 - Due 06/14

Final Project – Deliverables



- Project ideation – Friday May 28st
 - 5 points

- Project proposal – Sunday June 6th
 - 9 points

- Project presentations – Monday June 14th
 - 6 points

- Project submissions – Friday June 18th
 - 15 points

- http://coms2710.barnard.edu/final_project

Project Presentations – Monday June 14th



5 minute presentations by each group

Format:

- Research Question
- Motivation
 - Why should we care?
- Data Collected
 - Where did the data come from? How did you collect it?
 - What filtering was done?
 - Resulting corpus:
 - How many documents? Average size of documents? Vocabulary size?
- Results (preliminary)
 - figures

Goal:

- Publicizing your research is important
- You all to see what everyone else is working on



“Congratulations, offer letter from Google” Spam

“Congratulations, you won the lottery” Not Spam

Which mistake is worse?



Recall

- When we do not want false negatives

Precision

- When we do not want false positives



Logistic Regression

Summary of Logistic Regression



- Optimizes $P(Y | X)$ directly
- Define the **features**
- Learn a vector of **weights** for each label $y \in Y$
 - Gradient descent, update weights based on error
- Multiple feature vector by weight vector
- Output is $P(Y = y | X)$ after normalizing
- Choose the most probable Y



- Optimizes $P(Y | X)$ directly
- Define the **features**
- Learn a vector of **weights** for each label $y \in Y$
 - Gradient descent, update weights based on error
- **Multiple feature vector by weight vector**
- Output is $P(Y = y | X)$ after normalizing
- Choose the most probable Y

Scoring one document - dot product



$$\begin{aligned} [f_1, f_2, f_3] \cdot [w_1, &= (f_1 \times w_1) + (f_2 \times w_2) + (f_3 \times w_3) \\ w_2, & \\ w_3] &= \sum (f_i \times w_i) \end{aligned}$$

Score two documents



$$\begin{bmatrix} [f_{1,1}, f_{1,2}, f_{1,3}] \\ [f_{2,1}, f_{2,2}, f_{2,3}] \end{bmatrix} \cdot \begin{bmatrix} w_1, \\ w_2, \\ w_3 \end{bmatrix}$$

$$= \begin{bmatrix} (f_{1,1} \times w_1) + (f_{1,2} \times w_2) + (f_{1,3} \times w_3), \\ (f_{2,1} \times w_1) + (f_{2,2} \times w_2) + (f_{2,3} \times w_3) \end{bmatrix}$$



We can multiply two matrices A and B if
number of columns in A = number of rows in B

The size of the resulting matrix is
number of rows in A & the number of columns in B

Matrix Multiplication



$$\begin{bmatrix} 1 & 7 \\ 2 & 4 \end{bmatrix} \cdot \begin{bmatrix} 3 & 3 \\ 5 & 2 \end{bmatrix}$$

A

B

[Khan Academy](https://www.khanacademy.org)

Matrix Multiplication



$$\begin{array}{c} \vec{a}_1 \rightarrow \\ \vec{a}_2 \rightarrow \end{array} \begin{array}{c} \\ \\ \end{array} \begin{bmatrix} 1 & 7 \\ 2 & 4 \end{bmatrix} \cdot \begin{array}{c} \vec{b}_1 \quad \vec{b}_2 \\ \downarrow \quad \downarrow \\ \begin{bmatrix} 3 & 3 \\ 5 & 2 \end{bmatrix} \end{array}$$

$A \qquad B$

Matrix Multiplication



$$\begin{array}{c} \vec{a}_1 \rightarrow \\ \vec{a}_2 \rightarrow \end{array} \begin{array}{c} \vec{b}_1 \quad \vec{b}_2 \\ \downarrow \quad \downarrow \end{array} \begin{array}{c} \left[\begin{array}{cc} 1 & 7 \\ 2 & 4 \end{array} \right] \cdot \left[\begin{array}{cc} 3 & 3 \\ 5 & 2 \end{array} \right] = \left[\begin{array}{cc} \vec{a}_1 \cdot \vec{b}_1 & \vec{a}_1 \cdot \vec{b}_2 \\ \vec{a}_2 \cdot \vec{b}_1 & \vec{a}_2 \cdot \vec{b}_2 \end{array} \right] \\ A \quad B \quad C \end{array}$$

Matrix Multiplication Example



$$\begin{matrix} A \\ \begin{bmatrix} 1, & 2, & 3 \\ 4, & 5, & 6 \end{bmatrix} \end{matrix} \cdot \begin{matrix} B \\ \begin{bmatrix} 7, & 8, \\ 9, & 10, \\ 11, & 12 \end{bmatrix} \end{matrix} = \begin{matrix} C \end{matrix}$$

Question: What are the dimension of C?
2 rows x 2 columns



Document—Term Matrix



DMT:

- Rows represent a document
- Columns represent a word
- Values represent some feature of word w_i in document d_j

	w_1	w_2	w_3	w_4	w_v
d_1									
d_1									
...									
d_n									10

Properties of Document-Term Matrix



- Sparse matrix
 - Most values are 0
- Very large
 - Many, many, many columns
- Noisy

	w_1	w_2	w_3	w_4	w_v
d_1									
d_1									
...									
d_n									



- Make values in each cell more meaningful
- Reduce the size of the matrix
 - Dimensionality reduction
- Remove noise



— Matrix Factorization/ Dimensionality Reduction —

Abstract thought



<https://www.youtube.com/watch?v=dROx9Djr7mk>

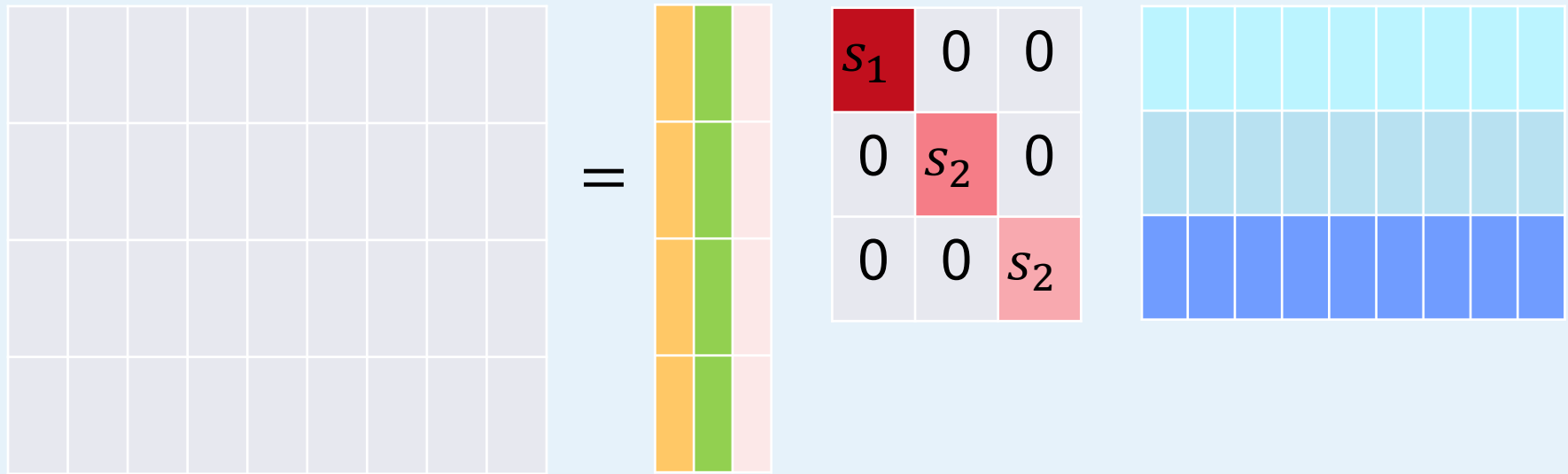


Singular Value Decomposition



$$\begin{array}{c} \mathbf{M} \\ n \times v \end{array} = \begin{array}{c} \mathbf{U} \\ n \times k \end{array} \begin{array}{c} \mathbf{S} \\ k \times k \end{array} \begin{array}{c} \mathbf{V}^T \\ k \times v \end{array}$$

Singular Value Decomposition



$$\begin{matrix} \mathbf{M} \\ n \times v \end{matrix} = \begin{matrix} \mathbf{U} \\ n \times k \end{matrix} \begin{matrix} \mathbf{S} \\ k \times k \end{matrix} \begin{matrix} \mathbf{V} \\ k \times v \end{matrix}$$



- Applying SVD to the DTM is called Latent Semantic Analysis

- The name of LDA is based on this



- Monday 06/07 - Matrix Factorization
- Tuesday 06/08 – Word Embeddings
- Wednesday 06/09 – Guest Lecture
 - Attendance required
- Thursday 06/10 – ngrams & phrases