

Practical Text Analytics: Latent Semantic Analysis

Fan Dai

Iowa State University

October 14, 2019

Recall: Text analytics process

Text analytics process

- Planning the text analytics projects
- Preparing and preprocessing text
- **Analyzing data**
 - ① *Latent Semantic Analysis (LSA)*
(Chapter 6, Anandarajan et al.(2019))
- Interpreting results

LSA: Motivating Example



Source: <https://kids.nationalgeographic.com/>

- What a computer can learn from text describing "cheetah" by looking at word frequency, proximity. . . ?
 - "cheetah" and "cat" are semantically related.
 - "cheetah" and "fastest" are more closely related than "cat" and "fastest".
- The computer makes the connection that *cheetah is the fastest cat*.

LSA: Definition

- The underlying idea

Extract and reveal information that conveyed from how words co-occur with each other across documents

- 1 Reflect a shared latent concept, e.g., {"fastest", "cat"} \Rightarrow "cheetah".
- 2 Classify the documents, e.g., {"dog", "cat", "apple", "blueberry", "orange"} \Rightarrow group documents into "animal" and "fruit"

- The analysis object: tokens

- 1 Meaning of words
- 2 Relationships with other tokens

LSA: singular value decomposition (SVD)

How to do when two words are related through a third word only?

e.g., Two words from Iowa's winery blogs: "Ackerman", "Tassel"

- Rarely together within the same document.
- Related through their frequent shared co-occurrence with other terms like "producer" or "price".

Solution:

Singular value decomposition (SVD) on weighted the term-document matrix (TDM)

TDM

- Rows: words/terms. Columns: documents
- Entries: the number of times that the i th term appears in the j th documents.

LSA: singular value decomposition (SVD)

Local weight for word i in document j		Global weight for word i	
Raw	$tf_{i,j}$	None	1
Binary	$\begin{cases} tf_{i,j} \geq 1: 1 \\ tf_{i,j} = 0: 0 \end{cases}$	IDF	$1 + \log_2 \left(\frac{n}{df(i)} \right)$
Log	$\log_2(tf_{i,j} + 1)$	Entropy	$1 + \sum_{j=document}^{corpus} \frac{\frac{tf_{i,j}}{gf_i} \times \log_2 \left(\frac{tf_{i,j}}{gf_i} \right)}{\log_2(n)}$
$tf_{i,j}$	Term frequency: number of times word i appears in the document j		
df_i	Document frequency: Number of documents word i appears in at least once in the corpus		
gf_i	Global frequency: Number of times word i appears across the entire corpus		
n	Number of documents in the corpus		
$Final\ weight = local\ weight \times global\ weight$			

Figure: Weighting for TDM

Source: **David Gefen; et al. (2017)**

tfidf-weighted TDM is preferred in LSA

LSA: singular value decomposition (SVD)

How does SVD work?

- identify underlying factors in the weighted TDM by transforming it into three matrices that represent terms, documents, and a matrix multiplier for reconstruction, respectively.

$$M = U \times \Sigma \times V^T$$

- M is the weighted TDM
- U contains the left singular vectors of terms
- Σ is a matrix with weight values on the diagonal
- V contains the right singular vectors of documents

LSA: singular value decomposition (SVD)

- in LSA, we apply a truncated SVD

$$M \approx A_k = U_k \times \Sigma_k \times V_k^T$$

A_k represents the LSA space out of the rank r matrix M

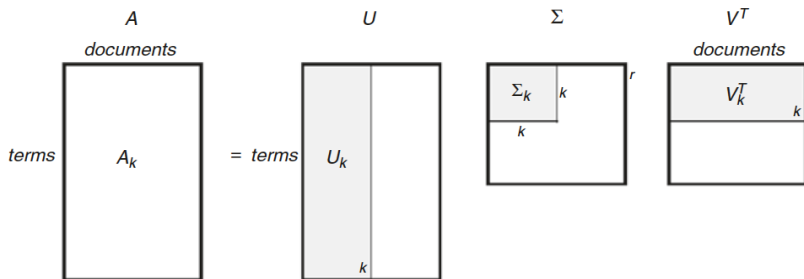


Figure: Truncated SVD process

Source: Martin and Berry (2007)

LSA: singular value decomposition (SVD)

- the k remaining singular vectors in U and V correspond to k "hidden concepts" where the terms and documents participate.
- k is too small: vectors that are related conceptually are not combined
- k is too large: redundant information is included (with those singular values that are "too small" and thus "negligible").
- determine k : e.g., scree plot showing variance explained by number of singular vectors.

- SVD v.s. Principal component analysis (PCA): *reference link*

LSA: singular value decomposition (SVD)

What can we do with the SVD in LSA?

Example: 10 respondents and descriptions of their dog (Anandarajan et al.(2019)):

- *Document 1: My Favorite Dog Is Fluffy and Tan.*
- Document 2: the dog is brown and cat is brown.*
- Document 3: My favorite coat is brown and hat is pink*
- Document 4: My dog has a hat and leash. ♡*
- ⋮*
- Document 10: MY fluffy dog has a white coat and hat!*

LSA: singular value decomposition (SVD)

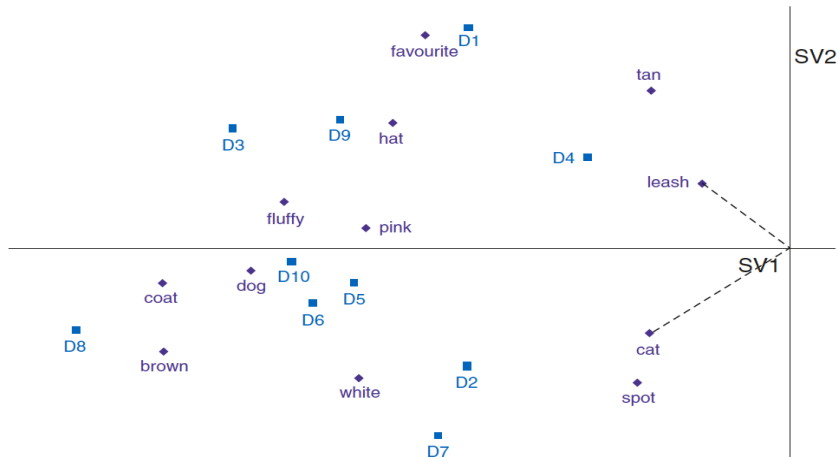
tfidf-weighted TDM

	D1	D2	D3	D4	D5	D6	D7	D8	D9	D10
<i>Brown</i>	0.0	3.0	1.5	0.0	1.5	3.0	1.5	1.5	1.5	0.0
<i>Cat</i>	0.0	4.3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
<i>Coat</i>	0.0	0.0	2.0	0.0	4.0	2.0	0.0	2.0	0.0	2.0
<i>Dog</i>	1.3	1.3	0.0	1.3	0.0	1.3	1.3	2.6	1.3	1.3
<i>Favorite</i>	2.7	0.0	2.7	0.0	0.0	0.0	0.0	0.0	2.7	0.0
<i>Fluffy</i>	1.7	0.0	0.0	0.0	1.7	1.7	0.0	1.7	1.7	1.7
<i>Hat</i>	0.0	0.0	2.3	2.3	0.0	0.0	0.0	0.0	2.3	2.3
<i>Leash</i>	0.0	0.0	0.0	4.3	0.0	0.0	0.0	0.0	0.0	0.0
<i>Pink</i>	0.0	0.0	3.3	0.0	0.0	0.0	0.0	3.3	0.0	0.0
<i>Spot</i>	0.0	0.0	0.0	0.0	0.0	0.0	4.3	0.0	0.0	0.0
<i>Tan</i>	4.3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
<i>White</i>	0.0	0.0	0.0	0.0	0.0	0.0	2.7	2.7	0.00	2.7

Source: **Practical text analytics: maximizing the value of text data**

LSA: singular value decomposition (SVD)

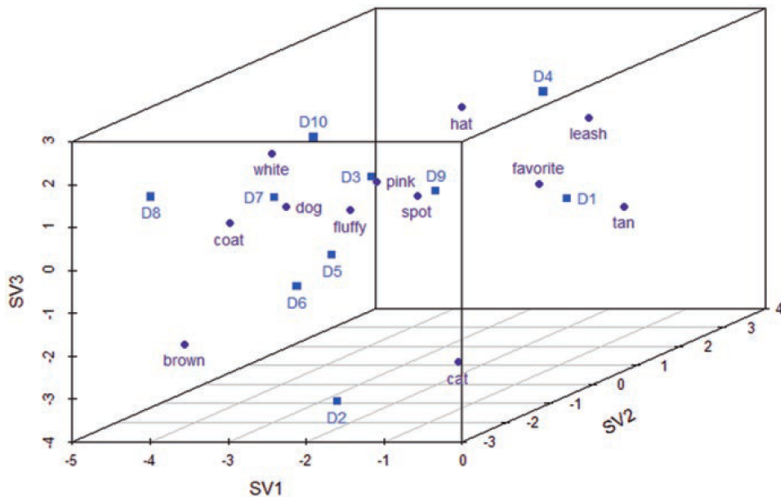
Apply truncated SVD with $k = 2$ and plot each terms and documents according to the 2 singular vectors in U and V , respectively.



Source: **Practical text analytics: maximizing the value of text data**

LSA: singular value decomposition (SVD)

Apply truncated SVD with $k = 3$.



Source: Practical text analytics: maximizing the value of text data

LSA: cosine similarity

Measure the similarity between two vectors in the LSA space.

- Cosine similarity: $\cos(\mathbf{v}_1, \mathbf{v}_2) = \frac{\mathbf{v}_1^\top \mathbf{v}_2}{\|\mathbf{v}_1\| \|\mathbf{v}_2\|}$, which can be applied to terms, documents (or both), and queries.

Example: Term-term cosine similarity measures via using the row vectors from the LSA space, A_3 , corresponding to each of the terms.

	<i>brown</i>	<i>cat</i>	<i>coat</i>	<i>dog</i>	<i>favorite</i>	<i>fluffy</i>	<i>hat</i>	<i>leash</i>	<i>pink</i>	<i>spot</i>	<i>tan</i>	<i>white</i>
<i>brown</i>	0.0											
<i>cat</i>	0.8	0.0										
<i>coat</i>	0.9	0.3	0.0									
<i>dog</i>	0.8	0.3	1.0	0.0								
<i>favorite</i>	0.3	0.0	0.5	0.5	0.0							
<i>fluffy</i>	0.8	0.3	0.9	0.9	0.8	0.0						
<i>hat</i>	0.3	-0.4	0.7	0.7	0.8	0.8	0.0					
<i>leash</i>	-0.3	-0.8	0.2	0.3	0.4	0.3	0.8	0.0				
<i>pink</i>	0.7	0.1	1.0	1.0	0.6	1.0	0.9	0.5	0.0			
<i>spots</i>	0.3	0.0	0.5	0.5	-0.5	0.2	0.1	0.2	0.4	0.0		
<i>tan</i>	0.2	0.0	0.2	0.2	1.0	0.5	0.6	0.2	0.4	-0.7	0.0	
<i>white</i>	0.5	-0.1	0.8	0.8	-0.1	0.5	0.5	0.4	0.7	0.9	-0.4	0.0

Source: **Practical text analytics: maximizing the value of text data**

LSA: cosine similarity

- Query: a scaled, weighted sum of the component term vectors.

$$\text{query} = q^{\top} U_k \Sigma_k^{-1}$$

- LSA uses the cosine measures to find documents that are similar to words that designated as query terms

Example: Cosine values between the query (brown, pink, tan) and documents.

<i>brown</i>	<i>cat</i>	<i>coat</i>	<i>dog</i>	<i>favorite</i>	<i>fluffy</i>	<i>hat</i>	<i>leash</i>	<i>pink</i>	<i>spot</i>	<i>tan</i>	<i>white</i>
1	0	0	0	0	0	0	0	1	0	1	0

Document	Cosine
6	0.81
5	0.78
9	0.73
2	0.71
1	0.69
3	0.66
8	0.24
10	-0.08
7	-0.30
4	-0.30

LSA: summaries

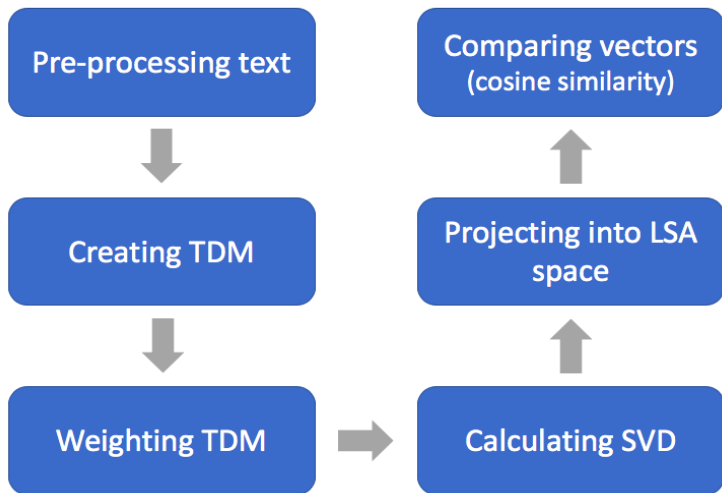


Figure: Main Steps for latent semantic analysis