# TEXT PROCESSING: LATENT SEMANTIC INDEXING -Applied Multivariate Analysis-

Lecturer: Darren Homrighausen, PhD

# A REPEATED THEME

Often, dimension reduction can be 'layered' between our original representations (in this case bag of words) and techniques that operate on that representation (in this case, finding relevant documents)

ORIGINAL REPRESENTATION [dimension reduction] STATISTICAL METHOD

How does this relate to text processing?

# PCA AND FACTOR ANALYSIS INTERPRETATION

• PCA: Find the directions of greatest variance. This doesn't on its face seem like it maintains correlations, but observe:

$$var(aX_1 + bX_2) = a^2 Var(X_1) + b^2 Var(X_2) + 2abCov(X_1, X_2)$$

If we standardize the matrix, then this reduces to

$$var(aX_1 + bX_2) = a^2 + b^2 + 2abCov(X_1, X_2)$$

This gets maximized over  $a^2 + b^2 = 1$ .

- If  $Cov(X_1, X_2) \approx 0$ , then this gets maximized by any  $a^2 + b^2 = 1$  (it doesn't matter)
- If  $Cov(X_1, X_2) \approx 1$ , then this gets maximized by setting  $a = b = 1/\sqrt{2}$
- FACTOR ANALYSIS: Defined by maintaining correlations.

So, in either case, we are really maintaining correlations

#### GRAPHICAL EXAMPLE OF THE PHENOMENON



FIGURE: Left: sig = 0. Right: sig =  $(999)^{1} + (399)^{2}$ 

Think about the document-term matrix  $\mathbb{X}$ .

The columns correspond to the words. If two words  $w_1, w_2$  commonly appear together then the  $w_1^{th}$  and  $w_2^{th}$  columns of X are correlated

When we have a large number of documents that are about a topic, it is common to have some, or most, of the documents using related, but not identical words

Therefore, if we were to search  $\mathbb X$  with a query, we would miss some of the important documents

#### AN EXAMPLE

Suppose we have a corpus of documents

We wish to search for documents containing agriculture

We can query Y = ( "agriculture" )

However, "agriculture" is not regularly explicitly mentioned in articles about agriculture

This is where correlations come in. Whenever agriculture is mentioned, it will occur very frequently along with many synonyms ("farming", for instance)

This is where latent semantic indexing comes in

#### AN EXAMPLE, GIVEN TO US BY AN INVISIBLE HAND

To see why it is called latent semantic *indexing*, observe the following

When a book is written, a list of terms (or topics) is written down and an index is formed saying where these terms appear. For example, here is the start to the entry for "Agriculture" in the index to *The Wealth of Nations* 

> AGRICULTURE, the labour of, does not admit of such subdivisions as manufactures, 6; this impossibility of separation, prevents agriculture from improving equally with manufactures, 6; natural state of, in a new colony, 92; requires more knowledge and experience than most mechanical professions, and yet is carried on without any restrictions, 127; the terms of rent, how adjusted between landlord and tenant, 144; is extended by good roads and navigable canals, 147; under what circumstances pasture land is more valuable than arable, 149; gardening not a very gainful employment, 152–3; vines the most profitable article of culture, 154; estimates of profit from projects, very fallacious, ib; cattle and tillage mutually improve each other, 220; ...

It is asking a lot for a computer to do this.

However, if we only want to get the pages where "agriculture" is the topic (like, 6, 92, 152–3, 220..), then we can make a document-term matrix out of the pages of the book.

This approach will fail if we search this document-term matrix directly.

However, asking for pages that contain highly correlated words (like "rent") should work very well

# LATENT SEMANTIC INDEXING (LSI)

If we have our document-term matrix  $\mathbb{X},$  then we write  $\mathbb{X} = UDV^{\top},$  where

- The matrix *U* is the concept-document matrix (and maps into the document space)
- The matrix V is the term-concept matrix (and maps into the term space)
- V is the matrix of loadings of the original words

If we have our query Y, we can map it into the document space by thinking of it as a new row in  $\mathbb X$ 

$$\mathbb{X} = UDV^{\top}$$
 is the same as  $\mathbb{X}VD^{-1} = U$ 

Which means we transform Y as

$$YVD^{-1}$$

#### TO CENTER OR NOT TO CENTER?

If we think about LSI as performing PCA, then we should technically do the SVD like

 $\mathbb{X} - \overline{\mathbb{X}} = UDV^{\top}$ 

However, observe the following example

$$A = \left[ \begin{array}{rrrr} a & b & 0 \\ d & 0 & f \\ 0 & h & i \end{array} \right]$$

(日)

10

What happens if we column center A?

#### TO CENTER OR NOT TO CENTER?

If we think about LSI as performing PCA, then we should technically do the SVD like

 $\mathbb{X} - \overline{\mathbb{X}} = UDV^{\top}$ 

However, observe the following example

$$A = \left[ \begin{array}{rrrr} a & b & 0 \\ d & 0 & f \\ 0 & h & i \end{array} \right]$$

What happens if we column center A?

We lose sparsity!

#### TO CENTER OR NOT TO CENTER?

If we think about LSI as performing PCA, then we should technically do the SVD like

 $\mathbb{X} - \overline{\mathbb{X}} = UDV^{\top}$ 

However, observe the following example

$$A = \left[ \begin{array}{rrrr} a & b & 0 \\ d & 0 & f \\ 0 & h & i \end{array} \right]$$

What happens if we column center A?

We lose sparsity!

The consensus is to column center if your data is small enough, otherwise, don't worry about it

#### EXAMPLE DATASET

Let's look at D = 20 Reuters news articles about crude oil production and importation.

#### The corpus has 860 words

"Diamond Shamrock Corp said that\neffective today it had cut its contract prices for crude oil by\n1.50 dlrs a barrel.\n The reduction brings its posted price for West Texas\nIntermediate to 16.00 dlrs a barrel, the copany said.\n \"The price reduction today was made in the light of falling\noil product prices and a weak crude oil market,\" a company\nspokeswoman said.\n Diamond is the latest in a line of U.S. oil companies that\nhave cut its contract, or posted, prices over the last two days\nciting weak oil markets.\n Reuter"

# Example dataset: R

Let's look at

```
mydtm = as.matrix(dtm.stem)
```

```
out.pca = prcomp(mydtm)
```

```
out.lsi = out.pca$rotation
```

```
signif(sort(out.lsi[,1],decreasing=TRUE)[1:24],2)
signif(sort(out.lsi[,1],decreasing=FALSE)[1:24],2)
```

# PC LOADINGS, FOR REUTERS DOCUMENTS

> signif(sort(out.lsi[,1],decreasing=TRUE)[1:24],2)

january	cubic :		fiscales	petroliferos	yacimi	entos	billion
0.5000	0.3200		0.3200	0.3200	0	.3200	0.2500
argentine	ne gas		metrers	metres	nat	tural p	produced
0.1600	0.1600		0.1600	0.1600	0	.1600	0.1600
totalled	ba	arrels	added	production		pct	mln
0.1600	.1600 0.0940		0.0860	0.0850	0	.0780	0.0770
output	ł	oudget	riyals	abdul-aziz	expend	iture	revenue
0.0630	(	0.0061	0.0061	0.0051	0	.0051	0.0041
> signif(sor	t(out.ls	si[,1],decı	reasing=F <i>l</i>	ALSE) [1:24],2]	)		
posted	canada	canadian	west	z power	bbl	lowered	texaco
-0.050	-0.044	-0.044	-0.041	L -0.040	-0.040	-0.036	-0.033
texas	brings	effective	dlrs	s grade	sweet	contract	ship
-0.033	-0.032	-0.032	-0.031	L -0.031	-0.031	-0.031	-0.030
price	changed	pay	postings	s decrease	company	benchmark	feb
-0.030	-0.028	-0.028	-0.028	-0.027	-0.027	-0.026	-0.026

These are the 24 largest and smallest loadings on the first PC

# PC LOADINGS, FOR REUTERS DOCUMENTS

> signif(sort(out.lsi[,1],decreasing=TRUE)[1:24],2)

	january		cubic	fiscales	petroliferos	yacimi	entos	billion
	0.5000	0.3200		0.3200	0.3200	0.3200 0.32		0.2500
	argentine		gas	metrers	metres	na	tural p	produced
	0.1600	600 0.1600		0.1600	0.1600	0	.1600	0.1600
	totalled	ba	arrels	added	production		pct	mln
	0.1600	0.1600 0.0940		0.0860	0.0850	0	.0780	0.0770
	output	ł	oudget	riyals	abdul-aziz	expend	iture	revenue
	0.0630	(	0.0061	0.0061	0.0051	0	.0051	0.0041
>	signif(sort(out.lsi[,1],decreasing=FALSE)[1:24],2)							
	posted	canada	canadian	west	z power	bbl	lowered	texaco
	-0.050	-0.044	-0.044	-0.041	L -0.040	-0.040	-0.036	-0.033
	texas	brings	effective	dlrs	s grade	sweet	contract	ship
	-0.033	-0.032	-0.032	-0.031	L -0.031	-0.031	-0.031	-0.030
	price	changed	pay	postings	s decrease	company	benchmark	feb
	-0.030	-0.028	-0.028	-0.028	3 -0.027	-0.027	-0.026	-0.026

These are the 24 largest and smallest loadings on the first PC

- · Large loadings correspond to things related to the international market
- Negative loadings correspond to the American/Canadian market

### PC LOADINGS, FOR TMNT DOCUMENTS

> signif(sort(out.lsi[,1],decreasing=TRUE)[1:24],2)

	cool	rude	dude	part	i rap	hael	-foot	x
	5.0e-01	5.0e-01	5.0e-01	5.0e-0	1 3.6	ie-02 -8.	1e-05 -	8.1e-05
	acclaim	acknowledg	add	adeq	ı admi	nist a	erial	ahead
	-8.1e-05	-8.1e-05	-8.1e-05	-8.1e-0	5 -8.1	e-05 -8.	1e-05 -	8.1e-05
	albiera	albrecht	alessandra	ama	z ambros	iana a	mount	analys
	-8.1e-05	-8.1e-05	-8.1e-05	-8.1e-0	5 -8.1	e-05 -8.	1e-05 -	8.1e-05
>	signif(so	ort(out.lsi	[,1],decrea	sing=FALSE	)[1:24],	2)		
	turtl	comic d	onatello	paint	ninja	florenc	leonardo	mutant
	-0.0180	-0.0100	-0.0085	-0.0081	-0.0076	-0.0073	-0.0070	-0.0063
	seri	statu	art	voic	anim	tmnt	game	brother
	-0.0061	-0.0047	-0.0047	-0.0045	-0.0044	-0.0043	-0.0039	-0.0038
	charact	shredder	duomo	medici	casey	splinter	teenag	penni
	-0.0037	-0.0036	-0.0035	-0.0034	-0.0032	-0.0032	-0.0032	-0.0032

These are the 24 largest and smallest projections onto the first PC

### PC LOADINGS, FOR TMNT DOCUMENTS

> signif(sort(out.lsi[,1],decreasing=TRUE)[1:24],2)

	cool	rude	dude	part	i rap	hael	-foot	x
	5.0e-01	5.0e-01	5.0e-01	5.0e-0	1 3.6	ie-02 -8.	1e-05 -	8.1e-05
	acclaim	acknowledg	add	adeq	u admi	nist a	erial	ahead
	-8.1e-05	-8.1e-05	-8.1e-05	-8.1e-0	5 -8.1	e-05 -8.	1e-05 -	8.1e-05
	albiera	albrecht	alessandra	ama	z ambros	iana a	mount	analys
	-8.1e-05	-8.1e-05	-8.1e-05	-8.1e-0	5 -8.1	e-05 -8.	1e-05 -	8.1e-05
>	signif(so	ort(out.lsi	[,1],decrea	sing=FALSE	)[1:24],	2)		
	turtl	comic d	onatello	paint	ninja	florenc	leonardo	mutant
	-0.0180	-0.0100	-0.0085	-0.0081	-0.0076	-0.0073	-0.0070	-0.0063
	seri	statu	art	voic	anim	tmnt	game	brother
	-0.0061	-0.0047	-0.0047	-0.0045	-0.0044	-0.0043	-0.0039	-0.0038
	charact	shredder	duomo	medici	casey	splinter	teenag	penni
	-0.0037	-0.0036	-0.0035	-0.0034	-0.0032	-0.0032	-0.0032	-0.0032

These are the 24 largest and smallest projections onto the first PC

What story can be told here?

• There is a substantial amount of overlap on the first PC

Let's look at a plot

# PLOT OF PC SCORES FOR TMNT



# PC loadings, for TMNT documents, second component

<pre>&gt; signif(sort(out.lsi[,2],decreasing=TRUE)[1:24],2)</pre>									
floren	c st	tatu	paint	duomo	art	basilic	a medici		
0.220	0 O.	.150	0.140	0.130	0.110	0.08	3 0.080		
muse	o equesti	rian ma	adonna	execut	del	bron	z firenz		
0.074	4 0.	.073	0.071	0.070	0.069	0.06	3 0.062		
cosimo	o jud	dith lo	orenzo	prophet	penni	vasar	i centuri		
0.059	9 0.	.059	0.057	0.057	0.055	0.05	4 0.054		
> signif(	sort(out)	.lsi[,2],d	lecreasin	g=FALSE) [:	1:24],2)				
turtl	comic	ninja	mutant	seri	voic	anim	tmnt game		
-0.420	-0.240	-0.180	-0.150	-0.140	-0.100	-0.100	-0.100 -0.092		
charact	brother	shredder	casey	teenag	splinter	foot	sai mirag		
-0.087	-0.083	-0.081	-0.077	-0.074	-0.074	-0.067	-0.066 -0.065		
movi	mikey	episod	mutat	clan	raph				
-0.063	-0.061	-0.056	-0.049	-0.046	-0.045				

These are the 24 largest and smallest projections onto the first PC

# PC loadings, for TMNT documents, second component

> signif(s	<pre>&gt; signif(sort(out.lsi[,2],decreasing=TRUE)[1:24],2)</pre>								
florenc	st	tatu	paint	duomo	art	basilio	ca m	edici	
0.220	0.	150	0.140	0.130	0.110	0.08	33	0.080	
museo	equesti	rian m	adonna	execut	del	bror	nz f	irenz	
0.074	0.	.073	0.071	0.070	0.069	0.06	33	0.062	
cosimo	jud	lith l	orenzo	prophet	penni	vasai	ri ce	nturi	
0.059	0.	059	0.057	0.057	0.055	0.05	54	0.054	
> signif(s	ort(out.	lsi[,2],	decreasin	g=FALSE) [:	1:24],2)				
turtl	comic	ninja	mutant	seri	voic	anim	tmnt	game	
-0.420	-0.240	-0.180	-0.150	-0.140	-0.100	-0.100	-0.100	-0.092	
charact	brother	shredder	casey	teenag	splinter	foot	sai	mirag	
-0.087	-0.083	-0.081	-0.077	-0.074	-0.074	-0.067	-0.066	-0.065	
movi	mikey	episod	mutat	clan	raph				
-0.063	-0.061	-0.056	-0.049	-0.046	-0.045				

16

These are the 24 largest and smallest projections onto the first PC

- · Positive values are related to the renaissance artists
- Negative values are related to the TMNTs

### DISTANCE TO QUERY USING LSI

ORIGINAL REPRESENTATION [dimension reduction] STATISTICAL METHOD

- Form the (normalized) document-term matrix  ${\mathbb X}$
- Compute its LSI  $X = UDV^{\top}$
- Get Y into LSI via  $\tilde{Y} = YVD^{-1}$
- Choose a K
- Find distances for documents  $d = 1, \dots, D$

$$\operatorname{distance}(d, \tilde{Y}) = ||U_{d,1:k} - \tilde{Y}||_2$$

### DISTANCE TO QUERY USING LSI

(tmnt leo) 0.9711807

- (tmnt rap) 0.7696525
- (tmnt mic) 0.7669749
- (tmnt don) 0.9718710
- (real leo) 0.9711512
- (real rap) 0.9709391
- (real mic) 0.9709492
- (real don) 0.9734319

query 0.000000

# WRAP-UP

Latent semantic indexing seeks to use correlations to help with document queries

This can be accomplished by forming the SVD of the document-term matrix  $\mathbb X.$  Alternatively, factor analysis is commonly used

The documents and the query get projected into a lower dimensional space where distances are computed

Some complications

- How do the normalization techniques we talked about affect this reduction?
- What happens if we want to make a new query, do we need to reform the entire SVD?
- How does factor analysis compare to PCA?