

HW: Latent Semantic Analysis (LSA)

- Given a sparse term-doc matrix
 - E.g., 4 terms and 3 docs

	Doc		
Term	2.3	0.0	4.2
	0.0	1.3	2.2
	3.8	0.0	0.5
	0.0	0.0	0.0



Row #Tem	Col. # Doc	Nonzero entries
4	3	6
2		2 nonzero entries at Col 0
0	2.3	Col 0, Row 0
2	3.8	Col 0, Row 2
1		1 nonzero entry at Col 1
1	1.3	Col 1, Row 1
3		3 nonzero entry at Col 2
0	4.2	Col 2, Row 0
1	2.2	Col 2, Row 1
2	0.5	Col 2, Row 2

- Each entry is weighted by *TFxIDF* score

- Perform SVD to obtain corresponding term and doc vectors represented in the latent semantic space
- Evaluate the information retrieval capability of the LSA approach by using varying sizes (e.g., 100, 200, .., 600 etc.) of LSA dimensionality

HW: Latent Semantic Analysis (cont.)

- Example: term-docmatrix

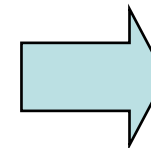
Indexing Term no.	Doc no.	Nonzero entries
51253	2265	218852
77		
508	7.725771	
596	16.213399	
612	13.080868	
709	7.725771	
713	7.725771	
744	7.725771	
1190	7.725771	
1200	16.213399	
1259	7.725771	
.....		

- SVD command (IR_svd.bat)

`svd -r st -o LSA100 -d 100 Term-Doc-Matrix`

sparse matrix input prefix of output files No. of reserved eigenvectors name of sparse matrix input

output



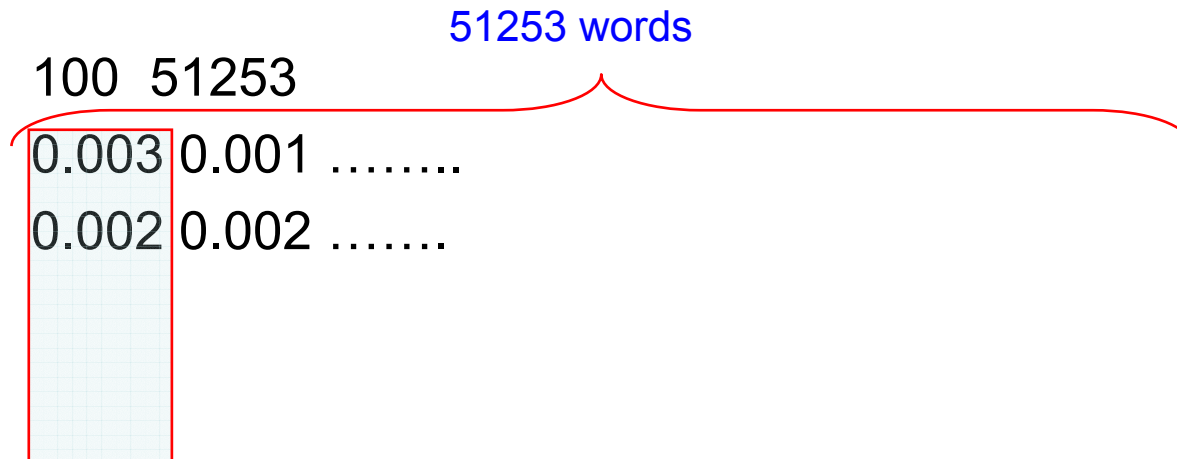
LSA100-Ut

LSA100-S

LSA100-Vt

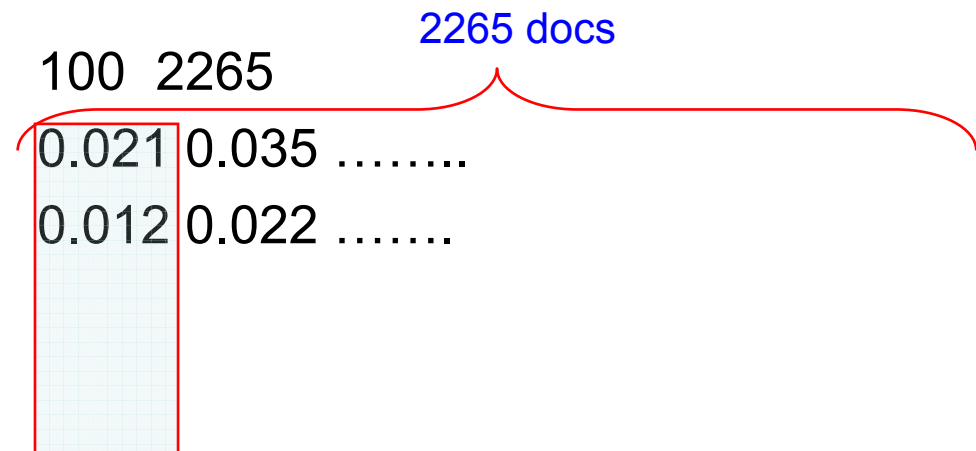
HW: Latent Semantic Analysis (cont.)

- **LSA100-U_t**



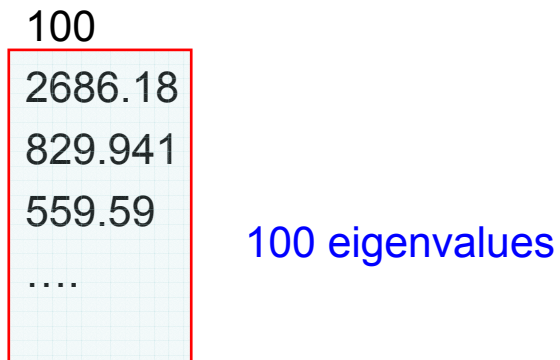
word vector (u^T): 1x100

- **LSA100-V_t**



doc vector (v^T): 1x100

- **LSA100-S**



HW: Latent Semantic Analysis (cont.)

- Fold-in a new $m \times 1$ query vector

$$\hat{q}_{1 \times k} = \left(q^T \right)_{1 \times m} U_{m \times k} \Sigma^{-1}_{k \times k}$$

Just like a row of V

Query represented by the weighted sum of its constituent term vectors

TFxIDF weighted beforehand

The separate dimensions are differentially weighted

- Cosine measure between the query and doc vectors in the latent semantic space

$$\text{sim}(\hat{q}, \hat{d}) = \text{coine}(\hat{q}\Sigma, \hat{d}\Sigma) = \frac{\hat{q}\Sigma^T \hat{d}}{\|\hat{q}\Sigma\| \|\hat{d}\Sigma\|}$$

SVDLIBC

- Doug Rohde's SVD C Library version 1.3 is based on the [SVDPACKC](#) library
- Download it at <http://tedlab.mit.edu/~dr/>