

Lab 09

文本與詞彙的向量表徵

Document-Term Matrix, Latent Semantic Analysis & Word Vectors

表徵文本

(Representing documents)

How to represent documents numerically?

doc1: “I baked the cake and the muffin”

doc2: “I loved the cake”

doc3: “I wrote the book”

How to represent documents numerically?

doc1: “I baked the **cake** and the **muffin**”
doc2: “I loved the **cake**”
doc3: “I wrote the **book**”

Document-Term Matrix

doc1: “I baked the **cake** and the **muffin**”

doc2: “I loved the **cake**”

doc3: “I wrote the **book**”

	I	baked	loved	wrote	the	and	cake	muffin	book
doc1	1	1	0	0	2	1	1	1	0
doc2	1	0	1	0	1	0	1	0	0
doc3	1	0	0	1	1	0	0	0	1

Document-Term Matrix

doc1: “I baked the **cake** and the **muffin**”
doc2: “I loved the **cake**”
doc3: “I wrote the **book**”

	I	baked	loved	wrote	the	and	cake	muffin	book
doc1	1	1	0	0	2	1	1	1	0
doc2	1	0	1	0	1	0	1	0	0
doc3	1	0	0	1	1	0	0	0	1

Representing documents as vectors

doc1: “I baked the **cake** and the **muffin**”
doc2: “I loved the **cake**”
doc3: “I wrote the **book**”

	I	baked	loved	wrote	the	and	cake	muffin	book	
doc1	1	1	0	0	2	1	1	1	0	vector 1
doc2	1	0	1	0	1	0	1	0	0	vector 2
doc3	1	0	0	1	1	0	0	0	1	vector 3

Vectors quantify similarity b/t documents

Euclidean Distance $d(\vec{p}, \vec{q}) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_n - q_n)^2}$

Cosine Similarity $\cos(\theta) = \frac{\vec{p} \cdot \vec{q}}{\|\vec{p}\| \|\vec{q}\|}$

	I	baked	loved	wrote	the	and	cake	muffin	book	
doc1	1	1	0	0	2	1	1	1	0	vector 1
doc2	1	0	1	0	1	0	1	0	0	vector 2
doc3	1	0	0	1	1	0	0	0	1	vector 3

Vectors quantify similarity b/t documents

Euclidean Distance $d(\vec{p}, \vec{q}) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_n - q_n)^2}$

Cosine Similarity $\cos(\theta) = \frac{\vec{p} \cdot \vec{q}}{\|\vec{p}\| \|\vec{q}\|}$

	l	cake	
doc1	1	1	vector 1
doc2	1	1	vector 2
doc3	1	0	vector 3

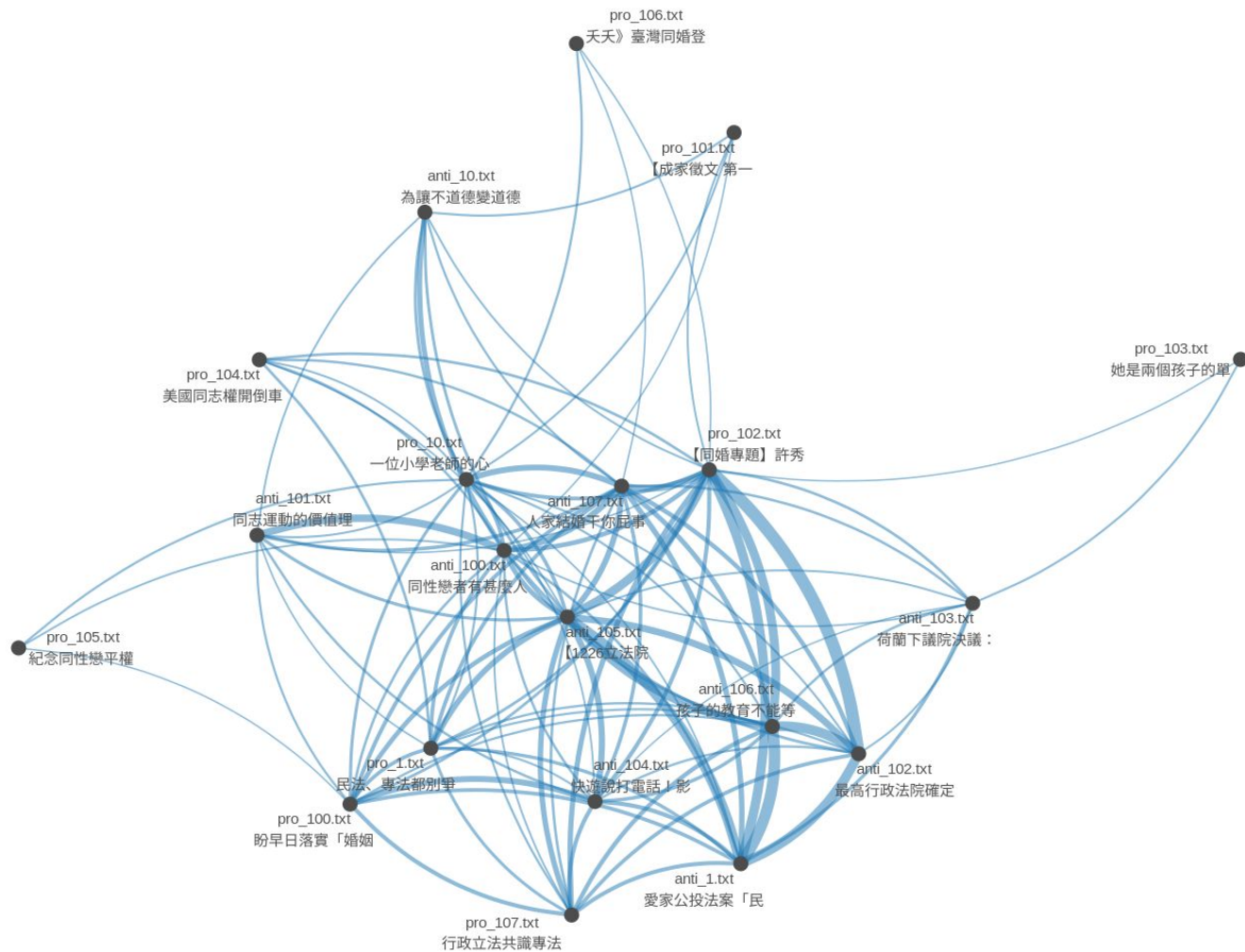
	$\cos(\theta)$	$d(\vec{p}, \vec{q})$
doc1, doc2	1	0
doc1, doc3	$1/\sqrt{2}$	1
doc2, doc3	$1/\sqrt{2}$	1

Selecting & Weighting Terms

tf-idf
weighting

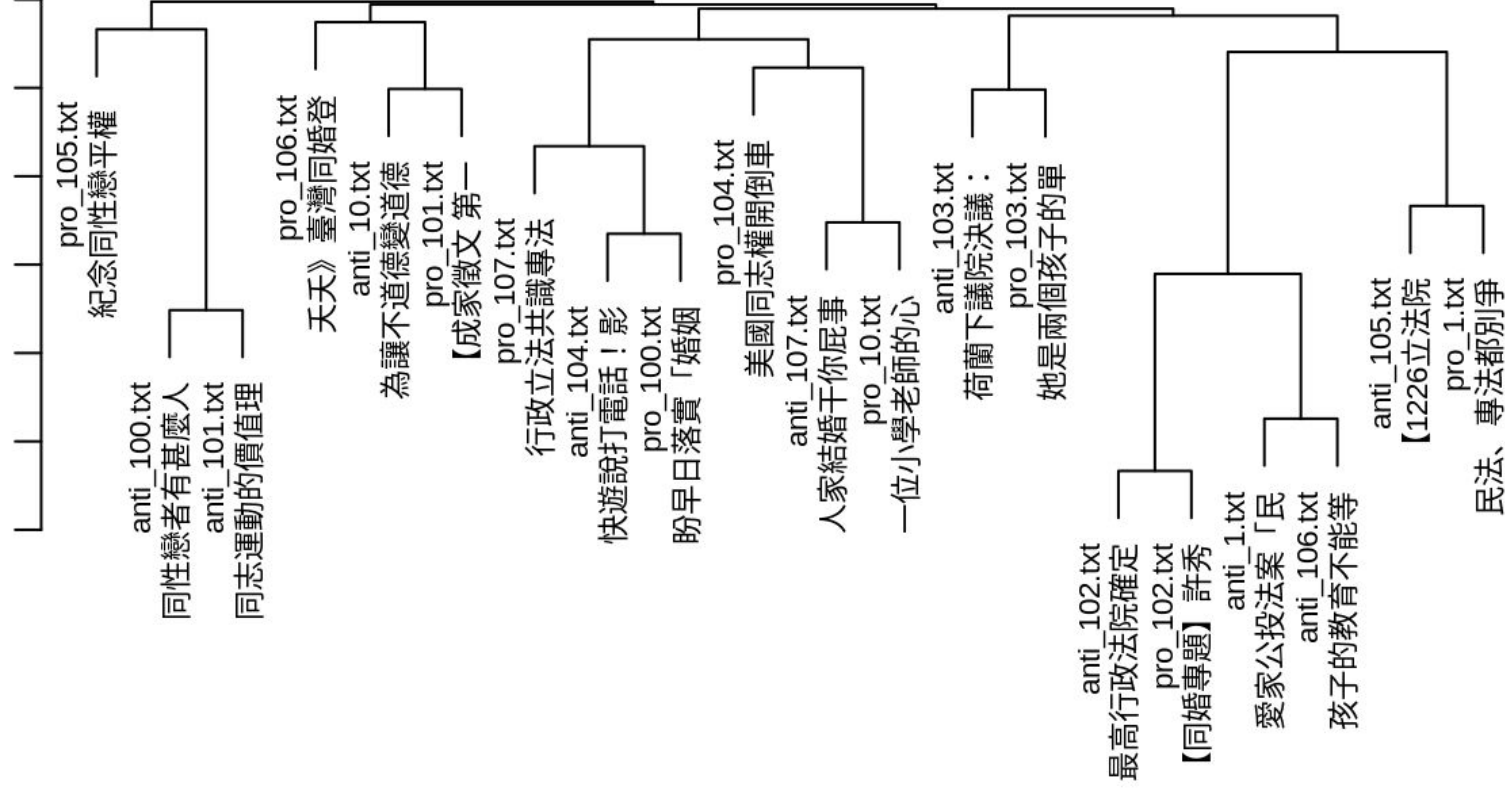
stopword

	I	baked	loved	wrote	the	and	cake	muffin	book
doc1	1	1	0	0	2	1	1	1	0
doc2	1	0	1	0	1	0	1	0	0
doc3	1	0	0	1	1	0	0	0	1



Height

0.70 0.80 0.90 1.00



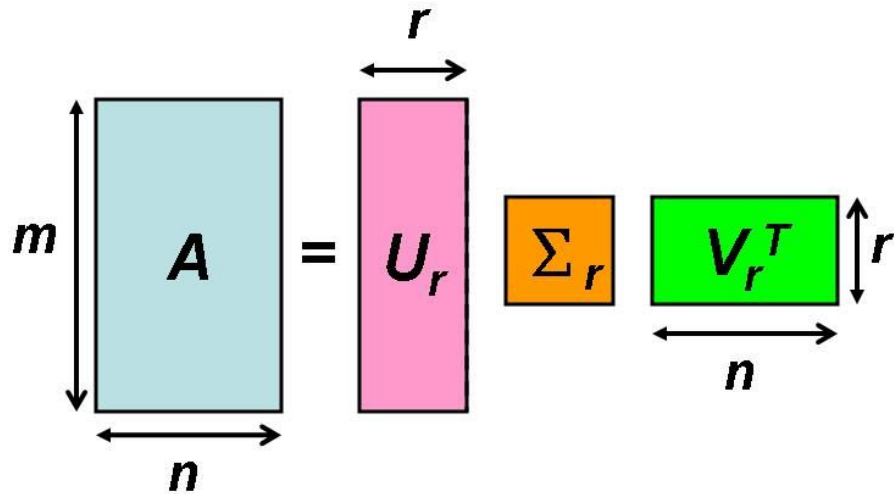
20 Selected Posts

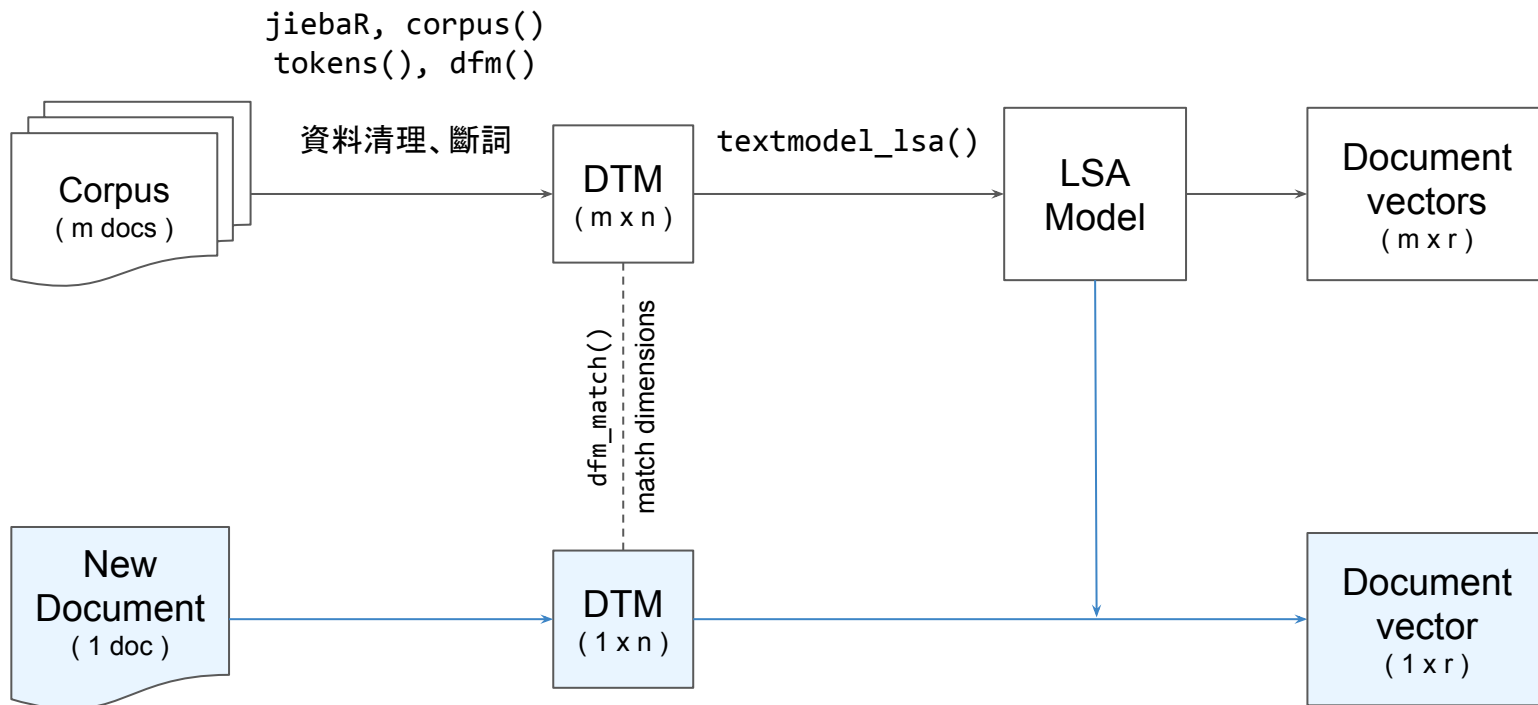
Dimensionality Reduction (SVD / LSA)

doc2: “I loved the cake”
doc4: “He liked my muffin”

	I	He	liked	loved	the	my	cake	muffin
doc2	1	0	0	1	1	0	1	0
doc4	0	1	1	0	0	1	0	1

Dimensionality Reduction (SVD / LSA)





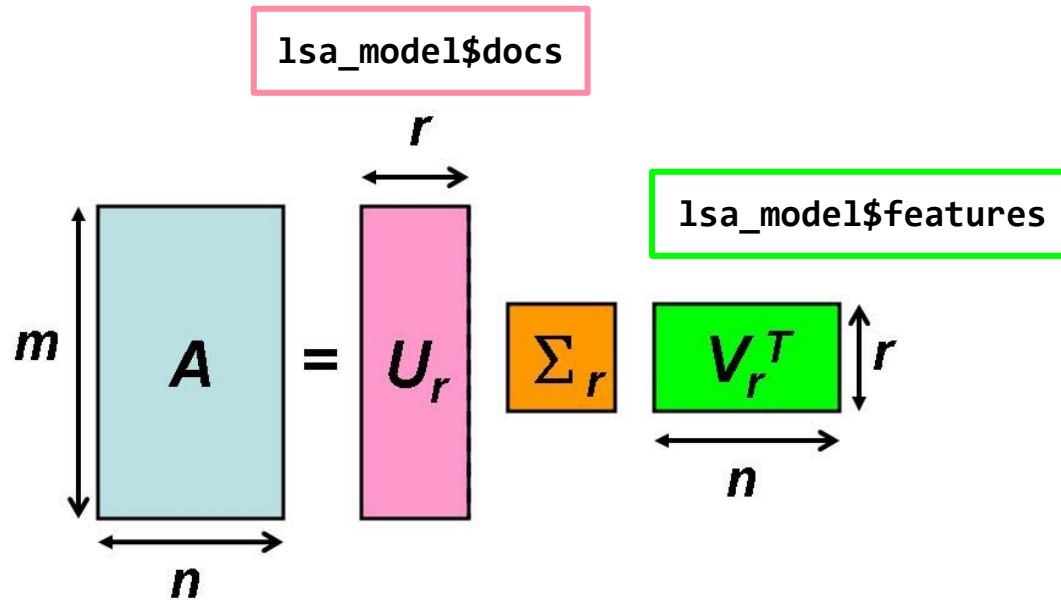
表徵詞彙

(Representing words)

Document-Term Matrix

	I	baked	loved	wrote	the	and	cake	muffin	book
doc1	1	1	0	0	2	1	1	1	0
doc2	1	0	1	0	1	0	1	0	0
doc3	1	0	0	1	1	0	0	0	1

LSA Word Vectors



Word Vectors Trained From Large Raw Text

- [PPMI](#) (count-based)
- [GloVe](#) (count-based)
- [fasttext](#) (neural network)
- [word2vec](#) (neural network)