

Exercises: Text Classification and Clustering

September 27, 2016

Task 1: Preprocessing

Preprocess the documents in Table 1 and construct the corresponding document-term matrix.

Table 2 contains the stopwords list to be used. Apply *suffix-s* stemming.

docID	content
1	The King's Speech
2	The Lord of the Rings: The Return of the King
3	Street Kings
4	The Scorpion King
5	The Lion King

Table 1: Document collection consisting of movie titles.

a	as	by	into	not	such	then	this	with
an	at	for	is	of	that	there	to	
and	be	if	it	on	the	these	was	
are	but	in	no	or	their	they	will	

Table 2: Standard English stopwords list.

Solution

DocID														
1														
2														
3														
4														
5														

Table 3: Document-term matrix.

Task 2: K-Nearest Neighbor Classification

Find $K = 2$ nearest neighbors of document 1 using the Jaccard coefficient and the cosine similarity.

$$\cos(d_1, d_2) = \frac{\sum_t n(t, d_1)n(t, d_2)}{\sqrt{\sum_t n(t, d_1)^2 \sum_t n(t, d_2)^2}} \quad (1)$$

Solution

Documents	Jaccard	Cosine
1 vs. 2		
1 vs. 3		
1 vs. 4		
1 vs. 5		

Table 4: Similarity between document 1 and documents 2–5.

Task 3: Naive Bayes Classification

Assuming a binary classification task, where the documents are labeled as shown in Table 5, how would the document $d = \text{“The Mummy King Returns”}$ be classified using a Naive Bayes classifier? Use Laplace smoothing for computing term probabilities.

docID	target (Oscar?)
1	Yes
2	Yes
3	No
4	No
5	Yes

Table 5: Target labels for the document collection

		$n(t, c)$											
Class	N_c												SUM
Oscar=Yes													
Oscar=No													

Table 6: Term counts.

$$P(c|d) \propto P(c) \prod_{t \in d} P(t|c)^{n(t,d)} \quad (2)$$

where

$$P(c) = \frac{N_c}{N} \quad P(t|c) = \frac{n(t, c) + 1}{\sum_t n(t, c) + |V|} \quad (3)$$

Solution

$$P(d|Yes) =$$

$$P(d|No) =$$

Task 4: K-Nearest Neighbors Clustering

Perform the first 2 iterations of K-Nearest Neighbors clustering with $K = 2$ for the documents in Table 1. Use Jaccard similarity. Take documents 1 and 5 as the initial centroids.

Solution

Round numbers to the first digit.

DocID	Term vector	Iteration 1		Iteration 2	
		Distance to centroid		Distance to centroid	
		of Cluster 1 ()	of Cluster 2 ()	of Cluster 1 ()	of Cluster 2 ()
1	()				
2	()				
3	()				
4	()				
5	()				

Table 7: K-Nearest Neighbors Clustering

Task 5: Hierarchical Agglomerative Clustering

Perform the first step of the Hierarchical Agglomerative Clustering method for the documents in Table 1. Use cosine similarity and the group average cluster proximity.

Solution

Init: Compute the proximity matrix

Each document corresponds to a (singleton) cluster.

	doc 1	doc 2	doc 3	doc 4	doc 5
doc 1	0				
doc 2		0			
doc 3			0		
doc 4				0	
doc 5					0

Table 8: Initial proximity matrix.

First iteration: Merge the two most similar clusters and update the proximity matrix

	doc	doc	doc	doc
doc	0			
doc		0		
doc			0	
doc				0

Table 9: Proximity matrix after the first iteration.