

Index Compression + Vector Space Model (Chapter 5+6)

Exercise 6/1

	car	truck	van
weight	1.5	2.5	1.5
length	4.5	5.5	4.5
height	1.5	2.5	1.5

Definition 4 (Sparse Inverse Sparse)

Given document D (matrix) D of size $n \times m$

$$D = \begin{bmatrix} d_{11} & \dots & d_{1m} \\ \vdots & \ddots & \vdots \\ d_{n1} & \dots & d_{nm} \end{bmatrix}$$

Let r be the number of all features and k the number of features (1 to the number of features for each).

Definition 5 (k-SIFT weighting scheme)

In the k-SIFT weighting scheme, a term t in a document d has weight

$$w_{dt} = \frac{d_{dt}}{k_d}$$

where k_d is the number of values (the term frequency) t in a document d .

Consider the frequency table of the words of three documents. Calculate the k-SIFT weight of the terms car, auto, insurance, and tent for each document. k-SIFT values of terms are in the table.

	car	auto	insurance	tent
car	1.5	0	0	0
auto	0	1.5	0	0
insurance	0	0	1.5	0
tent	0	0	0	1.5

Exercise 6/2

Consider document representations as normalized Euclidean weight vectors for each document from the previous exercise. Each vector has four components, one for each term.

→ $\vec{v}_1 = \begin{bmatrix} 1.5 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\vec{v}_2 = \begin{bmatrix} 0 \\ 1.5 \\ 0 \\ 0 \end{bmatrix}$ $\vec{v}_3 = \begin{bmatrix} 0 \\ 0 \\ 1.5 \\ 0 \end{bmatrix}$ $\vec{v}_4 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1.5 \end{bmatrix}$

→ $\| \vec{v}_1 \| = \sqrt{1.5^2} = 1.5$ $\| \vec{v}_2 \| = \sqrt{1.5^2} = 1.5$ $\| \vec{v}_3 \| = \sqrt{1.5^2} = 1.5$ $\| \vec{v}_4 \| = \sqrt{1.5^2} = 1.5$

→ $\vec{u}_1 = \frac{\vec{v}_1}{\| \vec{v}_1 \|} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\vec{u}_2 = \frac{\vec{v}_2}{\| \vec{v}_2 \|} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ $\vec{u}_3 = \frac{\vec{v}_3}{\| \vec{v}_3 \|} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ $\vec{u}_4 = \frac{\vec{v}_4}{\| \vec{v}_4 \|} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

Exercise 6/3

Based on the weights from the last exercise, return the underlying matrix (outer product) of the three documents for the query "car insurance".

	car	auto	insurance	tent
car	1.5	0	0	0
auto	0	1.5	0	0
insurance	0	0	1.5	0
tent	0	0	0	1.5

The each of the two weighting schemes of term weight is 1 if the query contains the word and 0 otherwise.

$Q = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ $D = \begin{bmatrix} 1.5 & 0 & 0 & 0 \\ 0 & 1.5 & 0 & 0 \\ 0 & 0 & 1.5 & 0 \\ 0 & 0 & 0 & 1.5 \end{bmatrix}$

→ $Q \cdot D = \begin{bmatrix} 1.5 & 0 & 0 & 0 \\ 0 & 1.5 & 0 & 0 \\ 0 & 0 & 1.5 & 0 \\ 0 & 0 & 0 & 1.5 \end{bmatrix}$

Exercise 6/4

Compare the Laplacian matrix L and D .

Write down the matrix of distances between all pairs as required in Algorithm 4.

	car	auto	insurance	tent
car	1.5	0	0	0
auto	0	1.5	0	0
insurance	0	0	1.5	0
tent	0	0	0	1.5



Algorithm 4 (Laplacian matrix)

Input: A graph $G = (V, E)$ with n vertices and m edges. Let D be the degree matrix and A the adjacency matrix of G .

Output: The Laplacian matrix $L = D - A$.

	car	auto	insurance	tent
car	1.5	0	0	0
auto	0	1.5	0	0
insurance	0	0	1.5	0
tent	0	0	0	1.5

	car	auto	insurance	tent
car	1.5	0	0	0
auto	0	1.5	0	0
insurance	0	0	1.5	0
tent	0	0	0	1.5

	car	auto	insurance	tent
car	1.5	0	0	0
auto	0	1.5	0	0
insurance	0	0	1.5	0
tent	0	0	0	1.5

where L_{ij} is the Laplacian matrix element L_{ij} between i and j , and L_{ii} is the degree of vertex i .