DEPARTMENT OF COMPUTER SCIENCE - THIRD YEAR LICENCE (ISIL)

# TD 02

## Exercise 01

Consider the following documents collection

d1 = "Big cats are nice and funny"
d2 = "Small dogs are better than big dogs"
d3 = "Small cats are afraid of small dogs"
d4 = "Big cats are not afraid of small dogs"
d5 = "Funny cats are not afraid of small dogs"

(a) Compute the tokens for each document.
(b) Normalize the tokens with respect to plurals and upper/lower case.
(c) Compute the dictionary relative to the documents collection.

- For each of the questions above write a pseudo code of the algorithm.

## Exercise 02

Assume the following documents

Doc1 "italy is world champion 2006"
Doc2 "germany and italy played each other in the semifinal"
Doc3 "germany was in the semifinal 2006"
Doc4 "germany won the semifinal in italy 1990"

Assume that the following terms are stop words: **is, and, in, the, was, each, other**. Construct an inverted index.

## Exercise 03

Starting from the document collection of exercise 01, build the documents-terms incidence matrix as required by the Boolean model.

## Exercise 04

Starting from the documents collection of Exercise 1 consider a Boolean model.

(a) Answer the query q1 = funny AND dog
(b) Answer the query q2 = nice OR dog
(c) Answer the query q3 = big AND dog AND NOT funny

# Exercise 05

Draw the term-document incidence matrix and the inverted index representation for the following document collection:

Doc 1 breakthrough drug for schizophrenia
Doc 2 new schizophrenia drug
Doc 3 new approach for treatment of schizophrenia
Doc 4 new hopes for schizophrenia patients

For the document collection above, what are the returned results for these queries:

– schizophrenia AND drug
– for AND NOT(drug OR approach)

## Term-document incidence matrix

| Term | Doc 1 | Doc 2 | Doc 3 | Doc 4 |
|---|---|---|---|---|
| breakthrough | 1 | 0 | 0 | 0 |
| drug | 1 | 1 | 0 | 0 |
| for | 1 | 0 | 1 | 1 |
| schizophrenia | 1 | 1 | 1 | 1 |
| new | 0 | 1 | 1 | 1 |
| approach | 0 | 0 | 1 | 0 |
| treatment | 0 | 0 | 1 | 0 |
| of | 0 | 0 | 1 | 0 |
| hopes | 0 | 0 | 0 | 1 |
| patients | 0 | 0 | 0 | 1 |

## Inverted index

- breakthrough → 1
- drug → 1, 2
- for → 1, 3, 4
- schizophrenia → 1, 2, 3, 4
- new → 2, 3, 4
- approach → 3
- treatment → 3
- of → 3
- hopes → 4
- patients → 4

## Queries

**schizophrenia AND drug**
- schizophrenia → {1, 2, 3, 4}
- drug → {1, 2}
- Result: 1, 2

**for AND NOT(drug OR approach)**
- for → {1, 3, 4}
- drug OR approach → {1, 2, 3}
- NOT(drug OR approach) → {4}
- for AND {4} → Result: 4