

# Norwegian University of Science and Technology

---

## **Assignment Title**

Assignment 1

-

Basic Definitions, IR Models, Term Weighting

## **Course**

TDT4117 Information Retrieval

## **Semester**

FALL 2018

## **Date Submitted**

30/09/2018

## **Submitted by**

Dusan Jakovic

---

## Task 1: Basic Definitions

**Explain the main differences between:**

### **1. Information Retrieval vs Data Retrieval**

Information retrieval: the information is about a subject or topic, semantics is frequently loose, small errors are tolerated, produces multiple results with ranking (partial match is allowed).

Data retrieval (database management system): which documents that contains a set of keywords, well defined semantics, a single erroneous object implies failure, and it produces exact results or no results if no exact match is found.

### **2. Structured Data vs Unstructured Data**

Structured Data: Clearly defined data types whose pattern makes them easily searchable. Used for relational databases and data warehouses, e.g: phone numbers, addresses, dates, credit card numbers, etc..

Unstructured Data: Data that is usually not as easily searchable. it is essentially everything else, e.g.: text files, website, media (images, video files, audio files), reports, etc..

## Task 2: IR Models

Given the following document collection containing words from the set  $O = \{\text{Autumn}, \text{Winter}, \text{Spring}\}$ , answer the questions in subtasks 1.1 and 1.2

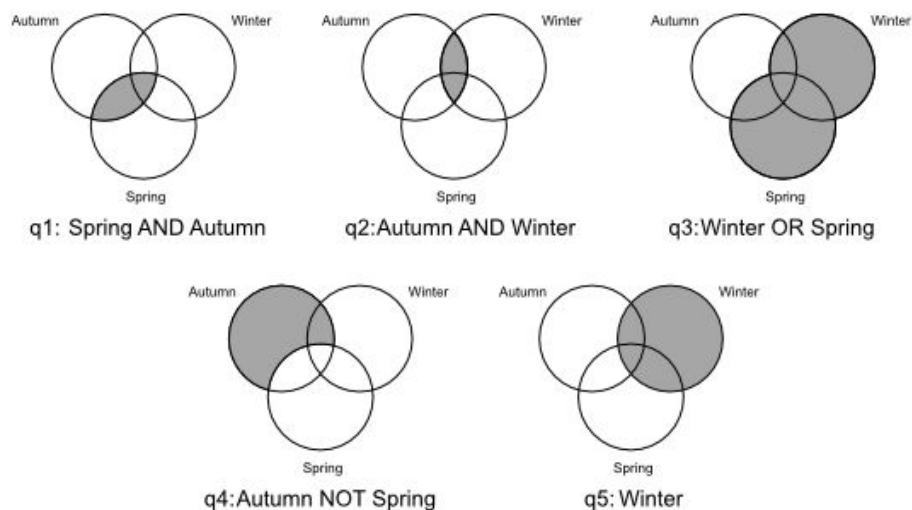
$\text{doc}_1 = \{\text{Winter Autumn Spring Spring}\}$   
 $\text{doc}_2 = \{\text{Spring Autumn}\}$   
 $\text{doc}_3 = \{\text{Winter Winter Autumn}\}$   
 $\text{doc}_4 = \{\text{Autumn Autumn Winter Autumn Spring Winter Spring Autumn}\}$   
 $\text{doc}_5 = \{\text{Winter}\}$   
 $\text{doc}_6 = \{\text{Spring Spring}\}$   
 $\text{doc}_7 = \{\text{Winter Autumn Spring}\}$   
 $\text{doc}_8 = \{\text{Spring Spring Autumn}\}$   
 $\text{doc}_9 = \{\text{Spring Spring Winter}\}$   
 $\text{doc}_{10} = \{\text{Winter Autumn Winter Spring}\}$

### SubTask 1.1: Boolean Model and Vector Space Model

Given the following queries:

$q_1 = \text{"Spring AND Autumn"}$   
 $q_2 = \text{"Autumn AND Winter"}$   
 $q_3 = \text{"Winter OR Spring"}$   
 $q_4 = \text{"Autumn NOT Spring"}$   
 $q_5 = \text{"Winter"}$

1. Which of the documents will be returned as the result for the above queries using the Boolean model? Explain your answers and draw a figure to illustrate.



Following table will illustrate each which of the documents that will be returned.

'1' = returned

'0' = not returned

	doc <sub>1</sub>	doc <sub>2</sub>	doc <sub>3</sub>	doc <sub>4</sub>	doc <sub>5</sub>	doc <sub>6</sub>	doc <sub>7</sub>	doc <sub>8</sub>	doc <sub>9</sub>	doc <sub>10</sub>
q1	1	1	0	1	0	0	1	1	0	1
q2	1	0	1	1	0	0	1	0	0	1
q3	1	1	1	1	1	1	1	1	1	1
q4	0	0	1	0	0	0	0	0	0	0
q5	1	0	1	1	1	0	1	0	1	1

2. What is the dimension of the vector space representing this document collection when you use the vector model and how is it obtained?

There are three unique terms in the document collection, therefore the dimension is 3.

3. Calculate the weights for the documents and the terms using tf and idf weighting. Put these values into a document-term-matrix. (Tip: use the equations in the book and state which one you used.)

Term frequency tf:  $tf(1 + \log_2(fi,j))$

#	term	f <sub>i,1</sub>	f <sub>i,2</sub>	f <sub>i,3</sub>	f <sub>i,4</sub>	f <sub>i,5</sub>	f <sub>i,6</sub>	f <sub>i,7</sub>	f <sub>i,8</sub>	f <sub>i,9</sub>	f <sub>i,10</sub>	TF <sub>i,1</sub>	TF <sub>i,2</sub>	TF <sub>i,3</sub>	TF <sub>i,4</sub>	TF <sub>i,5</sub>	TF <sub>i,6</sub>	TF <sub>i,7</sub>	TF <sub>i,8</sub>	TF <sub>i,9</sub>	TF <sub>i,10</sub>
1	Autumn	1	1	1	4	0	0	1	1	0	1	1	1	1	3	-	-	1	1	-	1
2	Spring	2	1	0	2	0	2	1	2	2	1	2	1	-	2	-	2	1	2	2	1
3	Winter	1	0	2	2	1	0	1	0	1	2	1	-	2	2	1	-	1	-	1	2
Doc length		4	2	3	8	1	2	3	3	3	4										

Inverse Document Frequency:  $idf(\log_2(N/df_i))$ :

#	term	n <sub>i</sub>	IDF = $\log_2(N/n_i)$
1	Autumn	7	0.515
2	Spring	8	0.322
3	Winter	7	0.515

Term Frequency - Inverse Document Frequency:  $TF-IDF = tf * idf$

#	Term	d <sub>1</sub>	d <sub>2</sub>	d <sub>3</sub>	d <sub>4</sub>	d <sub>5</sub>	d <sub>6</sub>	d <sub>7</sub>	d <sub>8</sub>	d <sub>9</sub>	d <sub>10</sub>
1	Autumn	0.515	0.515	0.515	1.545	-	-	0.515	0.515	-	0.515
2	Spring	0.644	0.322	-	0.644	-	0.644	0.322	0.644	0.644	0.322
3	Winter	0.515	-	1.03	1.03	0.515	-	0.515	-	0.515	1.03

4. Study the documents 1, 2, 4 and 10 and compare them to document 5. Calculate the similarity between document 5 and these four documents according to Euclidean distance. (Use tf-idf weights for your computations).

$$\text{Euclidean distance } (d(p,q)) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + (p_3 - q_3)^2}$$

$$d(d_5, d_1) = \sqrt{(0 - 0.515)^2 + (0 - 0.644)^2 + (0.515 - 0.515)^2} = 0.825$$

$$d(d_5, d_2) = \sqrt{(0 - 0.515)^2 + (0 - 0.322)^2 + (0.515 - 0)^2} = 0.796$$

$$d(d_5, d_4) = \sqrt{(0 - 1.545)^2 + (0 - 0.644)^2 + (0.515 - 1.03)^2} = 0.972$$

$$d(d_5, d_{10}) = \sqrt{(0 - 0.515)^2 + (0 - 0.322)^2 + (0.515 - 1.03)^2} = 0.796$$

5. Rank the documents for query q5 using cosine similarity.

$q = q_5 = \{\text{Winter}\}$

$$\text{sim}(d_j, q) = \frac{\vec{d_j} \cdot \vec{q}}{|\vec{d_j}| \times |\vec{q}|}$$

Calculation example for the cosine similarity ( $d_1, q$ ):

$$\text{sim}(d_1, q) = \frac{(1 * 0) + (2 * 0) + (1 * 1)}{\sqrt{1^2 + 2^2 + 1^2} * \sqrt{0^2 + 0^2 + 1^2}} = 0.408$$

$\text{sim}(d_j, q) : \downarrow$

	d <sub>1</sub>	d <sub>2</sub>	d <sub>3</sub>	d <sub>4</sub>	d <sub>5</sub>	d <sub>6</sub>	d <sub>7</sub>	d <sub>8</sub>	d <sub>9</sub>	d <sub>10</sub>
sim(d <sub>j</sub> , q)	0.408	0	0.894	0.408	1	0	0.577	0	0.447	0.816

#	d	#	d	#	d	#	d	#	d	#	d	#	d	#	d	#	d	#	d	#	d
1	d <sub>5</sub>	2	d <sub>3</sub>	3	d <sub>10</sub>	4	d <sub>7</sub>	5	d <sub>9</sub>	6	d <sub>1</sub>	7	d <sub>4</sub>	8	d <sub>2</sub>	9	d <sub>6</sub>	10	d <sub>8</sub>		

## SubTask 1.2: Probabilistic Models

1. What are the main differences between BM25 and the probabilistic model introduced by Robertson-Jones?

### Robertson-Jones Probabilistic Model:

The probability that the user will find the document  $d_j$  relevant.

The model has some flaws:

- Relevance to the user might be affected by variables *outside* the system. Thus, the ideal answer set produced using information available to the system might not be ideal from the point of view of the user.
- The principle does not state explicitly how to compute the probabilities of relevance.
- Taking the *odds of relevance* as the rank minimizes the probability of an erroneous judgement

### BM25 (Best Match 25)

The BM25 is a combination of the *BM11* and *BM15* ranking formulas. Contrary to the original conceptions in the classic probabilistic model, BM25 formula uses parameters and can be computed without any relevance information provided by the user, i.e., in fully automatic fashion.

2. Rank the documents using the BM25 model. Set the parameters to  $k = 1.2$ ,  $b = 0.75$  (Here we assume relevance information is not provided.)

Given the following queries:

$q_1$  = "Winter Spring"

$q_2$  = "Autumn"

$$score(D, Q) = \sum_{i=1}^n IDF(q_i) * \frac{f(q_i, D) * (k_1 + 1)}{f(q_i, D) + k_1 * ((1 - b + b * (\frac{|D|}{L_{ave}}))}$$

$L = 33, L_{ave} = 3.3$ .

Calculation example using BM25 model ( $d_1, q$ ):

$$score(d_1, q_1) + score(d_1, q_2) = 0.515 * \frac{1 * (1.2 + 1)}{1 + (1.2 * (1 - 0.75) + 0.75 * \frac{4}{3.3})} + 0.322 * \frac{2 * (1.2 + 1)}{2 + (1.2 * (1 - 0.75) + 0.75 * \frac{4}{3.3})}$$
$$= 0.892$$

$$score(d_1, q_2) = 0.515 * \frac{1 * (1.2 + 1)}{1 + (1.2 * (1 - 0.75) + 0.75 * \frac{4}{3.3})} = 0.474$$

	d <sub>1</sub>	d <sub>2</sub>	d <sub>3</sub>	d <sub>4</sub>	d <sub>5</sub>	d <sub>6</sub>	d <sub>7</sub>	d <sub>8</sub>	d <sub>9</sub>	d <sub>10</sub>
sim(d <sub>r</sub> ,q <sub>1</sub> )	0.892	0.384	0.535	0.936	0.720	0.498	0.867	0.534	0.989	0.770

#	d	#	d	#	#	#	d	#	d	#	d	#	d	#	d	#	d	#	d
1	d <sub>9</sub>	2	d <sub>4</sub>	3	d <sub>1</sub>	4	d <sub>7</sub>	5	d <sub>10</sub>	6	d <sub>5</sub>	7	d <sub>3</sub>	8	d <sub>8</sub>	9	d <sub>6</sub>	10	d <sub>2</sub>

	d <sub>1</sub>	d <sub>2</sub>	d <sub>3</sub>	d <sub>4</sub>	d <sub>5</sub>	d <sub>6</sub>	d <sub>7</sub>	d <sub>8</sub>	d <sub>9</sub>	d <sub>10</sub>
sim(d <sub>r</sub> ,q) <sub>2</sub>	0.474	0.614	0.727	0.506	0	0	0.534	0.454	0	0.668

#	d	#	d	#	#	#	d	#	d	#	d	#	d	#	d	#	d	#	d
1	d <sub>3</sub>	2	d <sub>10</sub>	3	d <sub>2</sub>	4	d <sub>7</sub>	5	d <sub>4</sub>	6	d <sub>1</sub>	7	d <sub>8</sub>	8	d <sub>5</sub>	9	d <sub>6</sub>	10	d <sub>9</sub>

## Task 3: Term Weighting

Explain:

1. Term Frequency (tf)  
The number of times a term occurring in the given document
2. Document Frequency  
The number of documents in the collection that contains a term
3. Inverse Document Frequency (idf)  
How important a term is.
4. Why *idf* is important for term weighting  
It diminishes the weight of terms that occur very frequently in the document set and increases the weight of terms that occur rarely.