

CS5560 Knowledge Discovery and Management

Problem Set 3

June 19 (T), 2017

Name: *Revanth Chakilam*

Class ID: *02*

Information Retrieval (Text Mining) with TF-IDF

Consider the following three short documents

Doc #1:

The researchers will focus on computational phenotyping and will produce disease prediction models from machine learning and statistical tools.

Doc #2:

The researchers will develop tools that use Bayesian statistical information to generate causal models from large and complex phenotyping datasets.

Doc #3:

The researchers will build a computational information engine that uses machine learning to combine gene function and gene interaction information from disparate genomic data sources.

- First remove stop words and punctuation; detect manually multi-word terms (using N-Gram or POS Tagging/Chunking); parse manually the documents and select the terms from the given 3 documents and created the dictionary (list of terms).
- Create the document vectors by computing TF-IDF weights. Show how to compute the TF-IDF weights for terms. For each form of weighting list the document vectors in the following format:

	Term1	Term2	Term3	Term4	Term5	Term6	Term7	Term8 ...
DOC1	0	3	1	0	0	2	1	0
DOC2	5	0	0	0	3	0	0	2
DOC3	3	0	4	3	4	0	0	5

9) Stop words:

Basically the stop words are those which do not contain important significance to be used in search queries. Usually these words are filtered out from queries because they return huge amount of data.

→ Removal of stop words / punctuation:

Doc #1: Dr. researches focus computational phenotyping produce disease prediction models machine learning statistical tools.

Doc #2: Researcher develop cross Bayesian statistical information generate causal models large complex phenotypic datasets.

Doc #3: researchers build computational information engine uses machine learning combine gene function gene interaction information disparate genomic data sources.

→ N-gram approach:

In this scenario, I am considering N-gram size as 3 ($N=3$)

Doc #1: The researcher will
 researcher will focus
 will focus on
 focus on computational
 can computational phenotyping
 computational phenotyping and
 phenotyping and will
 and will produce.
 so on
 and statistical tools.

the output for $N=3$ after removal stop words for doc#1 is

researcher focus computational

focus computational phenotyping

computational phenotyping produce

phenotyping produce disease

produce disease prediction

disease prediction model

prediction models machine

models machine learning

machine learning statistical

learning statistical tools

The N-gram
value - 3

Doc #2:

researcher develop tools

develop tools to Bayesian

tools Bayesian statistical

Bayesian statistical information

statistical information generate

generate causal models

causal models large

models large complex

large complex phenotyping

complex phenotyping datasets

N-gram - 3

Doc #3:

researcher build computational

build computational information

computational information engine

information engine uses

engine uses machine

uses machine learning

machine learning combine

learning combine gene

combine gene function

gene function gene

function gene interaction

N-gram
value - 3

interaction information -disparate
information disparate genomic
disparate genomic data
genomic data sources

→ parse manually the documents by select the terms from given 3
docs:

	#doc 1	doc #2	doc #3	Count in all docs
researcher	1	1	1	3
focus	1	0	0	1
Computational	1	0	1	2
phenotyping	1	1	0	2
proteins	1	0	0	1
dataset	1	0	0	1
prediction	1	0	0	2
models	1	1	1	2
machine	1	0	1	2
learning	1	0	0	2
Statistical	1	1	0	1
tools	1	1	0	1
develop	0	1	0	2
Bayesian	0	1	1	1
information	0	1	0	1
generate	0	1	0	1
Causal	0	1	0	1
large	0	1	0	1
Complex	0	1	0	1
datasets	0	1	1	1
build	0	0	1	1
seem	0	0	1	1
engine	0	0	1	1
genomic	0	0	1	1

(4)

b) Term Frequency: It means how frequently a term occurs in a document since every document is different in length. It is possible that a term would appear much more times.

$$T(t)/t = \frac{\text{No. of times 't' appears in doc}}{\text{Total no. of terms in the doc}}$$

Inverse Document Frequency:

$$IDF(t) = \log_e \left(\frac{\text{Total no. of documents}}{\text{No. of documents with term 't' on it}} \right)$$

TFIDF: Its weight often used in information retrieval and text mining.

★ The matrix shows how many times a term occurs in document is already shown in previous page (3) ★
Based on the matrix, each word count in all docs

TF-IDF values for each term in Documents:
Doc #1:

Researcher 8-

$$TF = 1/12 \\ = 0.083$$

$$IDF = \log_e \left(\frac{3}{12} \right)$$

$$TF-IDF = \frac{1}{12} \times \log_e \left(\frac{3}{12} \right) \\ = 0.0146$$

Fans:

$$TF = 1/12$$

$$IDF = \log_e \left(\frac{3}{1} \right) = 0.477$$

$$TF-IDF = \frac{1}{12} \times 0.477 = 0.039$$

Computational:

$$TF = 1/12, IDF = \log_e \left(\frac{3}{2} \right) = 0.176$$

$$TF-IDF = 0.0146$$

Phenotyping:

$$TF = 1/12, IDF = \log_e \left(\frac{3}{2} \right) = 0.176$$

$$TF-IDF = 0.0146$$

Produce:

$$TF = 1/12, IDF = \log_e \left(\frac{3}{1} \right) = 0.477$$

$$TF-IDF = \frac{1}{12} \times 0.477 = 0.039$$

Disease:

$$TF = 1/12 \quad IDF = \log(3/1) = 3 \quad TF-IDF = \frac{1}{12} \times 0.477 = 0.039$$

Prediction:

$$TF = 1/12, \quad IDF = \log(3/1) = 0.477, \quad TF-IDF = \frac{1}{12} \times 0.477 = 0.039$$

Models:

$$TF = 1/12 \quad IDF = 0.176 \quad TF-IDF = 0.0146$$

Malware:

$$TF = 1/12 \quad IDF = \log(3/2) = 0.176, \quad TF-IDF = 0.0146$$

For learning:

$$TF = 1/12, \quad IDF = \log(3/2) = 0.176, \quad TF-IDF = 0.0146$$

Statistical:

$$1/12, \quad 0.176, \quad 0.0146$$

tools:

$$1/12, \quad 0.176, \quad 0.0146.$$

The other words in Doc #1 have $TF=0, \quad TF-IDF=0.$

Doc #2:

Develop:

$$TF = 1/13, \quad IDF = \log(3/1) = 0.477, \quad TF-IDF = 0.036$$

Bayesian:

$$1/13, \quad 0.477, \quad 0.036$$

Information:

$$1/13, \quad \log(3/2) = 0.176, \quad 0.0135$$

Researcher:

$$1/13, \quad \log(3/2), \quad 0$$

tools:

$$1/13, \quad \log(3/2), \quad 0.0135$$

Statistical:

$$1/13, \quad \log(3/2), \quad 0.0135$$

Generate, large, datasets, Casual, Complex:

⑥

$$TF = 1/18, \quad IDF = \log(3/1), \quad TF-IDF = 0.036$$

Phenotyping, models

$$TF = 1/18, \quad IDF = \log(3/2), \quad TF-IDF = 0.0135$$

Doc #3 :-

$$\text{Researchers: } TF = 1/18, \quad IDF = \log(3/3), \quad TF-IDF = 0$$

$$\text{Build: } TF = 1/18, \quad IDF = \log(3/1), \quad 0.026$$

$$\text{Computational: } TF = 1/18, \quad IDF = \log(3/2), \quad 0.978$$

$$\text{Information: } TF = 2/18, \quad IDF = \log(3/2), \quad 0.0195$$

Engine, uses, machine, combine, function, interaction, data, sensors,
genomic, disparate?

$$TF = 1/18, \quad IDF = \log(3/1), \quad TF-IDF = 0.026$$

$$\text{Gene: } TF = 2/18, \quad IDF = \log(3/1), \quad TF-IDF = 0.0529$$

$$\text{Learning: } TF = 1/18, \quad IDF = \log(3/2), \quad TF-IDF = 0.978$$

— END —