

# Starfish

Connecting the Docs

Finding implicitly related items based on semantic similarities and metadata in a non-hierarchical network of documents

---

## Authors

R. van Ginkels

J. Peters

L. Weerts

Project commissioned by *Perceptum B.V.*

## Supervisors

Academic Supervisor	Raquel Fernandez
Company Supervisor	Robrecht Jurriaans
	Sander Latour
	Wijnand Baretta

June 23, 2014  
Universiteit van Amsterdam

## Abstract

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Product overview</b>	<b>4</b>
2.1	Vectorizer . . . . .	4
2.1.1	Text-based transformation . . . . .	4
2.1.2	Tag-based transformations . . . . .	4
2.2	Distance . . . . .	5
2.2.1	StarFish specific adaptations: Bayesian weighting . . . . .	5
2.3	Threshold value . . . . .	6
<b>3</b>	<b>Method</b>	<b>6</b>
3.1	Data . . . . .	6
3.2	Text vectorization . . . . .	6
3.2.1	Textvectorizer . . . . .	6
3.2.2	Weighted textvectorizer . . . . .	6
3.3	Tag vectorization . . . . .	6
3.3.1	Simple tag vectorizer . . . . .	6
3.3.2	Tag smoothing . . . . .	7
3.3.3	Glossaries of tags . . . . .	7
3.4	Distance metrics . . . . .	7
3.5	Bayesian weighting . . . . .	7
3.6	Thresholds . . . . .	7
<b>4</b>	<b>Experiments</b>	<b>7</b>
<b>5</b>	<b>Conclusion</b>	<b>7</b>
<b>6</b>	<b>Future Work &amp; Recommendations</b>	<b>8</b>

# 1 Introduction

This report describes the results of the Second Year's project of the Perceptum team. The project focused on creating a *document link recommender system* to the StarFish website.

StarFish, one of the projects of Perceptum, is a website that aims to share knowledge about the education domain by means of a connected graph. People from all around the world should get access to this knowledge graph in a simple, personalized manner. The nodes in this graph are documents and they are connected with links. These documents can be of all sorts of types - e.g. a good practice, information, a question. Each document has a set of tags associated with it, which describe the different aspects of educational innovation. StarFish is community-driven: both the content of the documents as the links between documents are determined by the users of StarFish.

The drawback of a community driven knowledge graph is that not all the users know the entire document base. Especially when the knowledge base grows it becomes impossible for a user, since one does not know the existence of one or more linkable documents. A possible solution could be to make use of administrators, which can devote more time in getting to know all the documents, but that approach has two main drawbacks. First of all, this would mean that some central authority determines whether or not two documents should be linked. This is not in line with the idea of a community-driven knowledge base. Secondly, if the knowledge base grows even further, it becomes impossible also for an administrator to keep track of all documents. Imagine one person having to link all pages on Wikipedia - an impossible job.

In order to overcome the problem of linking documents in a large knowledge base, this process should be automated. This project therefore focuses on automating making the connections between documents. Though ideally these connections should be made completely automatic, a first step would be to create a recommendation system. When a user adds a new document, he or she can choose from a list of proposed documents the documents he or she deems relevant. This means that the recommender system does not have to work perfect, but should work reasonably well enough. Defining 'well enough', however, is also a part of this project. Thus, the product vision of the system can be described in the following concise way:

## Product vision:

---

**For** StarFish users

**who** search for and edit knowledge in starfish

**the** starfish document linker

**is** a starfish core system addition

**that** finds related documents

**unlike** moderated or individual linking

**our** product uses algorithms and data to suggest document links

Within the time span of this project multiple ways of recommending links between documents have been explored. The results of these explorations will be discussed in this report.

## 2 Product overview

The product created in this project is a python program that takes a set of documents and a new document and returns the subset of documents that should be linked with the new document. For this, a descriptor-based approach was used, which consists of three steps. First, each of the documents is transformed into a descriptor: a vector containing numerical values that in some way describes the document (hence the term 'descriptor'). Creating these descriptors is not trivial and during the project several techniques have been explored. Secondly, a ranking is made of all documents based on the similarity of the document descriptors and the descriptor of the new added document. To compare the descriptors the Nearest Neighbour algorithm was implemented, including five different distance metrics that determine how near two vectors are. Thirdly, an algorithm chooses the proper amount of proposed links that must be returned.

```
python documentlinker.py -vectorizer <vectorizername>
-distance <distance metric>
-bayes <true/false>
-threshold <'auto' or a fixed number>
```

We will now discuss each of these parameters, since these will give more insight into the approach that was chosen to solve the problem. For the performance of the different parameters we refer to the evaluation section.

### 2.1 Vectorizer

The first step is to create document descriptors, which is done by algorithms that we call *vectorizers*. Two main paths have been explored: transformation based on text and transformation based on tags.

#### 2.1.1 Text-based transformation

**Textvectorizer** The text-based vectorizers use the textual content of the documents and are therefore generally applicable to other systems. The textual content is first transformed into a *bag of words*. Then, based on all the documents in the knowledge base, the *TF-IDF* value is calculated for each of the words in the bag of words. *TF-IDF* stands for Term Frequency-Inverse Document Frequency and is a number that represents the importance of a word to a document in a bigger set of documents. Thus, the document descriptor consists of a vector with all TF-IDF values for that document of all words in the corpus.

**Weighted\_textvectorizer** The weighted textvectorizer is implemented as an extension of the textvectorizer. Besides the descriptor of a document itself, this method also adds the vectors of documents linked to it with some weight. This captures the idea that if a new document resembles some of the documents that are linked to one particular document, it is more likely to be linked to this particular document.

#### 2.1.2 Tag-based transformations

**Simple\_tag\_similarity** The tag-based transformations are more StarFish specific, since they make use of the tags that are assigned to the documents. A tag is a keyword that describes a topic/term that is important for that document. For example, 'Online Support and Online Assessment for Teaching and Learning Chemistry' is tagged with 'chemistry', 'e-learning' and 'assessment'. The simple tag similarity vectorizer creates a vector where each value indicates whether or not one particular tag is assigned to the document.

**Tag\_smoothing** The tag smoothing vectorizer uses the co-occurrence of tags in estimating document similarity. Even though tags might not co-occur on any document in the data set, they can still provide information about each other. For example, the dataset exists of documents with associated tags like  $\{\{t_1, t_2\}, \{t_1, t_3\}\}$ . From the co-occurrence it does not follow that  $t_2$  and  $t_3$  are related, however by transitivity with  $t_1$  we want to create a small implicit link between  $t_2$  and  $t_3$ . The tag smoothing method does this based on work from Zhou et al. (2011).

**Glossaries\_of\_tags** Another way of capturing tag similarity is by using tag Glossaries. Most of the tags have a Glossary - a special type of Document which holds an explanation of a tag. Though glossaries are documents, they cannot be assigned as a link since this should be done by assigning a tag. The glossaries can still be used by applying a text-based transformation on the glossaries to indicate the similarity between tags. Thus, glossaries\_of\_tags can be seen as a hybrid form of the tag and text-based approaches, where the glossary of a tag is turned into a TF-IDF bag of words. The document descriptor consists of the sum of vectors of each of it's tags.

**Weighted\_tag\_vectorizer** This is an extension of glossaries of tags, where a weight is assigned to the tag vectors.

## 2.2 Distance

The nearest neighbour algorithm loops through all available document descriptors and compares these with the descriptor of the new document. The closer related the descriptors are, the higher their ranking will be. The distance metrics define the closeness of the descriptors - a lower distance means a closer relation. The following were implemented:

**Euclidian**

**Cosine**

**Bhattacharyya**

**Correlation**

**Intersection**

### 2.2.1 StarFish specific adaptations: Bayesian weighting

Both the tag-based and text-based approaches uses some kind of 'semantic similarity' - the similarity of tags or text. However, except for the weighted text vectorizers, no information about possible links is used. For example, the text on a person's profile might be similar to other persons, but within StarFish a person is almost never linked to another person. In the Bayesian Weigthed Vectorizer this is captured by weighting the vectors with the probability that two documents are linked together:

$$P(D_a \rightarrow D_b|t)$$

Thus, the weight of a tag within a vector is equal to the chance that given this particular vector, a document of type a (the type of the newly added document) and a document of type b (equal to the type of proposed link) are linked together.

## 2.3 Threshold value

In the end, an algorithm chooses the proper amount of proposed links that must be returned. If the threshold value is set to a fixed number, e.g. 10, than only the proposed links of rank 1-10 are returned. If the threshold is set to *'auto'*, the number of returned links is based on the gradient of the distances. For example, if the algorithm is only certain that the first two links are correct, it returns no more than two links. For user-friendliness a maximum of 15 links is returned.

## 3 Method

### 3.1 Data

- How big is the dataset
- Types of documents
- Properties of documents (title, author ect)
- Tags: meaning, aliases and glossaries
- Linking between documents: probabilities & standard linking of authors

### 3.2 Text vectorization

#### 3.2.1 Textvectorizer

- Motivation: straightforward, standard approach
- Implementation: using bag of words and TF-IDF and adding stopwords
- Evaluation: is good with questions (since text has many keywords), bad in persons because person descriptions are similar. Bad at texts in different languages, cannot use images (which will become more prevalent in starfish). Quite slow.

#### 3.2.2 Weighted textvectorizer

- Motivation: use more of the links in starfish
- Implementation
- Evaluation: works better than textvectorizer (find out why!!!) but still has the same problems. Quite slow but not much slower than textvectorizer.

### 3.3 Tag vectorization

#### 3.3.1 Simple tag vectorizer

- Motivation: use more starfish specific things
- Implementation: binary vectors
- Evaluation: works surprisingly well, only works bad on questions because no proper tagging, works extremely fast

### 3.3.2 Tag smoothing

- Motivation: use connection between tags
- Implementation: based on the paper
- Evaluation: very slow, similar results as the simple tag vectorizer, 'magic' so cannot really see why something happens

### 3.3.3 Glossaries of tags

- Motivation: find underlying network between tags by using their glossaries
- Implementation: hybrid form of text and tag vectorizer
- Evaluation: does not work well (find out why!)

## 3.4 Distance metrics

- Euclidian
- Cosine
- Bhattacharyya
- Correlation
- Intersection

## 3.5 Bayesian weighting

Explain motivation, implementation and short results

## 3.6 Thresholds

Explain motivation, implementation and short results

# 4 Experiments

Overview of performance over the entire pipeline

# 5 Conclusion

- Which vectorizers do not work
- Which vectorizers do work
- How well does the bayesian layer perform
- How well does the thresholds perform

## 6 Future Work & Recommendations

### References

Zhou, J., Nie, X., Qin, L., and Zhu, J. (2011). Web clustering based on tag set similarity. *Journal of computers*, 6(1):59–66.