

# Fast Cheap and Synthetic Oracle (FACSO) Proximity Measures to Capture Expert Knowledge in the "Bibliome"

Luis M. Rocha<sup>†‡</sup> and Andreas Rechtsteiner<sup>†</sup>

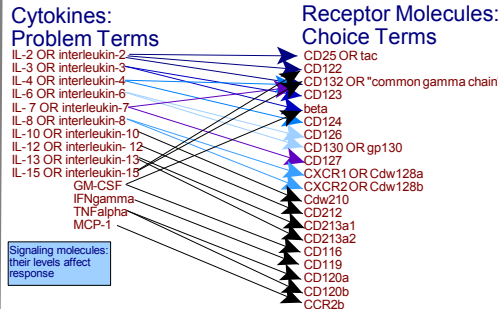
<sup>†</sup>CCS3 - Modeling, Algorithms, and Informatics, Los Alamos National Laboratory, MS B256, Los Alamos, NM 87545, USA

<sup>‡</sup>Instituto Gulbenkian de Ciência, 2781-901, Oeiras, Portugal

### Abstract

The present-day existence of very large electronic collections of documents, particularly in Biomedicine, allows simple information retrieval (IR) statistical methods to automatically capture expert knowledge with very good accuracy. Here we present a search utility (FACSO) built at Los Alamos, which queries the *Altavista*<sup>™</sup> search engine to determine associations between two classes of biological objects in Immunology (though applicable to any domain). This search utility uses measures of proximity built from the co-occurrence of keywords in documents. We present the accuracy of this utility, and contrast it to the Pointwise Mutual Information - IR method. A new version of this tool using PubMed is now under development.

## 1. Example: Expert Knowledge



## 2. Co-Occurrence Proximity

Given a binary relation  $A$  between sets of keywords  $K$  and documents  $D$  we extract a co-occurrence proximity measure:  $KDP(k_i, k_j)$  is the probability that both keywords  $k_i$  and  $k_j$  co-occur in the same document  $d \in D$ .

$$kdp(k_i, k_j) = \frac{\sum_{k=1}^m (a_{i,k} \wedge a_{j,k})}{\sum_{k=1}^m (a_{i,k} \vee a_{j,k})} = \frac{N_{\cap}(k_i, k_j)}{N_{\cup}(k_i, k_j)}$$

(Keyword Document Proximity)

## 3. Altavista Queries

### 3.1 Proximity Queries

$$prox_1(\text{problem}, \text{choice}_i) = \frac{\text{hits}(\text{problem AND choice}_i)}{\text{hits}(\text{problem OR choice}_i)}$$

$$prox_2(\text{problem}, \text{choice}_i) = \frac{\text{hits}(\text{problem NEAR choice}_i)}{\text{hits}(\text{problem OR choice}_i)}$$

**NEAR:** co-occurrence within 10 words

$$prox_3(\text{problem}, \text{choice}_i) = \frac{\text{hits}(\text{problem NEAR choice}_i \text{ NEAR context})}{\text{hits}(\text{problem OR choice}_i)}$$

**context** = {receptor}

### 3.2 Conditional Probability Queries

Turney, P.D. (2001). "Mining the Web for Synonyms: PMI-IR versus LSA on TOEFL." *Proceedings of the Twelfth European Conference on Machine Learning (ECML-2001)*, Freiburg, Germany, pp.491-502. <http://extractor.iit.nrc.ca/reports/ECML2001.html>

$$score_1(\text{choice}_i) = \frac{\text{hits}(\text{problem AND choice}_i)}{\text{hits}(\text{choice}_i)}$$

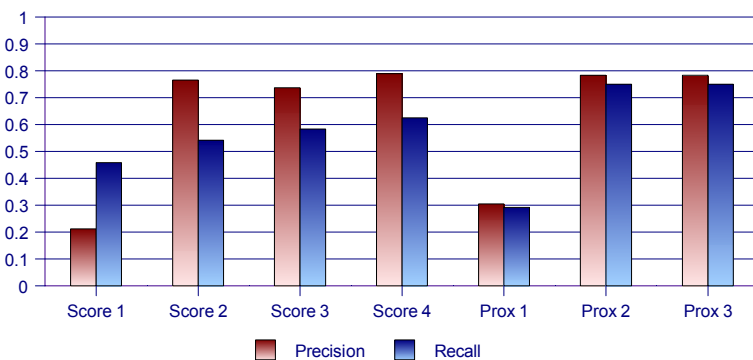
$$score_2(\text{choice}_i) = \frac{\text{hits}(\text{problem NEAR choice}_i)}{\text{hits}(\text{choice}_i)}$$

$$score_3(\text{choice}_i) = \frac{\text{hits}(\text{problem NEAR choice}_i) \text{ AND NOT}(\text{problem OR choice}_i \text{ NEAR "not"})}{\text{hits}(\text{choice}_i \text{ AND NOT}(\text{choice}_i \text{ NEAR "not"})}$$

$$score_4(\text{choice}_i) = \frac{\text{hits}(\text{problem NEAR choice}_i \text{ NEAR context AND NOT}(\text{problem OR choice}_i \text{ NEAR "not"})}{\text{hits}(\text{choice}_i \text{ NEAR context AND NOT}(\text{choice}_i \text{ NEAR "not"})}$$

## 4. Relevant Associations

- Normalize proximity/score of choice terms for each problem term (Column Normalization)
- Calculate  $\alpha$ -cut ( $\alpha = 0.5$ )
  - ▶ Choice terms  $\geq \alpha$  are considered relevant associations



**Precision:** probability that an identified association is relevant

$$precision = \frac{|\{relevant\} \cap \{retrieved\}|}{|\{retrieved\}|}$$

**Recall:** probability that an association has been identified given that it is relevant

$$recall = \frac{|\{retrieved\} \cap \{relevant\}|}{|\{relevant\}|}$$

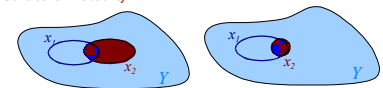
{relevant} = {associations defined in 1}

$$KDP(k_i, k_j) = \frac{1}{\frac{1}{P_K(k_j|k_i)} + \frac{1}{P_K(k_i|k_j)} - 1}$$

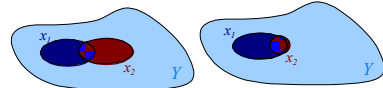
Where

$$P_K(k_i|k_j) = \frac{\sum_{k=1}^m (a_{i,k} \wedge a_{j,k})}{\sum_{k=1}^m (a_{j,k})} = \frac{N_{\cap}(k_i, k_j)}{N(k_j)}$$

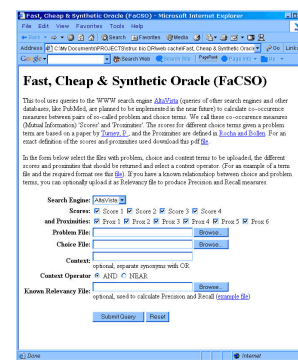
Conditional Probability



Proximity



## 5. Web Utility



# Boolean Score/proximity after using alpha-cut

```
1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0
1 0 1 0 0 0 0 0 1 1 0 0 0 0 0 0
0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 1 0 0 0 0 0 1 0 0 0 0 0 0 0
0 0 1 0 0 0 0 0 1 0 0 0 0 0 0 0
0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0
```

# Number of hits in the numerator of the score/proximity

```
1238 3 31 34 8 11 32 15 7 7 8 0 1 8
108 2 1 0 0 0 0 0 0 0 10 0 0 0 0
93 2 24 0 12 0 0 0 0 4 7 0 0 0 0
9 48 1 0 0 0 1 1 1 0 2 3 0 0 0 0
3 9 35 2 1 0 0 0 0 4 0 0 0 0 0 0
2 0 2 32 0 0 0 0 0 0 0 0 0 0 0 0
7 19 4 66 1 15 8 11 11 2 0 20 0 1 1
2 0 1 9 31 0 3 0 0 0 0 0 0 0 0 0
2 0 1 0 1 172 0 0 0 1 0 0 0 0 0 4
2 1 0 3 0 201 1 0 1 0 1 0 0 0 0 11
0 0 0 0 0 0 0 0 2 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 3 2 1 0 0 0 0 0
1 0 2 0 0 0 0 2 6 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 6 1 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 24 0 0 0 0
0 0 2 0 0 0 0 0 0 2 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0
0 0 0 0 0 0 3 0 0 0 0 0 0 0 0 15
```

- Discovers Relevant Associations
- Retrieves Documents substantiating the associations
- Being Developed for PubMed