# Usage of Semantic Similarity Measurements

## Tufts University Department of Computer Science

Thomas Schaffner

21 May 2016

# Abstract

Study and comparison of protein function is an important research topic in modern biology and bioinformatics. Better understanding of protein function aids in targeting medical and pharmacological research. Ontologies of functional terms organize and give structure to possible protein functions, while annotation corpuses apply functional labels to specific proteins. Many methods exist to compare protein functional annotations. These methods range from simply counting the number of overlapping functional labels to more complex methods that make use of the structure of ontologies. Specifically, we look at the Resnik semantic similarity measurement. Resnik scores make use of both the structure of an ontology and the distribution of functional labels throughout an annotation corpus. In this thesis, we see that incomplete data can lead to erroneous low Resnik values, while high Resnik values are likely to be more meaningful.

Using the Gene Ontology Consortium's ontology (GO) and annotation corpuses from UniProt Swiss-Prot and the Saccharomyces Genome Database (SGD), we analyze Resnik scores. We create matrices of Resnik scores for each species, representing the Resnik values between all pairs of proteins within a species. Using these matrices, we show that even high quality datasets such as SGD and UniProt Swiss-Prot do not completely label their proteins, leaving many proteins labeled with very general functions. We go on to discuss methods for identifying high and low Resnik values. We also show that matrix completion methods do not appear effective in predicting functional similarity between two proteins within a species.

1

# Contents

# Chapter 1

# Protein Comparison

## 1.1 Introduction

### 1.1.1 Proteins

Proteins perform many distinct functions within cells, from metabolism to regulation and production of other proteins. Any given species may have tens of thousands of unique proteins, each with their own roles in the life and activity of a cell. Understanding how proteins operate and interact is vital to modern medical and biological research.

By studying specific proteins, researchers can identify mechanisms at work in cancers and infectious diseases. By finding proteins involved in ailments, researchers can better target medical research or drug development. Additionally, by leveraging the genetic similarities of different species, researchers can find proteins that fill the same niche across multiple species. While developing a drug that targets a human protein may be extremely expensive, it can be cheaper and faster to experiment with proteins in a model organism that perform a similar function.

With modern technology, biologists are able to determine protein interactions and functions at unprecedented speeds. As the size of available data grows, both the possible benefits and the complexity of interpretation grow as well. Computational techniques make use of standardized semantics for discussing proteins in order to

compare proteins and even predict protein function.

## 1.1.2   Functional Annotations

In order to discuss the functions of proteins in large data sets and across species, researchers have established conventions for identifying and labeling protein functions. These standards are broken into two pieces: functional ontologies and annotation corpuses. Ontologies represent the set of all possible functional labels, as well as the relationships between these labels. Annotation corpuses apply specific labels to individual proteins. Several different ontologies exist, and annotation corpuses differ between datasets.

### Ontologies

GO, from the Gene Ontology Consortium [1], is a commonly used protein functional ontology [2]. GO provides a structured set of labels for protein function. The labels in GO are organized as a directed acyclic graph. Each node represents a functional label, while edges between nodes represent different relationships. GO contains data for six relationships between terms: "is-a", "part-of", "regulates", "positively-regulates", "negatively-regulates", and "has-part". For this thesis, we have only examined "is-a" relationships; if a functional label $l_i$ is a specification of another term $l_j$, an edge exists from $l_j$ to $l_i$. Note that any label may have multiple children and multiple parents. GO is partitioned into three separate domains. Each of these domains consists of its own directed acyclic graph. The root terms of these DAGs, or domain roots, are "cellular component", "biological process", and "molecular function".

Because different child terms of a single ancestor term may differ in their specificity, the "depth" of a term in the GO graph is not necessarily meaningful. For example, the GO terms for "single-organism organelle organization" and "transposition" are both children of the GO term "single-organism cellular process", and both have a depth of 4 in the GO graph. However, within the SGD annotation corpus

used in this thesis, "single-organism organelle organization" is used to label more than ten times the number of proteins that "transposition" labels (279 compared to 20). Even though both terms share the same depth in GO, and are in fact siblings, they are not equally common functional labels. GO is only able to provide explicit information on the relationships between a parent and a child. It is always true that a child is more specific than a parent, but it is not necessarily true that a node is equally specific as its siblings.

The terms specified in GO are not very useful by themselves. Protein functions are more meaningful when associated with actual proteins. Annotation corpuses contain the known associations of proteins with their functions. These associations are determined either through biological experimentation or through computational prediction.

**Annotation Corpuses**

Different groups may publish their own annotation corpuses, based on their own experimentation or functional prediction. Different groups may also compile others' data. Some experiments may be more error-prone than others, and computational prediction of protein function also comes with less than complete confidence. Because of these potential errors or sources of noise, annotation corpus providers may require different levels of confidence before approving a functional annotation and adding it to the corpus. All these differences lead to many slightly different datasets, potentially focusing on different species, different diseases, or proteins with specific functions.

Annotation corpuses are also far from complete. Even if all functional annotations were one hundred percent accurate, many proteins have never been tested for function. The biological experiments to determine function take time and resources, and therefore many proteins are currently overlooked and untested. Even when proteins have labels in an annotation corpus, those labels are not always terribly informative. For 1148 of 7014 proteins in the Saccharomyces Genome Database

(SGD), the most specific functional label is a domain root of GO. For roughly one seventh of the proteins in this dataset, we can provide no function more specific than "molecular function", "biological process", or "cellular component".

In this thesis, we will primarily discuss two annotation corpuses: the SGD annotation corpus mentioned above, and the Swiss-Prot annotation corpus provided by UniProt. Swiss-Prot is a multi-species annotation corpus. It contains labels for proteins that belong to various species. Swiss-Prot includes only experimentally-supported functional labels. SGD provides annotations solely for *Saccharomyces cerevisiae*. The labels in SGD are also solely experimentally-derived.

### 1.1.3   Data Sources

**Gene Ontology**

All calculations were performed using the GO release of June 13, 2015. The file used was go.obo, downloaded from geneontology.org

**Annotation Corpuses**

Two annotation corpuses were used in this thesis.

The UniProt Swiss-Prot annotation corpus was downloaded from uniprot.org on June 23, 2015. The SGD annotation corpus was downloaded from yeastgenome.org on December 27, 2015.

### 1.1.4   Functional Prediction

Predicting the functional labels of experimentally untested proteins is an important research problem. Predictions can be used to identify potentially interesting proteins for medical or biological research, among other uses. Many algorithms exist for protein functional prediction, often using additional biological data. Various methods can be used to append predicted functional labels within an annotation corpus, such as simple majority voting [3] or more complex algorithms that make

use of local neighborhoods within the network [4, 5], clustering within the network [6, 7, 8], and more [9, 10, 11, 12].

## 1.2 Techniques

In order to compare proteins, researchers look at several different aspects of each protein depending on their methodology. Proteins can be compared based on genetic sequence, structure, or the function of a protein. As discussed above, data regarding protein interactions can also be used to compare proteins and find similarities.

BLAST [13] is a tool that compares gene DNA sequence or protein amino acid sequence. BLAST produces a score measuring the similarity between sequences without examining protein functional labels, and functional annotations typically transfer if the BLAST similarity score is exceeds some threshold.

Another common tool is Pfam, which compares the secondary structure of proteins using machine learning [14]. Pfam uses hidden Markov models to classify proteins into groups with similar secondary structures.

Other techniques for comparing the function of proteins rely on existing information about the function of the proteins being compared. The simplest methods compare the exact labels of multiple proteins. For example, Jaccard GO [15, 16, 17] compares the number of functional terms shared by both proteins to the number of proteins associated with either protein. If $L(p)$ is the set of functional labels associated with protein "p", then:

$$JaccardGO(p_i, p_j) = \frac{|L(p_i) \cap L(p_j)|}{|L(p_i) \cup L(p_j)|}$$

Other techniques compare functional labels directly, relying on the structure of GO itself. Two methods $SimUI$ and $SimPE$ [18] make use of induced subgraphs of GO. They define the induced subgraph of a term $V(t)$ to be all nodes and edges present in all paths from the root to the $t$. For example, the highlighted nodes and edges in figure 1.1 represent the induced subgraph for GO:0044700. Note that figure

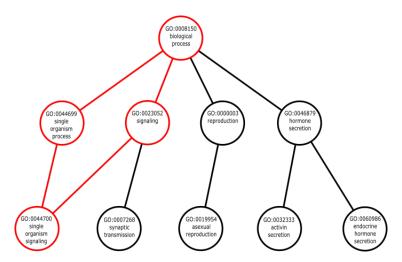1.1 refers to a partial representation of the GO DAG.



Figure 1.1: Induced Subgraph of GO:0044700 on partial GO DAG

$SimUI$ evaluates the subgraphs based on node count

$$SimUI(t_1, t_2) = \frac{|V(t_1) \cap V(t_2)|_n}{|V(t_1) \cup V(t_2)|_n}$$

while $SimPE$ evaluates the subgraphs based on edge count

$$SimPE(t_1, t_2) = \frac{|V(t_1) \cap V(t_2)|_e}{|V(t_1) \cup V(t_2)|_e}$$

Some techniques even make use of the distribution of functional labels within an annotation corpus. Annotation Corpuses are often used to generate information content (IC) measurements for specific terms, based on how well-represented terms are within the annotation corpus. These information content values are used in the Resnik semantic similarity comparison, as well as all of the modifications and adaptations of the Resnik comparison. IC is also widely used in other comparisons.

## 1.3   Information Content Comparisons

Information content is a measurement, originally from the field of natural language processing [19], used to indicate the specificity of a given functional term in GO. Information content is calculated based on a specific annotation corpus; two IC values are not comparable unless they are both calculated using the same annotation corpus.

### 1.3.1   Probability

The probability of a functional label, denoted $P(l)$, is equal to the portion of proteins in the annotation corpus that perform the function of $l$.

$$P(l) = \frac{|proteins(l)|}{|annotation\_corpus|}$$

Note that different curators may apply labels differently; given experimental evidence that a protein $p_i$ performs function $l_i$, one curator may choose to add the label $l_i$ by itself to the annotation corpus, while another curator may also add all ancestors of $l_i$. Either would be technically accurate, it is a matter of convention. The structure of the GO DAG means that $l_i$ is a specification of all of the ancestors of $l_i$. For example, if a protein performs a function related to synaptic transmission (GO:0007268), its function must also relate to signaling (GO:0023052).

In order to maintain consistency and accuracy when calculating probability, we define $proteins(l)$ as a function that returns all proteins in the annotation corpus that contain either $l$ or any of the descendants of $l$ as a functional label.

### 1.3.2   Information Content

If a term is close to its domain root, it is more likely to have a high probability value. Nearly all proteins in a given annotation corpus will contain the functional label "biological process" or one of biological process' descendants. Essentially, the

lower the probability of a functional label, the more specific it is likely to be. Therefore, information content is defined as

$$IC(l) = -ln(P(l))$$

The least-represented GO terms in a given annotation corpus will have a low probability and consequently a higher information content value. Because the probability of a term takes into account proteins labeled with that term's descendants, the probability of a term is always at least as large as any of the probabilities of that term's children. Therefore, a term in GO always has an information content that is at least as large as its parent term's information content value.

## 1.4 The Resnik Score

A Resnik similarity score is an indication of the level of similarity between two functional labels within the context of a common ontology and annotation corpus [19]. Given two functional labels $l_i$ and $l_j$, the maximum informative common ancestor (MICA) of those terms is the common ancestor of $l_i$ and $l_j$ with the highest information content.

All common ancestors of two functional labels represent functions shared by both terms. The MICA of two terms, therefore, represents the most specific or most informative shared function between the two functional terms being compared. The Resnik score of $l_i$ and $l_j$ is equal to the information content of the maximum informative common ancestor of $l_i$ and $l_j$.

$$simRes(l_i, l_j) = IC(MICA(l_i, l_j))$$

### 1.4.1 Resnik Scores Across GO Domains

Resnik scores rely on the MICA of two terms. If two functional labels are from different domains of GO, then they are part of separate DAGs. Therefore, they have

no common ancestors, and no MICA. In this case, GO is treated as having a true root node with three direct children: "molecular function", "biological process", and "cellular component" (the domain roots of GO). All functional labels in GO are descendants of this dummy root. Therefore, the probability

$$P(ROOT_{dummy}) = 1$$

The information content of the dummy root consequently equals zero. This dummy root allows us to treat the Resnik score between two functional labels from different GO domains as 0.

## 1.4.2 Resnik Scores and Distance

Note that Resnik scores between two GO terms are symmetric. However, a Resnik score is a measurement of similarity, not a distance metric. Resnik scores themselves do not obey the triangle inequality. For example, consider the terms "GO:0016075" (rRNA catabolic process), "GO:0007483" (genital disc morphogenesis), and "GO:0061558" (cranial ganglion maturation).

$$simRes(\text{GO:0016075}, \text{GO:0007483}) = 0.037477152772432756$$

$$simRes(\text{GO:0007483}, \text{GO:0061558}) = 3.7497179567995427$$

$$simRes(\text{GO:0016075}, \text{GO:0061558}) = 0.037477152772432756$$

This makes sense, because the Resnik comparison measures similarity rather than dissimilarity. However, even taking the multiplicative inverse of Resnik scores

does not provide a true distance metric. The reciprocals of Resnik scores between the functional labels "GO:0033471" (GDP-L-galactose metabolic process), "GO:1901805" (beta-glucoside catabolic process), and "GO:1901699" (cellular response to nitrogen compound) do not obey the triangle inequality.

$$\frac{1}{simRes(\text{GO:0033471}, \text{GO:1901805})} = 0.6011489018392736$$

$$\frac{1}{simRes(\text{GO:1901805}, \text{GO:1901699})} = 26.682923488669466$$

$$\frac{1}{simRes(\text{GO:0033471}, \text{GO:1901699})} = 1.908992864684634$$

Therefore, Resnik scores are not trivially interchangeable with other distance metrics between terms.

### 1.4.3 Modifications to Resnik Scores

Several semantic similarity measurements tweak the Resnik comparison slightly in order to normalize or adjust sensitivity. One measurement, $simLin$ [20], uses the information contents of the terms being compared to normalize Resnik scores.

$$simLin(l_i, l_j) = \frac{simRes(l_i, l_j)}{IC(l_i) + IC(l_j)}$$

### 1.4.4 Relative Specificity Similarity

Wu et al. have developed several other GO-based measurements of similarity between functional labels. Relative Specificity Similarity (RSS) is a similarity measurement that does not require an annotation corpus [21]. RSS relies on the most recent common ancestor (MRCA) of two terms. Wu et al. also define $\alpha$, $\beta$, and $\gamma$.

For terms $l_i$ and $l_j$

$$\alpha = max\left(|m \cap n|_n\right)$$

for all paths $m$ from root to $l_i$ and all paths $n$ from root to $l_j$

$$\beta = (gen(l_i), gen(l_j))$$

where $gen(l)$ is the smallest number of edges between $l$ and any descendant of $l$ that is a leaf (has no descendants of its own).

$$\gamma = dist(\text{MRCA}(l_i, l_j), l_i) + dist(\text{MRCA}(l_i, l_j), l_j)$$

$\alpha$ is the relative specificity of $l_i$ and $l_j$, $\beta$ is the relative generality of $l_i$ and $l_j$, and $\gamma$ is the sum of distances between $\text{MRCA}(l_i, l_j)$ and $l_i$ and $l_j$. Based on these definitions, RSS is defined as

$$RSS(l_i, l_j) = \frac{maxDepth^{GO}}{maxDepth^{GO} + \gamma} * \frac{\alpha}{\alpha + \beta}$$

### 1.4.5 Hybrid Relative Specificity Similarity

Wu et al. also modified their RSS measurement to make use of information content. They created Hybrid Relative Specificity Similarity (HRSS) using both the most informative leaf (MIL) and MICA of terms, as well as a Resnik comparison [22]. They redefine $\alpha_{IC}$, $\beta_{IC}$, and $\gamma_{IC}$ for terms $l_i$ and $l_j$ as follows.

$$\alpha_{IC} = simRes(l_i, l_j)$$

$$\beta_{IC} = \frac{(IC(MIL(l_i)) - IC(l_i)) + (IC(MIL(l_j)) - IC(l_j))}{2}$$

$$\gamma_{IC} = (IC(l_i) - simRes(l_i, l_j)) + (IC(l_j) - simRes(l_i, l_j))$$

## 1.5 Protein Comparisons

Proteins in the annotation corpus often have more than one label, and can have multiple labels that are very dissimilar functionally. However, many of these semantic similarity measurements only compare individual terms. In order to apply these measurements to protein functional comparison, several "mixing methods" are commonly used. The simplest method for comparing two proteins $p_i$ and $p_j$, and the method used in this thesis, is to simply take the maximum similarity score of all pairs of $(l_i, l_j)$, where $l_i$ is a functional label of $p_i$ and $l_j$ is a functional label of $p_j$. Another simple method is to take the average of all pairwise comparisons.

## 1.6 Discussion

### 1.6.1 Conditional Information Content

The relationship between the information content of a functional label and the information content of one of its descendants has a special meaning. A "conditional information content" score can indicate the similarity between an ancestor functional label and one of its descendants. Bayes' theorem shows that the difference in information content between an ancestor term and a descendant term is based on a conditional probability, defined as follows.

$$P(l_i \mid l_j) = \frac{|proteins(l_i) \cap proteins(l_j)|}{|proteins(l_j)|}$$

With this definition, it is possible to consider conditional information content values, defined below.

$$IC(l_i|l_j) = -ln(P(l_i \mid l_j))$$

When comparing two functional labels $l_i$ and $l_j$ where one label is an ancestor of the other, the intersection $proteins(l_i) \cap proteins(l_j)$ is guaranteed to be non-empty

as long as both $proteins(l_i)$ and $proteins(l_j)$ are both non-empty. If $l_i$ is the ancestor of $l_j$, note that every protein in $proteins(l_j)$ is also in $proteins(l_i)$ by definition. Therefore, the conditional probability $P(ancestor \mid descendant)$ will always be 1. According to Bayes' theorem

$$P(desc. \mid anc.) = \frac{P(anc. \mid desc.)P(desc.)}{P(anc.)}$$

Therefore

$$\text{IC}\,(\text{desc.} \mid \text{anc.}) = -\ln\,(\text{P}\,(\text{desc.} \mid \text{anc.}))$$

$$= -\ln\left(\frac{\text{P}\,(\text{anc.} \mid \text{desc.})\,\text{P}\,(\text{desc.})}{\text{P}\,(\text{anc.})}\right)$$

$$= -1 * (\ln\,(\text{P}\,(\text{anc.} \mid \text{desc.})) + \ln\,(\text{P}\,(\text{desc.})) - \ln\,(\text{anc.}))$$

$$= -\ln\,(1) - \ln\,(\text{P}\,(\text{desc.})) + \ln\,(\text{P}\,(\text{anc.}))$$

$$= \text{IC}\,(\text{desc.}) - \text{IC}\,(\text{anc.})$$

$$(1.1)$$

This shows that the difference between the information content values of two functional labels has a special meaning if one label is the ancestor of the other. The difference in information content values represents the relative information content of the descendant with respect to the ancestor. If the difference $IC\,(desc.) - IC\,(anc.)$ is high, then the proportion $\dfrac{|proteins\,(desc.)|}{|proteins\,(anc.)|}$ must be low. Therefore, this conditional information content score can be used to differentiate between functional terms that are significantly more specific than their parent, and functional terms that provide very little additional information compared to their parents.

## 1.6.2 Resnik Scores

Resnik scores are often discussed in relation to functional similarity problems. New measurements are often compared to the Resnik semantic similarity comparison method, which is one of the simplest information content-based measures. However, some of the issues of Resnik scores are rarely discussed.

It is important to know that Resnik scores are not true metrics. Because Resnik scores do not obey the triangle inequality (and cannot be trivially manipulated to obey the triangle inequality) they cannot be substituted for common distance metrics such as shortest path. Data quality issues are also rarely mentioned in depth. The difference in meaning and meaningfulness of high Resnik score values and low Resnik score values is often neglected. However, we have seen that many proteins in up to date annotation corpuses are only labeled with general functional terms. It is important to remember, when using Resnik scores as a comparison, that low scores do not necessarily indicate dissimilarity.

## 1.6.3 Resnik Scores and Data Quality

When discussing functional comparisons between proteins, the data quality issues of annotation corpuses discussed earlier impact the meaningfulness of semantic similarity measurements. In the case of Resnik scores (the primary semantic similarity measurement for the rest of this thesis), incomplete annotation corpuses lead to erroneously low estimates of the Resnik scores between some proteins.

Several situations can result in a low Resnik score. Two proteins that are truly dissimilar will always have a low Resnik score, as long as they have at least a single correct functional label. If these similar proteins are labeled with their most specific functional terms but their MICA is a domain root, they will have an extremely low Resnik score. However, incomplete labelings of similar proteins can also lead to a low Resnik score. If a protein has multiple functions but is only labeled with one function, that protein will appear dissimilar (low Resnik score) to other proteins that share the missing label as a true function. Recall that many proteins can have

very general functional terms as their sole labels. Even if a function is not strictly missing, it is still possible that a protein is labeled with a more general version of its function. This can lead to lower Resnik scores, regardless of the true functional similarity of proteins.

Incomplete annotation corpuses cannot, however, lead to falsely high Resnik scores under the maximum mixing method. The Resnik score between two proteins under the maximum mixing method is equal to the highest Resnik score between any pair of terms across the proteins. Therefore, if two proteins have a high Resnik score, they must be labeled with two functionally similar terms. Added functional labels cannot reduce the Resnik score of two proteins. Only false annotations can create a falsely high Resnik score. Because the annotation corpuses used in this thesis exclusively use functional labels backed by experimental evidence, all functional labels that we used are more likely to be accurate than labels from annotation corpuses that allow computationally predicted annotations.

# Chapter 2

# Resnik as a Matrix

We have calculated Resnik scores for all yeast (Saccharomyces cerevisiae), mouse (Mus musculus), and human (Homo sapiens) proteins. Mouse and human Resnik scores were calculated using the UniProt Swiss-Prot annotation corpus, while yeast Resnik scores were calculated using the SGD annotation corpus.

For each species, we created a square matrix of Resnik scores. Table 2 shows the number of proteins and total Resnik scores for each species. Each row and column within a Resnik score matrix represents a single protein, with element $(i; j)$ representing the Resnik score between protein $i$ and protein $j$. Because $simRes(l_i; l_j) = simRes(l_j; l_i)$, these matrices are symmetric. However, not all pairs of proteins have functional labels from the same domain of GO. Consequently, table 2 shows that a number of Resnik scores within each species are cross-domain, resulting in Resnik scores of 0.

| Species | Number of Proteins | Total Entries | Cross-Domain Scores |
|---------|-------------------|---------------|---------------------|
| Human | 20205 | 408242025 | 65409425 |
| Mouse | 16711 | 279257521 | 19243345 |
| yeast | 7014 | 49196196 | 3876 |

## 2.1 New Properties

When examining a matrix of Resnik scores, it is possible to find new patterns and properties. The distribution of values can provide a basis for differentiating "high" and "low" Resnik scores. As discussed above, low Resnik scores can arise from several situations while high Resnik scores are more likely to be accurate. Distinguishing between high and low scores in the matrix can identify high-confidence Resnik comparisons. Namely, by examining the distribution of Resnik scores, it is possible to identify the high-value Resnik scores as high confidence.

### 2.1.1 Diagonals in Resnik Score Matrix

Other properties of the matrix can be used to find poorly-labeled proteins and low-confidence Resnik scores. A simple indicator of a protein's labeling quality is the diagonal of the matrix. A diagonal at index $i$ represents the Resnik score between $protein_i$ and itself. Under the maximum mixing method, this is equal to the information content of the most informative term of $protein_i$. Essentially, the diagonal at index $i$ is equal to $max\left(IC(l)\right)$ for all functional labels $l$ associated with $protein_i$. If a diagonal at index $i$ has a low value, the information content of the most informative functional label of $protein_i$ must be low. $protein_i$ cannot, therefore, be labeled with any specific functional labels. For reference, figure 2.1 shows histograms of the diagonal values for each species' Resnik matrix.

**Diagonals as Upper Bound on Resnik Scores**

The MICA of two functional labels $l_i$ and $l_j$ must be an ancestor of both $l_i$ and $l_j$. The information content of $MICA(l_i, l_j)$ must be less than or equal to $min\left(IC(l_i), IC(l_j)\right)$. Therefore, the Resnik similarity between two terms is bounded by the less informative term. A diagonal at index $i$ in a Resnik matrix is an upper bound on values in column $i$ and row $i$. If $protein_i$ does not have any highly informative functional labels, then $simRes(protein_i, protein_j)$ will always be low for any $j$. $protein_i$ can therefore be considered poorly labeled, and all comparisons with $protein_i$ (all entries
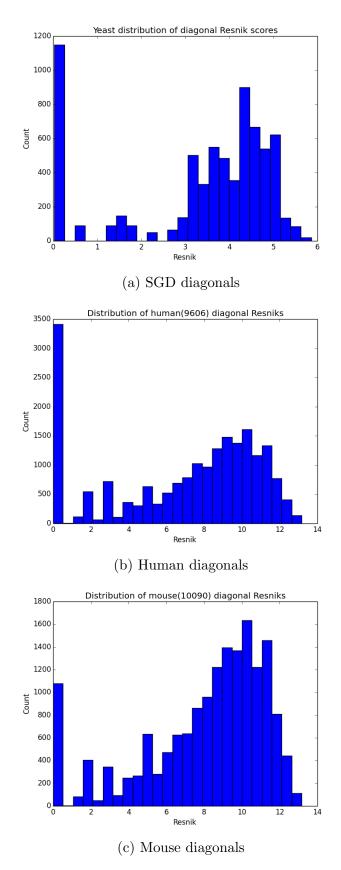
(a) SGD diagonals



(b) Human diagonals



(c) Mouse diagonals

Figure 2.1: Histograms of diagonal values in Resnik matrices

in row $i$ or column $i$ in the matrix) are likely not meaningful.

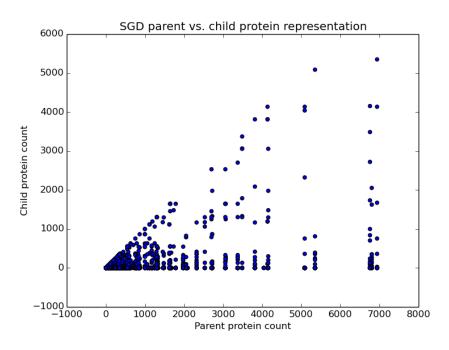## 2.1.2 Distribution of Specificity

Beyond diagonals, we can use statistics to determine low or high Resnik values. In the SGD Resnik matrix (generated using the SGD annotation corpus), Resnik values greater than 3.043 are two standard deviations above the mean. In the human matrix, the cutoff is 5.4435 while in mouse the cutoff is 5.9412. Alternately, we can classify some of the lowest values as less meaningful by comparing matrix entries to the IC values of highly-represented functional labels in the annotation corpus. The graphs in figure 2.2 show the size of $proteins(label)$ and $proteins(label) \cap proteins(label_{child})$ for each child $label_{child}$ of $label$. There is one data point on the graph for each parent-child relationship within GO.

In order to make use of these graphs for preliminary testing, we chose all data points with $|proteins(parent) \cap proteins(child)| \geq 3000$ for SGD, and $|proteins(parent) \cap proteins(child)| \geq 250000$ for UniProt Swiss-Prot. This leads to IC cutoff values of roughly 0.85 and 0.74 respectively.
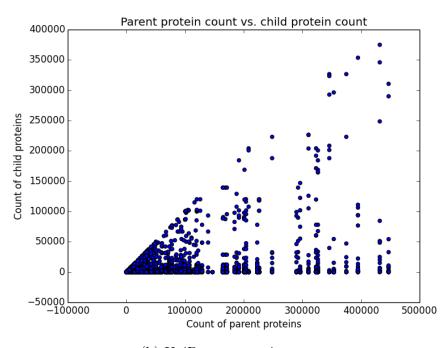
## 2.2 High Value Density

Looking at the SGD Resnik score matrix, we examined the density of high Resnik score values within each row of the Resnik score matrix. For this evaluation, a "high" Resnik score is a value greater than or equal to 3.043. For each row in the Resnik score matrix, we calculated the percentage of entries that were "high" values. Figure 2.3 shows a histogram with the distribution. For this histogram, we removed all percentages equal to 0 in order to make the rest of the graph more readable.

We noticed several outliers on the right of the graph. A few rows (corresponding to proteins) had exceptionally high percentages of high Resnik scores. The proteins five proteins with the highest percent of high Resnik scores were YBR160W, YER125W, YPL204W, YBR279W, and YER133W. According to yeastgenome.org,

(a) SGD annotation corpus



(b) UniProt annotation corpus

Figure 2.2: $|proteins(parent)|$ vs. $|proteins(parent) \cap proteins(child)|$

all of these proteins have regulatory functions. This similarity between the outliers may be useful for identifying regulatory proteins in the future.
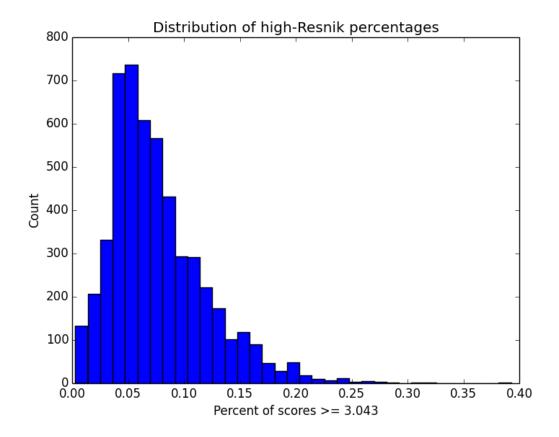


Figure 2.3: Distribution of high Resnik percentages in matrix rows

## 2.3 Resnik Score Cutoff

We have used a couple simple methods to classify Resnik scores as low or high in value. However, we have not performed experiments to determine the efficacy of the methods discussed in this thesis. In future work, we hope to refine and further study these methods. These methods would allow researchers to classify Resnik scores as high- or low-confidence with more accuracy and precision.

# Chapter 3

# Matrix Completion

## 3.1 Introduction

The problem of matrix completion, especially for matrices of low rank, has become widely studied recently. Matrix completion is applicable in several situations. Images can be represented as matrices, and matrix completion can be used to identify image subjects or components. In a more abstract setting, matrix completion can be used in the context of the Netflix prize; it can predict user preferences or ratings of movies.

## 3.2 Methods

Essentially, the matrix completion problem takes as input a matrix with some entries missing. A solution to the matrix completion problem returns a matrix with estimates in place of missing entries. We attempt to apply an existing matrix completion solution (LMaFit, translated to Python code by Professor Mark Crovella of Boston University) [23] to the Resnik matrices. LMaFit uses the nuclear norms method [24] for matrix completion. We treat low Resnik scores in the Resnik matrices as missing data entries. We then run LMaFit to estimate the missing values. Unfortunately, the preliminary results for matrix completion on Resnik matrices (discussed below in more detail) are not promising.

## 3.3 Results

We ran LMaFit on the SGD matrix using an IC cutoff value of 0.85 (as discussed in chapter 2). We first set all entries in the SGD matrix below 0.85 to 0. We then considered all remaining non-zero values to be "known" values, essentially high-enough confidence to keep. For each trial, we split the known values in half to run a 2-fold cross validation. Note, because the Resnik score matrices are symmetric, entries $(i, j)$ and $(j, i)$ were always partitioned into the same fold. For each fold, we set the opposite fold's values to zero and used LMaFit to predict non-zero values for all zeroed values in the matrix. The maximum Resnik score in the SGD matrix is roughly 5.8599, and we found an average error of 1.738 for all predicted matrix entry values using LMaFit.

For comparison, we also randomly selected entry values between 0 and the maximum value of the matrix (5.8599) to fill all zeroed entries. The average error for random completion was 1.771. Given the large error and negligible difference between LMaFit matrix completion results and random prediction results, LMaFit does not seem to accurately predict Resnik scores based on a Resnik score matrix.

## 3.4 Discussion

We only have preliminary results for matrix completion of Resnik score matrices. However, these results are not promising. With high error values roughly equivalent to randomly predicting missing entries, there is no evidence that matrix completion methods can predict functional similarity between proteins with any confidence. It may be worthwhile to investigate other methods for matrix completion, but there seems to be a low probability of success.

# Bibliography

[1] Michael Ashburner et al. "Gene Ontology: tool for the unification of biology". In: *Nature genetics* 25.1 (2000), pp. 25–29.

[2] Jiajie Peng et al. "Extending gene ontology with gene association networks". In: *Bioinformatics* (2015), btv712.

[3] Benno Schwikowski, Peter Uetz, and Stanley Fields. "A network of protein–protein interactions in yeast". In: *Nature biotechnology* 18.12 (2000), pp. 1257–1261.

[4] Hon Nian Chua, Wing-Kin Sung, and Limsoon Wong. "Exploiting indirect neighbours and topological weight to predict protein function from protein–protein interactions". In: *Bioinformatics* 22.13 (2006), pp. 1623–1630.

[5] Haretsugu Hishigaki et al. "Assessment of prediction accuracy of protein function from protein–protein interaction data". In: *Yeast* 18.6 (2001), pp. 523–531.

[6] Vicente Arnau, Sergio Mars, and Ignacio Marin. "Iterative cluster analysis of protein interaction data". In: *Bioinformatics* 21.3 (2005), pp. 364–378.

[7] Gary D Bader and Christopher WV Hogue. "An automated method for finding molecular complexes in large protein interaction networks". In: *BMC bioinformatics* 4.1 (2003), p. 1.

[8] Jimin Song and Mona Singh. "How and when should interactome-derived clusters be used to predict functional modules and protein function?" In: *Bioinformatics* 25.23 (2009), pp. 3143–3150.

[9] Elena Nabieva et al. "Whole-proteome prediction of protein function via graph-theoretic analysis of interaction maps". In: *Bioinformatics* 21.suppl 1 (2005), pp. i302–i310.

[10] Ulas Karaoz et al. "Whole-genome annotation by using evidence integration in functional-linkage networks". In: *Proceedings of the National Academy of Sciences of the United States of America* 101.9 (2004), pp. 2888–2893.

[11] Alexei Vazquez et al. "Global protein function prediction from protein-protein interaction networks". In: *Nature biotechnology* 21.6 (2003), pp. 697–700.

[12] Roded Sharan, Igor Ulitsky, and Ron Shamir. "Network-based prediction of protein function". In: *Molecular systems biology* 3.1 (2007), p. 88.

[13] Stephen F Altschul et al. "Basic local alignment search tool". In: *Journal of molecular biology* 215.3 (1990), pp. 403–410.

[14] Alex Bateman et al. "The Pfam protein families database". In: *Nucleic acids research* 30.1 (2002), pp. 276–280.

[15] Paul Jaccard. *Etude comparative de la distribution florale dans une portion des Alpes et du Jura*. Impr. Corbaz, 1901.

[16] Paul Jaccard. "The distribution of the flora in the alpine zone." In: *New phytologist* 11.2 (1912), pp. 37–50.

[17] Daniele Merico et al. "Enrichment map: a network-based method for gene-set enrichment visualization and interpretation". In: *PloS one* 5.11 (2010), e13984.

[18] Yuan-Peng Li and Bao-Liang Lu. "Semantic Similarity Definition over Gene Ontology by Further Mining of the Information Content." In: (2008), pp. 155–164.

[19] Philip Resnik. "Using information content to evaluate semantic similarity in a taxonomy". In: *arXiv preprint cmp-lg/9511007* (1995).

[20] Pietro H Guzzi et al. "Semantic similarity analysis of protein data: assessment with biological features and issues". In: *Briefings in bioinformatics* 13.5 (2012), pp. 569–585.

[21] Xiaomei Wu et al. "Prediction of yeast protein–protein interaction network: insights from the Gene Ontology and annotations". In: *Nucleic acids research* 34.7 (2006), pp. 2137–2150.

[22] Xiaomei Wu et al. "Improving the measurement of semantic similarity between gene ontology terms and gene products: insights from an edge-and IC-based hybrid method". In: *PloS one* 8.5 (2013), e66745.

[23] Zaiwen Wen, Wotao Yin, and Yin Zhang. "Solving a low-rank factorization model for matrix completion by a nonlinear successive over-relaxation algorithm". In: *Mathematical Programming Computation* 4.4 (2012), pp. 333–361.

[24] Emmanuel J Candès and Benjamin Recht. "Exact matrix completion via convex optimization". In: *Foundations of Computational mathematics* 9.6 (2009), pp. 717–772.

# Appendix A

# SemSimCalculator

## A.1    semsimcalc.py

I have produced a python class that can be used to calculate several semantic simi-
larity measurements. The code is included as Appendix A, and the full README
can be found on github, under TuftsBCB/semsimcalc. Given a file representing GO
and an annotation corpus, the SemSimCalculator class can calculate probability and
information content values, Resnik scores, conditional probabilities, conditional in-
formation content scores, and a couple other semantic similarity measurements. The
SemSimCalculator class also contains two mixing methods for protein comparison:
maximum and average. The relationship between proteins and GO terms is stored
internally, so there is no need to look up labels of proteins separately.

## A.2    Code

```python
1   #!/usr/bin/python
2
3   # See http://bib.oxfordjournals.org/content/13/5/569.full
4   # For definitions
5
6   import sys
7   import time
8   import networkx as nx
9   import math
10  import pickle
11  import numpy
12
13  # Helper functions
14  def announce(message):
15          """ Timestamped output to stdout """
16          print time.strftime('%H:%M'),message
17          sys.stdout.flush()
18
19  def open_or_abort(filename, option='r'):
20          """ Output error message to stderr if file opening failed """
21          try:
22                  newfile = open(filename, option)
23          except IOError:
24                  sys.stderr.write("Could not open {} -- Aborting\n".format(filename))
25                  raise IOError
26          return newfile
27
28
29  # NOTE(tfs): Accepted GO file format:
30  #
31  #       ! comments
32  #
```

```
33  #          [Term]
34  #          id: GO_term
35  #          ...
36  #          is_a: GO_term
37  #          is_a: GO_term
38  #
39  #          [Term]
40  #          ...
41  #
42  #           [Typedef]
43  #
44  # The [Typedef] tag signals end of GO terms.
45  # It is necessary in the current implementation
46  #
47  def parse_go_file(go_file_name):
48          """ Parses and returns (does not natively store) GO data """
49
50          go_file = open_or_abort(go_file_name)
51
52          # Setup
53          go_file.seek(0)
54          go_graph = nx.DiGraph()
55
56          alt_ids = {}
57
58          go_term = ''
59          parents = []
60          is_obsolete = False
61
62          # Don't start paying attention until we see '[Term]'
63          valid_to_read = False
64
65          # Main parsing loop
66          for line in go_file:
67
68                  # Only if we're within a '[Term]' header
69                  if valid_to_read:
70                          if line.startswith('alt_id:'):
71                                  alt_id = line.strip()[8:]
72                                  alt_ids[alt_id] = go_term
73
74                          elif line.startswith('id:'):
75
76                                  # Only log if the entry is valid
77                                  if not is_obsolete:
78                                          if go_term != '':
79
80                                                  # Only add connected node
81                                                  if len(parents) > 0:
82                                                          go_graph.add_node(go_term)
83
84                                                          for parent in parents:
85                                                                  if parent != '':
86                                                                          go_graph.add_edge(parent, go_term)
87
88                                  # Reset regardless of logging status
89                                  parents = []
90                                  is_obsolete = False
91
92                                  go_term = line.strip()[4:]
93                          elif line.startswith('is_a:'):
94
95                                  # Store is_a as a parent
96                                  parents.append(line.split('!')[0].strip()[6:])
97
98                          elif line.startswith('is_obsolete: true'):
99
100                                 # Do not store the data under this '[Term]' header
101                                 is_obsolete = True
102
103                         elif line.startswith('[Typedef]'):
104                                 # Write if the previous entries were valid
105
```

29

```
106                                    # Only log if the entry is valid
107                                    if not is_obsolete:
108                                        if go_term != '':
109
110                                            # Only add connected node
111                                            if len(parents) > 0:
112
113                                                go_graph.add_node(go_term)
114
115                                                for parent in parents:
116                                                    if parent != '':
117                                                        go_graph.add_edge(parent, go_term)
118
119                                    # Reset regardless of logging status
120                                    parents = []
121                                    is_obsolete = False
122
123                                    go_term = '' # No valid ID to reset with under a '[Typedef]' header
124
125                                    # Stop paying attention
126                                    valid_to_read = False
127                        else:
128                            if '[Term]' in line:
129
130                                    # Start paying attention
131                                    valid_to_read = True
132
133            go_file.close()
134
135            return (go_graph, alt_ids)
136
137
138    # NOTE(tfs): Accepted AC file format:
139    #
140    #        -
141    #        protein_name
142    #        GO_term
143    #        GO_term
144    #        GO_term
145    #          -
146    #
147    def parse_annotation_corpus(ac_file_name, alt_ids=None):
148        """
149            Parses annotation corpus. Returns a dictionary of { gene: [terms] }.
150            If a term is a key in alt_ids, saves the associated value instead (if provided).
151        """
152
153        ac_file = open_or_abort(ac_file_name)
154
155        # Setup
156        prot_to_gos = {}
157        go_to_prots = {}
158
159        ac_file.seek(0)
160
161        curr_prot = ''
162        curr_gos = []
163        new_entry = True
164
165        for line in ac_file:
166
167                # Start information from new entry
168                if line.startswith('-'):
169
170                        # Only update if we have enough information for the last entry
171                        if curr_prot != '' and len(curr_gos) > 0:
172
173                                # Update prot_to_gos
174                                if curr_prot in prot_to_gos:
175                                        prot_to_gos[curr_prot] = prot_to_gos[curr_prot] + curr_gos
176                                else:
177                                        prot_to_gos[curr_prot] = curr_gos
178
```

```
179                             # Update go_to_prots
180                             for go in curr_gos:
181                                     if go in go_to_prots:
182                                             go_to_prots[go].append(curr_prot)
183                                     else:
184                                             go_to_prots[go] = [curr_prot]
185
186                         # Reset, regardless of whether or not we updated
187                         curr_prot = ''
188                         curr_gos = []
189                         new_entry = True
190
191                 # If we've just started looking at a new entry, parse as protein name
192                 # DON'T do this if we're still on the delimiter line ('-')
193                 elif new_entry:
194                         curr_prot = line.strip().strip(';')
195                         new_entry = False
196                 # Otherwise, parse as GO term
197                 else:
198                         if ("GO:" in line):
199                                 new_go = line.strip().strip(';')
200                                 if alt_ids is not None:
201                                         if new_go in alt_ids:
202                                                 new_go = alt_ids[new_go]
203                                 curr_gos.append(line.strip().strip(';'))
204
205         ac_file.close()
206
207         return (prot_to_gos, go_to_prots)
208
209
210 # Load a saved SemSimCalculator
211 def load_semsimcalc(saved_path):
212         """
213                 Loads (unpickles) a saved SemSimCalculator
214         """
215
216         return pickle.load(open(saved_path, 'rb'))
217
218
219
220
221
222 ###############################
223 ### SemSim_Calculator class ###
224 ###############################
225
226
227 class SemSimCalculator():
228         """
229                 Stores GO and annotation corpus data internally.
230                 Calculates different semantic similarity metrics.
231         """
232
233         def __init__(self, go_file_name, ac_file_name):
234                 """ Initialize using GO and annotation corpus files (pass in file name, not file object) """
235
236                 self._go_graph, self._alt_list = parse_go_file(go_file_name)
237                 self._prot_to_gos, self._go_to_prots = parse_annotation_corpus(ac_file_name, self._alt_list)
238                 self._proteins = [x[0] for x in self._prot_to_gos.items()]
239                 self._num_proteins = len(self._proteins)
240                 self._ic_vals = {} # For memoizing IC values (they are unchanging given an ontology and annotation corpus)
241
242                 self._go_terms = self._go_graph.nodes()
243
244                 self._mica_store = None
245
246         def link_mica_store(self, mica_store):
247                 """ Stores a reference to a MicaStore instance """
248
249                 self._mica_store = mica_store
250
251         def unlink_mica_store(self):
```

```
252                         """ Removes link to a MicaStore instance (sets to None) """
253
254                         self._mica_store = None
255
256                 def save(self, filepath):
257                         """
258                                 Saves (pickles) to filepath
259
260                                 NOTE: Does not save reference to MicaStore instance (as this will likely be broken on load)
261                         """
262
263                         # Do not store reference to MicaStore instance
264                         temp = self._mica_store
265                         self._mica_store = None
266
267                         pickle.dump(self, open(filepath, 'wb'))
268
269                         # Restore _mica_store reference
270                         self._mica_store = temp
271
272                 def get_go_graph(self):
273                         """ Return nx graph for GO """
274
275                         return nx.DiGraph(self._go_graph)
276
277                 def get_alt_list(self):
278                         """ Return alt_list """
279
280                         return dict(self._alt_list)
281
282                 def get_ptg(self):
283                         """ Return copy of prot_to_gos """
284
285                         return dict(self._prot_to_gos)
286
287                 def get_gtp(self):
288                         """ Return copy of go_to_prots """
289
290                         return dict(self._go_to_prots)
291
292                 def get_proteins(self):
293                         """ Return copy of proteins """
294
295                         return list(self._proteins)
296
297                 def get_num_proteins(self):
298                         """ Return number of proteins """
299
300                         return int(self._num_proteins)
301
302                 def get_ic_vals(self):
303                         """
304                                 Return all stored ic_vals.
305                                 Not all values are guaranteed to exist.
306                                 Consider running precompute_ic_vals first.
307                         """
308
309                         return dict(self._ic_vals)
310
311                 def get_go_terms(self):
312                         """ Return list of GO terms """
313
314                         return list(self._go_terms)
315
316                 def get_mica_store(self):
317                         """ Returns copy of mica_store """
318
319                         return self._mica_store
320
321                 def calc_term_prob(self, term):
322                         """ Probability of term or desc(term) to occur as a label within the annotation corpus """
323
324                         if term == None or (not term in self._go_graph):
```

```
325                     return None
326
327             # Find all descendants of term, including term
328             terms = nx.algorithms.dag.descendants(self._go_graph, term)
329             terms.add(term)
330
331             annotated_proteins = {}
332
333             # Mark any protein labeled with term or a descendant of term
334             for term in terms:
335                     if term in self._go_to_prots:
336                             for prot in self._go_to_prots[term]:
337                                     annotated_proteins[prot] = True
338
339             prob = float(len(annotated_proteins.items())) / float(self._num_proteins)
340
341             return prob
342
343     def calc_conditional_prob(self, term, condition):
344             """
345                     Probability that term or desc(term) appears
346                     as label in annotaiton corpus,
347                     given that condition appears as a term.
348             """
349
350             if term == None or (not term in self._go_graph):
351                     return None
352
353             # Find all descendants of condition, including condition
354             cond_terms = nx.algorithms.dag.descendants(self._go_graph, condition)
355             cond_terms.add(condition)
356
357             # Find all descendants of term, including term
358             terms = nx.algorithms.dag.descendants(self._go_graph, term)
359             terms.add(term)
360
361             conditional_proteins = {}
362             for cond_term in cond_terms:
363                     if cond_term in self._go_to_prots:
364                             for prot in self._go_to_prots[cond_term]:
365                                     conditional_proteins[prot] = True
366
367             restricted_term_proteins = {}
368             for r_term in terms:
369                     if r_term in self._go_to_prots:
370                             for prot in self._go_to_prots[r_term]:
371                                     if prot in conditional_proteins.keys():
372                                             restricted_term_proteins[prot] = True
373
374             if len(conditional_proteins.items()) == 0:
375                     return None
376             else:
377                     prob = float(len(restricted_term_proteins.items()))
378                     prob = prob / float(len(conditional_proteins.items()))
379                     return prob
380
381     def IC(self, term):
382             """ Information content: IC(c) = -log(p(c)) """
383
384             # Check if IC has been computed for term already
385             if not (term in self._ic_vals):
386                     # If not seen before, compute IC
387                     prob = self.calc_term_prob(term)
388
389                     if prob == 0 or prob == None:
390                             self._ic_vals[term] = None
391                             return None
392                     else:
393                             ic = (-1) * math.log(prob)
394                             self._ic_vals[term] = ic # Memoize IC value
395                             return ic
396             else:
397                     # If seen before, return memoized value
```

```
398                    return self._ic_vals[term]
399
400         def conditional_IC(self, term, condition):
401             """ Conditional Information Content: cIC(t | c) = -log(p(t | c)) """
402
403             # Too many values to memoize
404             cond_prob = self.calc_conditional_prob(term, condition)
405
406             if cond_prob == 0 or cond_prob == None:
407                 return None
408             else:
409                 cic = (-1) * math.log(cond_prob)
410                 return cic
411
412         def precompute_ic_vals(self):
413             """ Compute and store IC values for all ontology terms """
414
415             for term in self._go_graph.nodes():
416                 self.IC(term)
417
418         def MICA(self, left, right):
419             """
420                 Maximum Informative Common Ancestor:
421                 MICA(t1, t2) = arg max, IC(tj)
422                                        tj in ancestors(t1, t2)
423
424                 (returns a term, common ancestor of left and right)
425
426                 NOTE: If a MicaStore instance is linked, first try querying the stored instance
427             """
428
429             if not left in self._go_terms:
430                 if left in self._alt_list:
431                     left = self._alt_list[left]
432                 else:
433                     return None
434
435             if not right in self._go_terms:
436                 if right in self._alt_list:
437                     right = self._alt_list[right]
438                 else:
439                     return None
440
441             # Attempt lookup in linked MicaStore instance
442             if (self._mica_store != None):
443                 mica = self._mica_store.mica_lookup(left, right)
444
445                 if (mica != None) and (mica != '') and (mica != 'None'):
446                     return mica
447                     #if (mica == ''):
448                         # MICA is stored, but does not exist (None is a possible MICA value)
449                     #     return None
450                     #else:
451                     #     return mica
452
453             # Fall through and calculate MICA
454
455             # Find common ancestors as intersection of two ancestor sets
456             # NOTE(tfs): Python sets are very slow. List comprehensions are faster
457             left_ancs = nx.algorithms.dag.ancestors(self._go_graph, left)
458             left_ancs.add(left)
459             right_ancs = nx.algorithms.dag.ancestors(self._go_graph, right)
460             right_ancs.add(right)
461             ancestors = [a for a in left_ancs if a in right_ancs]
462
463             # Edge case where left and right are the same. Treat left and right as a common ancestor
464             #if left == right:
465             #     ancestors.append(left)
466
467             max_term = None
468             max_IC = 0
469
470             # Calculate IC for all ancestors; store maximum IC value and term
```

```
471                    for ancestor in ancestors:
472                            anc_IC = self.IC(ancestor)
473                            if anc_IC != None and anc_IC > max_IC:
474                                    max_IC = anc_IC
475                                    max_term = ancestor
476
477                    return max_term
478
479        def simRes(self, left, right):
480                """
481                        simRes(t1, t2) = IC[MICA(t1, t2)]
482                        Returns a value (IC result)
483                """
484
485                return self.IC(self.MICA(left, right))
486
487        def simLin(self, left, right):
488                """
489                        simLin(t1, t2) = [IC[MICA(t1, t2)]] / [IC(t1) + IC(t2)]
490                        Returns a value
491                        Currently untested
492                """
493
494                leftIC = self.IC(left)
495                rightIC = self.IC(right)
496
497                if leftIC == None or rightIC == None:
498                        return None
499                else:
500                        return self.IC(self.MICA(left, right)) / (leftIC + rightIC)
501
502
503        def simJC(self, left, right):
504                """
505                        simJC(t1, t2) = 1 - IC(t1) + IC(t2) - 2xIC[MICA(t1, t2)]
506                        Returns a value
507                        Currently untested
508                """
509
510                leftIC = self.IC(left)
511                rightIC = self.IC(right)
512
513                if leftIC == None or rightIC == None:
514                        return None
515                else:
516                        return 1 - self.IC(left) + self.IC(right) - (2*self.IC(self.MICA(left, right)))
517
518        def pairwise_average_term_comp(self, lefts, rights, metric):
519                """
520                        Compares each pair of terms in two sets or lists of terms.
521                        Returns the average of these comparison scores.
522                        Uses metric(left, right) to make each comparison.
523                        metric must take in two ontology terms (left and right) and return a numeric score.
524                """
525
526                total_score = 0
527                num_scores = 0
528                for left in lefts:
529                        for right in rights:
530                                new_score = metric(left, right)
531
532                                # Count a new_score of None in the denominator, but treat it as a value of 0
533                                # This mimics a dummy root node if there are multiple roots in the GO DiGraph
534                                if new_score != None:
535                                        total_score += new_score
536                                num_scores += 1
537
538                if total_score == 0:
539                        return None
540                else:
541                        return total_score / num_scores
542
543        def pairwise_max_term_comp(self, lefts, rights, metric):
```

```
544                    """
545                            Compares each pair of terms in two sets or lists of terms.
546                            Returns the maximum score found in these comparisons.
547                            Uses metric(left, right) to make each comparison.
548                            metric must take in two ontology terms (left and right) and return a numeric score.
549                    """
550
551                    if (len(lefts) == 0) or (len(rights) == 0):
552                            return None
553
554                    max_score = 0
555
556                    for left in lefts:
557                            for right in rights:
558                                    temp_score = metric(left, right)
559                                    if temp_score != None and temp_score > max_score:
560                                            max_score = temp_score
561
562                    return max_score
563
564        def average_protein_comp(self, left_prot, right_prot, metric):
565                    """
566                            Looks up all go terms for left_prot and right_prot.
567                            Uses pairwise_average_term_comp to compare the above sets of terms.
568                            metric must take in two ontology terms (left and right) and return a numeric score.
569                    """
570
571                    left_terms = self._prot_to_gos[left_prot]
572                    right_terms = self._prot_to_gos[right_prot]
573
574                    return self.pairwise_average_term_comp(left_terms, right_terms, metric)
575
576        def max_protein_comp(self, left_prot, right_prot, metric):
577                    """
578                            Looks up all terms for left_prot and right_prot.
579                            Uses pairwise_max_term_comp to compare the above sets of terms.
580                            metric must take in two ontology terms (left and right) and return a numeric score.
581                    """
582
583                    left_terms = []
584                    right_terms = []
585
586                    if (left_prot in self._prot_to_gos):
587                            left_terms = self._prot_to_gos[left_prot]
588
589                    if (right_prot in self._prot_to_gos):
590                            right_terms = self._prot_to_gos[right_prot]
591
592                    return self.pairwise_max_term_comp(left_terms, right_terms, metric)
593
594
595
596 ####################################
597 ###  End SemSim_Calculator class ###
598 ####################################
599
600
601
602
603
604
605 ########################
606 ### MicaStore class ###
607 ########################
608
609
610 class MicaStore():
611        """
612            Loads a matrix of MICA scores (and a list of GO term indices),
613            Provides accessors for MICA score lookup
614        """
615
616        def __init__(self, matrix_filename, ordering_filename):
```

```
617                     """
618                             Loads the .npy numpy array,  matrix_filename,
619                             Stores the indices for each GO term in ordering_filename
620                     """
621                     orderfile = open_or_abort(ordering_filename)
622
623                     self._micas = numpy.load(matrix_filename)
624                     self._go_to_index = {}
625
626                     index = 0
627                     for line in orderfile:
628                             self._go_to_index[line.strip()] = index
629                             index += 1
630
631                     orderfile.close()
632
633             def get_micas(self):
634                     """
635                             Returns reference to numpy matrix of MICA values.
636                             NOTE: This is a large matrix
637                     """
638
639                     return self._micas
640
641             def get_ordering(self):
642                     """
643                             Returns copy of the dictionary mapping
644                             GO terms to indices in the _micas matrix
645                     """
646
647                     return dict(self._go_to_index)
648
649             def get_index(self, term):
650                     """
651                             Returns the index of a GO term in the ordering of _micas (using _go_to_index)
652                             Returns None if term is not in _go_to_index
653                     """
654
655                     if (term in self._go_to_index):
656                             return self._go_to_index[term]
657                     else:
658                             return None
659
660             def mica_lookup(self, left, right):
661                     """
662                             If a MICA value can be found in _micas, return that MICA
663                             Else, return None
664                     """
665
666                     left_index = self.get_index(left)
667                     right_index = self.get_index(right)
668
669                     if (left_index != None) and (right_index != None):
670                             mica = self._micas[left_index, right_index]
671                     else:
672                             mica = None
673
674                     if (mica == ''):
675                             # Indicates that the mica was found, but does not exist (None is a valid MICA value)
676                             mica = ''
677
678                     return mica
679
680
681     ###########################
682     ### End MicaStore class ###
683     ###########################
```