

# The *GOSim* package

Holger Fröhlich

October 11, 2007

## 1 Introduction

The Gene Ontology (GO) has become one of the most widespread systems for systematically annotating gene products within the bioinformatics community and is developed by the Gene Ontology Consortium <sup>?</sup>. It is specifically intended for describing gene products with a controlled and structured vocabulary. GO terms are part of a Directed Acyclic Graph (DAG), covering three orthogonal taxonomies or "aspects": *molecular function*, *biological process* and *cellular component*. Two different kinds of relationship between GO terms exist: the "is-a" relationship and the "part-of" relationship. Providing a standard vocabulary across any biological resources, the GO enables researchers to use this information for automated data analysis.

The *GOSim* package provides the researcher with various information theoretic similarity concepts for GO terms <sup>???????</sup>. It additionally implements different methods for computing functional similarities between gene products based on the similarities between the associated GO terms. This can, for instances, be used for clustering genes according to their biological function <sup>??</sup> and thus may help to get a better understanding of the biological aspects covered by a set of genes.

## 2 Usage of *GOSim*

To elucidate the usage of *GOSim* we show an example workflow and explain the employed similarity concepts. We create a character vector of Entrez gene IDs:

```
> library(GOSim)
> genes = c("207", "208", "596", "901", "780", "3169", "9518",
+          "2852", "26353", "8614", "7494")
```

Next we investigate the GO annotation within the current ontology (which is *biological process* by default):

```
> getGOInfo(genes)
```

## 2.1 Term Similarities

Let us examine the similarity of the GO terms for genes "8614" and "2852" in greater detail:

```
> getTermSim(c("GO:0007166", "GO:0007267", "GO:0007584", "GO:0007165",
+             "GO:0007186"), method = "Resnik", verbose = FALSE)
```

	GO:0007166	GO:0007267	GO:0007584	GO:0007165	GO:0007186
GO:0007166	1.0000000	0.3032191	0.3032191	0.3125535	0.3607165
GO:0007267	0.3032191	1.0000000	0.3032191	0.3032191	0.3032191
GO:0007584	0.3032191	0.3032191	1.0000000	0.3032191	0.3032191
GO:0007165	0.3125535	0.3032191	0.3032191	1.0000000	0.3125535
GO:0007186	0.3607165	0.3032191	0.3032191	0.3125535	1.0000000

This calculates Resnik's pairwise similarity between GO terms ??:

$$sim(t, t') = IC_{ms}(t, t') := \max_{\hat{t} \in Pa(t, t')} IC(\hat{t}) \quad (1)$$

Here  $Pa(t, t')$  denotes the set of all common ancestors of GO terms  $t$  and  $t'$ , while  $IC(t)$  denotes the information content of term  $t$ . It is defined as (e.g. ?)

$$IC(\hat{t}) = -\log P(\hat{t}) \quad (2)$$

i.e. as the negative logarithm of the probability of observing  $t$ . The information content of each GO term is already precomputed for each ontology based on the empirical observation, how many times a specific GO term or any of its direct or indirect offsprings appear in the annotation of the GO with gene products. The association between gene products and GO identifiers is reported regularly by the NCBI.

```
> data("ICsBPall")
> IC[c("GO:0007166", "GO:0007267", "GO:0007584", "GO:0007165",
+      "GO:0007186")]
```

GO:0007166	GO:0007267	GO:0007584	GO:0007165	GO:0007186
5.238172	6.624280	9.867617	4.538770	5.607354

This loads the information contents of all GO terms within "biological process". Likewise, the data files ICsMFA11 and ICsCCA11 contain the information contents of all GO terms within "molecular function" and "cellular component". If only GO terms having evidence codes "IMP" (inferred from mutant phenotype), "IGI", (inferred from genetic interaction), "IDA" (inferred from direct assay), "IEP" (inferred from expression pattern) or "IPI" (inferred from physical interaction) are wanted, one can use the data files ICs-BPIMP\_IGI\_IDA\_IEP\_IPI, ICsMFIMP\_IGI\_IDA\_IEP\_IPI and ICsCCIMP\_IGI\_IDA\_IEP\_IPI,

respectively. The information contents for GO terms filtered with respect to different evidence codes must be calculated explicitly using the function `calcICs`. Please refer to the manual pages for details.

For the similarity computation in (Eq.: 1) normalized information contents are used, which are obtained by dividing the raw information contents by its maximal value:

```
> IC[c("GO:0007166", "GO:0007267", "GO:0007584", "GO:0007165",
+      "GO:0007186")]/max(IC[IC != Inf])
```

```
GO:0007166 GO:0007267 GO:0007584 GO:0007165 GO:0007186
0.3607165  0.4561681  0.6795141  0.3125535  0.3861395
```

To continue our example from above, let us also calculate Jiang and Conrath's pair-wise similarity between GO terms, which is the default, for comparison reasons ?:

```
> getTermSim(c("GO:0007166", "GO:0007267", "GO:0007584", "GO:0007165",
+             "GO:0007186"), verbose = FALSE)
```

```
GO:0007166 GO:0007267 GO:0007584 GO:0007165 GO:0007186
GO:0007166 1.0000000 0.7895537 0.5662077 0.9518371 0.9745770
GO:0007267 0.7895537 1.0000000 0.4707560 0.8377166 0.7641307
GO:0007584 0.5662077 0.4707560 1.0000000 0.6143706 0.5407847
GO:0007165 0.9518371 0.8377166 0.6143706 1.0000000 0.9264140
GO:0007186 0.9745770 0.7641307 0.5407847 0.9264140 1.0000000
```

Jiang and Conrath's similarity measure is defined as

$$sim(t, t') = 1 - \min(1, IC(t) - 2IC_{ms}(t, t') + IC(t')) \quad (3)$$

i.e. the similarity between  $t$  and  $t'$  is 0, if their normalized distance is at least 1.

Likewise, we can also compute Lin's pairwise similarity between GO terms ?:

```
> getTermSim(c("GO:0007166", "GO:0007267", "GO:0007584", "GO:0007165",
+             "GO:0007186"), method = "Lin", verbose = FALSE)
```

```
GO:0007166 GO:0007267 GO:0007584 GO:0007165 GO:0007186
GO:0007166 1.0000000 0.7423794 0.5829845 0.9284641 0.9659600
GO:0007267 0.7423794 1.0000000 0.5339859 0.7888919 0.7199725
GO:0007584 0.5829845 0.5339859 1.0000000 0.6112872 0.5690764
GO:0007165 0.9284641 0.7888919 0.6112872 1.0000000 0.8946806
GO:0007186 0.9659600 0.7199725 0.5690764 0.8946806 1.0000000
```

It is defined as:

$$sim(t, t') = \frac{2IC_{ms}(t, t')}{IC(t) + IC(t')} \quad (4)$$

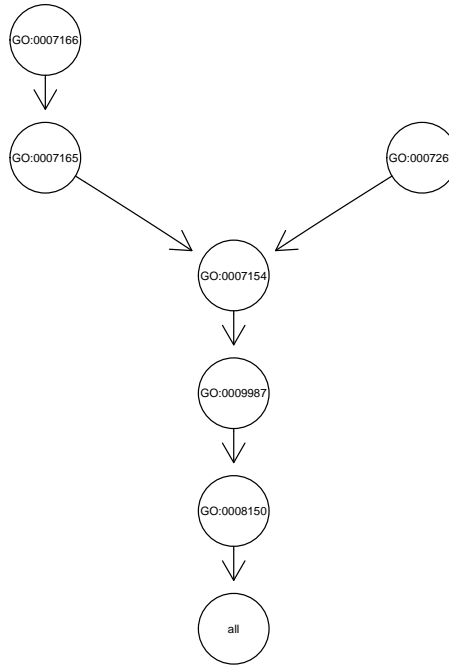


Figure 1: Example of a GO graph starting with leaves GO:0007166 and GO:0007267.

Resnik's, Jiang-Conraths's and Lin's term similarities all refer to  $IC_{ms}(t, t')$ , the information content of the minimum subsumer of  $t$  and  $t'$ , i.e. of the lowest common ancestor in the hierarchy. For illustration let us plot the GO graph with leaves GO:0007166 and GO:0007267 and let us compute their minimum subsumer (see Fig. 1):

```

> library(Rgraphviz)
> G = getGOGraph(c("GO:0007166", "GO:0007267"))
> plot(G)

> getMinimumSubsumer("GO:0007166", "GO:0007267")

[1] "GO:0007154"

```

In contrast to the above defined similarity measures Couto et al. [?] introduced a concept, which is not based on the minimum subsumer, but on the set of all disjunctive common ancestors. Roughly speaking, the idea is not to consider the common ancestor having the highest information content only, but also others, if they are somehow "separate" from each other, i.e. there is a path to  $t$  and  $t'$  not passing any other of the disjunctive common ancestors.

```

> getDisjCommAnc("GO:0007166", "GO:0007267")

```

[1] "GO:0007154"

In this case the set of disjunctive common ancestors only consists of the minimum subsumer, because any path from the other ancestors to GO:0007166 and GO:0007267 would have to pass the minimum subsumer (see Fig. 1).

Based on the notion of disjunctive common ancestors Resnik's similarity concept can be extended by defining:

$$sim(t, t') = IC_{share}(t, t') = \frac{1}{|DisjCommAnc|} \sum_{t \in DisjCommAnc} IC(t) \quad (5)$$

Likewise, Jiang-Conraths's and Lin's measures can be extended as well by replacing  $IC_{ms}(t, t')$  by  $IC_{share}(t, t')$ .

```
> getTermSim(c("GO:0007166", "GO:0007267", "GO:0007584", "GO:0007165",  
+ "GO:0007186"), method = "CoutoResnik", verbose = FALSE)
```

	GO:0007166	GO:0007267	GO:0007584	GO:0007165	GO:0007186
GO:0007166	1.0000000	0.3032191	0.232439	0.3032191	0.3125535
GO:0007267	0.3032191	1.0000000	0.232439	0.3032191	0.3032191
GO:0007584	0.2324390	0.2324390	1.0000000	0.2324390	0.2324390
GO:0007165	0.3032191	0.3032191	0.232439	1.0000000	0.3032191
GO:0007186	0.3125535	0.3032191	0.232439	0.3032191	1.0000000

Finally, it should be mentioned that also the depth and density enriched term similarity by Couto et al. [?] has been integrated into *GOSim*:

```
> setEnrichmentFactors(alpha = 0.5, beta = 0.3)  
> getTermSim(c("GO:0007166", "GO:0007267", "GO:0007584", "GO:0007165",  
+ "GO:0007186"), method = "CoutoEnriched", verbose = FALSE)
```

	GO:0007166	GO:0007267	GO:0007584	GO:0007165	GO:0007186
GO:0007166	1.0000000	0.1438407	0.1791987	0.1094908	0.1372791
GO:0007267	0.1438407	1.0000000	0.2089535	0.1287481	0.1514393
GO:0007584	0.1791987	0.2089535	1.0000000	0.1613989	0.1880462
GO:0007165	0.1094908	0.1287481	0.1613989	1.0000000	0.1154679
GO:0007186	0.1372791	0.1514393	0.1880462	0.1154679	1.0000000

## 2.2 Functional Gene Similarities

The special strength of *GOSim* lies in the possibility not only to calculate similarities for individual GO terms, but also for genes based on their complete GO annotation. For this purpose three basic ideas have been implemented:

1. Maximum and average pairwise GO term similarity

2. Computation of a so-called *optimal assignment* of terms from one gene to those of another one ?.
3. Embedding of each gene into a feature space defined by the gene's similarity to certain prototype genes ??. Within this feature space similarities naturally arise as dot products between the feature vectors. These dot products can be understood as so-called *kernel functions* ?, as used in e.g. Support Vector Machines ?.

## 2.3 Maximum and Average Pairwise GO Term Similarity

The idea of the maximum pairwise GO term similarity is straight forward. Given two genes  $g$  and  $g'$  annotated with GO terms  $t_1, \dots, t_n$  and  $t'_1, \dots, t'_m$  we define the functional similarity between  $g$  and  $g'$  as

$$sim_{gene}(g, g') = \max_{\substack{i = 1, \dots, n \\ j = 1, \dots, m}} sim(t_i, t'_j) \quad (6)$$

where  $sim$  is some similarity measure to compare GO terms  $t_i$  and  $t'_j$ . The resulting value is then further normalized to account for an unequal number of GO terms for both genes:

$$sim_{gene}(g, g') \leftarrow \frac{sim_{gene}(g, g')}{\sqrt{sim_{gene}(g, g) sim_{gene}(g', g')}} \quad (7)$$

Instead of computing the maximum pairwise GO term similarity one may also take the average here.

### 2.3.1 Optimal Assignment Gene Similarities

To elucidate the idea of the optimal assignment, consider the GO terms associated with gene "8614" on one hand and gene "2852" on the other hand:

```
> getGOInfo(c("8614", "2852"))
```

```
$`8614`
```

```
$`8614`$`GO:0007166`
```

```
GOID: GO:0007166
```

```
Term: cell surface receptor linked signal transduction
```

```
Ontology: BP
```

```
Definition: Any series of molecular signals initiated by the binding of  
an extracellular ligand to a receptor on the surface of the target  
cell.
```

```
$`8614`$`GO:0007267`
```

```
GOID: GO:0007267
```

Term: cell-cell signaling

Ontology: BP

Definition: Any process that mediates the transfer of information from one cell to another.

Synonym: cell-cell signalling

\$`8614`\$`G0:0007584`

G0ID: G0:0007584

Term: response to nutrient

Ontology: BP

Definition: A change in state or activity of a cell or an organism (in terms of movement, secretion, enzyme production, gene expression, etc.) as a result of a nutrient stimulus.

Synonym: response to nutrients

Synonym: nutritional response pathway

\$`2852`

\$`2852`\$`G0:0007165`

G0ID: G0:0007165

Term: signal transduction

Ontology: BP

Definition: The cascade of processes by which a signal interacts with a receptor, causing a change in the level or activity of a second messenger or other downstream target, and ultimately effecting a change in the functioning of the cell.

Synonym: signaling

Synonym: signalling

\$`2852`\$`G0:0007186`

G0ID: G0:0007186

Term: G-protein coupled receptor protein signaling pathway

Ontology: BP

Definition: The series of molecular signals generated as a consequence of a G-protein coupled receptor binding to its physiological ligand.

Synonym: G protein coupled receptor protein signaling pathway

Synonym: G protein coupled receptor protein signalling pathway

Synonym: G-protein coupled receptor protein signalling pathway

Synonym: G-protein-coupled receptor protein signaling pathway

Synonym: G-protein-coupled receptor protein signalling pathway

Synonym: GPCR protein signaling pathway

Synonym: GPCR protein signalling pathway

Given a similarity concept *sim* to compare individual GO terms, the idea is now to assign each term of the gene having fewer annotation to exactly one term of the other gene such that the overall similarity is maximized. More formally this can be stated as follows: Let  $\pi$  be some permutation of either an  $n$ -subset of natural numbers  $\{1, \dots, m\}$  or an  $m$ -subset of natural numbers  $\{1, \dots, n\}$  (this will be clear from context). Then we are looking for the quantity

$$sim_{gene}(g, g') = \begin{cases} \max_{\pi} \sum_{i=1}^n sim(t_i, t'_{\pi(i)}) & \text{if } m > n \\ \max_{\pi} \sum_{j=1}^m sim(t_{\pi(j)}, t'_j) & \text{otherwise} \end{cases} \quad (8)$$

The computation of (8) corresponds to the solution of the classical maximum weighted bipartite matching (optimal assignment) problem in graph theory and can be carried out in  $O(\max(n, m)^3)$  time ?. To prevent that larger lists of terms automatically achieve a higher similarity we should further normalize  $sim_{gene}$  according to (Eq. 7)

In our example, using Lin's GO term similarity measure the following assignments are found:

$$GO : 0007165 \rightarrow GO : 0007267 \quad (9)$$

$$GO : 0007186 \rightarrow GO : 0007166 \quad (10)$$

The resulting similarity matrix is:

```
> getGeneSim(c("8614", "2852"), similarity = "OA", similarityTerm = "Lin",
+ verbose = FALSE)
```

```
      8614      2852
8614 1.0000000 0.7164153
2852 0.7164153 1.0000000
```

Note the difference to a gene similarity that is just based on the maximum GO term similarity:

```
> getGeneSim(c("8614", "2852"), similarity = "max", similarityTerm = "Lin",
+ verbose = FALSE)
```

```
      8614      2852
8614 1.00000 0.96596
2852 0.96596 1.00000
```

### 2.3.2 Feature Space Embedding of Gene Products

To calculate the feature vectors for each gene we can either define certain prototype genes a priori or we use one of the heuristics implemented in the function `selectPrototypes`. The default behavior is to select the 250 best annotated genes, i.e. which have been annotated with GO terms most often:

```
> proto = selectPrototypes(verbose = FALSE)
```

We now calculate for each gene  $g$  feature vectors  $\phi(g)$  by using their similarity to all prototypes  $p_1, \dots, p_n$ :

$$\phi(g) = (\text{sim}'(g, p_1), \dots, \text{sim}'(g, p_n))^T \quad (11)$$

Here  $\text{sim}'$  by default is the maximum pairwise GO term similarity. Alternatively, one can use the optimal assignment similarity for  $\text{sim}'$  as well. Both similarity measures can by itself again be combined with arbitrary GO term similarity concepts. The default is the Jiang-Conrath term similarity.

Because the feature vectors are very high-dimensional we usually perform a principal component analysis (PCA) to project the data into a lower dimensional subspace:

```
> PHI = getGeneFeaturesPrototypes(genes, prototypes = proto, verbose = FALSE)
```

This uses the above define prototypes to calculate feature vectors and performs a PCA afterwards. The number of principal components is chosen such that at least 95% of the total variance in feature space can be explained (this is a relatively conservative criterion).

We can now plot our genes in the space spanned by the first 2 principal components to get an impression of the relative "position" of the genes to each other in the feature space (see Fig. 2). The feature vectors are normalized to Euclidian norm 1 by default:

```
> x = seq(min(PHI$features[, 1]), max(PHI$features[, 1]), length.out = 100)
> y = seq(min(PHI$features[, 2]), max(PHI$features[, 2]), length.out = 100)
> plot(x, y, xlab = "principal component 1", ylab = "principal component 2",
+      type = "n")
> text(PHI$features[, 1], PHI$features[, 2], labels = genes)
```

Finally, we can directly calculate the similarities of the genes to each other, this time using the Resnik's GO term similarity concept. These similarities may then be used to cluster genes with respect to their function:

```
> sim = getGeneSimPrototypes(genes, prototypes = proto, similarityTerm = "Resnik",
+   verbose = FALSE)
> h = hclust(as.dist(1 - sim$similarity), "ward")
> plot(h, xlab = "")
```

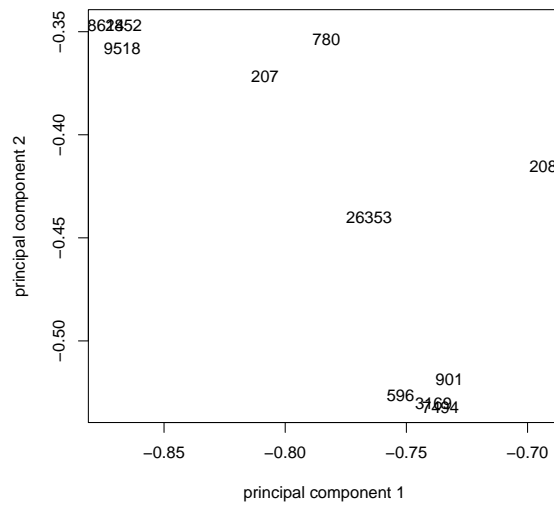


Figure 2: Embedding of the genes into the feature space spanned by the first 2 principal components

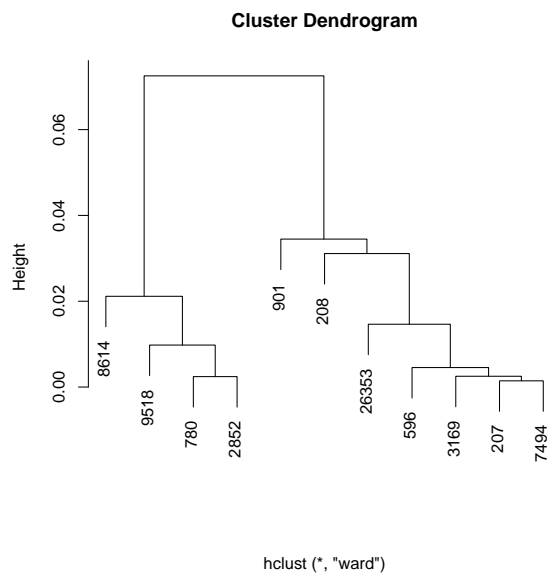


Figure 3: Possible functional clustering of the genes using Ward's method.

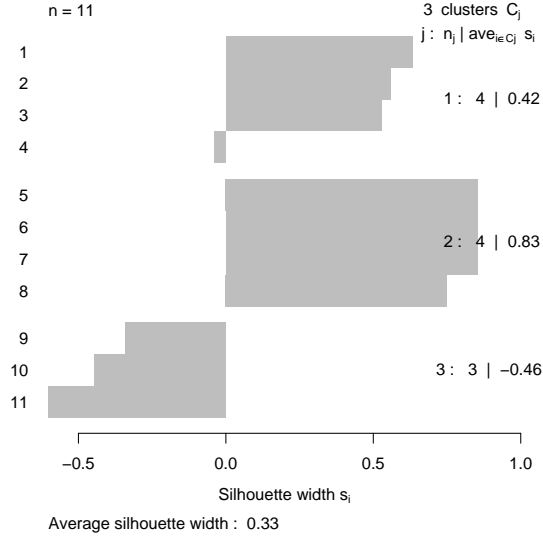


Figure 4: Silhouette plot of a possible given grouping of genes.

This produces a hierarchical clustering of all genes using Ward's method (see Fig. 3).

It should be mentioned that up to now all similarity computations were performed within the ontology "biological process". One could imagine to combine functional similarities between gene products with regard to different taxonomies. An obvious way for doing so would be to consider the sum of the respective similarities:

$$sim_{total}(g, g') = sim_{Ontology1}(g, g') + sim_{Ontology2}(g, g') \quad (12)$$

Of course, one could also use a weighted averaging scheme here, if desired.

## 2.4 Cluster Evaluations

*GOSim* has the possibility to evaluate a given clustering of genes or terms by means of their GO similarities. Supposed, based on other experiments (e.g. microarray), we have decided to put genes "8614", "9518", "780", "2852" in one group, genes "3169", "207", "7494", "596" in a second and the rest in a third group. Then we can ask ourselves, how similar these groups are with respect to their GO annotations:

```
> ev = evaluateClustering(c(2, 3, 2, 3, 1, 2, 1, 1, 3, 1, 2), sim$similarity)
> plot(ev$clustersil, main = "")
```

A good indication of the clustering quality can be obtained by looking at the cluster silhouettes ? (see Fig. 4). This shows that clusters 1 and 2 are relatively homogenous

with respect to the functional similarity of the genes contained in it, while the genes in cluster 3 are more dissimilar.