

# GO-terms Semantic Similarity Measures

Guangchuang Yu

April 21, 2009

## 1 Introduction

Functional similarities of gene products can be estimated by controlled biological vocabularies, such as Gene Ontology (GO). GO comprises of three orthogonal ontologies, molecular function (MF), biological process (BP), and cellular component (CC).

Four methods proposed by Resnik[Philip, 1999], Jiang[Jiang and Conrath, 1997], Lin[Lin, 1998] and Schlicker[Schlicker et al., 2006] respectively have presented to determine the semantic similarities of two GO terms based on the annotation statistics of their common ancestor terms. Wang [Wang et al., 2007] proposed a new method to measure the similarities based on the graph structure of GO. Each of these methods has its own strengths and weaknesses. The *GOSemSim* package implemented all these five methods.

## 2 Semantic Similarity Measures

The *GOSemSim* package contains functions to estimate similarity scores of GO terms. Details about Wang's method can be seen in [Wang et al., 2007], details about Rel method can be seen in [Schlicker et al., 2006] and the details about Resnik, Lin, and Jiang's methods can be seen in [Lord et al., 2003]. Resnik, Lin, Rel, and Jiang's methods based on the information content of the GO terms while Wang's method based on both the locations of these terms in the GO graph and their relations with their ancestor terms.

The method proposed by [Wang et al., 2007] is based on the graph structure of each term.

Formally, a GO term A can be represented as  $DAG_A = (A, T_A, E_A)$  where  $T_A$  is the set of GO terms in  $DAG_A$ , including term A and all of its ancestor terms in the GO graph, and  $E_A$  is the set of edges connecting the

GO terms in  $DAG_A$ .

To encode the semantics of a GO term in a measurable format to enable a quantitative comparison of two term's semantics. Firstly define the semantic value of term A as the aggregate contribution of all terms in  $DAG_A$  to the semantics of term A. Terms closer to term A in  $DAG_A$  contribute more to its semantics. Thus, define the contribution of a GO term  $t$  to the semantics of GO term A as the S-value of GO term  $t$  related to term A. For any of term  $t$  in  $DAG_A = (A, T_A, E_A)$ , its S-value related to term A.  $S_A(t)$  is defined as:

$$\begin{cases} S_A(A) = 1 \\ S_A(t) = \max\{w_e \times S_A(t') | t' \in \text{childrenof}(t)\} \text{ if } t \neq A \end{cases}$$

where  $w_e$  is the semantic contribution factor for edge  $e \in E_A$  linking term  $t$  with its child term  $t'$ . Term A contribute to its own is defined as one. After obtaining the S-values for all terms in  $DAG_A$ , the semantic value of GO term A,  $SV(A)$ , is calculated as:

$$SV(A) = \sum_{t \in T_A} S_A(t)$$

Given two GO terms A and B, the semantic similarity between these two terms,  $GO_{A,B}$ , is defined as:

$$S_{GO}(A, B) = \sum_{t \in T_A \cap T_B} \frac{S_A(t) + S_B(t)}{SV(A) + SV(B)}$$

where  $S_A(t)$  is the S-value of GO term  $t$  related to term A and  $S_B(t)$  is the S-value of GO term  $t$  related to term B.

Details about this method can be seen in [Wang et al., 2007]. This method determines the semantic similarity of two GO terms based on both the locations of these terms in the GO graph and their relations with their ancestor terms.

The GOSemSim package implemented four other methods which are based on information content were proposed by Resnik[Philip, 1999], Jiang[Jiang and Conrath, 1997], Lin[Lin, 1998] and Schlicker[Schlicker et al., 2006] respectively.

Information content is defined as frequency of each term occurs in the GO corpus. We used Bioconductor package `org.Hs.eg.db`, `org.Dm.eg.db`,

org.Mm.eg.db, org.Rn.eg.db, org.Sc.sgd.db to calculate the information content of human, fly, mouse, rat and yeast species respectively. The information content will update regularly.

Given the information content, we applied the four measures to estimate the semantic similarity between terms.

As GO allows multiple parents for each concept, two terms can share parents by multiple paths. We take the minimum  $p(t)$ , where there is more than one shared parents. The  $p_{ms}$  is defined as :

$$p_{ms}(t1, t2) = \min_{t \in S(t1, t2)} \{p(t)\}$$

where  $S(t1, t2)$  is the set of parent terms shared by  $t1$  and  $t2$ .

The first method Resnik[Philip, 1999] is defined as:

$$sim(t1, t2) = -\ln p_{ms}(t1, t2)$$

The second method Lin[Lin, 1998] is defined as:

$$sim(t1, t2) = \frac{2 \times \ln(p_{ms}(t1, t2))}{\ln p(t1) + \ln p(t2)}$$

The third method Rel[Schlicker et al., 2006] combine Resnik's and Lin's method is defined as:

$$sim = \frac{2 \times \ln p_{ms}(t1, t2)}{\ln p(t1) + \ln p(t2)}$$

The last method Jiang[Jiang and Conrath, 1997] define a semantic distance as:

$$d(t1, t2) = \ln p(t1) + \ln p(t2) - 2 \times \ln p_{ms}(t1, t2)$$

and the corresponding similarity measure for  $d(t1, t2)$  is given by:

$$sim(t1, t2) = 1 - \min(1, d(t1, t2))$$

The semantic similarity of one GO term  $go$  and a GO terms set  $GO = \{go_1, go_2 \dots go_k\}$  is defined as:

$$Sim(go, GO) = \max_{1 \leq i \leq k} (S_{GO}(go, GO_i))$$

Therefore, given two GO terms sets  $GO_1 = \{go_{11}, go_{12} \cdots go_{1m}\}$  and  $GO_2 = \{go_{21}, go_{22} \cdots go_{2n}\}$ , the semantic similarity between them is defined as:

$$Sim(GO_1, GO_2) = \frac{\sum_{1 \leq i \leq m} Sim((go_{1i}), (GO_2)) + \sum_{1 \leq j \leq n} Sim((go_{2j}), (GO_1))}{m+n}$$

```
> library(GOSemSim)
> goSim("GO:0004022", "GO:0005515", ont = "MF", measure = "Wang")

[1] 0.252
```

The function goSim generates the semantic similarity score for a pair of GO terms.

```
> go1 = c("GO:0004022", "GO:0004024", "GO:0004174")
> go2 = c("GO:0009055", "GO:0005515")
> mgoSim(go1, go2, ont = "MF", measure = "Wang")

[1] 0.299
```

The function mgoSim generates the similarity score of two GO terms lists.

```
> geneSim("241", "2561", ont = "MF", organism = "human", measure = "Wang")

$geneSim
[1] 0.29

$G01
[1] "GO:0005488" "GO:0008047"

$G02
[1] "GO:0004890"
```

The function geneSim estimate two genes's semantic similarity. The mapping from Gene IDs to GO IDs can be restricted based on evidence codes. It supports five species, which are "human", "rat", "mouse", "fly", and "yeast".

### 3 Functional Clustering

Given GO based similarity scores, gene products may be clustered by their function. *GOsemSim* package provides a function, `mgeneSim`, that returns pairwise similarities scores for a list of genes. It can be used by other functions to perform clustering.

```
> sim <- mgeneSim(c("835", "5261", "241", "934"), ont = "MF", organism = "human",
+   measure = "Wang")
> sim

      835  5261  241
835  1.000  0.149  0.531
5261  0.149  1.000  0.238
241  0.531  0.238  1.000

> library(cluster)
> pamCluster <- pam(as.dist(1 - sim[complete.cases(sim), complete.cases(sim)]),
+   2)
> pamCluster$clustering

 835 5261  241
  1   2   1
```

We also implemented two functions for estimating similarities among gene clusters. *clusterSim* for calculating semantic similarity between two gene clusters and *mclusterSim* for calculating pairwise similarities of a set of gene clusters. For calculate two gene clusters similarities, we first calculate pairwise similarities among genes, and the average similarity between all gene products was taken since all genes contribute to the gene cluster.

```
> cluster1 <- c("snR67", "snR40", "snR48", "snR17a", "snR8")
> cluster2 <- c("YOR251C", "YPR137C-B", "YPR010C", "YPR072W")
> cluster3 <- c("YNL133C", "YOL041C", "YOL018C", "YOR236W", "YOR179C",
+   "YOR230W")
> clusterSim(cluster1, cluster2, ont = "MF", organism = "yeast",
+   measure = "Wang")

[1] 0.215

> clusters <- list(a = cluster1, b = cluster2, c = cluster3)
> mclusterSim(clusters, ont = "MF", organism = "yeast", measure = "Wang")
```

