

GO-terms Semantic Similarity Measures

Guangchuang Yu

May 11, 2010

1 Introduction

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

Four methods have been presented to determine the semantic similarity of two GO terms based on the annotation statistics of their common ancestor terms (Resnik[Philip, 1999], Jiang[Jiang and Conrath, 1997], Lin[Lin, 1998] and Schlicker[Schlicker et al., 2006]). Wang [Wang et al., 2007] proposed a new method to measure the similarity based on the graph structure of GO. Each of these methods has its own advantages and weaknesses. The **GOSemSim** package [Yu et al., 2010] is developed to compute semantic similarity among GO terms, sets of GO terms, gene products, and gene clusters, providing both five methods mentioned above.

2 Semantic Similarity Measurement Based on GO

The **GOSemSim** package contains functions to estimate semantic similarity of GO terms based on Resnik's, Lin's, Jiang and Conrath's, Rel's and Wang's method. Details about Resnik's, Lin's, and Jiang and Conrath's methods can be seen in [Lord et al., 2003], details about Rel's method can be seen in [Schlicker et al., 2006], details about Wang's method can be seen in [Wang et al., 2007].

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 between two term's semantics, we 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' . We defined term A contributes to its own 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.

The method described above is proposed by Wang[Wang et al., 2007]. The Wang's 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 other four methods proposed by Resnik[Philip, 1999], Jiang[Jiang and Conrath, 1997], Lin[Lin, 1998] and Schlicker[Schlicker et al., 2006] are information content (IC) based, which depend on the frequencies of two GO terms involved and that of their closest common ancestor term in a specific corpus of GO annotations. Information content is defined as frequency of each term occurs in the corpus. At present, **GOSemSim** supports analysis on many species. We used Bioconductor package `org.At.tair.db`, `org.Ag.eg.db`, `org.Bt.eg.db`, `org.Cf.eg.db`, `org.Gg.eg.db`, `org.Pt.eg.db`, `org.Sco.eg.db`, `org.EcK12.eg.db`, `org.EcSakai.eg.db`, `org.Dm.eg.db`, `org.Hs.eg.db`, `org.Pf.plasmo.db`,

org.Mm.eg.db, org.Ss.eg.db, org.Rn.eg.db, org.Mmu.eg.db, org.Ce.eg.db, org.Xl.eg.db, org.Sc.sgd.db and org.Dr.eg.db to calculate the information content of Arabidopsis, Anopheles, Bovine, Canine, Chicken, Chimp, Coelicolor, E coli strain K12 and strain Sakai, Fly, Human, Malaria, Mouse, Pig, Rat, Rhesus, Worm, Xenopus, Yeast and Zebrafish respectively. The information content will update regularly.

As GO allow 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 similarity of Resnik's method is defined as:

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

The similarity of Lin's is defined as:

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

The similarity of Schlicker's method combine Resnik's and Lin's method is defined as:

$$sim(t1, t2) = \frac{2 \times \ln p_{ms}(t1, t2)}{\ln p(t1) + \ln p(t2)} \times (1 - p_{ms}(t1, t2))$$

The Jiang and Conrath's method define a semantic similarity as:

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

where

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

In **GOSemSim**, on the basis of semantic similarity between GO terms, we can also compute semantic similarity among sets of GO terms, gene products, and gene clusters. We implemented four methods which called *max*, *average*, *rcmax*, and *rcmax.avg* to combine semantic similarity scores of multiple GO terms. The similarities among gene products and gene clusters which annotated by multiple GO terms were also calculated by the same combine methods mentioned above.

Given two GO terms sets $GO_1 = \{go_{11}, go_{12} \cdots go_{1m}\}$ and $GO_2 = \{go_{21}, go_{22} \cdots go_{2n}\}$, method *max* calculate the maximum semantic similarity score over all pairs

of GO terms between these two sets, method *average* calculate the average semantic similarity score over all pairs of GO terms.

Similarities between GO terms form a matrix, and method *rcmax* use the maximum of RowScore and ColumnScore as the similarity, where RowScore (or ColumnScore) is the average of maximum similarities on each row (or column).

And method *rcmax.avg* calculate the average of all maximum similarities on each row and column, and defined as:

$$Sim(GO1, GO2) = \frac{\sum_{1 \leq i \leq m} \max(Sim((go_{1i}), (GO_2))) + \sum_{1 \leq j \leq n} \max(Sim((go_{2j}), (GO_1)))}{m+n}$$

3 Functions

Six functions are provided by **GOSemSim** package. The function *goSim*, *mgoSim*, *geneSim*, *clusterSim* can be used to compute the semantic similarity among GO terms, sets of GO terms, GO descriptions of gene products and GO descriptions of gene clusters respectively. The function *mgeneSim* and *mclusterSim* are designed to calculate the similarity scores matrix of a set of genes and gene clusters.

For all these six functions, the ontology of GO used in measurement can be restricted by assigning the corresponding parameter to "BP" (biological process), "MF" (molecular function) and "CC" (cellular component).

Users must set parameter *measure* to specify which method to be used to measure the similarity and set parameter *combine* to specify which method for combining semantic similarity scores of multiple GO terms associated with protein or protein clusters.

The function *goSim* gives the semantic similarity score for a pair of GO terms. The output value of *goSim* is between 0 and 1. The higher the value obtained the more similar between them. For example:

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

[1] 0.152
```

The function *mgoSim* generates the similarity score of two GO terms lists. For example:

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

[1] 0.225
```

The function *geneSim* estimate semantic similarity of two genes. The mapping from Gene IDs to GO IDs can be restricted based on evidence codes. Gene IDs and species are needed for the function. For yeast, the Gene ID is refer to ORF identifiers from Saccharomyces Genome Database (SGD), and for arabidopsis, the Gene ID is refer to TAIR gene identifiers from The Arabidopsis Information Resource (TAIR) and for other supported species, the Gene IDs are all refer to Entrez Gene ID. For example:

```
> geneSim("241", "251", ont = "MF", organism = "human", measure = "Wang",
+         combine = "rcmax.avg")
```

```
$geneSim
[1] 0.088
```

```
$G01
[1] "GO:0047485" "GO:0050544"
```

```
$G02
[1] "GO:0004035"
```

For reducing the dependences of **GOSemSim**, it only depends on species specific annotation package org.Hs.eg.db for human. When calculated semantic similarities of other species, **GOSemSim** will auto load the corresponding annotation package if it was installed or download and install the package automatically if it was not installed.

Given GO based similarity scores, gene products may be clustered by their function. **GOSemSim** package provides *mgeneSim* function to compute pairwise similarity scores for a list of genes. The *mgeneSim* function will automatically remove the genes without annotations. Gene IDs and species are needed for the function. The *mgeneSim* function can be used for large-scale analysis, such as gene clustering. For example:

```
> sim <- mgeneSim(as.character(100:200), ont = "MF", organism = "human",
+               measure = "Wang", combine = "rcmax.avg")
> sim[1:5, 1:5]
```

```

      100    101    102    103    104
100 1.000 1.000 0.438 0.848 0.351
101 1.000 1.000 0.649 0.149 0.351
102 0.438 0.649 1.000 0.136 0.178
103 0.848 0.149 0.136 1.000 0.076
104 0.351 0.351 0.178 0.076 1.000

> library(cluster)
> pamCluster <- pam(as.dist(1 - sim[complete.cases(sim), complete.cases(sim)]),
+   5)
> pamCluster$clustering

100 101 102 103 104 105 107 108 116 118 119 120 124 125 126 127 128 130 131 141
   1   1   1   2   3   3   1   2   1   4   4   2   1   2   2   2   2   2   1   2
142 143 148 150 153 154 155 156 158 159 162 163 164 165 166 175 176 177 178 181
   1   1   5   5   1   5   1   1   2   1   1   1   2   3   1   2   1   1   2   1
182 183 185 186 189 190 191 196 197 199
   1   1   1   1   1   3   1   1   1   4

```

Two functions have also been implemented to estimate similarities among gene clusters. The *clusterSim* function is developed for calculating semantic similarity between two gene clusters. The *clusterSim* function first calculate all pair-wise similarities between gene products from different clusters, and then combine the similarity matrix by one of the combine method *max*, *average*, *rcmax* and *rcmax.avg* metioned above.

The *mclusterSim* function is developed for calculating pair-wise similarities of a set of gene clusters, and its kernel algorithm is just the same as the *clusterSim* function.

Gene IDs, species are needed for these functions. Here we use the clustering result of the previous example to show the calculation of clusters' similarities:

```

> g <- pamCluster$clustering
> cluster1 <- names(g[g == 1])
> cluster2 <- names(g[g == 2])
> cluster3 <- names(g[g == 3])
> clusters <- list(a = cluster1, b = cluster2, c = cluster3)
> clusterSim(cluster1, cluster2, ont = "MF", organism = "human",
+   measure = "Wang", combine = "rcmax.avg")

[1] 0.496

```


James Z Wang, Zhidian Du, Rapeeporn Payattakool, Philip S Yu, and Chin-Fu Chen. A new method to measure the semantic similarity of go terms. *Bioinformatics (Oxford, England)*, 23: 1274–81, May 2007. ISSN 1460-2059. doi: btm087. URL <http://www.ncbi.nlm.nih.gov/pubmed/17344234>. PMID: 17344234.

Guangchuang Yu, Fei Li, Yide Qin, Xiaochen Bo, Yibo Wu, and Shengqi Wang. Gosemsim: an r package for measuring semantic similarity among go terms and gene products. *Bioinformatics*, 26:976–978, 2010. ISSN 1367-4803. doi: 10.1093/bioinformatics/btq064. URL <http://bioinformatics.oxfordjournals.org/cgi/content/abstract/26/7/976>. PMID: 20179076.