

Bioconductor项目简介

及其在生物信息学中的应用

陈钢 chen.gang1983@gmail.com

中南大学计算机系

2008年12月13日

- 1 提纲
- 2 Bioconductor项目简介
 - Bioconductor中软件包的主要分类
 - Bioconductor中的主要软件包
 - Bioconductor的安装
- 3 蛋白质相互作用网络和GO语义相似性
 - 蛋白质相互作用网络
 - Gene Ontology
 - 蛋白质的GO语义相似性
- 4 实验结果
 - 试验数据
 - 试验规模
 - 代码
 - 试验结果
- 5 与Java实现的简单比较

- 网址: <http://www.bioconductor.org>





- 网址: <http://www.bioconductor.org>
- **Bioconductor** is an open source and open development software project for the analysis and comprehension of genomic data.



- 网址: <http://www.bioconductor.org>
- **Bioconductor** is an open source and open development software project for the analysis and comprehension of genomic data.
- Bioconductor 2.3: 294个软件包, 针对R2.8, 新加入36个软件包, 加入对新一代测序技术的支持。

- 软件（Software）：

- 软件（Software）：
 - LIMMA：基因芯片数据线性建模

- 软件（Software）：
 - LIMMA：基因芯片数据线性建模
 - affyio：Affymetrix公司芯片数据的处理

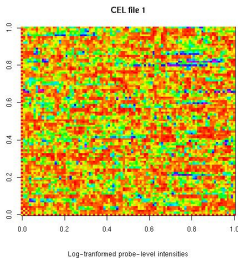
- 软件（Software）：
 - LIMMA：基因芯片数据线性建模
 - affyio：Affymetrix公司芯片数据的处理
- 注释数据（AnnotationData）：

- 软件（Software）：
 - LIMMA：基因芯片数据线性建模
 - affyio：Affymetrix公司芯片数据的处理
- 注释数据（AnnotationData）：
 - GO.db：Gene Ontology注释数据

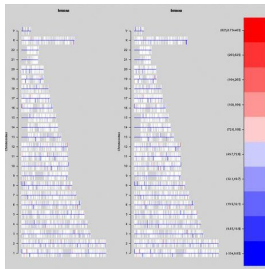
- 软件（Software）：
 - LIMMA：基因芯片数据线性建模
 - affyio：Affymetrix公司芯片数据的处理
- 注释数据（AnnotationData）：
 - GO.db：Gene Ontology注释数据
 - hgu133a2.db：Affymetrix Human Genome U133A 2.0芯片的注释数据

- 软件（Software）：
 - LIMMA：基因芯片数据线性建模
 - affyio：Affymetrix公司芯片数据的处理
- 注释数据（AnnotationData）：
 - GO.db：Gene Ontology注释数据
 - hgu133a2.db：Affymetrix Human Genome U133A 2.0芯片的注释数据
- 实验数据（ExperimentData）：

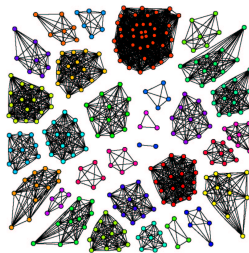
- 软件（Software）：
 - LIMMA：基因芯片数据线性建模
 - affyio：Affymetrix公司芯片数据的处理
- 注释数据（AnnotationData）：
 - GO.db：Gene Ontology注释数据
 - hgu133a2.db：Affymetrix Human Genome U133A 2.0芯片的注释数据
- 实验数据（ExperimentData）：
 - hapmapsnp5：人类单体型计划SNP数据



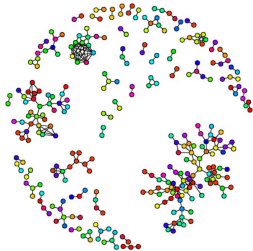
affy: Image of Affymetrix probe-level intensities, *image()* function.



geneplotter: Comparing mean expression levels between two groups.

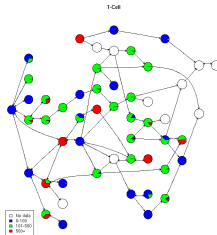


graph and **Rgraphviz**: Completely connected subgraphs for clusters of genes with similar expression profiles. Created with **graph** and plotted using **Rgraphviz**.



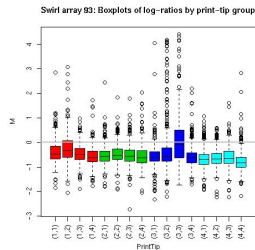
graph and Rgraphviz:

Graph of interacting protein pairs with proteins (nodes) colored according to cell cycle expression profile cluster membership. Created with **graph** and plotted using **Rgraphviz**.



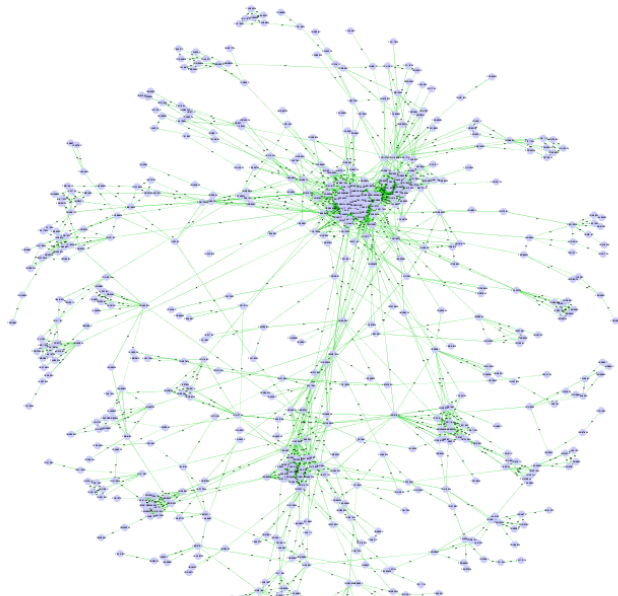
graph and Rgraphviz:

Graph representation of the integrin mediated cell adhesion pathway with nodes containing pie charts displaying expression levels for the corresponding gene over a set of samples.



marrayPlots: boxplots of log-ratios $\log_2 R/G$ by print-tip group, *maBoxplot()* function.

```
source("http://bioconductor.org/biocLite.R")  
biocLite()
```

酵母的部分蛋白质相互作用

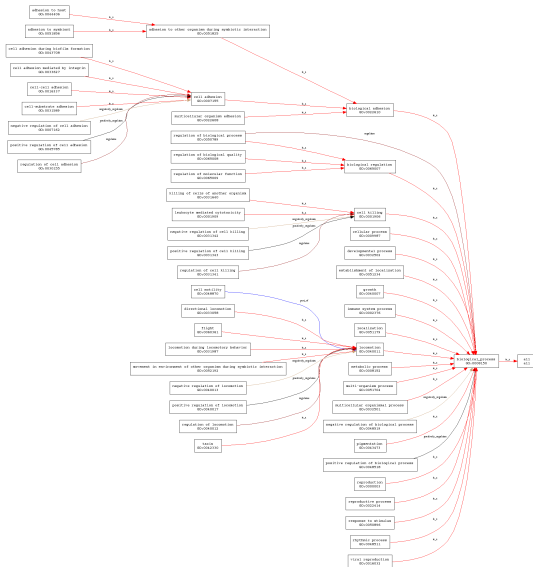
- 点表示蛋白质
- 边表示蛋白质之间的相互作用

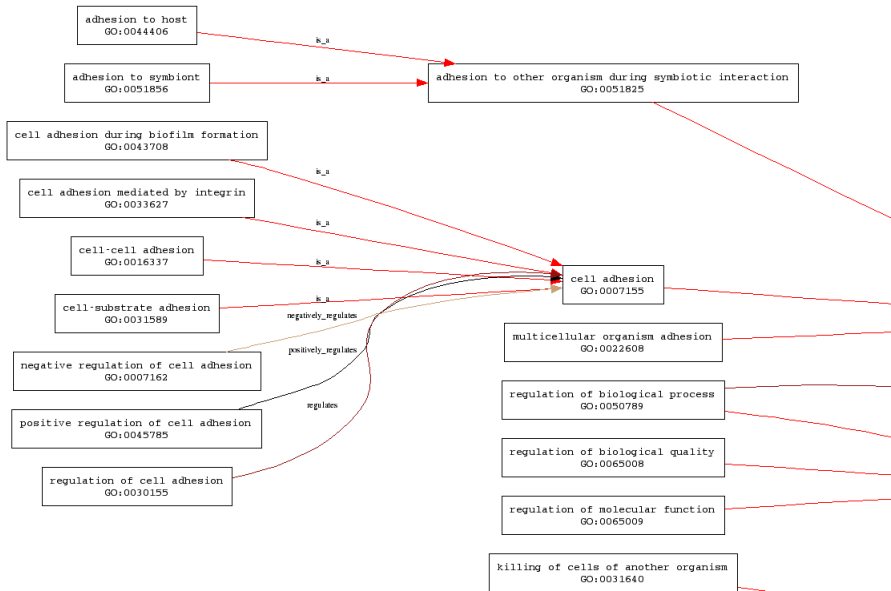
提纲

Bioconductor项目简介
蛋白质相互作用网络和GO语义相似性
实验结果
与Java实现的简单比较

Gene Ontology

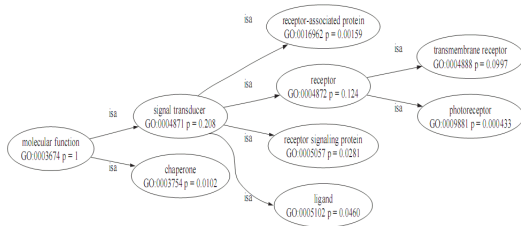
蛋白质的GO语义相似性





GO Term之间相似性的定义:

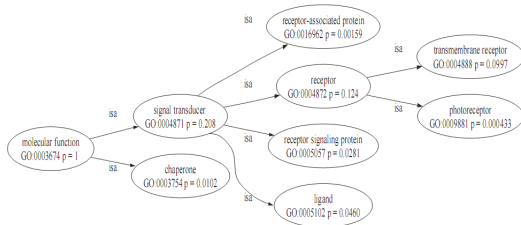
- P. W. Lord



$$Sim_{Lord}(term1, term2) = \min_{term \in (term1, term2)} \{p(term)\}$$

GO Term之间相似性的定义:

- P. W. Lord



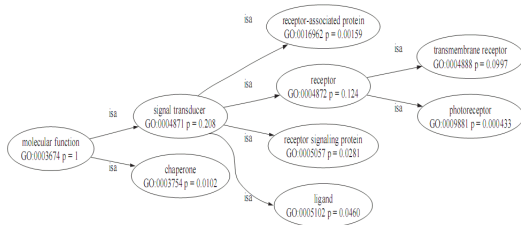
$$Sim_{Lord}(term1, term2) = \min_{term \in (term1, term2)} \{p(term)\}$$

- Resnik

$$Sim_{Resnik}(term1, term2) = -\ln Sim_{Lord}(term1, term2)$$

GO Term之间相似性的定义:

- P. W. Lord



$$Sim_{Lord}(term1, term2) = \min_{term \in (term1, term2)} \{p(term)\}$$

- Resnik

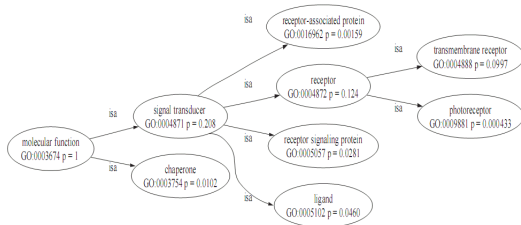
$$Sim_{Resnik}(term1, term2) = -\ln Sim_{Lord}(term1, term2)$$

- Jiang

$$Sim_{Jiang}(term1, term2) = 1 - \min(1, IC(term1) - 2IC_{ms} + IC(term2))$$

GO Term之间相似性的定义:

- P. W. Lord



$$Sim_{Lord}(term1, term2) = \min_{term \in (term1, term2)} \{p(term)\}$$

- Resnik

$$Sim_{Resnik}(term1, term2) = -\ln Sim_{Lord}(term1, term2)$$

- Jiang

$$Sim_{Jiang}(term1, term2) = 1 - \min(1, IC(term1) - 2IC_{ms} + IC(term2))$$

蛋白质之间GO语义相似性的定义：

- 平均值

蛋白质之间GO语义相似性的定义：

- 平均值
- 最大值

蛋白质之间GO语义相似性的定义：

- 平均值
- 最大值
- 匹配得最好的GO术语的相似性的平均值

蛋白质之间GO语义相似性的定义：

- 平均值
- 最大值
- 匹配得最好的GO术语的相似性的平均值
- 匹配得最好的GO术语的相似性的最大值

- 数据来源：人类蛋白质参考数据库（Human Protein Reference Database）

- 数据来源：人类蛋白质参考数据库（Human Protein Reference Database）
- 蛋白质数量：25,661

- 数据来源：人类蛋白质参考数据库（Human Protein Reference Database）
- 蛋白质数量：25,661
- 蛋白质相互作用数量：37,107

- GO term的分类：3种

- GO term的分类：3种
- GO Term之间相似性定义：7种

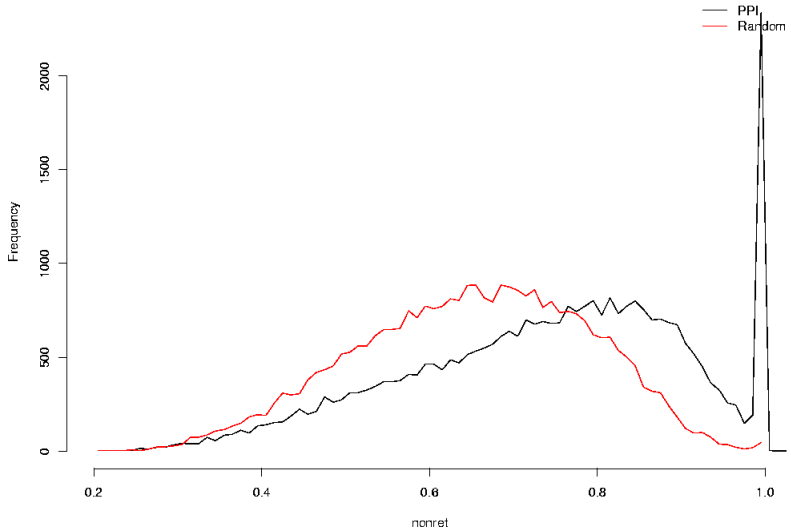
- GO term的分类：3种
- GO Term之间相似性定义：7种
- Gene Product之间的GO相似性定义：4种

- GO term的分类：3种
- GO Term之间相似性定义：7种
- Gene Product之间的GO相似性定义：4种
- 试验量： $3 \times 7 \times 4 = 84$

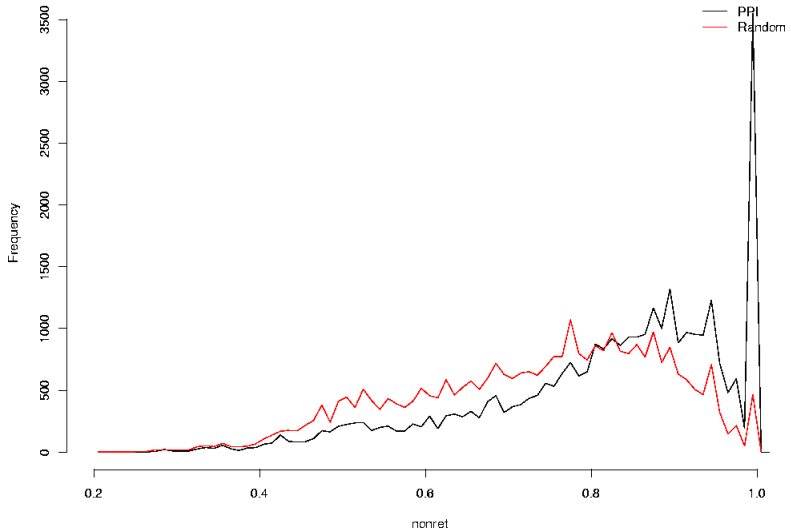
```
filterGenePairs <- function(genePairList) {  
  newGeneList <- matrix(ncol=2,nrow=0)  
  
  num <- 1  
  for(i in 1:(length(genePairList)/2)){  
    if(length(filterGO(genePairList[i,1])) != 0  
    && length(filterGO(genePairList[i,2])) != 0 )  
      newGeneList <- rbind(newGeneList, genePairList  
    )  
  }  
  return(newGeneList)  
}
```

```
getGenePairSim <- function(geneA, geneB, similarity = "funS  
similarityTerm = "Lin", verbose = FALSE) {  
  if(geneA & FALSE || geneB & FALSE) return(0)  
  if(geneA == geneB) return (1)  
  
  temp <- getGeneSim(c(geneA, geneB),  
                      similarity = similarity,  
                      similarityTerm = similarityTerm)  
  
  ret <- temp[2,1]  
  return(ret)  
}
```

Histogram of nonret



Histogram of nonret



	基于Bioconductor的实现	基于Java的实现
代码量	120行	约2500行
运行时间*	20小时	48小时内未完成计算
开发时间**	4小时	约1个月

如何参与？

- 使用、宣传

如何参与？

- 使用、宣传
- 报告Bug

如何参与？

- 使用、宣传
- 报告Bug
- 参软件包的开发，贡献自己的代码

如何参与？

- 使用、宣传
- 报告Bug
- 参软件包的开发，贡献自己的代码
- 贡献自己的软件包

如何参与？

- 使用、宣传
- 报告Bug
- 参软件包的开发，贡献自己的代码
- 贡献自己的软件包
- 进入core team？

所面临的问题:

- R代码的并行？ 分布式？

所面临的问题:

- R代码的并行？ 分布式？
- R代码的管理？ 团队协作？

所面临的问题:

- R代码的并行？ 分布式？
- R代码的管理？ 团队协作？
- R的组内培训？

谢谢