

Relatório Final de Atividades

Integração de dados

Inferência de GRNs a partir da integração dos dados da *Arabidopsis thaliana*

Euler Angelo de Menezes Junior

Bolsista CNPq

Tecnologia em Análise de Desenvolvimento de Sistemas

Data de ingresso no programa: 08/2014

Prof. Dr. Fabrício Martins Lopes

Bioinformática

CAMPUS CORNÉLIO PROCÓPIO, 2015

Integração de dados

Euler Angelo de Menezes Junior, Fabrício Martins Lopes

Universidade Tecnológica Federal do Paraná – UTFPR
Campus Cornélio Procópio
Avenida Alberto Carazzai, 1640 – Centro, Cornélio Procópio – PR, 86300-000

euler@alunos.utfpr.edu.br, fabricio@utfpr.edu.br

Palavras-chave: Integração de dados, *Arabidopsis thaliana*, Bioinformática

INTRODUÇÃO

Na era da informação, a quantidade de dados produzida é tão grande que supera a quantidade que conseguimos processar e analisar. Gerenciar e processar dados que crescem exponencialmente em tamanho em um curto período de tempo [1], é um dos maiores desafios da informática.

Isso não é diferente no contexto de dados biológicos, onde a tecnologia de aquisição das expressões gênicas vêm evoluindo rapidamente, sendo possível obter a expressão gênica de milhares de genes simultaneamente e em múltiplos instantes de tempo.

Com a grande quantidade de dados gerados, eles foram espalhados por vários bancos de dados públicos, como TAIR, KEGG e NCBI, porém muitos em formatos diferentes, dificultando a aplicação algoritmos computacionais.

Sendo assim, viu-se a necessidade de integrar todos esses dados em um único banco de dados, categorizar de uma forma mais simplificada e indexar as informações com identificadores únicos e numéricos para uma possível aplicação de fórmulas matemáticas e estatísticas.

Na bioinformática, a inferência de redes de regulação gênica (GRNs) a partir de seus perfis de expressão é um problema a ser solucionado, sendo necessário o desenvolvimento de métodos alternativos para montar as redes de forma mais adequada e precisa.

Tomamos como estudo de caso a planta *Arabidopsis thaliana*, onde é feita a integração dos dados biológicos obtidos dos bancos de dados públicos, e disponibilização dos mesmos para visualização em uma interface WEB, tão como uma API onde a partir dela é possível montar e inferir as GRNs.

DESENVOLVIMENTO

Esse trabalho consiste em integrar os dados em um único banco, eliminando informações duplicadas e indexar de uma forma que seja possível descobrir os inter-relacionamentos entre os genes.

O projeto foi dividido entre as seguintes etapas:

1. Identificação das informações a serem extraídas
2. Extração dos dados
3. Montagem do banco de dados
4. Eliminação de informações duplicadas e indexação dos dados
5. Disponibilização dos dados para pesquisa

Depois da seleção das informações a serem extraídas para cada banco, foi feito o trabalho de extração dos dados do banco de dados TAIR e KEGG para cada gene (identificado pelo seu Locus), no caso de existir o NCBI GeneID de algum gene nas informações obtidas do KEGG, executou-se a extração dos dados daquele gene no NCBI. A tecnologia utilizada foi NodeJS [2], com os módulos request [3], para baixar o HTML, e cheerio [4], para extrair os dados navegando pelo DOM da página. As informações foram salvas no formato JSON [5], contendo o locus identificador e um vetor de informações para cada categoria de informação.

Para o banco de dados, foi utilizado um banco relacional, PostgreSQL [6], contendo uma tabela com todos os genes, uma tabela para cada categoria de informação e, por fim, várias tabelas de relacionamento entre os genes e cada categoria de informação. Em todas as tabelas foram colocados chaves primárias numéricas (identificadores). A escolha de um banco de dados relacional com identificadores numéricos se deu pela facilidade de execução de algoritmos sobre esses identificadores, podendo assim gerar saídas com somente números para cálculos e estatísticas. Foram também inseridas VIEWS para listar os genes e suas informações relacionadas.

Na inserção, foi executado um script, também em NodeJS, para ler os dados no formato JSON, verificar a existência de duplicados e inserir no banco criado em PostgreSQL.

Com o banco devidamente montado e populado, deu-se início ao desenvolvimento de uma interface WEB e uma API JSON, utilizando o framework Express [7] para NodeJS, juntamente com pg-query [8] para conexão com o banco de dados. Nelas são possíveis a listagem e visualização de todos os dados e utilização dos mesmos em outros sistemas e ambientes computacionais.

RESULTADOS E DISCUSSÃO

Na extração dos dados, foi observada a dificuldade de manter o software desenvolvido, pois, sempre que ocorrerem mudanças nas estruturas HTML dos bancos de dados, terá que ser revisado e muitas vezes corrigido o algoritmo de extração, uma solução é pegar os dados de algum outro formato, como XML, porém não existem muitas opções disponíveis em todos os bancos de dados.

Links utilizados para extração dos dados:

- TAIR: <http://arabidopsis.org/servlets/TairObject?type=locus&name=LOCUS>
- KEGG: http://www.genome.jp/dbget-bin/www_bget?ath:LOCUS
- NCBI: <http://www.ncbi.nlm.nih.gov/gene/NCBIGENEID>

De cada campo, foram selecionadas algumas informações, separadas por categorias:

- TAIR: description, gene_model_type, other_names, located_in, functions_in, has;
- KEGG: definition, module, other_dbs, pathway;
- NCBI: official_symbol, official_full_name, see_related, gene_symbol, gene_description, primary_source, locus_tag, gene_type, rna_name, refseq_status, organism, lineage, also_known_as, summary, sequence, location, id_protein.

No Apêndice A temos o código para extração dos dados, montado como um módulo para o NodeJS, contendo uma função para cada banco de dados onde, para o TAIR e KEGG, passamos o LOCUS e a função para tratamento dos dados extraídos e, para o NCBI, passamos o NCBI-GeneId e a função para tratamento dos dados. Os HTMLs baixados foram salvos caso necessário o uso posterior. Caso nenhuma função de tratamento definida, é executada uma função para salvar no formato JSON.

Os arquivos criados seguiram uma estrutura predefinida de subdiretórios, sendo arquivos/BANCO/html/LOCUS.html (ex.: arquivos/kegg/html/AT1G01010.html) e arquivos/BANCO/json/LOCUS.json (ex.: arquivos/tair/json/AT1G01010.json), no caso do NCBI, substituindo o locus pelo geneid.

Abaixo o código do TAIR, na linha 9 é definido 2 parâmetros, o locus e a função de callback, nela executamos o módulo request com o link para o banco de dados.

```
9.     TAIR: function(locus, callback) {
10.         request('http://arabidopsis.org/servlets/TairObject?type=locus&name=' + locus, function (error, response, html) {
11.             var obj = { db : 'tair',
12.                 locus : locus,
13.                 description : [],
14.                 gene_model_type : [],
15.                 other_names : [],
16.                 located_in : [],
17.                 functions_in : [],
18.                 has : []};
```

Como resultado obtemos o conteúdo HTML, junto com os códigos de resposta ou algum erro, para cada informação a ser extraída, foi inicializado um vetor vazio e salvo no objeto JSON.

```
19. if (typeof callback === 'undefined') {
20.     callback = save;
21. }
```

Caso nenhuma função callback seja passada, é utilizada uma predefinida (save).

```
22. if (!error && response.statusCode == 200) {
23.     var $ = cheerio.load(html);
24.     var split = {};
25.
26.     split.description = '';
27.     split.other_names = '';
28.     split.annotations = '';
29.
30.     var tmp = '';
31.     $('<clearfix>').first().contents().find('tr').each(function(i, elem)
32.     {
33.         var aux = $(this).find('th').text().trim().toLowerCase().split('
34.         ').join('_').replace(':', '');
35.         if(obj[aux]) {
36.             obj[aux] = $(this).find('th').next().next().text().trim().split('
37.             ').join('_');
38.             if(obj[aux2]) {
39.                 obj[aux2] = $(this).find('td').first().next().next().text
39.                 ().trim().toLowerCase().split(' ').join('_');
```

```

40.         var aux3 = $(this).find('td').first().next().next().
next().next().text().trim().split(split.annotations);
41.         obj[aux2] = obj[aux2].concat(aux3);
42.     }
43.     });
44. }
45. });
46.     callback(null, obj);
47. } else {
48.     callback(error, obj);
49. }

```

Caso nenhum erro ocorrido, navegamos pelo DOM do HTML a fim de capturar as informações, onde cada informação encontrada é inserida no vetor definido anteriormente.

Por fim, chamamos a função callback passando o erro junto com um objeto vazio, ou nulo quando nenhum erro juntamente com o objeto populado.

O arquivo HTML é salvo automaticamente pelo request utilizando seu método pipe.

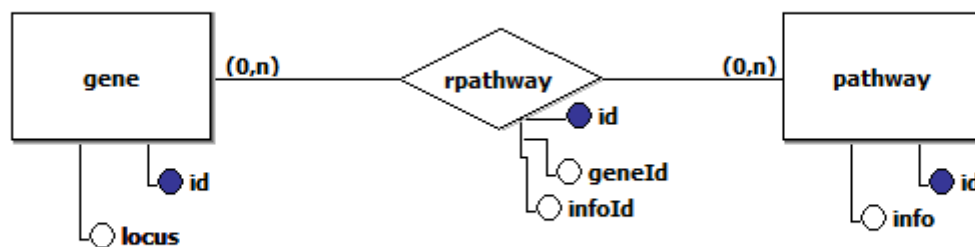
```

50.
51.     }).pipe(fq.createWriteStream('arquivos/html/tair/'+locus+'.html'));

```

No banco de dados, os resultados foram satisfatórios, apesar de não ter nenhum problema a ser corrigido, o ideal é a utilização de algum padrão de banco para dados biológicos, ou, melhor estruturação desse mesmo.

Nele foi criada uma tabela para todos os genes, contendo uma chave primária numérica juntamente com o LOCUS, uma tabela para cada categoria de informação, também com uma chave primária mais a string da informação, e, para cada informação, uma tabela de relacionamento. Como exemplo, diagrama dos genes com as informações “pathway” extraídas do KEGG na imagem abaixo:



A existência de dados duplicados foi verificada antes da inserção, dada uma informação qualquer de um gene, foi dado em SELECT em sua respectiva tabela a procura dela, inserida caso inexistente, criando também o relacionamento com o gene e, caso encontrada, foi dado outro SELECT na tabela de relacionamentos, onde era inserido caso inexistente.

O script completo para criação das tabelas encontra-se no Apêndice B, e para as views no Apêndice C.

Os genes foram inseridos manualmente utilizando um script SQL com os códigos de inserção (INSERT INTO gene (locus) VALUES ('AT1G01010')), já o script para inserção das informações pode ser visualizado no Apêndice D.

No banco, não há diferenciação das informações para o TAIR, KEGG e NCBI além de cada categoria ter uma tabela e um relacionamento separado, nesse quesito, uma melhoria a ser feita, é identificar a duplicidade de informações entre tabelas diferentes, pois como a

estrutura dos dados foi mantida a mesma, pode haver informações duplicadas em formatos diferentes.

A interface WEB para visualização dos dados é satisfatória, podendo abrir espaço para mais funcionalidades como gerar GRNs.

A API funciona como o esperado, via web, podendo ser utilizada para gerar as tabelas, em qualquer linguagem que consiga ler objetos JSON, tornando os dados facilmente acessíveis.

Tanto a interface WEB como a API oferecem os seguintes formatos:

- Lista de genes
- Lista de informações por categoria
- Lista de informações compartilhadas entre genes
- Lista de informações por gene
- Lista de informações compartilhadas entre dois ou mais genes

Como resultado dos trabalhos, os identificadores das informações foram extraídos para uma tabela em texto, sendo utilizados no desenvolvimento de um artigo (Apêndice E), aceito no CIARP 2015.

AGRADECIMENTOS

Agradeço meu orientador Fabrício Martins Lopes, pelo apoio e orientação nas pesquisas, a Universidade Tecnológica Federal do Paraná (UTFPR) e ao órgão CNPq.

REFERÊNCIAS

- [1] Ian Gorton, Paul Greenfield, Alex Szalay, Johns Hopkins. **Data-Intensive Computing in the 21st Century**.
<<http://www.computer.org/csdl/mags/co/2008/04/mco2008040030.html#bibmco20080400301>>
- [2] Node.js Foundation. **Node.JS** <<https://nodejs.org/en/>>
- [3] Lalit Kapoor, Loïc Mahieu, Mikeal Rogers, Maciej Małecki, James Nysten, Sean Hagstrom, simo, eriksm. **Request** <<https://github.com/request/request>>
- [4] Matt Mueller. **Cheerio** <<https://github.com/cheeriojs/cheerio>>
- [5] Ecma International. **JSON** <<http://www.ecma-international.org/publications/files/ECMA-ST/ECMA-404.pdf>>
- [6] The PostgreSQL Global Development Group. **PostgreSQL**
<<http://www.postgresql.org/>>
- [7] StrongLoop. **Express** <<http://expressjs.com/>>
- [8] Brian Carlson. **Node pg-query** <<https://github.com/brianc/node-pg-query>>

Apêndice A – Mecanismo de busca e extração das informações

```
1. var request = require('request');
2. var cheerio = require('cheerio');
3. var fs = require('fs');
4. var FileQueue = require('filequeue');
5. var fq = new FileQueue(100);
6.
7. module.exports = {
8.
9.   TAIR: function(locus, callback) {
10.     request('http://arabidopsis.org/servlets/TairObject?type=locus&name=' + locus,
11.     function (error, response, html) {
12.       var obj = { db : 'tair',
13.                 locus : locus,
14.                 description : [],
15.                 gene_model_type : [],
16.                 other_names : [],
17.                 located_in : [],
18.                 functions_in : [],
19.                 has : []};
20.       if (typeof callback === 'undefined') {
21.         callback = save;
22.       }
23.       if (!error && response.statusCode == 200) {
24.         var $ = cheerio.load(html);
25.         var split = {};
26.         split.description = '; '
27.         split.other_names = ', '
28.         split.annotations = ', '
29.
30.         var tmp = '';
31.         $('.clearfix').first().contents().find('tr').each(function(i, elem) {
32.           var aux = $(this).find('th').text().trim().toLowerCase().split(' ')
33.           .join('_').replace(':', '');
34.           if(obj[aux]) {
35.             obj[aux] = $(this).find('th').next().next().text().trim().split
36.             (split[aux]);
37.           }
38.           if(aux == 'annotations_category_relationship_type_keyword') {
39.             $(this).find('tr').each(function(i, elem) {
40.               var aux2 = $(this).find('td').first().next().next().text().
41.               trim().toLowerCase().split(' ').join('_');
42.               if(obj[aux2]) {
43.                 var aux3 = $(this).find('td').first().next().next().next
44.                 t().next().text().trim().split(split.annotations);
45.                 obj[aux2] = obj[aux2].concat(aux3);
46.               }
47.             });
48.           }
49.         });
50.         callback(null, obj);
51.       } else {
52.         callback(error, obj);
53.       }
54.     }).pipe(fq.createWriteStream('arquivos/html/tair/'+locus+'.html'));
55.   },
56.
57.   KEGG: function(locus, callback) {
58.     request('http://www.genome.jp/dbget-
59.     bin/www_bget?ath:' + locus, function (error, response, html) {
60.       var obj = { db : 'kegg',
```

```

57.         locus : locus,
58.         definition : [],
59.         module : [],
60.         other_dbs : [],
61.         pathway : []];
62.     if (typeof callback === 'undefined') {
63.         callback = save;
64.     }
65.     if (!error && response.statusCode === 200) {
66.         var $ = cheerio.load(html);
67.         var split = {};
68.
69.         split.definition = ',';
70.         split.module = ' ';
71.         split.other_dbs = ' ';
72.         split.pathway = ' ';
73.
74.         var tmp = '';
75.         $('>.th10').each(function(i, elem) {
76.             var aux = $(this).text().toLowerCase().trim().split(' ').join('_');
77.
78.             var aux2 = $(this).next();
79.             if (obj[aux]) {
80.                 if(aux == 'other_dbs') {
81.                     aux2.find('div').each(function(i, elem) {
82.                         if(i % 2 != 0) {
83.                             var odb = {
84.                                 name : $(this).text().trim().replace(':', ' '),
85.                                 id : $(this).next().text().trim()
86.                             };
87.                             obj.other_dbs = obj.other_dbs.concat(odb);
88.                         }
89.                     });
90.                 } else if(aux == 'pathway' || aux == 'module') {
91.                     aux2.find('tr').each(function(i, elem) {
92.                         obj[aux] = obj[aux].concat($(this).text().trim().replace(
93.                             e( /\s\n\r]+/g, ' ' ));
94.                     });
95.                 } else {
96.                     aux2.find('div').each(function(i, elem) {
97.                         obj[aux] = obj[aux].concat($(this).text().trim());
98.                     });
99.                 }
100.            });
101.            var aux = $(this).text().toLowerCase().trim().split(' ').join('_');
102.
103.            var aux2 = $(this).next();
104.            if (obj[aux]) {
105.                if(aux == 'other_dbs') {
106.                    aux2.find('div').each(function(i, elem) {
107.                        if(i % 2 != 0) {
108.                            var odb = {
109.                                name : $(this).text().trim().replace(':', ' '),
110.                                id : $(this).next().text().trim()
111.                            };
112.                            obj.other_dbs = obj.other_dbs.concat(odb);
113.                        }
114.                    });
115.                } else if(aux == 'pathway' || aux == 'module') {
116.                    aux2.find('tr').each(function(i, elem) {
117.                        obj[aux] = obj[aux].concat($(this).text().trim().replace(
118.                            e( /\s\n\r]+/g, ' ' ));

```



```

117.         });
118.     } else {
119.         aux2.find('div').each(function(i, elem) {
120.             obj[aux] = obj[aux].concat($(this).text().trim());
121.         });
122.     }
123. }
124. });
125.
126.     callback(null, obj);
127. } else {
128.     callback(error, obj);
129. }
130. }).pipe(fq.createWriteStream('arquivos/html/kegg/'+locus+'.html'));
131. },
132.
133. NCBI: function(geneid, callback) {
134.     request('http://www.ncbi.nlm.nih.gov/gene/' + geneid, function (error, response
, html) {
135.         var obj = { db : 'ncbi',
136.             geneid : geneid,
137.             official_symbol : [],
138.             official_full_name : [],
139.             see_related : [],
140.             gene_symbol : [],
141.             gene_description : [],
142.             primary_source : [],
143.             locus_tag : [],
144.             gene_type : [],
145.             rna_name : [],
146.             refseq_status : [],
147.             organism : [],
148.             lineage : [],
149.             also_known_as : [],
150.             summary : [],
151.             sequence : [],
152.             location : [],
153.             id_protein : []};
154.         if (typeof callback === 'undefined') {
155.             callback = save;
156.         }
157.         if (!error && response.statusCode == 200) {
158.             var $ = cheerio.load(html);
159.             var split = {};
160.
161.             split.official_full_name = ', '
162.             split.see_related = '; ';
163.             split.lineage = '; ';
164.             split.also_known_as = '; ';
165.
166.             $('#summaryDl').find('dt').each(function(i, elem) {
167.                 var aux = $(this).text().trim().toLowerCase().split(' ').join('_');
168.
169.                 var aux2 = [];
170.                 aux2 = $(this).next().text().trim().split(split[aux]);
171.                 if(obj[aux])
172.                     obj[aux] = aux2;
173.             });
174.
175.             $('.gc_cont').find('dt').each(function(i, elem) {
176.                 var aux = $(this).text().trim().toLowerCase().replace(':', '').split
(' ').join('_');
177.                 var aux2 = $(this).next().text().trim().split(split[aux]);
178.                 if(obj[aux])
179.                     obj[aux] = aux2;
180.             });

```

```
180.
181.         callback(null, obj);
182.     } else {
183.         callback(error, obj);
184.     }
185.     }).pipe(fq.createWriteStream('arquivos/html/ncbi/'+geneid+'.html'));
186. },
187.
188.     save: function(error, obj) {
189.         var name;
190.         var path = 'arquivos/json/' + obj.db + '/';
191.         if (obj.db === 'ncbi') {
192.             name = obj.geneid;
193.         } else {
194.             name = obj.locus;
195.         }
196.         if(error) {
197.             console.log(error);
198.         } else {
199.             fq.writeFile(path + name + '.json', JSON.stringify(obj, null, 4), function(
200. err) {
201.                 if(err) {
202.                     console.log(err);
203.                 }
204.             });
205.         }
206.     }
```

Apêndice B – Script para gerar o banco de dados

```
1. CREATE TABLE gene
2. (
3.     id serial NOT NULL,
4.     name character varying(20) NOT NULL,
5.     ncbigeneid integer DEFAULT 0,
6.     CONSTRAINT gene_pkey PRIMARY KEY (id),
7.     CONSTRAINT gene_name_key UNIQUE (name)
8. );
9.
10. CREATE TABLE also_known_as
11. (
12.     id serial NOT NULL,
13.     info text NOT NULL,
14.     CONSTRAINT also_known_as_pkey PRIMARY KEY (id),
15.     CONSTRAINT also_known_as_info_key UNIQUE (info)
16. );
17.
18. CREATE TABLE definition
19. (
20.     id serial NOT NULL,
21.     info text NOT NULL,
22.     CONSTRAINT definition_pkey PRIMARY KEY (id),
23.     CONSTRAINT definition_info_key UNIQUE (info)
24. );
25.
26. CREATE TABLE description
27. (
28.     id serial NOT NULL,
29.     info text NOT NULL,
30.     CONSTRAINT description_pkey PRIMARY KEY (id),
31.     CONSTRAINT description_info_key UNIQUE (info)
32. );
33.
34. CREATE TABLE gene_description
35. (
36.     id serial NOT NULL,
37.     info text NOT NULL,
38.     CONSTRAINT gene_description_pkey PRIMARY KEY (id),
39.     CONSTRAINT gene_description_info_key UNIQUE (info)
40. );
41.
42. CREATE TABLE gene_model_type
43. (
44.     id serial NOT NULL,
45.     info text NOT NULL,
46.     CONSTRAINT gene_model_type_pkey PRIMARY KEY (id),
47.     CONSTRAINT gene_model_type_info_key UNIQUE (info)
48. );
49.
50. CREATE TABLE gene_symbol
51. (
52.     id serial NOT NULL,
53.     info text NOT NULL,
54.     CONSTRAINT gene_symbol_pkey PRIMARY KEY (id),
55.     CONSTRAINT gene_symbol_info_key UNIQUE (info)
56. );
57.
58. CREATE TABLE gene_type
59. (
60.     id serial NOT NULL,
61.     info text NOT NULL,
62.     CONSTRAINT gene_type_pkey PRIMARY KEY (id),
63.     CONSTRAINT gene_type_info_key UNIQUE (info)
```

```

64. );
65.
66. CREATE TABLE id_protein
67. (
68.     id serial NOT NULL,
69.     info text NOT NULL,
70.     CONSTRAINT id_protein_pkey PRIMARY KEY (id),
71.     CONSTRAINT id_protein_info_key UNIQUE (info)
72. );
73.
74. CREATE TABLE lineage
75. (
76.     id serial NOT NULL,
77.     info text NOT NULL,
78.     CONSTRAINT lineage_pkey PRIMARY KEY (id),
79.     CONSTRAINT lineage_info_key UNIQUE (info)
80. );
81.
82. CREATE TABLE location
83. (
84.     id serial NOT NULL,
85.     info text NOT NULL,
86.     CONSTRAINT location_pkey PRIMARY KEY (id),
87.     CONSTRAINT location_info_key UNIQUE (info)
88. );
89.
90. CREATE TABLE locus_tag
91. (
92.     id serial NOT NULL,
93.     info text NOT NULL,
94.     CONSTRAINT locus_tag_pkey PRIMARY KEY (id),
95.     CONSTRAINT locus_tag_info_key UNIQUE (info)
96. );
97.
98. CREATE TABLE module
99. (
100.    id serial NOT NULL,
101.    info text NOT NULL,
102.    CONSTRAINT module_pkey PRIMARY KEY (id),
103.    CONSTRAINT module_info_key UNIQUE (info)
104. );
105.
106. CREATE TABLE official_full_name
107. (
108.    id serial NOT NULL,
109.    info text NOT NULL,
110.    CONSTRAINT official_full_name_pkey PRIMARY KEY (id),
111.    CONSTRAINT official_full_name_info_key UNIQUE (info)
112. );
113.
114. CREATE TABLE official_symbol
115. (
116.    id serial NOT NULL,
117.    info text NOT NULL,
118.    CONSTRAINT official_symbol_pkey PRIMARY KEY (id),
119.    CONSTRAINT official_symbol_info_key UNIQUE (info)
120. );
121.
122. CREATE TABLE organism
123. (
124.    id serial NOT NULL,
125.    info text NOT NULL,
126.    CONSTRAINT organism_pkey PRIMARY KEY (id),
127.    CONSTRAINT organism_info_key UNIQUE (info)
128. );
129.

```

```

130. CREATE TABLE other_dbs
131. (
132.   id serial NOT NULL,
133.   name text NOT NULL,
134.   info text NOT NULL,
135.   CONSTRAINT other_dbs_pkey PRIMARY KEY (id)
136. );
137.
138. CREATE TABLE other_names
139. (
140.   id serial NOT NULL,
141.   info text NOT NULL,
142.   CONSTRAINT other_names_pkey PRIMARY KEY (id),
143.   CONSTRAINT other_names_info_key UNIQUE (info)
144. );
145.
146. CREATE TABLE pathway
147. (
148.   id serial NOT NULL,
149.   info text NOT NULL,
150.   CONSTRAINT pathway_pkey PRIMARY KEY (id),
151.   CONSTRAINT pathway_info_key UNIQUE (info)
152. );
153.
154. CREATE TABLE primary_source
155. (
156.   id serial NOT NULL,
157.   info text NOT NULL,
158.   CONSTRAINT primary_source_pkey PRIMARY KEY (id),
159.   CONSTRAINT primary_source_info_key UNIQUE (info)
160. );
161.
162. CREATE TABLE refseq_status
163. (
164.   id serial NOT NULL,
165.   info text NOT NULL,
166.   CONSTRAINT refseq_status_pkey PRIMARY KEY (id),
167.   CONSTRAINT refseq_status_info_key UNIQUE (info)
168. );
169.
170. CREATE TABLE rna_name
171. (
172.   id serial NOT NULL,
173.   info text NOT NULL,
174.   CONSTRAINT rna_name_pkey PRIMARY KEY (id),
175.   CONSTRAINT rna_name_info_key UNIQUE (info)
176. );
177.
178. CREATE TABLE see_related
179. (
180.   id serial NOT NULL,
181.   info text NOT NULL,
182.   CONSTRAINT see_related_pkey PRIMARY KEY (id),
183.   CONSTRAINT see_related_info_key UNIQUE (info)
184. );
185.
186. CREATE TABLE sequence
187. (
188.   id serial NOT NULL,
189.   info text NOT NULL,
190.   CONSTRAINT sequence_pkey PRIMARY KEY (id),
191.   CONSTRAINT sequence_info_key UNIQUE (info)
192. );
193.
194. CREATE TABLE summary
195. (

```

```

196. id serial NOT NULL,
197. info text NOT NULL,
198. CONSTRAINT summary_pkey PRIMARY KEY (id),
199. CONSTRAINT summary_info_key UNIQUE (info)
200. );
201.
202. CREATE TABLE located_in
203. (
204. id serial NOT NULL,
205. info text NOT NULL,
206. CONSTRAINT located_in_pkey PRIMARY KEY (id),
207. CONSTRAINT located_in_info_key UNIQUE (info)
208. );
209.
210. CREATE TABLE functions_in
211. (
212. id serial NOT NULL,
213. info text NOT NULL,
214. CONSTRAINT functions_in_pkey PRIMARY KEY (id),
215. CONSTRAINT functions_in_info_key UNIQUE (info)
216. );
217.
218. CREATE TABLE has
219. (
220. id serial NOT NULL,
221. info text NOT NULL,
222. CONSTRAINT has_pkey PRIMARY KEY (id),
223. CONSTRAINT has_info_key UNIQUE (info)
224. );
225.
226. CREATE TABLE r_also_known_as
227. (
228. id serial NOT NULL,
229. gid integer NOT NULL,
230. cid integer NOT NULL,
231. CONSTRAINT r_also_known_as_pkey PRIMARY KEY (id),
232. CONSTRAINT r_also_known_as_cid_fkey FOREIGN KEY (cid)
233. REFERENCES also_known_as (id) MATCH SIMPLE
234. ON UPDATE CASCADE ON DELETE CASCADE,
235. CONSTRAINT r_also_known_as_gid_fkey FOREIGN KEY (gid)
236. REFERENCES gene (id) MATCH SIMPLE
237. ON UPDATE CASCADE ON DELETE CASCADE
238. );
239.
240. CREATE TABLE r_definition
241. (
242. id serial NOT NULL,
243. gid integer NOT NULL,
244. cid integer NOT NULL,
245. CONSTRAINT r_definition_pkey PRIMARY KEY (id),
246. CONSTRAINT r_definition_cid_fkey FOREIGN KEY (cid)
247. REFERENCES definition (id) MATCH SIMPLE
248. ON UPDATE CASCADE ON DELETE CASCADE,
249. CONSTRAINT r_definition_gid_fkey FOREIGN KEY (gid)
250. REFERENCES gene (id) MATCH SIMPLE
251. ON UPDATE CASCADE ON DELETE CASCADE
252. );
253.
254. CREATE TABLE r_description
255. (
256. id serial NOT NULL,
257. gid integer NOT NULL,
258. cid integer NOT NULL,
259. CONSTRAINT r_description_pkey PRIMARY KEY (id),
260. CONSTRAINT r_description_cid_fkey FOREIGN KEY (cid)
261. REFERENCES description (id) MATCH SIMPLE

```

```

262.         ON UPDATE CASCADE ON DELETE CASCADE,
263.     CONSTRAINT r_description_gid_fkey FOREIGN KEY (gid)
264.         REFERENCES gene (id) MATCH SIMPLE
265.         ON UPDATE CASCADE ON DELETE CASCADE
266. );
267.
268. CREATE TABLE r_gene_description
269. (
270.     id serial NOT NULL,
271.     gid integer NOT NULL,
272.     cid integer NOT NULL,
273.     CONSTRAINT r_gene_description_pkey PRIMARY KEY (id),
274.     CONSTRAINT r_gene_description_cid_fkey FOREIGN KEY (cid)
275.         REFERENCES gene_description (id) MATCH SIMPLE
276.         ON UPDATE CASCADE ON DELETE CASCADE,
277.     CONSTRAINT r_gene_description_gid_fkey FOREIGN KEY (gid)
278.         REFERENCES gene (id) MATCH SIMPLE
279.         ON UPDATE CASCADE ON DELETE CASCADE
280. );
281.
282. CREATE TABLE r_gene_model_type
283. (
284.     id serial NOT NULL,
285.     gid integer NOT NULL,
286.     cid integer NOT NULL,
287.     CONSTRAINT r_gene_model_type_pkey PRIMARY KEY (id),
288.     CONSTRAINT r_gene_model_type_cid_fkey FOREIGN KEY (cid)
289.         REFERENCES gene_model_type (id) MATCH SIMPLE
290.         ON UPDATE CASCADE ON DELETE CASCADE,
291.     CONSTRAINT r_gene_model_type_gid_fkey FOREIGN KEY (gid)
292.         REFERENCES gene (id) MATCH SIMPLE
293.         ON UPDATE CASCADE ON DELETE CASCADE
294. );
295.
296. CREATE TABLE r_gene_symbol
297. (
298.     id serial NOT NULL,
299.     gid integer NOT NULL,
300.     cid integer NOT NULL,
301.     CONSTRAINT r_gene_symbol_pkey PRIMARY KEY (id),
302.     CONSTRAINT r_gene_symbol_cid_fkey FOREIGN KEY (cid)
303.         REFERENCES gene_symbol (id) MATCH SIMPLE
304.         ON UPDATE CASCADE ON DELETE CASCADE,
305.     CONSTRAINT r_gene_symbol_gid_fkey FOREIGN KEY (gid)
306.         REFERENCES gene (id) MATCH SIMPLE
307.         ON UPDATE CASCADE ON DELETE CASCADE
308. );
309.
310. CREATE TABLE r_gene_type
311. (
312.     id serial NOT NULL,
313.     gid integer NOT NULL,
314.     cid integer NOT NULL,
315.     CONSTRAINT r_gene_type_pkey PRIMARY KEY (id),
316.     CONSTRAINT r_gene_type_cid_fkey FOREIGN KEY (cid)
317.         REFERENCES gene_type (id) MATCH SIMPLE
318.         ON UPDATE CASCADE ON DELETE CASCADE,
319.     CONSTRAINT r_gene_type_gid_fkey FOREIGN KEY (gid)
320.         REFERENCES gene (id) MATCH SIMPLE
321.         ON UPDATE CASCADE ON DELETE CASCADE
322. );
323.
324. CREATE TABLE r_id_protein
325. (
326.     id serial NOT NULL,
327.     gid integer NOT NULL,

```

```

328.     cid integer NOT NULL,
329.     CONSTRAINT r_id_protein_pkey PRIMARY KEY (id),
330.     CONSTRAINT r_id_protein_cid_fkey FOREIGN KEY (cid)
331.         REFERENCES id_protein (id) MATCH SIMPLE
332.         ON UPDATE CASCADE ON DELETE CASCADE,
333.     CONSTRAINT r_id_protein_gid_fkey FOREIGN KEY (gid)
334.         REFERENCES gene (id) MATCH SIMPLE
335.         ON UPDATE CASCADE ON DELETE CASCADE
336. );
337.
338. CREATE TABLE r_lineage
339. (
340.     id serial NOT NULL,
341.     gid integer NOT NULL,
342.     cid integer NOT NULL,
343.     CONSTRAINT r_lineage_pkey PRIMARY KEY (id),
344.     CONSTRAINT r_lineage_cid_fkey FOREIGN KEY (cid)
345.         REFERENCES lineage (id) MATCH SIMPLE
346.         ON UPDATE CASCADE ON DELETE CASCADE,
347.     CONSTRAINT r_lineage_gid_fkey FOREIGN KEY (gid)
348.         REFERENCES gene (id) MATCH SIMPLE
349.         ON UPDATE CASCADE ON DELETE CASCADE
350. );
351.
352. CREATE TABLE r_location
353. (
354.     id serial NOT NULL,
355.     gid integer NOT NULL,
356.     cid integer NOT NULL,
357.     CONSTRAINT r_location_pkey PRIMARY KEY (id),
358.     CONSTRAINT r_location_cid_fkey FOREIGN KEY (cid)
359.         REFERENCES location (id) MATCH SIMPLE
360.         ON UPDATE CASCADE ON DELETE CASCADE,
361.     CONSTRAINT r_location_gid_fkey FOREIGN KEY (gid)
362.         REFERENCES gene (id) MATCH SIMPLE
363.         ON UPDATE CASCADE ON DELETE CASCADE
364. );
365.
366. CREATE TABLE r_locus_tag
367. (
368.     id serial NOT NULL,
369.     gid integer NOT NULL,
370.     cid integer NOT NULL,
371.     CONSTRAINT r_locus_tag_pkey PRIMARY KEY (id),
372.     CONSTRAINT r_locus_tag_cid_fkey FOREIGN KEY (cid)
373.         REFERENCES locus_tag (id) MATCH SIMPLE
374.         ON UPDATE CASCADE ON DELETE CASCADE,
375.     CONSTRAINT r_locus_tag_gid_fkey FOREIGN KEY (gid)
376.         REFERENCES gene (id) MATCH SIMPLE
377.         ON UPDATE CASCADE ON DELETE CASCADE
378. );
379.
380. CREATE TABLE r_module
381. (
382.     id serial NOT NULL,
383.     gid integer NOT NULL,
384.     cid integer NOT NULL,
385.     CONSTRAINT r_module_pkey PRIMARY KEY (id),
386.     CONSTRAINT r_module_cid_fkey FOREIGN KEY (cid)
387.         REFERENCES module (id) MATCH SIMPLE
388.         ON UPDATE CASCADE ON DELETE CASCADE,
389.     CONSTRAINT r_module_gid_fkey FOREIGN KEY (gid)
390.         REFERENCES gene (id) MATCH SIMPLE
391.         ON UPDATE CASCADE ON DELETE CASCADE
392. );
393.

```



```

394. CREATE TABLE r_official_full_name
395. (
396.   id serial NOT NULL,
397.   gid integer NOT NULL,
398.   cid integer NOT NULL,
399.   CONSTRAINT r_official_full_name_pkey PRIMARY KEY (id),
400.   CONSTRAINT r_official_full_name_cid_fkey FOREIGN KEY (cid)
401.     REFERENCES official_full_name (id) MATCH SIMPLE
402.     ON UPDATE CASCADE ON DELETE CASCADE,
403.   CONSTRAINT r_official_full_name_gid_fkey FOREIGN KEY (gid)
404.     REFERENCES gene (id) MATCH SIMPLE
405.     ON UPDATE CASCADE ON DELETE CASCADE
406. );
407.
408. CREATE TABLE r_official_symbol
409. (
410.   id serial NOT NULL,
411.   gid integer NOT NULL,
412.   cid integer NOT NULL,
413.   CONSTRAINT r_official_symbol_pkey PRIMARY KEY (id),
414.   CONSTRAINT r_official_symbol_cid_fkey FOREIGN KEY (cid)
415.     REFERENCES official_symbol (id) MATCH SIMPLE
416.     ON UPDATE CASCADE ON DELETE CASCADE,
417.   CONSTRAINT r_official_symbol_gid_fkey FOREIGN KEY (gid)
418.     REFERENCES gene (id) MATCH SIMPLE
419.     ON UPDATE CASCADE ON DELETE CASCADE
420. );
421.
422. CREATE TABLE r_organism
423. (
424.   id serial NOT NULL,
425.   gid integer NOT NULL,
426.   cid integer NOT NULL,
427.   CONSTRAINT r_organism_pkey PRIMARY KEY (id),
428.   CONSTRAINT r_organism_cid_fkey FOREIGN KEY (cid)
429.     REFERENCES organism (id) MATCH SIMPLE
430.     ON UPDATE CASCADE ON DELETE CASCADE,
431.   CONSTRAINT r_organism_gid_fkey FOREIGN KEY (gid)
432.     REFERENCES gene (id) MATCH SIMPLE
433.     ON UPDATE CASCADE ON DELETE CASCADE
434. );
435.
436. CREATE TABLE r_other_dbs
437. (
438.   id serial NOT NULL,
439.   gid integer NOT NULL,
440.   cid integer NOT NULL,
441.   CONSTRAINT r_other_dbs_pkey PRIMARY KEY (id),
442.   CONSTRAINT r_other_dbs_cid_fkey FOREIGN KEY (cid)
443.     REFERENCES other_dbs (id) MATCH SIMPLE
444.     ON UPDATE CASCADE ON DELETE CASCADE,
445.   CONSTRAINT r_other_dbs_gid_fkey FOREIGN KEY (gid)
446.     REFERENCES gene (id) MATCH SIMPLE
447.     ON UPDATE CASCADE ON DELETE CASCADE
448. );
449.
450. CREATE TABLE r_other_names
451. (
452.   id serial NOT NULL,
453.   gid integer NOT NULL,
454.   cid integer NOT NULL,
455.   CONSTRAINT r_other_names_pkey PRIMARY KEY (id),
456.   CONSTRAINT r_other_names_cid_fkey FOREIGN KEY (cid)
457.     REFERENCES other_names (id) MATCH SIMPLE
458.     ON UPDATE CASCADE ON DELETE CASCADE,
459.   CONSTRAINT r_other_names_gid_fkey FOREIGN KEY (gid)

```

```

460. REFERENCES gene (id) MATCH SIMPLE
461. ON UPDATE CASCADE ON DELETE CASCADE
462. );
463.
464. CREATE TABLE r_pathway
465. (
466. id serial NOT NULL,
467. gid integer NOT NULL,
468. cid integer NOT NULL,
469. CONSTRAINT r_pathway_pkey PRIMARY KEY (id),
470. CONSTRAINT r_pathway_cid_fkey FOREIGN KEY (cid)
471. REFERENCES pathway (id) MATCH SIMPLE
472. ON UPDATE CASCADE ON DELETE CASCADE,
473. CONSTRAINT r_pathway_gid_fkey FOREIGN KEY (gid)
474. REFERENCES gene (id) MATCH SIMPLE
475. ON UPDATE CASCADE ON DELETE CASCADE
476. );
477.
478. CREATE TABLE r_primary_source
479. (
480. id serial NOT NULL,
481. gid integer NOT NULL,
482. cid integer NOT NULL,
483. CONSTRAINT r_primary_source_pkey PRIMARY KEY (id),
484. CONSTRAINT r_primary_source_cid_fkey FOREIGN KEY (cid)
485. REFERENCES primary_source (id) MATCH SIMPLE
486. ON UPDATE CASCADE ON DELETE CASCADE,
487. CONSTRAINT r_primary_source_gid_fkey FOREIGN KEY (gid)
488. REFERENCES gene (id) MATCH SIMPLE
489. ON UPDATE CASCADE ON DELETE CASCADE
490. );
491.
492. CREATE TABLE r_refseq_status
493. (
494. id serial NOT NULL,
495. gid integer NOT NULL,
496. cid integer NOT NULL,
497. CONSTRAINT r_refseq_status_pkey PRIMARY KEY (id),
498. CONSTRAINT r_refseq_status_cid_fkey FOREIGN KEY (cid)
499. REFERENCES refseq_status (id) MATCH SIMPLE
500. ON UPDATE CASCADE ON DELETE CASCADE,
501. CONSTRAINT r_refseq_status_gid_fkey FOREIGN KEY (gid)
502. REFERENCES gene (id) MATCH SIMPLE
503. ON UPDATE CASCADE ON DELETE CASCADE
504. );
505.
506. CREATE TABLE r_rna_name
507. (
508. id serial NOT NULL,
509. gid integer NOT NULL,
510. cid integer NOT NULL,
511. CONSTRAINT r_rna_name_pkey PRIMARY KEY (id),
512. CONSTRAINT r_rna_name_cid_fkey FOREIGN KEY (cid)
513. REFERENCES rna_name (id) MATCH SIMPLE
514. ON UPDATE CASCADE ON DELETE CASCADE,
515. CONSTRAINT r_rna_name_gid_fkey FOREIGN KEY (gid)
516. REFERENCES gene (id) MATCH SIMPLE
517. ON UPDATE CASCADE ON DELETE CASCADE
518. );
519.
520. CREATE TABLE r_see_related
521. (
522. id serial NOT NULL,
523. gid integer NOT NULL,
524. cid integer NOT NULL,
525. CONSTRAINT r_see_related_pkey PRIMARY KEY (id),

```

```

526.     CONSTRAINT r_see_related_cid_fkey FOREIGN KEY (cid)
527.         REFERENCES see_related (id) MATCH SIMPLE
528.         ON UPDATE CASCADE ON DELETE CASCADE,
529.     CONSTRAINT r_see_related_gid_fkey FOREIGN KEY (gid)
530.         REFERENCES gene (id) MATCH SIMPLE
531.         ON UPDATE CASCADE ON DELETE CASCADE
532. );
533.
534. CREATE TABLE r_sequence
535. (
536.     id serial NOT NULL,
537.     gid integer NOT NULL,
538.     cid integer NOT NULL,
539.     CONSTRAINT r_sequence_pkey PRIMARY KEY (id),
540.     CONSTRAINT r_sequence_cid_fkey FOREIGN KEY (cid)
541.         REFERENCES sequence (id) MATCH SIMPLE
542.         ON UPDATE CASCADE ON DELETE CASCADE,
543.     CONSTRAINT r_sequence_gid_fkey FOREIGN KEY (gid)
544.         REFERENCES gene (id) MATCH SIMPLE
545.         ON UPDATE CASCADE ON DELETE CASCADE
546. );
547.
548. CREATE TABLE r_summary
549. (
550.     id serial NOT NULL,
551.     gid integer NOT NULL,
552.     cid integer NOT NULL,
553.     CONSTRAINT r_summary_pkey PRIMARY KEY (id),
554.     CONSTRAINT r_summary_cid_fkey FOREIGN KEY (cid)
555.         REFERENCES summary (id) MATCH SIMPLE
556.         ON UPDATE CASCADE ON DELETE CASCADE,
557.     CONSTRAINT r_summary_gid_fkey FOREIGN KEY (gid)
558.         REFERENCES gene (id) MATCH SIMPLE
559.         ON UPDATE CASCADE ON DELETE CASCADE
560. );
561.
562. CREATE TABLE r_located_in
563. (
564.     id serial NOT NULL,
565.     gid integer NOT NULL,
566.     cid integer NOT NULL,
567.     CONSTRAINT r_located_in_pkey PRIMARY KEY (id),
568.     CONSTRAINT r_located_in_cid_fkey FOREIGN KEY (cid)
569.         REFERENCES located_in (id) MATCH SIMPLE
570.         ON UPDATE CASCADE ON DELETE CASCADE,
571.     CONSTRAINT r_located_in_gid_fkey FOREIGN KEY (gid)
572.         REFERENCES gene (id) MATCH SIMPLE
573.         ON UPDATE CASCADE ON DELETE CASCADE
574. );
575.
576. CREATE TABLE r_functions_in
577. (
578.     id serial NOT NULL,
579.     gid integer NOT NULL,
580.     cid integer NOT NULL,
581.     CONSTRAINT r_functions_in_pkey PRIMARY KEY (id),
582.     CONSTRAINT r_functions_in_cid_fkey FOREIGN KEY (cid)
583.         REFERENCES functions_in (id) MATCH SIMPLE
584.         ON UPDATE CASCADE ON DELETE CASCADE,
585.     CONSTRAINT r_functions_in_gid_fkey FOREIGN KEY (gid)
586.         REFERENCES gene (id) MATCH SIMPLE
587.         ON UPDATE CASCADE ON DELETE CASCADE
588. );
589.
590. CREATE TABLE r_has
591. (

```

```
592. id serial NOT NULL,  
593. gid integer NOT NULL,  
594. cid integer NOT NULL,  
595. CONSTRAINT r_has_pkey PRIMARY KEY (id),  
596. CONSTRAINT r_has_cid_fkey FOREIGN KEY (cid)  
597. REFERENCES has (id) MATCH SIMPLE  
598. ON UPDATE CASCADE ON DELETE CASCADE,  
599. CONSTRAINT r_has_gid_fkey FOREIGN KEY (gid)  
600. REFERENCES gene (id) MATCH SIMPLE  
601. ON UPDATE CASCADE ON DELETE CASCADE  
602. );
```

Apêndice C – Script para gerar as views

1. **CREATE VIEW** vdescription **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** description i, r_description r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
2. **CREATE VIEW** vgene_model_type **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** gene_model_type i, r_gene_model_type r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
3. **CREATE VIEW** vother_names **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** other_names i, r_other_names r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
4. **CREATE VIEW** vlocated_in **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** located_in i, r_located_in r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
5. **CREATE VIEW** vfunctions_in **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** functions_in i, r_functions_in r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
6. **CREATE VIEW** vhas **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** has i, r_has r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
7. **CREATE VIEW** vdefinition **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** definition i, r_definition r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
8. **CREATE VIEW** vmodule **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** module i, r_module r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
9. **CREATE VIEW** vother_dbs **AS SELECT** i.name **AS** db, i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** other_dbs i, r_other_dbs r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
10. **CREATE VIEW** vpathway **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** pathway i, r_pathway r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
11. **CREATE VIEW** vofficial_symbol **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** official_symbol i, r_official_symbol r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
12. **CREATE VIEW** vofficial_full_name **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** official_full_name i, r_official_full_name r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
13. **CREATE VIEW** vsee_related **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** see_related i, r_see_related r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
14. **CREATE VIEW** vgene_symbol **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** gene_symbol i, r_gene_symbol r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
15. **CREATE VIEW** vgene_description **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** gene_description i, r_gene_description r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
16. **CREATE VIEW** vprimary_source **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** primary_source i, r_primary_source r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
17. **CREATE VIEW** vlocus_tag **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** locus_tag i, r_locus_tag r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
18. **CREATE VIEW** vgene_type **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** gene_type i, r_gene_type r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
19. **CREATE VIEW** vrna_name **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** rna_name i, r_rna_name r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
20. **CREATE VIEW** vrefseq_status **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** refseq_status i, r_refseq_status r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
21. **CREATE VIEW** vorganism **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** organism i, r_organism r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
22. **CREATE VIEW** vlineage **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** lineage i, r_lineage r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
23. **CREATE VIEW** valso_known_as **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** also_known_as i, r_also_known_as r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;

24. **CREATE VIEW** vsummary **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** summary i, r_summary r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
25. **CREATE VIEW** vsequence **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** sequence i, r_sequence r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
26. **CREATE VIEW** vlocation **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** location i, r_location r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;
27. **CREATE VIEW** vid_protein **AS SELECT** i.info **AS** info, i.id **AS** id, g.id **AS** gid, g.name **AS** locus **FROM** id_protein i, r_id_protein r, gene g **WHERE** i.id = r.cid **AND** g.id = r.gid;

Apêndice D – Script para inserção das informações no banco

```
1. var fs = require('fs');
2. var FileQueue = require('filequeue');
3. var fq = new FileQueue(100);
4. var anyDB = require('any-db');
5. var conn = anyDB.createPool('postgres://postgres:testebio@localhost/teste', {min: 1,
max: 1});
6.
7. function escape (str) {
8.     return str.replace(/[\0\x08\x09\x1a\n\r"'\%]/g, function (char) {
9.         switch (char) {
10.            case "\0":
11.                return "\\0";
12.            case "\x08":
13.                return "\\b";
14.            case "\x09":
15.                return "\\t";
16.            case "\x1a":
17.                return "\\z";
18.            case "\n":
19.                return "\\n";
20.            case "\r":
21.                return "\\r";
22.            case "\"":
23.                return "\\\"";
24.            case "%":
25.                return "\\%"+char;
26.        }
27.    });
28. }
29.
30. var tair = fs.readdirSync(__dirname + '/arquivos/json/tair/');
31. var kegg = fs.readdirSync(__dirname + '/arquivos/json/kegg/');
32. var ncbi = fs.readdirSync(__dirname + '/arquivos/json/ncbi/');
33. var files = {tair: tair, kegg: kegg, ncbi: ncbi}
34.
35. function addInfo (n1, n2, ni) {
36.     var infoArray = [{i: 'description', db: 'tair'},
37.                     {i: 'gene_model_type', db: 'tair'},
38.                     {i: 'other_names', db: 'tair'},
39.                     {i: 'located_in', db: 'tair'},
40.                     {i: 'functions_in', db: 'tair'},
41.                     {i: 'has', db: 'tair'},
42.                     {i: 'definition', db: 'kegg'},
43.                     {i: 'module', db: 'kegg'},
44.                     {i: 'other_dbs', db: 'kegg'},
45.                     {i: 'pathway', db: 'kegg'},
46.                     {i: 'official_symbol', db: 'ncbi'},
47.                     {i: 'official_full_name', db: 'ncbi'},
48.                     {i: 'see_related', db: 'ncbi'},
49.                     {i: 'gene_symbol', db: 'ncbi'},
50.                     {i: 'gene_description', db: 'ncbi'},
51.                     {i: 'primary_source', db: 'ncbi'},
52.                     {i: 'locus_tag', db: 'ncbi'},
53.                     {i: 'gene_type', db: 'ncbi'},
54.                     {i: 'rna_name', db: 'ncbi'},
55.                     {i: 'refseq_status', db: 'ncbi'},
56.                     {i: 'organism', db: 'ncbi'},
57.                     {i: 'lineage', db: 'ncbi'},
58.                     {i: 'also_known_as', db: 'ncbi'},
59.                     {i: 'summary', db: 'ncbi'},
60.                     {i: 'sequence', db: 'ncbi'},
61.                     {i: 'location', db: 'ncbi'},
```

```

62.         {i: 'id_protein', db: 'ncbi'}];
63.     var info = infoArray[ni].i;
64.     var db = infoArray[ni].db;
65.     var infoPrefixo = '';
66.     var relPrefixo = 'r_';
67.     var query1 = '';
68.     var query2 = '';
69.
70.     var data = fs.readFileSync(__dirname + '/arquivos/json/' + db + '/' + files[db][n
1]);
71.     data = JSON.parse(data);
72.
73.     var where = 'WHERE name = \'' + data.locus + '\'';
74.     if (db === 'ncbi') {
75.         where = 'WHERE ncbigeneid = \'' + data.geneid + '\'';
76.     }
77.     var x = data[info];
78.
79.     if (n2 == x.length) {
80.         if (++n1 < files[db].length) {
81.             if (x.length == 0)
82.                 do {
83.
84.                     d = fs.readFileSync(__dirname + '/arquivos/json/' + db + '/' + files[
db][++n1]);
85.                     var d = JSON.parse(d);
86.                     var z = d[info].length;
87.                     if (z > 0)
88.                         break;
89.
90.                 } while (n1 < files[db].length-1);
91.                 addInfo(n1, 0, ni);
92.             } else {
93.                 console.log(info + ' Populado!');
94.                 if (++ni < infoArray.length)
95.                     addInfo(0, 0, ni);
96.             }
97.         } else {
98.             if (info === 'other_dbs') {
99.                 query1 = 'SELECT id FROM ' + info + ' WHERE info = \'' + escape(x[n2].id
+ '\'' AND name = \'' + escape(x[n2].name) + '\'';
100.                query2 = 'INSERT INTO ' + info + ' (info, name) VALUES(\'' + escape(x[n2]
.id) + '\',\'' + escape(x[n2].name) + '\') RETURNING id';
101.            } else {
102.                query1 = 'SELECT id FROM ' + infoPrefixo + info + ' WHERE info = \'' + es
cape(x[n2]) + '\'';
103.                query2 = 'INSERT INTO ' + infoPrefixo + info + ' (info) VALUES(\'' + esca
pe(x[n2]) + '\') RETURNING id';
104.            }
105.            conn.query(query1, function (err, res) {
106.                if(err) {
107.                    console.log('1: '+err);
108.                } else {
109.                    if(res.rows.length == 0) {
110.                        conn.query(query2, function (err, res) {
111.                            if(err) {
112.                                console.log('2: '+err);
113.                            } else {
114.                                conn.query('INSERT INTO ' + relPrefixo + info + ' (gid, c
id) VALUES ((SELECT id FROM gene ' + where + '), \'' + res.rows[0].id + '\')', functi
on (err, res) {
115.                                    if(err)
116.                                        console.log('3: '+err);
117.                                });
118.                                n2++;
119.                                addInfo(n1, n2, ni);

```



```

120.         }
121.     });
122.     } else {
123.         conn.query('SELECT id FROM ' + relPrefixo + info + ' WHERE gid =
(SELECT id FROM gene ' + where + ') AND cid = \'' + res.rows[0].id + '\'', function (
err, resp) {
124.             if(err)
125.                 console.log('4: '+err);
126.             if (resp.rows.length == 0)
127.                 conn.query('INSERT INTO ' + relPrefixo + info + ' (gid, c
id) VALUES ((SELECT id FROM gene ' + where + '), \'' + res.rows[0].id + '\')', functi
on (err, res) {
128.                     if(err)
129.                         console.log('5: '+err);
130.                 });
131.                 n2++;
132.                 addInfo(n1, n2, ni);
133.             });
134.         }
135.     }
136. });
137. }
138. }
139.
140. addInfo(0, 0, 0);

```

Apêndice E - Artigo

A feature selection approach for evaluate the inference of GRNs through biological data integration - A case study on *A. thaliana*

Fábio F. R. Vicente^{1,2}, Euler Menezes¹, Gabriel Rubino¹, Juliana de Oliveira³
and Fabrício M. Lopes¹

¹ Federal University of Technology, Paraná, Brazil.

`fabiofernandes, fabricio@utfpr.edu.br`

² Institute of Mathematics and Statistics, University of São Paulo, São Paulo, Brazil.

³ Department of Biological Sciences Faculty of Sciences and Letters of Assis - FCLA
University of São Paulo State - UNESP Av. Dom Antonio, 2100, Parque Universitário
19806-900 - Assis, São Paulo, Brazil.

`juliana@assis.unesp.br`

Abstract. The inference of gene regulatory networks (GRNs) from expression profiles is a great challenge in bioinformatics due to the curse of dimensionality. For this reason, several methods that perform data integration have been developed to reduce the estimation error of the inference. However, it is not completely formulated how to use each type of biological information available. This work address this issue by proposing feature selection approach in order to integrate biological data and evaluate three types of biological information regarding their effect on the similarity of inferred GRNs. The proposed feature selection method is based on sequential forward floating selection (SFFS) search algorithm and the mean conditional entropy (MCE) as criterion function. An expression dataset was built as an additional contribution of this work containing 22746 genes and 1206 experiments regarding *A. thaliana*. The experimental results achieve 39% of GRNs improvement in average when compared to non-use of biological data integration. Besides, the results showed that the improvement is associated to a specific type of biological information: the cellular localization, which is a valuable and information for the development of new experiments and indicates an important insight for investigation.

Keywords: Gene regulatory networks, feature selection, data integration, bioinformatics, *Arabidopsis thaliana*.

1 Introduction

One of the main challenges in bioinformatics is to perform the reverse engineering of GRNs (gene regulatory networks) from expression profiles [2]. Since the genes and their corresponding proteins are only *part* of a whole biological system, it is important to discover the network of interactions between the cell components to better understand how the systems work. Furthermore, such knowledge can be useful to perform intervention and control of the system[23]. The development of

techniques such as DNA microarrays, SAGE and RNA-Seq [29] allowed studying the dynamic of the cell in a scale that was unfeasible with the previous methods. Since the interaction between the genes works as a network that produces the observed expression, the idea is to perform a reverse engineering in order to recover the network from the expression profiles.

The search space of possible networks is huge. Thus, the inference can be modeled as a feature selection problem, where a criterion function is used to evaluate a subset of predictor genes (features) that better classify the state of a given target gene based on the expression values [17,16]. In this way, the inference consists in searching for a subset of predictors for each target.

A GRN can be modeled as Bayesian Networks [12], Boolean Networks [8], Probabilistic Boolean Networks (PBN) [25] or Probabilistic Genetic Network (PGN) [3]. Also, distinct *criterion function* can be used for inference. Some criterion function as Pearson correlation [26] are limited to the evaluation of pairs of genes and the combinatorial regulation of multiple predictors cannot be assessed. The Coefficient of Determination (CoD) [13] is not limited to a number of variables and allows to model the multivariate nature of the regulation. Other criterion functions are based on information theory such as Mutual Information [22], Tsallis Entropy [19] and Mean Conditional Entropy (MCE) [18].

One of the main limitations of the inference is the high dimensionality. For this reason, it is important the development of new methods that take into account other *information* than the expression profiles in order to reduce the intrinsic estimation error of the inference and to recover plausible biological networks.

Distinct approaches have been proposed to implement data integration in GNs. Some works evaluate the data integration in biological networks [21, 28, 20]. Other methods perform a clustering analysis to identify groups of genes related to some particular property [7]. Other methods include known information about the topological properties of the network such as degree of connectivity, average path length and other local and global characteristics [16, 27].

However, a methodology to define how to select the biological information in order to reduce error on inference is still not completely formulated. Moreover, the study of distinct biological information can reveal the contribution of each one in the inference and favor the discovery of new biological knowledge such as biological information that are decisive to characterize the network together with the system dynamic. This work proposes a feature selection approach in order to evaluate a GRNs inference model based on a criterion function that encodes multivalued biological features applied on the *A. thaliana* organism.

2 Materials and Methods

2.1 Gold standard network

The gold standard network is the directed graph with the known physical direct interactions between each predictor (TF) and their target genes. The *Arabidopsis thaliana* gold standard network (AtRegNet) was obtained from AGRIS (The Arabidopsis Gene Regulatory Information Server) [30] on which each interaction

is verified in at least one experimental approach. The graph has 8154 genes of which 67 are predictors and 8131 are target genes, comprising 11481 edges.

Expression data All expression data were obtained from GEO [4] and only samples of the platform GPL198 (*Affymetrix Arabidopsis ATH1 Genome Array*) were selected. The chip contains 22810 *probe sets* that were mapped to genes through annotation data from TAIR [15], Gene Ontology [1] and TIGR [5]. A probe set is a collection of sequences (probes) used to identify a gene sequence and to measure its expression.

The GEO files containing the expression samples were obtained through the R Bioconductor , GEOquery and Biobase [6]. The files in SOFT format were acquired in three types: GDS (GEO Dataset, 13 files), GSE (GEO Series, 58 files) and GSM (GEO Samples, 109 files), each one with a set of expression experiments. The files were preprocessed to obtain the expression values and to filter out samples with missing data. Thus the 180 files resulted on a table with 22746 genes and 1206 experiments.

Biological Information A set of features associated to each probe set in the expression data were obtained from TAIR, KEGG [14] and NCBI [11]. For each probe set a corresponding *locus identifier* (TAIR) and a *NCBI gene id* were obtained. Then, the biological information associated to the NCBI gene id were obtained from the NCBI and features associated to the locus identifier were obtained from TAIR and KEGG. The dataset contains 23593 annotated elements and 15 features related to several biological aspects.

2.2 GRNs Inference

To infer the GRN it is necessary determine the best subset of predictors for each target gene. Thus, in a dataset with n transcription factors there are 2^n possible subset of predictors for each gene. Since the search space is huge, the sequential forward floating selection (SFFS)[24] algorithm was adopted to search for the best solution. The criterion function evaluates each subset by taking into account both the expression values and the biological information on the dataset. The expression value part is evaluated computing the mean conditional entropy (MCE) [18]. The MCE is computed as the average of the entropies of a target gene Y given the values of the predictors X (Eq. 1).

$$H(Y|X) = \sum_{x \in X} H(Y|x)P(x). \quad (1)$$

The MCE take into account only the expression values. Because the small number of samples it is common the occurrence of many ties between distinct set of predictors. Thus, only the expression data is commonly insufficient to determine a suitable solution. For this reason, the proposed criterion function includes a term related to another type of feature, i.e., a biological information.

$$F(Y, X) = \alpha[H(Y|X)] + \beta[D(Y, X)], \quad (2)$$

where $D(Y, X)$ is equals to 1 when the biological information of the gene Y is equals to X in the biological data set D . The parameters α and β defines the weight of each type of information where $\beta \in (0, 1)$ and $\alpha = 1 - \beta$. Thus, with the proposed feature selection approach the SFFS will search for the subset of predictors that not just minimize the entropy, but also for those that are coherent to the other biological aspects of the target gene.

3 Results and discussion

Data preprocessing The expression data are represented in a table with the variables (genes) in rows and the expression samples $s(i)$ in columns. The expression values of the genes of each sample $s(i)$, $s(i) \in \mathbb{R}$, were normalized by the *normal transformation*, defined as follows:

$$\eta[s(i)] = \frac{s(i) - E[s(i)]}{\sigma[s(i)]}, \quad (3)$$

where $\sigma[s(i)]$ is the standard deviation of $s(i)$ and $E[s(i)]$ is the expectation of $s(i)$. The normalized data were discretized into three levels $\{-1, 0, +1\}$ through a threshold mapping

$$s(i) = \begin{cases} -1 & \alpha < l \\ 0 & l \leq \alpha \leq h \\ +1 & \alpha > h \end{cases} \quad (4)$$

where, α is the expression value of the sample s , of the gene g .

Biological features The annotation of genes on public databases makes available several types of information, which refer to distinct aspects. They can be descriptive, some times adopting conventional terms to classify the gene as belonging to a category, for example, to a specific biological pathway, a known function, to cite but a few. Thus, in this work were adopted three features related to *physical* aspects and to the *activity* of the gene: (i) cellular location, (ii) pathway and (iii) function. Then, since all annotated features into the databases are attributes of type *nominal*, the value of each feature was indexed by an integer value resulting on a dataset with 20817 unique annotated pairs (*probe set*, *gene locus*) Table 1.

Table 1: Biological information available on public databases.

Database	Feature	Number of distinct values
TAIR	function	121
TAIR	located in	6
KEGG	pathway	128

Intersection between datasets To evaluate the proposed approach, the subset of genes in the intersection of the three datasets (gold standard network, expression and biological information) was adopted, resulting on 5974 Gene locus (Table 2).

Table 2: The validation dataset comprises 5974 genes.

Dataset	Description	Quantity
Gold standard network	Edges	8589
Expression data	Samples	1206
Biological features	Features	3

Evaluation of the inference In order to measure the effectiveness of the proposed approach it was adopted the *Similarity* between the gold standard and the inferred network, which was presented by [9] and is widely used for validation of methods of GNs inference. The validation is based on a confusion matrix, as described in Table 3. The *Similarity* is computed as presented in (Eq.5) and the normalized value relative to the inference based only on expression data ($\alpha = 1$), defined as follows:

Table 3: Confusion matrix. TP = true positive, FN = false negative, FP = false positive, TN = true negative.

Edge	Inferred	Not Inferred
Present	TP	FN
Absent	FP	TN

$$\text{Similarity} = \sqrt{\text{precision} \times \text{recall}} = \sqrt{\frac{TP}{TP + FP} \times \frac{TP}{TP + FN}}. \quad (5)$$

The experimental results have presented distinct performance for the three types of biological information (Figure 1(a)). In particular, the *localization* improves the similarity while *pathway* and *function* increased the error of the inference. For *cellular localization*, the curve increases from β varying from 0.1 to 0.7 and decreases for values greater than 0.7. This indicates that the expression combined to cellular localization can improve the similarity and that the use of only one type of data can limit the inference. The improvement varies according to the threshold, being higher for small values. The average improvement when the same weight is given to both expression and biological information (i.e. $\beta = 0.5$) is 39.1%.

The other two biological information presented a negative effect in similarity. A possible explanation is that the gold standard regulatory network refers to *physical* interactions and metabolic pathway and gene function refers to other type of relationship. The localization refers to a physical space on the cell where genes acts, metabolic pathway is more generic than localization and refers to the biological pathway the genes are related (not necessarily the genes are in the same local or interacts physically each other) and, an even more generic than metabolic pathway is the gene function. Thus, several proteins can have the same molecular function or acts in common metabolic pathways and can be an explanation for the decreasing in similarity. This points to the importance of the correct feature selection that can reduce the estimation error, specially in this type of multilevel data integration.

Figure 1(b) shows the improvement in similarity. Results show the similarity increases more and faster for lower thresholds and that both expression and

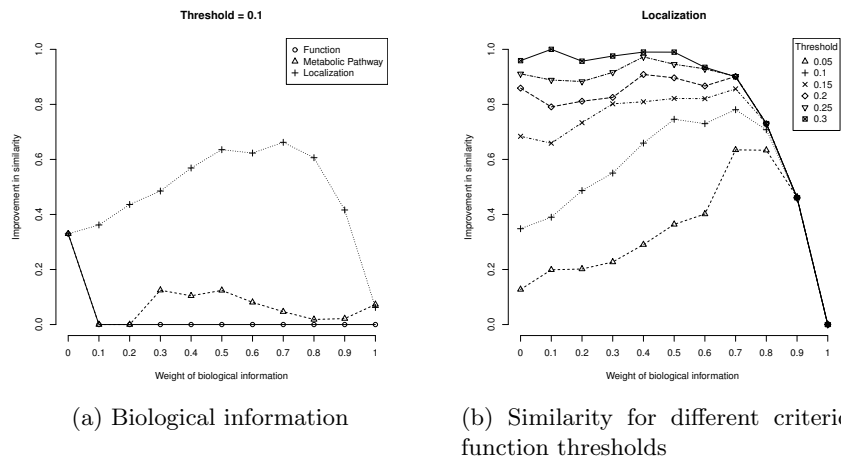


Fig. 1: Evaluation of the data integration of three types of biological information: Function, Metabolic pathway and Cellular localization. Weight = 0 means only expression. Weight = 1 means only biological information.

localization are important. For weight of biological information over 0.7 the similarity decreases faster despite of the threshold.

Another important issue is the validation of GNs inference. Here, we adopted the validation regarding to a know network, which useful to evaluate the *inferential* perspective as pointed by [10]. The validation of methods of data integration is still an problem to be solved in GNs inference.

4 Conclusions

This work presented a feature selection approach for integration of distinct datasets based on expression data and biological information for the inference of gene regulatory networks, which is based on SFFS search algorithm and MCE criterion function. An expression dataset was built as an additional contribution of this work containing 22746 genes and 1206 experiments regarding *A. thaliana*. The dataset was composed with a gold standard network, expression data and biological features were assembled and the proposed approach was applied on the composed dataset. The results showed the increasing on the similarity (average 39%) of the recovered network when the criterion function is based on cellular localization. Also, the results showed that the performance is better when expression and cellular localization are combined instead of the use of each one only, which is a valuable information for the development of new experiments and indicates an important insight for investigation.

The evaluation of GNs inference is commonly performed based on the know links of a gold-standard network. However, other biological information could be also used to validate the inference methods. Thus, as future work it is suggest the investigation of validation measures that take into account distinct biological data. Also, the inclusion of more types of biological information and to perform the proposed approach using these biological information can be considered as a

further work. Moreover, would be also important to evaluate the best set of the parameters (weight of biological information and the threshold of the criterion function).

Acknowledgements

This work was supported by FAPESP grant 2011/50761-2, UTFPR, CNPq, Fundação Araucária, CAPES and NAP eScience - PRP - USP.

References

1. Ashburner, M., Ball, C.A., Blake, J.A., Botstein, D., Butler, H., Cherry, J.M., Davis, A.P., Dolinski, K., Dwight, S.S., Eppig, J.T., Harris, M.A., Hill, D.P., Issel-Tarver, L., Kasarskis, A., Lewis, S., Matese, J.C., Richardson, J.E., Ringwald, M., Rubin, G.M., Sherlock, G.: Gene ontology: tool for the unification of biology. *Nat Genet* 25(1), 25–29 (2000)
2. Baralla, A., Mentzen, W.L., de la Fuente, A.: Inferring gene networks: dream or nightmare? *Annals of the New York Academy of Sciences* 1158, 246–56 (Mar 2009)
3. Barrera, J., Cesar-Jr, R.M., Martins-Jr, D.C., Vencio, R.Z.N., Merino, E.F., Yamamoto, M.M., Leonardi, F.G., Pereira, C.A.B., Portillo, H.A.: Constructing probabilistic genetic networks of *Plasmodium falciparum*, from dynamical expression signals of the intraerythrocytic development cycle. In: McConnell, P., Lin, S.M., Hurban, P. (eds.) *Meth. of Microarray Data Analysis*, pp. 11–26. Springer (2007)
4. Barrett, T., Wilhite, S.E., Ledoux, P., Evangelista, C., Kim, I.F., Tomashevsky, M., Marshall, K.a., Phillippy, K.H., Sherman, P.M., Holko, M., Yefanov, A., Lee, H., Zhang, N., Robertson, C.L., Serova, N., Davis, S., Soboleva, A.: NCBI GEO: Archive for functional genomics data sets - Update. *NAR* 41(D1), 991–995 (2013)
5. Childs, K.L., Hamilton, J.P., Zhu, W., Ly, E., Cheung, F., Wu, H., Rabinowicz, P.D., Town, C.D., Buell, C.R., Chan, A.P.: The tigr plant transcript assemblies database. *Nucleic Acids Research* 35(suppl 1), D846–D851 (2007)
6. Davis, S., Meltzer, P.: Geoquery: a bridge between the gene expression omnibus (geo) and bioconductor. *Bioinformatics* 14, 1846–1847 (2007)
7. De Haan, J., Piek, E., van Schaik, R., de Vlieg, J., Bauerschmidt, S., Buydens, L., Wehrens, R.: Integrating gene expression and go classification for pca by preclustering. *BMC Bioinformatics* 11(1), 158 (2010)
8. D’haeseleer, P., Liang, S., Somogyi, R.: Genetic network inference: from co-expression clustering to reverse engineering. *Bioinformatics* 16(8), 707–726 (2000)
9. Dougherty, E.R.: Validation of inference procedures for gene regulatory networks. *Current Genomics* 8(6), 351–359 (2007)
10. Dougherty, E.R.: Validation of gene regulatory networks: scientific and inferential. *Briefings in Bioinformatics* 12(3), 245–252 (2011)
11. Edgar, R., Domrachev, M., Lash, A.E.: Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. *Nucleic acids research* 30(1), 207–210 (2002)
12. Friedman, N., Linial, M., Nachman, I., Pe’er, D.: Using bayesian networks to analyze expression data. *Journal of Computational Biology* 7(3-4), 601–620 (2000)
13. Hashimoto, R.F., Kim, S., Shmulevich, I., Zhang, W., Bittner, M.L., Dougherty, E.R.: Growing genetic regulatory networks from seed genes. *Bioinformatics* 20(8), 1241–1247 (May 2004)
14. Kanehisa, M., Goto, S.: KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids Research* 28(1), 27–30 (2000)

15. Lamesch, P., Berardini, T.Z., Li, D., Swarbreck, D., Wilks, C., Sasidharan, R., Muller, R., Dreher, K., Alexander, D.L., Garcia-Hernandez, M., Karthikeyan, A.S., Lee, C.H., Nelson, W.D., Ploetz, L., Singh, S., Wensel, A., Huala, E.: The arabidopsis information resource (TAIR): improved gene annotation and new tools. *NAR* (2011)
16. Lopes, F.M., Jr., D.C.M., Barrera, J., Jr., R.M.C.: A feature selection technique for inference of graphs from their known topological properties: Revealing scale-free gene regulatory networks. *Information Sciences* 272(0), 1–15 (2014)
17. Lopes, F.M., Martins-Jr, D.C., Barrera, J., Cesar-Jr, R.M.: SFFS-MR: a floating search strategy for GRNs inference. In: *Pattern Recognition in Bioinformatics, Proceedings. LNCS*, vol. 6282, pp. 407–418. Springer Berlin / Heidelberg (2010)
18. Lopes, F.M., Martins-Jr, D.C., Cesar-Jr, R.M.: Feature selection environment for genomic applications. *BMC Bioinformatics* 9(1), 451 (October 2008)
19. Lopes, F.M., de Oliveira, E.A., Cesar-Jr, R.M.: Inference of gene regulatory networks from time series by Tsallis entropies. *BMC Systems Biology* 5(1), 61 (2011)
20. Lopes, F.M., Ray, S.S., Hashimoto, R.F., Jr., R.M.C.: Entropic biological score: a cell cycle investigation for GRNs inference. *Gene* 541(2), 129–137 (2014)
21. Lu, L.J., Xia, Y., Paccanaro, A., Yu, H., Gerstein, M.: Assessing the limits of genomic data integration for predicting protein networks. *Gen. Res.* 15(7), 945–53 (2005)
22. Margolin, A., Basso, K.N., Wiggins, C., Stolovitzky, G., Favera, R., Califano, A.: ARACNE: An algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context. *BMC Bioinformatics* 7(Suppl 1), S7 (2006)
23. Pavlopoulos, G.A., Secrier, M., Moschopoulos, C.N., Soldatos, T.G., Kossida, S., Aerts, J., Schneider, R., Bagos, P.G.: Using graph theory to analyze biological networks. *BioData mining* 4(1), 10 (jan 2011)
24. Pudil, P., Novovičová, J., Kittler, J.: Floating search methods in feature-selection. *Pattern Recognition Letters* 15(11), 1119–1125 (November 1994)
25. Shmulevich, I., Dougherty, E.R., Kim, S., Zhang, W.: Probabilistic boolean networks: a rule-based uncertainty model for gene regulatory networks. *Bioinformatics* 18(2), 261–274 (2002)
26. Stuart, J.M., Segal, E., Koller, D., Kim, S.K.: A gene-coexpression network for global discovery of conserved genetic modules. *Science* 302(5643), 249–255 (2003)
27. Vicente, F.F.R., Lopes, F.M.: SFFS-WS: A feature selection algorithm exploring the small-world properties of GNs. In: *Pattern Recognition in Bioinformatics, Proceedings. LNCS*, vol. 8626, pp. 60–71. Springer Berlin / Heidelberg (2014)
28. Vicente, F.F.R., Lopes, F.M., Hashimoto, R.F., Cesar-Jr, R.M.: Assessing the gain of biological data integration in gene networks inference. *BMC Genomics* 13(Suppl 6), S7 (2012)
29. Wang, Z., Gerstein, M., Snyder, M.: RNA-Seq: a revolutionary tool for transcriptomics. *Nat Rev Genet* 10(1), 57–63 (2009)
30. Yilmaz, A., Mejia-Guerra, M.K., Kurz, K., Liang, X., Welch, L., Grotewold, E.: Agris: the arabidopsis gene regulatory information server, an update. *Nucleic Acids Research* 39(suppl 1), D1118–D1122 (2011)