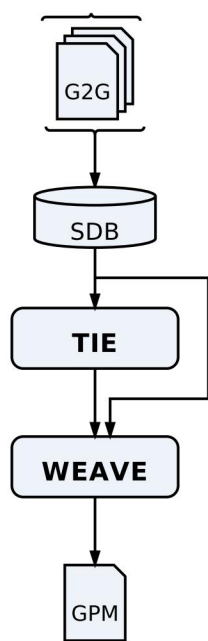


TIE & WEAVE outline



- G2G

Gene pairscore lists for different tracks.
Formatted as txt files where:
<string:gene> <string:gene> <float:value>
- SDB

Systems Database (PostgreSQL)
gene-gene pairs serve as the
primary key
- TIE

Integrate pairscore information
- WEAVE

Build multigraph from pairscore data
- GPM

Gene Pairscore Multigraph
Formatted as txt file where:
<string:gene> <string:gene> <int:track> <float:value>

BCB420: Systemikon TIE & WEAVE outline

Victor Kofia

03/10/2015

Input

Let $N = |LoG|$ where $|LoG|$ is the total # of distinct $\langle \text{string} : \text{gene} \rangle \langle \text{int} : \text{accession} \rangle$ pairs

INTERACT.txt # protein-protein interaction pairscores

1. $\langle \text{string} : \text{gene} \rangle \langle \text{string} : \text{gene} \rangle \langle \text{float} : \text{value} \rangle$
2. $\langle \text{string} : \text{gene} \rangle \langle \text{string} : \text{gene} \rangle \langle \text{float} : \text{value} \rangle$
...
 $N^2.$ $\langle \text{string} : \text{gene} \rangle \langle \text{string} : \text{gene} \rangle \langle \text{float} : \text{value} \rangle$

COEX.txt # co-expression pairscores

1. $\langle \text{string} : \text{gene} \rangle \langle \text{string} : \text{gene} \rangle \langle \text{float} : \text{value} \rangle$
2. $\langle \text{string} : \text{gene} \rangle \langle \text{string} : \text{gene} \rangle \langle \text{float} : \text{value} \rangle$
...
 $N^2.$ $\langle \text{string} : \text{gene} \rangle \langle \text{string} : \text{gene} \rangle \langle \text{float} : \text{value} \rangle$

PATHNE.txt # pathway neighbour pairscores

1. $\langle \text{string} : \text{gene} \rangle \langle \text{string} : \text{gene} \rangle \langle \text{float} : \text{value} \rangle$
2. $\langle \text{string} : \text{gene} \rangle \langle \text{string} : \text{gene} \rangle \langle \text{float} : \text{value} \rangle$
...
 $N^2.$ $\langle \text{string} : \text{gene} \rangle \langle \text{string} : \text{gene} \rangle \langle \text{float} : \text{value} \rangle$

GOSEM.txt # gene-ontology semantic similarity based

1. $\langle \text{string} : \text{gene} \rangle \langle \text{string} : \text{gene} \rangle \langle \text{float} : \text{value} \rangle$
2. $\langle \text{string} : \text{gene} \rangle \langle \text{string} : \text{gene} \rangle \langle \text{float} : \text{value} \rangle$
...
 $N^2.$ $\langle \text{string} : \text{gene} \rangle \langle \text{string} : \text{gene} \rangle \langle \text{float} : \text{value} \rangle$

Database Access

N/A

Algorithm

TIE

For database, use open source RDBMS: *PostgreSQL*.

Add pairscore information in **INTERACT.txt**, **COEX.txt**, **PATHNE.txt** and **GOSEM.txt** to database.

Database schema for *SDB*, the *Systemikon Database*:

genes(gid, gname, symbol, organism)

gene_pairs(pid, gid1, gid2)

tracks(tid, tname)

analysis(aid, pid, tid)

scores(aid, score)

```
INTEGRATE-SCORES(INTERACT.txt, COEX.txt, PATHNE.txt, GOSEM.txt) {
  initialize scoring matrix  $A$  where  $A = (a_{ij}) \in \mathbb{R}^{N \times N}$  and INTERACT.txt  $\rightarrow A$ 
  ... COEX.txt  $\rightarrow B$  # ... same as above
  ... PATHNE.txt  $\rightarrow C$ 
  ... GOSEM.txt  $\rightarrow D$ 
  initialize scoring matrix  $T$  where  $T = (t_{ij}) \in \mathbb{R}^{N \times N}$ 
  initialize txt file MULT.txt
  for  $i$  in  $[1, 2, \dots, N]$ 
    for  $j$  in  $[1, 2, \dots, N]$ 
       $t_{ij} = a_{ij} + b_{ij} + c_{ij} + d_{ij}$  # Concatenate scores
      add  $\langle i \rangle \ \langle j \rangle \ \langle 1 \rangle \ \langle a_{ij} \rangle$  to MULT.txt
      ...  $\langle i \rangle \ \langle j \rangle \ \langle 2 \rangle \ \langle b_{ij} \rangle$  to MULT.txt
      ...  $\langle i \rangle \ \langle j \rangle \ \langle 3 \rangle \ \langle c_{ij} \rangle$  to MULT.txt
      ...  $\langle i \rangle \ \langle j \rangle \ \langle 4 \rangle \ \langle d_{ij} \rangle$  to MULT.txt
      ...  $\langle i \rangle \ \langle j \rangle \ \langle 5 \rangle \ \langle t_{ij} \rangle$  to MULT.txt # 5 tracks in total
    return MULT.txt
}
```

Let MULT.txt = INTEGRATE-SCORES(INTERACT.txt, COEX.txt, PATHNE.txt, GOSEM.txt)

Add pairscore information in **MULT.txt** to database.

Output

SDB: Systemikon Database

MULT.txt # multigraph

-
- $\langle \text{string : gene} \rangle \ \langle \text{string : gene} \rangle \ \langle \text{int : track} \rangle \ \langle \text{float : value} \rangle$
 - $\langle \text{string : gene} \rangle \ \langle \text{string : gene} \rangle \ \langle \text{int : track} \rangle \ \langle \text{float : value} \rangle$

...

$4N^2$.*< string : gene > < string : gene > < int : track > < float : value >*

Tests & Verification

Input? Assumed to have been tested. Multi-graph quality is dependent on input.

Algorithm? Simple unit tests (perhaps with the assistance of an external library) should suffice.

Verification? Ensure that line numbers (of each of the text files) correspond to each other. So, for instance if $|LoG| = 10$ then INTERACT.txt, COEX.txt, PATHNE.txt and GOSEM.txt should all have 100 lines and MULT.txt should have 400 lines.

Database testing? This is tricky ...

From [Wikipedia](#):

Blackbox Testing

- Verify incoming data
- Verify outgoing data from query functions

Whitebox Testing

- Validate database tables & database schema
- Check referential integrity