

Using QuACN to Analyze Complex Biological Networks

Laurin Müller, Matthias Dehmer

December 2, 2010

Contents

1	Overview	1
2	Networks	1
3	Network Descriptors	2
3.1	Descriptors Based on Distances in a Graph	2
3.2	Descriptors Based on Other Graph-Invariants	4
3.3	Classical entropy based descriptors	5
3.4	Parametric Graph Entropy Measures	8
4	Session Info	10

1 Overview

This vignette provides an overview about the usage of QuACN.

Chapter 2 will give you an idea how to import already existing networks. In Chapter 3 a brief description of the implemented measures is presented, and it demonstrates how to call the related method in R.

2 Networks

```
> library("QuACN")
> set.seed(666)
> g <- randomGraph(1:8, 1:5, 0.36)
> plot(g, "neato")
> g
```

A *graphNEL* graph with undirected edges

Number of Nodes = 8

Number of Edges = 16

We generate a random graph with 8 nodes. This graph will be used to explain the implemented methods. To analyze a network the network has to be represented by a *graphNEL*-object, which is part of the Bioconductor *graph* package.

If you have already created networks that you want to analyze with QuACN, R offers several ways to import them. (It is important to know that networks have to be represented by *graphNEL*-objects.) Note that there is no general procedure to get your networks into an R-workspace. Some possibilities to import network data are listed below:

- **Adjacency matrix:** A representation of your network as an adjacency matrix can be easily imported and converted into a *graphNEL* object.
- **Node- and Edge-List:** With a list of nodes and Edges it is easy to create a *graphNEL*-object.
- **read.graph():** The `read.graph()` method of the *graph*-package offers the possibility to import graphs that are represented in different formats. For details see the manual of the *graph*-package.

- **System Biology Markup Language(SBML)** [1]: With the *RSBML*-package it is possible to import SBML-Models.
- **igraph-package**: Networks created with the *igraph*-package can be converted into graphNEL objects.

3 Network Descriptors

This section provides a overview of the network descriptors that are included in the QuACN package. Here we describe the respective descriptor and how to call it in R.

Note that every descriptor has at least two parameters, the *graphNEL*-object and the distance matrix representing the network. It is not necessary to pass the distance matrix to a function. If the parameters stays empty or is set to *NULL* the distance matrix will be estimated within each function. But if you want to calculate more than one descriptor, it is recommended to calculate the distance matrix separately and pass it to each method. Some of the methods need the degree of each node or the adjacency matrix to calculate their results. If they were calculated once they should have kept for later use. Specially with large networks it saves a lot of time, not to calculate these parameters for each descriptor again, and will enhance the performance of your script.

```
> mat.adj <- adjacencyMatrix(g)
> mat.dist <- distanceMatrix(g)
> vec.degree <- graph::degree(g)
> ska.dia <- diameter(g)
> ska.dia <- diameter(g, mat.dist)
```

3.1 Descriptors Based on Distances in a Graph

This section describes network measures based on distances in the network.

Wiener Index [2]:

$$W(G) := \frac{1}{2} \sum_{i=1}^{|N|} \sum_{j=1}^{|N|} d(v_i, v_j). \quad (1)$$

where $|N(G)| := |N|$ denotes the number of Nodes of the complex network. $d(v_i, v_j)$ stands for shortest distances between $v_i, v_j \in V$.

```
> wien <- wiener(g)
> wiener(g, mat.dist)
```

```
[1] 43
```

Hararay Index [3]:

$$H(G) := \frac{1}{2} \sum_{i=1}^{|N|} \sum_{j=1}^{|N|} (d(v_i, v_j))^{-1}, \quad i \neq j. \quad (2)$$

```
> harary(g)
```

```
[1] 21.16667
```

```
> harary(g, mat.dist)
```

```
[1] 21.16667
```

Balaban J Index [4]:

$$J(G) := \frac{|E|}{\mu + 1} \sum_{(v_i, v_j) \in E} [DS_i DS_j]^{-\frac{1}{2}}, \quad (3)$$

```
> balabanJ(g)
```

```
[1] 2.414364
```

```
> balabanJ(g, mat.dist)
```

```
[1] 2.414364
```

where $|E(G)| := |E|$ denotes the number of edges of the complex network, DS_i denotes the distance sum (row sum) of v_i and $\mu := |E| + 1 - |N|$ denotes the cyclomatic number.

Mean distance deviation [5]:

$$\Delta\mu(G) := \frac{1}{N} \sum_{i=1}^N |\mu(v_i) - \bar{\mu}|, \quad (4)$$

where

$$\mu(v_i) := \sum_{j=1}^N d(v_i, v_j), \quad (5)$$

and

$$\bar{\mu} := \frac{2W}{N}. \quad (6)$$

```
> meanDistanceDeviation(g)
```

```
[1] 1.6875
```

```
> meanDistanceDeviation(g, mat.dist)
```

```
[1] 1.6875
```

Compactness [6]:

$$C(G) := \frac{4W}{|N|(|N| - 1)}. \quad (7)$$

```
> compactness(g)
```

```
[1] 3.071429
```

```
> compactness(g, mat.dist)
```

```
[1] 3.071429
```

```
> compactness(g, mat.dist, wiener(g, mat.dist))
```

```
[1] 3.071429
```

Product of Row Sums index [7]:

$$\text{PRS}(G) = \prod_{i=1}^{|N|} \mu(v_i) \quad \text{or} \quad \log(\text{PRS}(G)) = \log \left(\prod_{i=1}^{|N|} \mu(v_i) \right). \quad (8)$$

```
> productOfRowSums(g, log = FALSE)
```

```
[1] 157464000
```

```
> productOfRowSums(g, log = TRUE)
```

```
[1] 27.23045
> productOfRowSums(g, mat.dist, log = FALSE)
[1] 157464000
> productOfRowSums(g, mat.dist, log = TRUE)
[1] 27.23045
```

Hyper-distance-path index [8]

$$D_P(G) := \frac{1}{2} \sum_{i=1}^{|N|} \sum_{j=1}^{|N|} d(v_i, v_j) + \frac{1}{2} \sum_{i=1}^{|N|} \sum_{j=1}^{|N|} \binom{d(v_i, v_j)}{2}. \quad (9)$$

```
> hyperDistancePathIndex(g)
[1] 60
> hyperDistancePathIndex(g, mat.dist)
[1] 60
> hyperDistancePathIndex(g, mat.dist, wiener(g, mat.dist))
[1] 60
```

3.2 Descriptors Based on Other Graph-Invariants

This section describes network measures based on other invariants than distances.

Index of total adjacency [9]:

$$A(G) := \frac{1}{2} \sum_{i=1}^{|N|} \sum_{j=1}^{|N|} a_{ij}. \quad (10)$$

```
> totalAdjacency(g)
[1] 17
> totalAdjacency(g, mat.adj)
[1] 17
```

Zagreb group indices [10]:

$$Z_1(G) := \sum_{i=1}^{|N|} k_{v_i}, \quad (11)$$

where k_{v_i} is the degree of the node v_i .

$$Z_2(G) := \sum_{(v_i, v_j) \in E} k_{v_i} k_{v_j}. \quad (12)$$

```
> zagreb1(g)
[1] 32
> zagreb1(g, vec.degree)
[1] 32
> zagreb2(g)
[1] 298
> zagreb2(g, vec.degree)
[1] 298
```

Randić connectivity index [11]:

$$R(G) := \sum_{(v_i, v_j) \in E} [k_{v_i} k_{v_j}]^{-\frac{1}{2}}. \quad (13)$$

```
> randic(g)
[1] 3.602215
> randic(g, vec.degree)
[1] 3.602215
```

The complexity index B [9]:

$$B(G) := \sum_{i=1}^{|N|} \frac{k_{v_i}}{\mu(v_i)}. \quad (14)$$

```
> complexityIndexB(g)
[1] 3.255556
> complexityIndexB(g, mat.dist)
[1] 3.255556
> complexityIndexB(g, mat.dist, vec.degree)
[1] 3.255556
```

Normalized edge complexity [9]:

$$E_N(G) := \frac{A(G)}{|N|^2}. \quad (15)$$

```
> normalizedEdgeComplexity(g)
[1] 0.265625
> normalizedEdgeComplexity(g, totalAdjacency(g, mat.adj))
[1] 0.265625
```

3.3 Classical entropy based descriptors

These measures are based on grouping the elements of an arbitrary graph invariant (vertices, edges, and distances etc.) using an equivalence criterion.

Topological information content [12, 13]:

$$I_{orb}^V(G) := - \sum_{i=1}^k \frac{|N_i^V|}{|N|} \log \left(\frac{|N_i^V|}{|N|} \right). \quad (16)$$

$|N_i^V|$ denotes the number of vertices belonging to the i -th vertex orbit.

```
> topologicalInfoContent(g)
[1] 2.25
> topologicalInfoContent(g, mat.dist)
[1] 2.25
> topologicalInfoContent(g, mat.dist, vec.degree)
[1] 2.25
```

Bonchev - Trinajstić indices [14]:

$$I_D(G) := -\frac{1}{|N|} \log \left(\frac{1}{|N|} \right) - \sum_{i=1}^{\rho(G)} \frac{2k_i}{|N|^2} \log \left(\frac{2k_i}{|N|^2} \right), \quad (17)$$

$$I_D^W(G) := W(G) \log(W(G)) - \sum_{i=1}^{\rho(G)} ik_i \log(i). \quad (18)$$

k_i is the occurrence of a distance possessing value i in the distance matrix of G .

```
> #I_D(G)
> bonchev1(g)

[1] 1.208931

> bonchev1(g, mat.dist)

[1] 1.208931

> #I^W_D(G)
> bonchev2(g)

[1] 170.3098

> bonchev2(g, mat.dist)

[1] 170.3098

> bonchev2(g, mat.dist, wiener(g))

[1] 170.3098
```

BERTZ complexity index [15]:

$$C(G) := 2N \log(|N|) - \sum_{i=1}^k |N_i| \log(|N_i|). \quad (19)$$

$|N_i|$ are the cardinalities of the vertex orbits as defined in Eqn. (16).

```
> bertz(g)

[1] 42

> bertz(g, mat.dist)

[1] 42

> bertz(g, mat.dist, vec.degree)

[1] 42
```

Radial centric information index [16]:

$$I_{C,R}(G) := \sum_{i=1}^k \frac{|N_i^e|}{|N|} \log \left(\frac{|N_i^e|}{|N|} \right). \quad (20)$$

$|N_i^e|$ is the number of vertices having the same eccentricity.

```
> radialCentric(g)

[1] 0.954434

> radialCentric(g, mat.dist)

[1] 0.954434
```

Vertex degree equality-based information index [16]:

$$I_{deg}(G) := \sum_{i=1}^{\bar{k}} \frac{|N_i^{k_v}|}{|N|} \log \left(\frac{|N_i^{k_v}|}{|N|} \right). \quad (21)$$

$|N_i^{k_v}|$ is the number of vertices with degree equal to i and $\bar{k} := \max_{v \in V} k_v$.

> `vertexDegree(g)`

[1] 2.25

> `vertexDegree(g, vec.degree)`

[1] 2.25

Balaban-like information indices [17]:

$$U(G) := \frac{|E|}{\mu + 1} \sum_{(v_i, v_j) \in E} [u(v_i)u(v_j)]^{-\frac{1}{2}}, \quad (22)$$

$$X(G) := \frac{|E|}{\mu + 1} \sum_{(v_i, v_j) \in E} [x(v_i)x(v_j)]^{-\frac{1}{2}}, \quad (23)$$

where

$$u(v_i) := - \sum_{j=1}^{\sigma(v_i)} \frac{j|S_j(v_i, G)|}{\mu(v_i)} \log \left(\frac{j}{\mu(v_i)} \right), \quad (24)$$

$$x(v_i) := -\mu(v_i) \log(d(v_i)) - y_i, \quad (25)$$

$$y_i := \sum_{j=1}^{\sigma(v_i)} j|S_j(v_i, G)| \log(j), \quad (26)$$

$$\mu(v_i) := \sum_{j=1}^{|N|} d(v_i, v_j) = \sum_{j=1}^{|N|} j|S_j(v_i, G)|. \quad (27)$$

> `#Balaban-like information index U(G)`

> `balabanlike1(g)`

[1] 8.831362

> `balabanlike1(g, mat.dist)`

[1] 8.831362

> `#Balaban-like information index X(G)`

> `balabanlike2(g)`

[1] 0.8436946

> `balabanlike2(g, mat.dist)`

[1] 0.8436946

Graph vertex complexity index [18]:

$$I_V(G) := \sum_{i=1}^N v_i^c, \quad (28)$$

where v_i^c is the so-called vertex complexity expressed by

$$v_i^c := \sum_{j=0}^{\sigma(v_i)} \frac{k_j^{v_i}}{N} \log \left(\frac{k_j^{v_i}}{N} \right). \quad (29)$$

$k_k^{v_i}$ is the number of distances starting from $V_i \in V$ equal to j .

```
> graphVertexComplexity(g)
```

```
[1] -12.08022
```

```
> graphVertexComplexity(g, mat.dist)
```

```
[1] -12.08022
```

3.4 Parametric Graph Entropy Measures

Measures of this group [19, 20] assign a probability value to each vertex of the network using a so-called information functional f which captures structural information of the network G . We yield [19],

$$I_f(G) := - \sum_{i=1}^{|N|} \frac{f(v_i)}{\sum_{j=1}^{|N|} f(v_j)} \log \left(\frac{f(v_i)}{\sum_{j=1}^{|N|} f(v_j)} \right), \quad (30)$$

where $I_f(G)$ represents a family of graph entropy [19] measures depending on the information functional. Further we implemented the following measurement[20]:

$$I_f^\lambda(G) := \lambda \left(\log(|N|) + \sum_{i=1}^{|N|} p(v_i) \log(p(v_i)) \right), \quad (31)$$

$$p(v_i) := \frac{f(v_i)}{\sum_{j=1}^{|N|} f(v_j)}, \quad (32)$$

where $p^V(v_i)$ are the vertex probabilities, $\lambda > 0$ a scaling constant. This measure can be interpreted as the distance between the entropy defined in equation 30 and maximum entropy ($\log(|N|)$).

We integrated 3 different information functionals:

1. An information functional using the j-spheres ("sphere"):

$$f^V(v_i) := c_1 |S_1(v_i, G)| + c_2 |S_2(v_i, G)| + \dots + c_{\rho(G)} |S_{\rho(G)}(v_i, G)|, \quad (33)$$

where $c_k > 0$.

2. An information functional using path lengths ("pathlength"):

$$f^{P_2}(v_i) := c_1 l(P(L_G(v_i, 1))) + c_2 l(P(L_G(v_i, 2))) + \dots + c_{\rho(G)} l(P(L_G(v_i, \rho(G)))), \quad (34)$$

where $c_k > 0$.

3. An information functional using vertex centrality("vertcent") :

$$f^{C_2}(v_i) := c_1 \beta^{L_G(v_i, 1)}(v_i) + c_2 \beta^{L_G(v_i, 2)}(v_i) + \dots + c_{\rho(G)} \beta^{L_G(v_i, \rho(G))}(v_i), \quad (35)$$

where $c_k > 0$.

We implemented 4 different settings (as example settings) of the weighting parameter c_i :

1. constant

$$c_1 := 1, c_2 := 1, \dots, c_{\rho(G)} := 1. \quad (36)$$

2. linear

$$c_1 := \rho(G), c_2 := \rho(G) - 1, \dots, c_{\rho(G)} := 1. \quad (37)$$

3. quadratic

$$c_1 := \rho(G)^2, c_2 := (\rho(G) - 1)^2, \dots, c_{\rho(G)} := 1. \quad (38)$$

4. exponential

$$c_1 := \rho(G), c_2 := \rho(G)e^{-1}, \dots, c_{\rho(G)} := \rho(G)e^{-\rho(G)+1}. \quad (39)$$

$\rho(G)$ represents the diameter of the network.

To call this type of network measure we provide the method *infoTheoreticGCM*. It has following input parameters:

- *g*: the network as a graphNEL object - it is the only mandatory parameter
- *dist*: the distance matrix of *g*
- *coeff*: specifies the weighting parameter: "const", "lin", "quad", "exp", "const" or "cust" are available constants. If it is set to "cust" you have to specify your customized weighting schema with the parameter *custCoeff*.
- *infofunct*: specifies the information functional: "sphere", "pathlength" or "vertcent" are available settings.
- *lamda*: scaling constant for the distance, default set to 1000.
- *custCoeff*: specifies the customized weighting schema. To use it you need to set *coeff*="const".

The method returns a list with following entries:

- *entropy*: contains the entropy, see formula 30
- *distance*: contains the distance described in formula 31
- *pis*: contains the probability distribution, see formula 32
- *fvi*: contains the values of the used information functional for each vertex v_i

```
> l1 <- infoTheoreticGCM(g)
> l2 <- infoTheoreticGCM(g, mat.dist, coeff = "lin", infofunct = "sphere",
+   lamda = 1000)
> l3 <- infoTheoreticGCM(g, mat.dist, coeff = "exp", infofunct = "sphere",
+   lamda = 1000)
> l4 <- infoTheoreticGCM(g, mat.dist, coeff = "const", infofunct = "pathlength",
+   lamda = 4000)
> l5 <- infoTheoreticGCM(g, mat.dist, coeff = "quad", infofunct = "vertcent",
+   lamda = 1000)
> l1

$entropy
[1] 2.990011

$distance
[1] 9.9892

$pis
      1      2      3      4      5      6      7      8
0.1376812 0.1304348 0.1304348 0.1376812 0.1376812 0.0942029 0.1159420 0.1159420

$fvis
  1  2  3  4  5  6  7  8
19 18 18 19 19 13 16 16

> l5

$entropy
[1] 2.924897

$distance
[1] 75.10322

$pis
      1      2      3      4      5      6      7
0.1666667 0.12851406 0.12851406 0.1666667 0.15160643 0.04518072 0.10642570
```

```

      8
0.10642570

$fviz
      1      2      3      4      5      6      7      8
55.33333 42.66667 42.66667 55.33333 50.33333 15.00000 35.33333 35.33333

```

4 Session Info

```

> sessionInfo()

R version 2.11.0 (2010-04-22)
x86_64-unknown-linux-gnu

locale:
 [1] LC_CTYPE=en_US.utf8      LC_NUMERIC=C
 [3] LC_TIME=en_US.utf8      LC_COLLATE=en_US.utf8
 [5] LC_MONETARY=C           LC_MESSAGES=en_US.utf8
 [7] LC_PAPER=en_US.utf8     LC_NAME=C
 [9] LC_ADDRESS=C           LC_TELEPHONE=C
[11] LC_MEASUREMENT=en_US.utf8 LC_IDENTIFICATION=C

attached base packages:
[1] grid      stats      graphics  grDevices  utils      datasets  methods
[8] base

other attached packages:
[1] QuACN_1.0      combinat_0.0-7  igraph_0.5.3    Rgraphviz_1.26.0
[5] RBGL_1.24.0    graph_1.26.0

loaded via a namespace (and not attached):
[1] tools_2.11.0

```

References

- [1] M. Hucka, A. Finney, H. M. Sauro, H. Bolouri, J. C. Doyle, H. Kitano, A. P. Arkin, B. J. Bornstein, D. Bray, A. Cornish-Bowden, A. A. Cuellar, S. Dronov, E. D. Gilles, M. Ginkel, V. Gor, I. I. Goryanin, W. J. Hedley, T. C. Hodgman, J.-H. Hofmeyr, P. J. Hunter, N. S. Juty, J. L. Kasberger, A. Kremling, U. Kummer, N. L. Novák, L. M. Loew, D. Lucio, P. Mendes, E. Minch, E. D. Mjolsness, Y. Nakayama, M. R. Nelson, P. F. Nielsen, T. Sakurada, J. C. Schaff, B. E. Shapiro, T. S. Shimizu, H. D. Spence, J. Stelling, K. Takahashi, M. Tomita, J. Wagner, J. Wang, and S. B. M. L. Forum, "The systems biology markup language (sbml): a medium for representation and exchange of biochemical network models." *Bioinformatics*, vol. 19, no. 4, pp. 524–531, Mar 2003.
- [2] H. Wiener, "Structural determination of paraffin boiling points," *Journal of the American Chemical Society*, vol. 69, no. 1, pp. 17–20, Jan. 1947. [Online]. Available: <http://dx.doi.org/10.1021/ja01193a005>
- [3] A. T. Balaban and O. Ivanciuc, "Historical development of topological indices," in *Topological Indices and Related Descriptors in QSAR and QSPAR*, J. Devillers and A. T. Balaban, Eds. Gordon and Breach Science Publishers, 1999, pp. 21–57, amsterdam, The Netherlands.
- [4] A. T. Balaban, "Highly discriminating distance-based topological index," *Chem. Phys. Lett.*, vol. 89, pp. 399–404, 1982.
- [5] V. A. Skorobogatov and A. A. Dobrynin, "Metrical analysis of graphs," *Commun. Math. Comp. Chem.*, vol. 23, pp. 105–155, 1988.
- [6] J. K. Doyle and J. E. Garver, "Mean distance in a graph," *Discrete Mathematics*, vol. 17, pp. 147–154, 1977.

- [7] H. P. Schultz, E. B. Schultz, and T. P. Schultz, "Topological organic chemistry. 4. graph theory, matrix permanents, and topological indices of alkanes," *Journal of Chemical Information and Computer Sciences*, vol. 32, no. 1, pp. 69–72, 1992.
- [8] R. Todeschini, V. Consonni, and R. Mannhold, *Handbook of Molecular Descriptors*. Wiley-VCH, 2002, weinheim, Germany.
- [9] D. Bonchev and D. H. Rouvray, *Complexity in Chemistry, Biology, and Ecology*, ser. Mathematical and Computational Chemistry. Springer, 2005, New York, NY, USA.
- [10] M. V. Diudea, I. Gutman, and L. Jäntschi, *Molecular Topology*. Nova Publishing, 2001, new York, NY, USA.
- [11] X. Li and I. Gutman, *Mathematical Aspects of Randić-Type Molecular Structure Descriptors*, ser. Mathematical Chemistry Monographs. University of Kragujevac and Faculty of Science Kragujevac, 2006.
- [12] A. Mowshowitz, "Entropy and the complexity of the graphs I: An index of the relative complexity of a graph," *Bull. Math. Biophys.*, vol. 30, pp. 175–204, 1968.
- [13] N. Rashevsky, "Life, information theory, and topology," *Bull. Math. Biophys.*, vol. 17, pp. 229–235, 1955.
- [14] D. Bonchev and N. Trinajstić, "Information theory, distance matrix and molecular branching," *J. Chem. Phys.*, vol. 67, pp. 4517–4533, 1977.
- [15] S. H. Bertz, "The first general index of molecular complexity," *Journal of the American Chemical Society*, vol. 103, pp. 3241–3243, 1981.
- [16] D. Bonchev, *Information Theoretic Indices for Characterization of Chemical Structures*. Research Studies Press, Chichester, 1983.
- [17] A. T. Balaban and T. S. Balaban, "New vertex invariants and topological indices of chemical graphs based on information on distances," *J. Math. Chem.*, vol. 8, pp. 383–397, 1991.
- [18] C. Raychaudhury, S. K. Ray, J. J. Ghosh, A. B. Roy, and S. C. Basak, "Discrimination of isomeric structures using information theoretic topological indices," *Journal of Computational Chemistry*, vol. 5, pp. 581–588, 1984.
- [19] M. Dehmer, "Information processing in complex networks: Graph entropy and information functionals," *Applied Mathematics and Computation*, vol. 201, pp. 82–94, 2008.
- [20] M. Dehmer, K. Varmuza, S. Borgert, and F. Emmert-Streib, "On entropy-based molecular descriptors: Statistical analysis of real and synthetic chemical structures," *J. Chem. Inf. Model.*, vol. 49, pp. 1655–1663, 2009.