

# Determinants of cross-regional R&D collaboration networks: an application of exponential random graph models

# 04.13

Tom Broekel and Matté Hartog

## **Impressum:**

Working Papers on Innovation and Space  
Philipps-Universität Marburg

Herausgeber:

Prof. Dr. Dr. Thomas Brenner  
Deutschhausstraße 10  
35032 Marburg  
E-Mail: [thomas.brenner@staff.uni-marburg.de](mailto:thomas.brenner@staff.uni-marburg.de)

Erschienen: 2013

# Determinants of cross-regional R&D collaboration networks: an application of exponential random graph models

Tom Broekel<sup>1</sup>

Institute of Economic and Cultural Geography, Leibniz University, Hannover.

Matté Hartog

Section of Economic Geography, Faculty of Geosciences, Utrecht University, Utrecht.

## Abstract:

This study investigates the usefulness of exponential random graph models (ERGM) to analyze the determinants of cross-regional R&D collaboration networks. Using spatial interaction models, most research on R&D collaboration between regions is constrained to focus on determinants at the node level (e.g. R&D activity of a region) and dyad level (e.g. geographical distance between regions). ERGMs represent a new set of network analysis techniques that have been developed in recent years in mathematical sociology. In contrast to spatial interaction models, ERGMs additionally allow considering determinants at the structural network level while still only requiring cross-sectional network data.

The usefulness of ERGMs is illustrated by an empirical study on the structure of the cross-regional R&D collaboration network of the German chemical industry. The empirical results confirm the importance of determinants at all three levels. It is shown that in addition to determinants at the node and dyad level, the structural network level determinant “triadic closure” helps in explaining the structure of the network. That is, regions that are indirectly linked to each other are more likely to be directly linked as well.

**Keywords:** cross-regional R&D collaboration, exponential random graph models, network analysis, chemical industry, Germany.

**JEL Classifications:** R11, O32, D85

---

<sup>1</sup> Corresponding Author: Tom Broekel, Institute of Economic and Cultural Geography, Leibniz University of Hannover, Schneiderberg 50, 30167 Hannover, Germany. E-Mail: broekel@wigeo.uni-hannover.de.

# 1 Introduction

There is growing scientific interest in the creation of knowledge and its diffusion among organizations. In the new growth theory, new knowledge is regarded as being pivotal to economic growth by generating increasing returns (Romer, 1990). In evolutionary economics, the re-combination of existing knowledge from different sources is argued to be crucial for new innovations to occur (Nelson and Winter, 1982). These theories and the according empirical evidence also impacted the policy level. For instances, one of the most well known policy instruments to stimulate knowledge diffusion and innovation are the Framework Programmes of the European Union. These programs have been in existence since 1984 and are used to fund thousands of collaborative research projects between organizations in the EU.

Such R&D collaboration networks, which are induced by policy, alter the spatial diffusion of knowledge. This put the investigation of their spatial structures on the agenda of regional economists and economic geographers (Autant-Bernard et al., 2007). While the geographical structures of inter-organizational collaboration networks are frequently analyzed from an organizational perspective (cf. Giuliani and Bell, 2005), these researchers rather employ a regional perspective and focus on cross-regional R&D collaboration networks (cf. Scherngell and Barber, 2009; 2011; Hoekman, et al., 2010). In order to investigate factors explaining the structure of cross-regional networks, most commonly spatial interaction models are used, which allow for considering factors at the node and dyad level. An example of a factor at the node level is the size of a region that matters as regions with more organizations are also more likely to have links to regions elsewhere. At the dyad level, most attention has been paid to the effect of increasing geographical distance that decreases the chances of research collaboration (cf. Ponds et al., 2007; Scherngell and Barber, 2009; Hoekman et al., 2009; 2010).

In addition to the node and dyad level, factors at the structural network level may also be important, though. That is, the creation of new links might not only depend on attributes of regions or region pairs. It may also be influenced by the existing structure of the cross-regional network. For instance, a key hypothesis in organizational network science is the tendency towards triadic closure (or transitivity), which implies in this context that regions, which are indirectly linked, are more likely to link themselves as well. However, factors at the structural network level cannot be included in spatial interaction models.

The paper presents exponential random graph models (ERGM) as an alternative empirical tool. These models have been developed in mathematical sociology in recent years (Snijders et al., 2006; Robins et al., 2006, 2007; Wang et al., 2012) and are increasingly used across scientific disciplines, for example in bioscience (Saul and Filkov, 2007), political science (Desmarais and Cranmer, 2012) and organization science (Uddin et al., 2012). The advantage of these models is that they allow for simultaneously estimating the effect of factors at the node, dyad, and structural network level for networks that are observed at one point in time.

We illustrate the usefulness of ERGMs by exemplarily investigating the structure and its determinants of the cross-regional R&D collaboration network in the German chemical industry between 2005 and 2010.

The study is structured as follows. The second section gives an overview of the literature on spatial structures of R&D collaboration networks and their determinants. This includes a brief discussion of factors at the node, dyad, and structural network level that may impact network structures. The third section elaborates on the exponential random graph model that we subsequently use to investigate the structure of the cross-regional network. We

present the empirical data in the fourth section. It is followed by the discussion of the results in the fifth section and some concluding remarks in the sixth section.

## 2 Determinants of cross-regional R&D collaboration

The structural determinants of cross-regional R&D collaboration networks can be distinguished at three different levels. These are the node level, the dyad level, and the structural network level. In the following, we elaborate about factors that become effective at these three different levels.

Node level factors are properties of network entities themselves. With respect to regional R&D collaboration networks, regions' size and their research intensity are particularly important. Firstly, large organizations are likely to have more ties than small organizations because their position in the industry is more prominent. They also have more resources at their disposal to create and maintain ties. For instance, Boschma and Ter Wal (2007) find that larger organizations are more central in the knowledge network of footwear producers in Barletta. At the regional level in general, regions with more organizations can be expected to have more ties because they have more collaboration opportunities. Secondly, research intensity may matter. Giuliani and Bell (2005) show that organizations with a more advanced knowledge base are more often approached by other organizations to exchange knowledge because they are perceived to be 'technological leaders'. The research intensity of a region is generally characterized by large numbers of R&D employees, many organizations being engaged in R&D-intensive activities, and by the presence of universities or other research institutes. All these characteristics are likely to increase the number of links a region has to other regions, i.e. a region's (degree) centrality in the cross-regional collaboration network.

Factors at the dyad level are characteristics of relationships between two entities (nodes) in a network. In the context of the paper it refers to the relation between two regions. A key idea in sociology is that entities are more likely to link when they have similar attributes, known as homophily effect (McPherson et al., 2001). For instance, regions with organizations that operate with similar routines and under comparable incentive mechanisms are more likely to be linked in R&D collaboration. Another example are universities, which are subject to different incentive frameworks than firms when it comes to knowledge creation and diffusion as they aim to publish new knowledge, whereas firms have an incentive to keep new knowledge secret. Hence, because of their institutional proximity (Metcalf, 1995), universities are more likely to collaborate with others and especially with other universities (cf. Broekel and Boschma 2012; Broekel and Hartog, 2013). This is likely to translate to the regional level as regions rarely house more than one university. Accordingly, university regions have a higher likelihood of being linked to each other.

In addition to institutional proximity, other forms of proximity may also be particularly relevant: geographical proximity, technological proximity, and social proximity. Many studies confirm that cross-regional R&D collaboration is more likely when regions are located close to one another in space (e.g. Magionni et al., 2007; Scherngell and Barber, 2009; Hoekman et al., 2009, 2010). This may be due to a variety of reasons, for instance geographical proximity facilitates face-to-face contact, which stimulates the diffusion of information about potential collaboration partners. The likelihood of cross-regional R&D collaboration is shown to increase when regions have similar technological profiles and specializations (Fischer et al., 2006; LeSage et al., 2007; Scherngell and Barber, 2009). A potential explanation is that organizations are more prone to collaborate with organizations with related knowledge assets. Similar technological profiles (technological proximity) ensure

that two organizations can easily communicate and learn from each other (Cohen and Levinthal, 1990; Nooteboom, 2000). Social proximity may also increase the likelihood of R&D collaboration (cf. Autant-Bernard et al., 2007). People already knowing each other find it easier to develop trust-based relations, which in turn facilitate knowledge exchange and ease interactions across regional boundaries (Maskell and Malmberg, 1999; Sobrero and Roberts, 2001; Breschi and Lissoni, 2009).

In addition to these factors at the node and dyad level, factors at the structural network level may also matter for the structure of cross-regional R&D collaboration. Such factors relate to properties of the entire network. Three factors are commonly put forward in this context: triadic closure (transitivity), multi-connectivity, and preferential attachment (cf. Glückler, 2010; Ter Wal and Boschma, 2009).

*Triadic closure* predicts that partners of organizations are likely to become partners themselves as well. As a result, a network will consist of many triangles, i.e. dense cliques of strongly interconnected organizations (Ter Wal, 2011). Such cliques can be regarded as a sign of social capital (Coleman, 1988) that may enhance trust and willingness among actors to invest in mutual goals, such as research collaboration. In contrast, *multi-connectivity* suggests that organizations will connect to others in multiple ways to decrease the dependency on a single link. It implies that multiple paths are formed amongst organizations leading to multiple reachability. Evidence for this is found in the creation of inter-firm alliances between US biotech firms (Powell et al., 2005). *Preferential attachment* means that organizations with many links are more likely to create or attract new links in the future. If a network is shaped by this factor, its degree distribution follows a power law (Barabasi and Albert, 1999). Gulati (1999) shows that in the case of multinational firms, the likelihood of creating new alliances increases the better organizations are connected in the network. Hence, the network of alliances among multinational firms is subject to preferential attachment processes.

If these processes involve organizations located in different regions their effects will naturally be translated to the regional level. This implies that multi-connectivity, preferential attachment, and triadic closure may also characterize cross-regional networks. Consequently, these factors at the structural network level need to be considered when analyzing such networks' structures.

To estimate the relative impact of the above factors on the structure of a network, they need to be simultaneously incorporated in the empirical model. This is not possible with the models most frequently used to investigate cross-regional collaboration: spatial interaction models in general and gravity models in particular (cf. Scherngell and Barber, 2009). These models can account for factors at the node and dyad level. However, they cannot be used to evaluate factors at the structural network level. In light of the theoretical relevance of factors at the structural network level, we therefore argue that network analysis modeling techniques represent a powerful alternative because they are able to simultaneously incorporate factors at all three levels.

When longitudinal data is available, a stochastic actor-based network approach can be used. It models the change of a network from one point in time to another as part of an iterative Markov chain process (see for technical details: Snijders et al., 2010). When it comes to the analysis of research collaboration networks of regions, however, such an approach is less useful. By aggregating collaboration data to the regional level and creating cross-regional networks, researchers generally are interested in approximating the relational interaction structures of regions and investigate their structures and determinants. Such networks are unlikely to drastically change within short time periods, though, as they are results of long-term social, regional, and industrial evolution processes. Hence, even when longitudinal data on these cross-regional networks structures are available, it is unlikely to cover a sufficiently long time period. It may include multiple time periods (years) and thereby principally allow

for employing longitudinal network analysis to study changes in the underlying cross-regional interaction structures.<sup>1</sup> However, the results generated with stochastic actor-based network approaches are unlikely to yield meaningful insights because the empirically observed changes in the network structures are dominated by short-term fluctuations that are of little interest to the researcher.

We therefore argue that exponential random graph models are the preferred option when investigating the structure of cross-regional interaction on the basis of data with a cross-sectional nature and factors at the structural network are to be considered. We elaborate on these models in the next section.

### 3 Exponential random graph models

Exponential random graph models are stochastic models that approach link creation as a time-continuous process. They regard a network observed at one point in time as one particular realization out of a set of multiple hypothetical networks with similar properties. This allows applying these models to purely cross-sectional network data.

The aim of exponential random graph models is to identify factors that maximize the probability of the emergence of a network with similar properties as the structure of the observed network. The general form of exponential random graph models is defined as follows (Robins et al., 2007):

$$\Pr(Y = y) = \left( \frac{1}{\kappa} \right) \exp \left\{ \sum_A \eta_A g_A(y) \right\} \quad (\text{eq. 1})$$

where  $\Pr(Y=y)$  is the probability that the network ( $Y$ ) generated by an exponential random graph is identical to the observed network ( $y$ ),  $\kappa$  is a normalizing constant to ensure that the equation is a proper probability distribution (summing up to 1),  $\eta_A$  is the parameter corresponding to network configuration  $A$ , and  $g_A(y)$  represents the network statistic. Network configurations can be factors at the node level, dyad level, and structural network level.

Estimation of the parameters is done with maximum pseudo likelihood or a Markov Chain Monte Carlo Maximum Likelihood Estimation procedure. The latter has been developed most recently and is regarded as the preferred procedure as it is often most accurate (Snijders, 2002; Van Duin et al., 2009). It is based on the generation of a distribution of random graphs by stochastic simulation from a starting set of parameter values, and subsequent refinement of those parameter values by comparing the obtained random graphs against the observed graph. This process is repeated until the parameter estimates stabilize. If they do not, the model might fail to converge and hence becomes unstable (see for technical details, Handcock 2003, and Hunter et al. 2008).

Checking whether the parameters predict the observed network well, i.e. evaluating a model's goodness of fit, is done by comparing the structure of the simulated networks to the structure of the observed network. According to Hunter et al. (2008), the comparison consists of the degree distribution, the distribution of edgewise shared partners (the number of links in which two organizations have exactly  $k$  partners in common, for each value of  $k$ ), and the

---

<sup>1</sup> The relational data derived from the 5<sup>th</sup>, 6<sup>th</sup>, and 7<sup>th</sup> EU-Framework Programmes are (currently) a good example in this respect. While they represent longitudinal data, it covers only a limited time-period (1998-2013). Of course, this may change when data on future programs will become available.

geodesic distribution (the number of pairs for which the shortest path between them is of length  $k$ , for each value of  $k$ ). The more the distributions of the simulated networks are in line with those of the observed network, the more accurate are the parameters of the exponential random graph model. In the next section, we construct an exponential random graph model to investigate the structure of the network of subsidized R&D collaboration in the German chemical industry.

## 4 Determinants of cross-regional R&D collaboration in the German chemical industry

### 4.1 Data

We analyze R&D collaboration that has been funded by the German federal government. As in most other advanced countries, the government actively supports public and private R&D activities with subsidies. While the Federal Ministry of Education and Research (BMBF) is the prime source of subsidies, the Federal Ministry of Economics and Technology (BMWi) and the Federal Ministry for the Environment, Nature Conservation and Nuclear Safety (BMU) contribute as well. The federal ministries publish comprehensive information about subsidized projects in the so-called “Förderkatalog” (subsidies catalog). This catalog contains detailed information on more than 150,000 individual subsidies that have been granted between 1960 and 2012. The catalog also includes information on the cooperative nature of projects. It specifically indicates if projects are joint projects realized by consortia of organizations. Participants in joint projects agree to a number of regulations that guarantee significant knowledge exchange between the partners. Accordingly, two organizations are defined to cooperate if they participate in the same joint project. Hence, the original network is a two-mode network (project-organizations links), which we transform into a one-mode projection of the network (organization-organization links). All organizations can be assigned to labor market regions allowing for regionalizing the network (see for more details on the data Broekel and Graf, 2012). The data is comparable to the EU Framework Programmes (EU-FP) data by and large, which is extensively used to model research collaboration networks (cf. Scherngell and Barber, 2009). In contrast to the EU-FP data, the data at hand exclusively covers collaboration between German organizations.

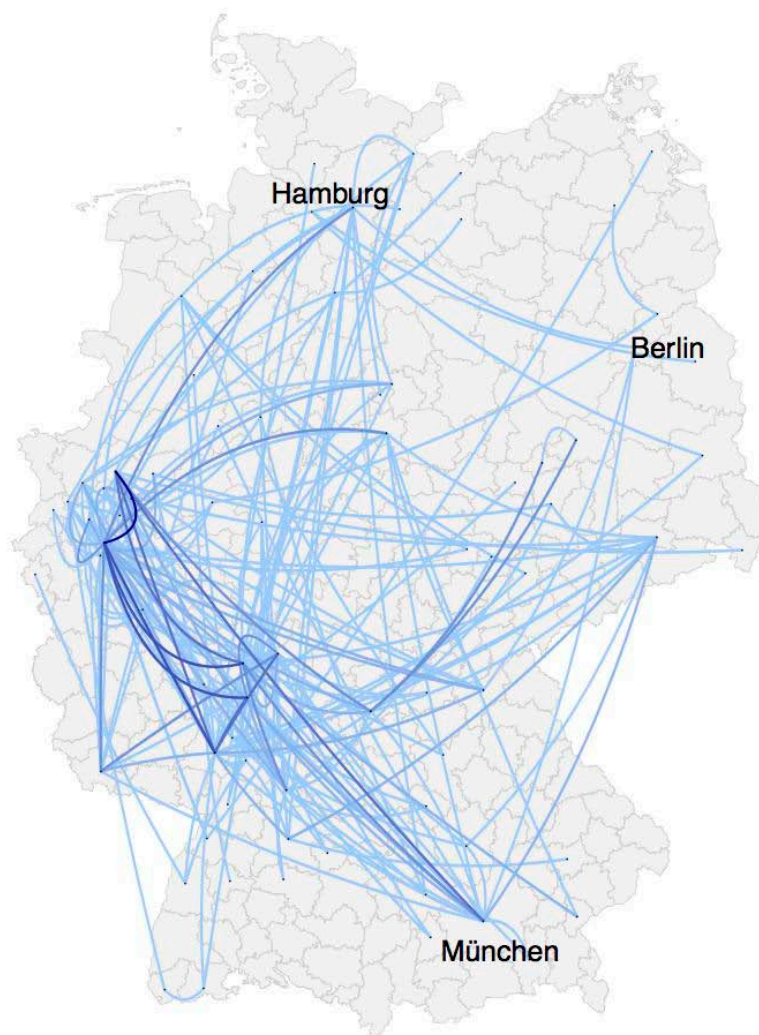
To construct the network of subsidized R&D collaboration in the German chemical industry, we first identify all firms in the data that are classified as being involved in the 2-digit NACE code C20 ‘Manufacture of chemicals and chemical products’. Subsequently, all cooperative projects are extracted in which at least one of these firms participates. On the basis of the joint appearance in a project, we construct the inter-organizational network among all chemical firms participating in these projects. We only consider links among firms: links to universities, research organizations, associations, and to firms belonging to other industries are excluded. We believe that this approach provides the most conservative picture of the (subsidized) R&D collaboration network in the chemical industry.<sup>2</sup> The corresponding inter-organizational undirected network is subsequently aggregated to the regional level using information on organizations’ location in the 270 German labor market regions. The 270 labor market regions are defined by the German Institute for Labor and Employment (e.g. Greif and Schmiedl, 2002). We construct the network that existed between 1 January 2005 and 31

---

<sup>2</sup> Alternatively one may consider all organizations active in joint projects in which at least one firm of the chemical industry is participating. However, such seems to be a very broad definition of an industry-specific network, which makes the definition of appropriate empirical variables more difficult.

December 2010. In this period, 775 projects were subsidized in which at least one firm of the chemical industry was involved. These projects are split into 975 individual funds allocated to 557 German firms belonging to the chemical industry.<sup>3</sup> 133 of the 775 projects are joint projects, which involve on average 2.8 firms. The resulting cross-regional R&D collaboration network is shown in Figure 1.

Figure 1: Network of subsidized R&D collaboration among firms in the German chemical industry (2005-2010)



The network is dichotomized, as we are only interested in whether or not a link exists between regions. The figure shows that the large agglomerations of the Ruhr Area, Frankfurt am Main, and Munich are important nodes in the network. In addition, a number of central regions are located along the Rhine River in the west. The region of Dresden is a central node in East Germany. All these regions are well-known centers of the chemical industry in

---

<sup>3</sup> This figure is based on the number of executing organizations („Ausführende Stelle“) as given in the data. Many of these organizations are part of larger organizations. This has however little relevance for the results as all data are aggregated to the regional level.

Germany. Some additional descriptive statistics of the network are presented in Table 1 in the Appendix.

## 4.2 Construction of empirical variables

### 4.2.1 Node level variables

The most important node-level factors are probably the intensity of regional R&D and innovation activities in the field of chemistry. Foremost, this is because undertaking R&D activities is necessary to receive R&D subsidies. Regions with large R&D activities are likely to host more organizations that are involved in R&D collaboration. Moreover, such regions may also be the location of the most successful innovators, which are preferred collaboration partners. We therefore consider the number of applied patents in chemistry by regional organizations as proxy for the intensity of regional R&D activities in this field. The regionalized data on patent applications are published in Greif and Schmiedl (2002) and Greif et al. (2006), which include applications to the German as well as to the European Patent Office, with a correction for double counts. The patents are assigned to labor market regions according to the inventor principle. The patent data is organized according to IPC-classes, which is matched to the 2-digit NACE industry using the concordance of Broekel (2007). Lacking the data for the years 2005-2010, we construct the first node-level variable as the summed number of patents of regional firms in the field of chemistry in the years 2001-2005.<sup>4</sup> The variable is denoted as PATS.

We take into account the effect of urbanization by including population density (POP) and the gross-domestic product (GDP) of a region in the year 2005. The corresponding data are obtained from the German Federal Institute for Research on Building.

Firms located in regions with strong public research infrastructure may also be more likely to link across regions. For instance, being co-located with public research institutes may induce knowledge spillovers and give better access to highly qualified personnel (e.g. Fritsch and Slavtchev, 2007). Accordingly, firms in these regions may be more prone to conduct R&D, engage in R&D collaboration, and be more successful in terms of innovation. In order to approximate this, we measure regions' public R&D infrastructure with three variables. The presence of universities in a region is modeled by counting their numbers of graduates in natural sciences in 2005 (UNI). Similarly, the analysis includes the number of employees working in regional research institutes of the Max Planck Society (MPG) and the Fraunhofer Society (FHG). More precise, only the numbers of employees working in the institutes' technological or natural science institutes in the year 2005 enter the analysis.<sup>5</sup>

### 4.2.1 Dyad level variables

We construct three variables at the dyad level. We measure geographical proximity with the physical distance between two regions' geographic centers. The variable is denoted as (GEO\_DIST). The chance of two regions being linked is expected to decrease with geographical distance. Geographical proximity frequently correlates with social proximity (Boschma, 2005), which needs to be considered in the interpretation.

We also include the variable SAME\_REG that has a value of 1 if both regions are located in the same federal state (i.e. NUTS 1 region), and 0 if not. SAME\_REG not only accounts for geographical proximity. It is likely to represent institutional proximity as well, as regions in the same federal state are probably similar in their R&D-related institutional

---

<sup>4</sup> The latest version of the „Patentatlas“ was published in 2006 and includes the patent data up to 2005. We use the aggregated numbers for 2001-2005 to minimize annual fluctuation.

<sup>5</sup> The employment numbers are relatively stable over time. Using data for a single year is therefore considered appropriate.

framework. The reason for this is the significant role the federal level is playing in the German R&D landscape. For instance, each federal state ("*Bundesland*") is responsible for its own resource endowment of universities and has its own R&D policies.

We also take into account that two regions with universities may be more likely to be linked. Firms in such regions are probably structurally more similar than two firms of which one is not located in a university region. It can be expected that firms in university regions are more R&D intensive and technologically more advanced as are more probable to benefit from knowledge spillovers (cf. Jaffe, 1998). To take this into account, we include the variable `UNI_1`, which has a value of one if both regions have a university and zero otherwise.

Notably, we do not construct a measure of technological similarity, which has been shown to make regions more likely to be linked (Scherngell and Barber, 2009). This is primarily motivated by data constraints. We analyze a network among firms of the same industry aggregated at the regional level. Hence, for the construction of a meaningful technological similarity measure we need information about the technological profiles of all regional firms in the chemical industry. Unfortunately, we miss such information and have to leave this issue to future research.

#### 4.2.3 Structural network level variables

We include five variables at the structural network level. Triadic closure (or transitivity) is captured by the geometrically weighted edgewise-shared partner statistic (GWESP-statistic: Snijders et al., 2006; Hunter et al., 2008). It measures the number of triangles in the network whilst taking into account the number of links that are involved in multiple triangles (multimodality) (see for details: Hunter et al., 2008). It thereby captures how frequently two nodes are connected by a direct link as well as by an indirect connection of length 2 (i.e. „two-path“) through another node (e.g. Hunter, 2007). If a positive coefficient is found for this statistic, there is a tendency towards triadic closure in the network.

We consider the geometrically weighted dyad shared partner statistic (GWDSP), which is an advanced version of the alternating k-two-path statistic put forward by Snijders et al. (2006). It measures the extent to which a network shows a tendency of nodes not directly linked to each other being at least indirectly linked. In other words, the statistic approximates whether multiple paths exist between such nodes. Accordingly, it captures multi-connectivity for nodes that are not tied directly.

Another variable at the network level is `EDGES`. It equals the number of links in the network and should always be included in exponential random graph models. Moreover, `EDGES` represents a so-called k-star(1) parameter. K-stars are essential configurations in networks. For instance, a k-star(2), or 2-star, corresponds to three nodes of which one is linked to each of the other two. Accordingly, a k-star(3) shows as four nodes with one node being linked to the other three. A triangle, i.e. three mutually connected nodes, logically includes three k-stars(2). This means that these configurations are hierarchically related (cf. Snijders et al., 2006, Hunter, 2007). While the `EDGES` parameter corresponds to a type of intercept parameter in the model, it is especially useful when considering the `GWDEGREE` statistic as well.

`GWDEGREE` is the geometrically weighted degree statistic, which helps modeling the observed network's degree distribution. Notably, the statistic can also be seen as an equivalent to the more traditional k-star statistic (Hunter, 2007). When being considered alongside the `EDGES` statistic, `GWDEGREE` (broadly) allows modeling preferential-attachment processes. More precise, if this statistic obtains a negative coefficient it signals the presence of preferential-attachment and a negative coefficient indicates anti-preferential attachment (Hunter, 2007).

For all three statistics, GWESP, GWDSP, and GWD, decay parameters have to be specified. Because few attempts have been made to systematically identify the best fitting parameter combinations (cf. Wright, 2010), researchers commonly rely on a manual iterative trial-and-error process of estimating varying model specifications. These specifications differ in terms of included variables as well as decay parameters of the GWDSP, GWESP and GWDEGREE statistics. This process ends when the best fitting model is identified. The best fitting model is a model that is stable and converges (when the Markov Chain Monte Carlo approach is used, the parameter traces should be horizontal) and provides the most appropriate goodness-of-fit statistics (matching degree, edgewise shared partners, and geodesic distributions) given the empirical data (observed network). In other words, the best fitting model most accurately predicts the structure of the observed network.

Once this model is identified the final goodness-of-fit statistics and MCMC trace plots are generated exclude all variables that are not significant at the 0.05 level in the original model. These variables are excluded because they represent noise that may distort the model and thereby bias the according statistics (cf. Wright, 2010). This “refined” model is used to generate all goodness-of-fit related statistics. We present the best fitting ERG-model for the cross-regional R&D collaboration network in the next section.

## 5 Results

Table 2 presents the results of the final, i.e. best fitting, model and those of its refined variant. Included are factors at the node, dyad, and structural network level. The model is stable and converges. Moreover, it is characterized by appropriate goodness-of-fit statistics (matching degree, edgewise shared partners, and geodesic distributions (Figure 2 in the Appendix) and horizontal parameter traces (Figure 3a-d in the Appendix).

Before we will discuss the variables with significant coefficients, it is also worthwhile to take a brief look at the insignificant ones. The insignificance of GDP implies that the economic prosperity of regions does not impact the structure of the cross-regional R&D collaboration network in the German chemical industry. The measure of the absolute physical distance (GEO\_DIST) between regions better captures the effect of geographic distance than when considering whether two regions are part of the same federal state (SAME\_REG), as the latter’s coefficient is insignificant while that of the first is not. The finding moreover questions the role of institutional proximity, which we argued to be reflected by SAME\_REG.

The measure of the network’s degree distribution (GWDEGREE) does not help in explaining the structure of the network. This means that we do not find evidence for preferential attachment processes, i.e. well-connected regions are not more prone to gain additional links than sparsely connected regions. The same applies to the GWDSP-statistic suggesting that two regions without a direct link are unlikely to be indirectly connected. Accordingly, we observe insignificant coefficients for variables at all three levels (node, dyad, and structural network level).

Table 2: Results of exponential random graph model with dyad level, node level and structural network level variables

	Main model				Refined model		
Variable	Estimate	Std. Error	p-value	Significance	Estimate	Std. Error	Significance
Node level							
PATS	0.00056	0.00013	< 1e-04	***	0.00028	0.00008	***
UNI	-0.00069	0.00017	< 1e-04	***	-0.00119	0.00015	***
POP_DEN	0.00009	0.00004	0.022735	*	0.00022	0.00001	***
GDP	-0.00113	0.00159	0.478296				
MPG	0.00037	0.00011	0.000882	***	0.00071	0.00009	***
FHG	0.00064	0.00026	0.013101	*	0.00135	0.00016	***
Dyad level							
GEO_DIST	-0.00164	0.00021	< 1e-04	***	-0.00072	0.00018	***
SAME_REG	0.07019	0.10950	0.521505				
Nodematch.UNI_1	0.30200	0.07094	< 1e-04	***	0.14760	0.07873	*
Structural network level							
EDGES	-4.36800	0.17230	< 1e-04	***	-7.24000	0.20440	***
GWESP, 0.69, fix	1.04400	0.06772	< 1e-04	***	2.02	0.00902	***
GWDEGREE	-2.86600	14.81000	0.846554				
GWDSP, 0.15, fix	0.02133	0.02736	0.435589				
Null Deviance: 50343.3 on 36315 degrees of freedom					50343.3 on 36315 degrees of freedom		
Residual Deviance: 1753.3 on 36302 degrees of freedom					1619.3 on 36305 degrees of freedom		
Deviance: 48589.0 on 13 degrees of freedom					48724.0 on 9 degrees of freedom		
AIC:	1779.3				1639.3		
BIC:	1889.8				1724.3		
* Significant at 90%; **Significant at 95%; *** Significant at 99%							

Now, we turn towards the significant variables reported in Table 2. As expected, regions with R&D intensive firms (PATS) tend to have more links. The same applies to urban regions (POP\_DEN) and regions in which institutes of the Max-Planck (MPG) and Fraunhofer (FHG) societies are located. The according coefficients of PATS, POP\_DEN, MPG, and FHG are all positive and significant. UNI obtains a negative significant coefficient suggesting that university regions tend to have fewer links. While this contradicts our expectations, it is essential to also consider the positive significant coefficient of the dyad-level variable UNI\_1 in the explanation. Accordingly, university regions generally have less links but they are more likely to link to other university regions. The latter is in line with our expectations and signals the presence of a homophily effect.

The dyad-level variable GEO\_DIST is characterized by a negative significant coefficient. Hence, geographical distance hampers link creation, which confirms existing empirical studies (cf. Magionni et al., 2007; Ponds et al., 2007; Scherngell and Barber, 2009; Hoekman et al., 2009; 2010; Broekel and Boschma, 2012).

We argued above that the main advantage of exponential random graph models is their ability to take into account factors at the structural network level in addition to factors at the node and dyad level. The significant coefficients of two variables at the structural network level empirically confirm this level's relevance. The coefficient of EDGES is negative and significant. By being similar to an intercept variable, EDGES represents the overall density of the network when all other effects are excluded. Its negative coefficient is a common feature of networks established by social processes indicating that such networks tend to be less dense than exponential random networks (cf. Varas, 2007).

In addition, we find a positive and significant coefficient of the GWESP-statistic. It means that triangles are a common feature of the network, which corresponding to the visual inspection of the network (see Figure 1). In other words, regions that are directly linked are also more likely to link through indirect connections. Hence, the result suggests that triadic closure is a driving force in the network formation processes. There might however be an alternative explanation. When constructing the empirical network, we transformed a bipartite network into a one-mode type. Such transformation more or less automatically increases the likelihood of triplets in the final one-mode network. Accordingly, the positive GEWSP-statistic might pick up this effect and act as a kind of control parameter for the one-mode projection procedure. However, we pointed out in Section 4.1 that on average less than three firms (2.8) are jointly participating in a cooperative project. Hence, it is most likely a combination of both effects that explains the statistic's significance. In any case, this structural network factors significantly helps in modeling the structure of the network.

In sum, we find that the structure of the network is best explained by factors at the node level, dyad level, and structural network level. This highlights the usefulness of exponential random graph models as a tool for analyzing structures of cross-regional collaboration networks.

## 6 Conclusion

The aim of this study was to discuss exponential random graph models (ERGM) as promising tools for the investigation of cross-regional collaboration networks. We pointed out that most existing studies focus on the evaluation of factors at the node and dyad level. However, network science suggests that factors at the structural network level may also be relevant in this respect. Such factors cannot be considered in methods commonly applied in this context. For instance, spatial interaction models allow only for factors at the node and dyad level. We argued that ERG-models represent a powerful alternative as they take into account factors at all three levels and require only cross-sectional network data.

We illustrated the application of ERGMs by analyzing the structure of the cross-regional R&D collaboration network in the German chemical industry between 2005 and 2010. By using an exponential random graph model, we considered factors at all three levels that might influence the network's structure. At the node level, it was shown that urban regions (reflected by population density) and regions with high research intensities are more likely to be linked to other regions. At the dyad level, we found regions to be more likely being linked when they have a university. Moreover, our results confirmed the negative impact of geographical distance on the likelihood of research collaboration. Finally, at the structural network level, evidence was provided for the existence of a triadic closure (transitivity) effect: regions that are indirectly linked to each other are likely to be directly linked as well.

Clearly, the study is only a first step towards understanding the role factors at the structural network level play for the formation of cross-regional collaboration networks. It

nevertheless underlines the usefulness of exponential random graph models for future research endeavors on this subject.

## References

Autant-Bernard, C., Billand, P., Frachisse, D., Massard, N. (2007) Social distance versus spatial distance in R&D co-operation: Empirical evidence from European collaboration choices in micro and nanotechnologies, *Papers in Regional Science*, 86, pp. 495-519.

Barabasi, A.L. and Albert, R. (1999) Emergence of Scaling in Random Networks, *Science*, 15, pp. 509-512.

Boschma, R.A. and Ter Wal, A.L.J. (2007) Knowledge networks and innovative performance in an industrial district: The case of a footwear district in the south of Italy, *Industry & Innovation*, 14, pp. 177-199.

Breschi, S. and Lissoni, F. (2009) Mobility of skilled workers and co-invention networks: an anatomy of localized knowledge flows, *Journal of Economic Geography*, 9, pp. 439-468.

Broekel, T. (2007). A concordance between industries and technologies - matching the technological fields of the Patentatlas to the German industry classification. Jenaer Economic Research Papers, 2007-013.

Broekel, T. and Boschma, R. (2012) Knowledge networks in the Dutch aviation industry: the proximity paradox, *Journal of Economic Geography*, 12, pp. 409-433.

Broekel, T. and Hartog, M. (2013), Explaining the structure of inter-organizational networks using exponential random graph models, *Industry and Innovation*, 20(3): forthcoming

Broekel, T. and Graf, H. (2012). Public research intensity and the structure of German R&D networks: A comparison of 10 technologies. *Economics of Innovation and New Technology*, 21(4):345–372.

Cohen, W.M. and Levinthal, D.A. (1990) Absorptive capacity: a new perspective on learning and innovation, *Administrative Science Quarterly*, 35, pp. 128-152.

Coleman, J.S. (1988) Social capital in the creation of human capital, *American Journal of Sociology*, 94 (supplement), pp. S95-S120.

Desmarais, B.A. and Cranmer, S.J. (2012) Micro-Level Interpretation of Exponential Random Graph Models with an Application to Estuary Networks, *Policy Studies Journal*, 40, DOI: 10.1111/j.1541-0072.2012.00459.x

Fischer, M.M., Scherngell, T. and Jansenberger, E. (2006), The geography of knowledge spillovers between high-technology firms in Europe. Evidence from a spatial interaction modeling perspective, *Geographical Analysis*, 38, pp. 288-309.

Fritsch, M. and Slavtchev, V. (2007), Universities and Innovation in Space, *Industry and Innovation*, 14, pp. 201–218.

- Giuliani, E. and Bell, M. (2005) The micro-determinants of meso-learning and innovation: evidence from a Chilean wine cluster, *Research Policy*, 34, pp. 47-68.
- Gulati, R. (1999), Network location and learning: the influence of network resources and firm capabilities on alliance formation, *Strategic Management Journal*, 20, pp. 397-420.
- Greif, S. and Schmiedl, D. (2002). Patentatlas 2002 Dynamik und Strukturen der Erfindungstätigkeit. Deutsches Patent- und Markenamt, München.
- Greif, S., Schmiedl, D., and Niedermeyer, G. (2006). Patentatlas 2006. Regionaldaten der Erfindungstätigkeit. Deutsches Patent- und Markenamt, München.
- Glückler, J. (2007) Economic geography and the evolution of networks, *Journal of Economic Geography*, 7, pp. 619-634.
- Handcock, M.S. (2003) Statistical Models for Social Networks: Degeneracy and Inference, in: R.L. Breiger, K.M. Carley & P. Pattison (Eds.), *Dynamic Social Network Modeling and Analysis: Workshop Summary and Papers*, pp. 229–240 (Washington, DC: National Academies Press).
- Hoekman, J., Frenken, K. and Van Oort, F. (2009) The geography of collaborative knowledge production in Europe, *The Annals of Regional Science*, 43, pp. 721-738.
- Hoekman, J., Frenken, K. & Tijssen, R.J.W. (2010) Research collaboration at a distance: Changing spatial patterns of scientific collaboration within Europe, *Research Policy*, 39, pp. 662-673.
- Hunter, D. R. (2007). Curved exponential family models for social networks. *Social Networks*, 29(2):216–230.
- Hunter, D.R., Goodreau, S.M. and Handcock, M.S. (2008) Goodness of Fit for Social Network Models, *Journal of the American Statistical Association*, 103, pp. 248-258.
- Jaffe, A. (1989). Real effects of academic research. *American Economic Review*, 79(5):957–970.
- LeSage, J., Fischer, M.M., Scherngell, T. (2007) Knowledge spillovers across Europe. Evidence from a Poisson spatial interaction model with spatial effects, *Papers in Regional Science*, 86, pp. 393-421.
- Maggioni, M.A., Nosvelli, M. and Uberti, T.E. (2007) Space versus networks in the geography of innovation: a European analysis, *Papers in Regional Science*, 86, pp. 471-493.
- Maskell, P. and Malmberg, A. (1999) Localized Learning and Industrial Competitiveness, *Cambridge Journal of Economics*, 23, pp. 167–186.
- McPherson, M., Smith-Lovin, L. and Cook, J.M. (2001) Birds of a Feather: Homophily in Social Networks, *Annual Review of Sociology*, 27, pp. 415-444.

- Metcalf, S. (1995) The Economic Foundations of Technology Policy: Equilibrium and Evolutionary Perspectives. In: Stoneman, P. (Ed), *Handbook of the Economics of Innovation and Technological Change*, pp. 409-512 (Oxford: Basil Blackwell).
- Nelson R. Winter S. (1982) *An Evolutionary Theory of Economic Change* (Cambridge MA: Belknap)
- Nooteboom, B. (2000) *Learning and innovation in organizations and economies* (Oxford: Oxford University Press).
- Ponds, R., Van Oort, F. And Frenken, K. (2007) The geographical and institutional proximity of research collaboration, *Papers in Regional Science*, 86, pp. 423-443.
- Powell, W. W., White, D. R., Koput, K. W., and Owen-Smith, J. (2005) Network dynamics and field evolution: The growth of interorganizational collaboration in the life sciences, *American Journal of Sociology*, 110, pp. 1132-1206.
- Robins, R., Snijders, T., Wang, P., Handcock, and Pattison, P. (2006) Recent developments in exponential random graph ( $p^*$ ) models for social networks, *Social networks*, 29, pp. 192-215.
- Robins, G., Pattison P., Kalish Y. and Lusher, D. (2007) An introduction to exponential random graph ( $p^*$ ) models for social networks, *Social Networks*, 29, pp 173-191.
- Romer, P.M. (1990) Endogenous technological change, *Journal of Political Economy*, 98, pp. 71-102.
- Saul, Z.M. and Filkov, V (2007) Exploring biological network structure using exponential random graph models, *Bioinformatics*, 23, pp. 2604- 2611.
- Scherngell, T. and Barber, M. J. (2009) Spatial interaction modeling of cross- region R&D collaboration. empirical evidence from the 5th EU framework programme, *Papers in Regional Science*, 88: 531–546.
- Scherngell, T. and Barber, M. J. (2011). Distinct spatial characteristics of industrial and public research collaborations: evidence from the fifth EU Framework Programme. *Annals of Regional Science*, 46:247–266.
- Snijders, T.A.B. (2002) Markov Chain Monte Carlo Estimation of Exponential Random Graph Models, *Journal of Social Structure* 3. Web journal available from: <<http://www2.heinz.cmu.edu/project/INSNA/joss/index1.html>>
- Snijders, T.A.B, Pattison P.E., Robins G. and Handcock, M.S. (2006) New specifications for exponential random graph models, *Sociological Methodology*, 36, pp. 99-153.
- Snijders, T.A.B., van de Bunt, G.G. and Steglich, C.E.G. (2010) Introduction to stochastic actor-based models for network dynamics, *Social Networks*, 32, pp. 44-60.
- Sobrero, M. and Roberts, E.B. (2001) The trade-off between efficiency and learning in interorganizational relationships for product development, *Management Science*, 47, pp. 493-511.

Ter Wal, A.L.J. and R.A. Boschma (2009) Applying social network analysis in economic geography: framing some key analytic issues, *Annals of Regional Science*, 43, pp. 739-756.

Ter Wal, A.L. J. (2011) The Dynamics of Inventor Networks in German Biotechnology: Geographical Proximity versus Triadic Closure. *Papers in Evolutionary Economic Geography* 11.02.

Uddin, S., Hamra, J. and Hossain L. (2012) Exploring communication networks to understand organizational crisis using exponential random graph models, *Computational and Mathematical Organization Theory*, DOI: 10.1007/s10588-011-9104-8

Van Duin, M.A.J., Gille, K.J. and Handcock, M.S. (2009) A framework for the comparison of maximum pseudo-likelihood and maximumlikelihood estimation of exponential family random graph models, *Social Networks*, 31, pp. 52-62.

Varas, M. L. L. (2007) Essays in social space: applications to Chilean communities on inter-sector social linkages, social capital, and social justice. Dissertation at University of Illinois, Urbana-Champaign, USA.

Wang, P., Pattison, P. and Robins, G. (2012) Exponential random graph model specifications for bipartite networks – A dependence hierarchy, *Social Networks*, <http://dx.doi.org/10.1016/j.socnet.2011.12.004>

Wright, D. (2010), Repression and network science: Tools in fight against terrorism. Dissertation at University of Michigan:  
[http://deepblue.lib.umich.edu/bitstream/handle/2027.42/77898/dewright\\_1.pdf?sequence=1](http://deepblue.lib.umich.edu/bitstream/handle/2027.42/77898/dewright_1.pdf?sequence=1)

## APPENDIX

<b>Table 1: Descriptives of empirical variables</b>								
<i>Variables</i>	<i>n</i>	<i>mean</i>	<i>st. deviation</i>	<i>median</i>	<i>min</i>	<i>max</i>	<i>skew</i>	<i>kurtosis</i>
PATS	270	69.55	199.12	12.48	0	1691.31	5.34	32.55
POP_DEN	270	825.35	1265.19	244.5	40	8523	3.06	11.44
GDP	270	40.46	33.58	26.75	14.1	296.9	3.66	19.83
UNI	270	101.51	244.73	0	0	1812	3.46	15.55
MPG	270	49.12	248.08	0	0	3438	10.20	128.50
FHG	270	30.81	123.52	0	0	978	5.22	29.24
GEO_DIST	72900	379.81	186.03	368.54	0	977.45	0.29	-0.52
SAME_REG	72900	0.11	0.31	0	0	1	2.49	4.22
UNI_1	72900	0.62	0.49	1	0	1	-0.47	-1.77

Figure 2: Goodness of fit of exponential random graph model with dyad level, node level + structural network level variables

### Goodness-of-fit diagnostics

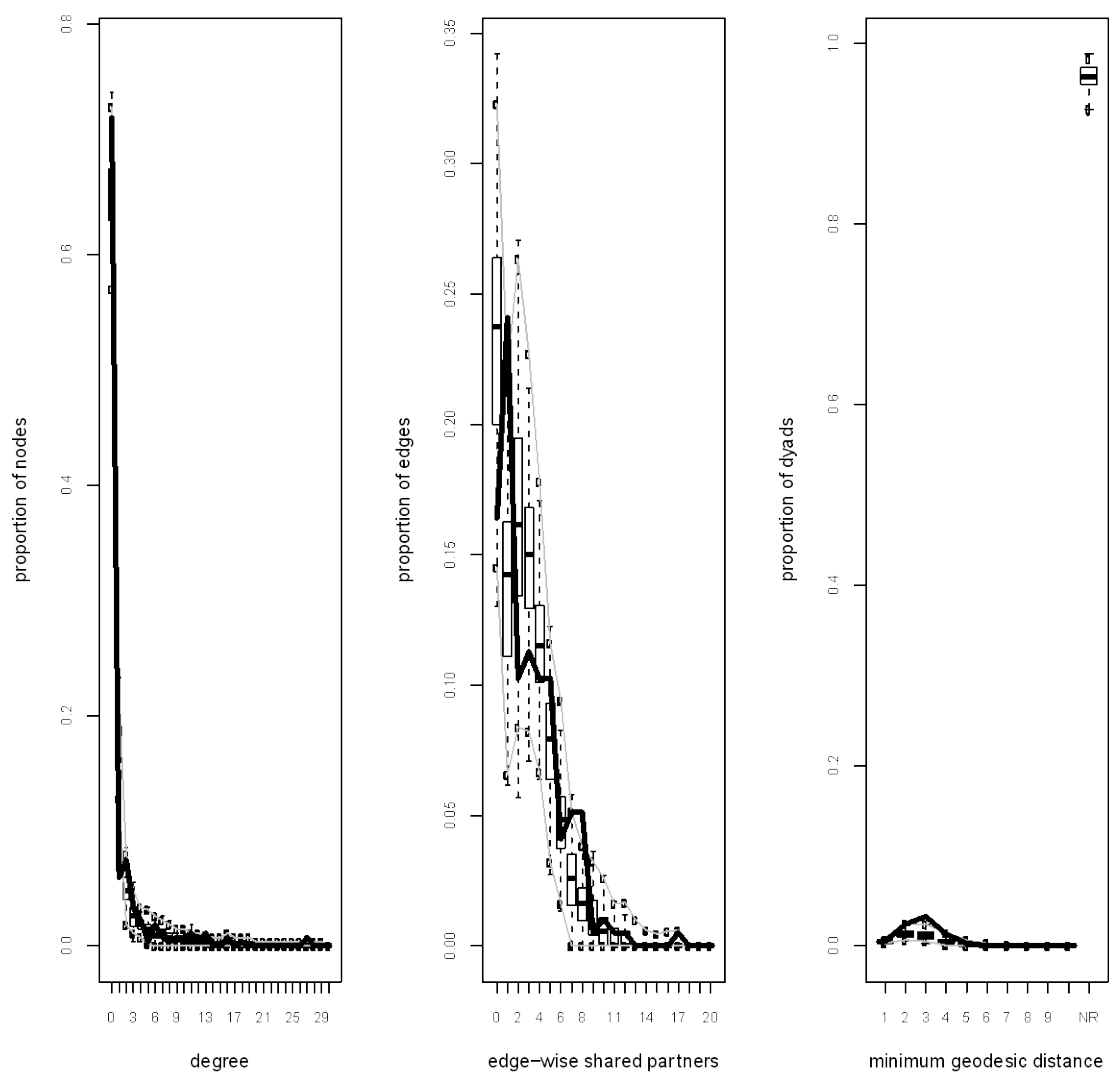


Figure 3 (a): MCMC-Statistics of exponential random graph model with dyad level, node level and structural network level variables

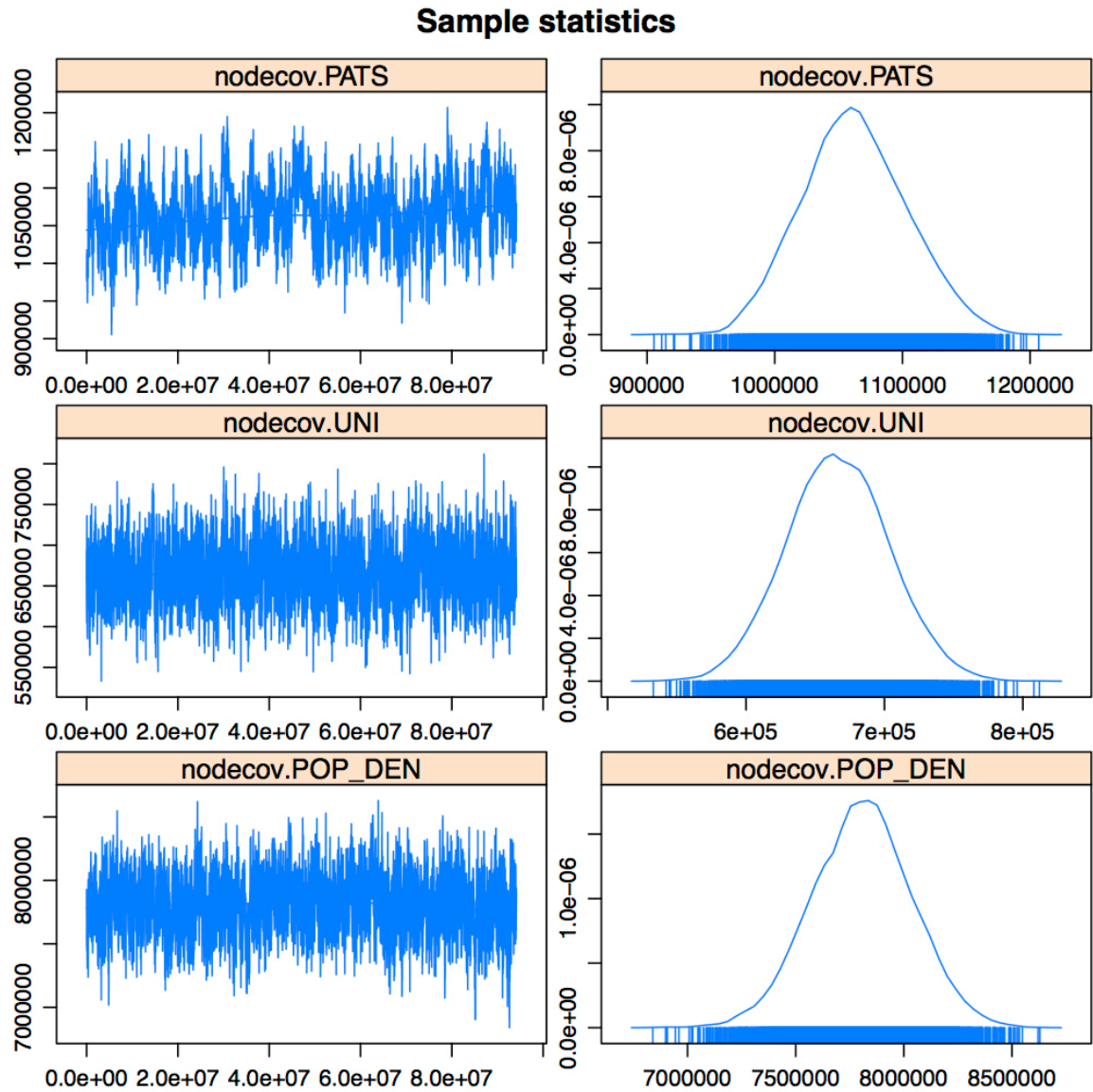


Figure 3 (b): MCMC-Statistics of exponential random graph model with dyad level, node level and structural network level variables

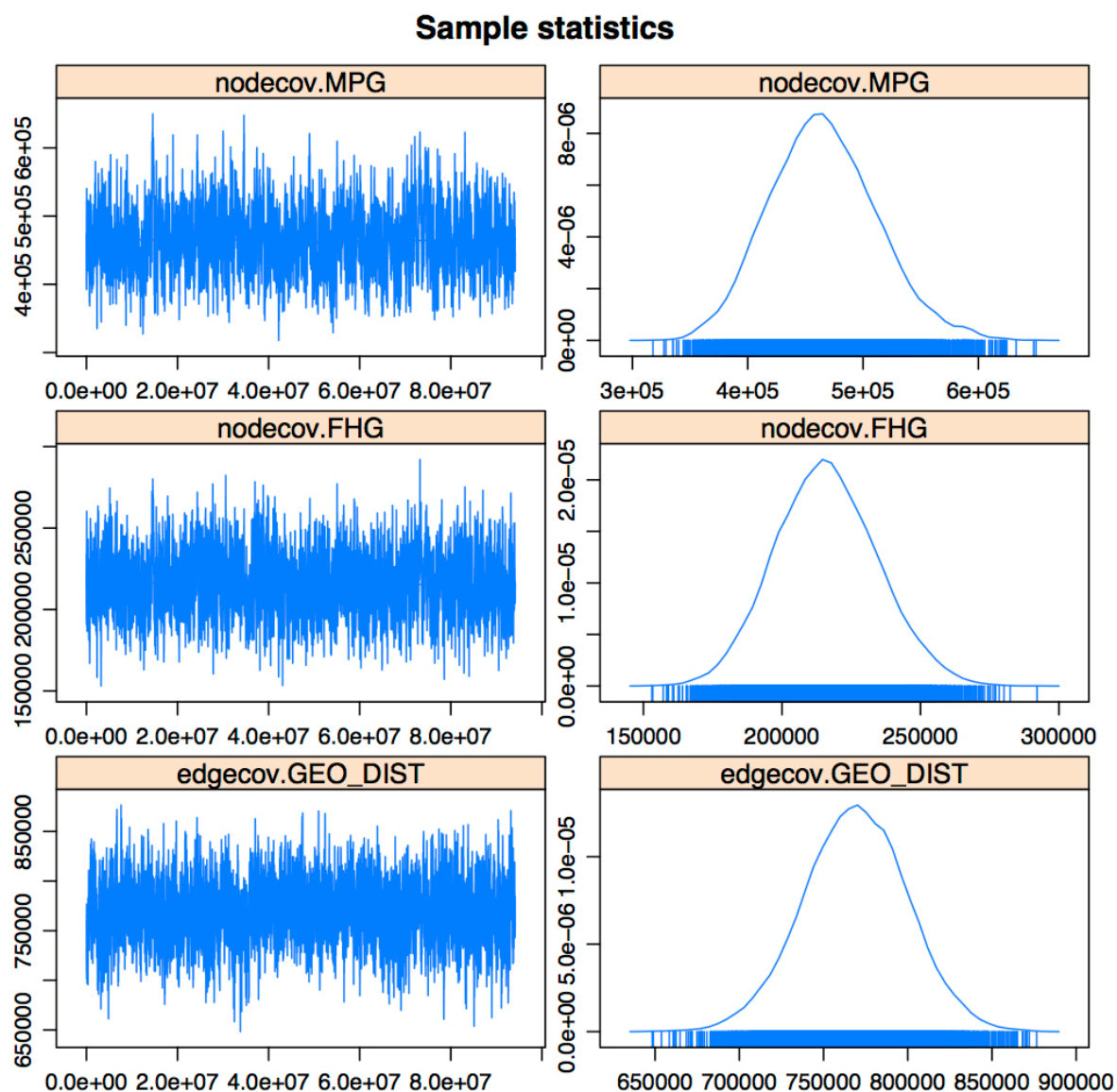


Figure 3 (c): MCMC-Statistics of exponential random graph model with dyad level, node level and structural network level variables

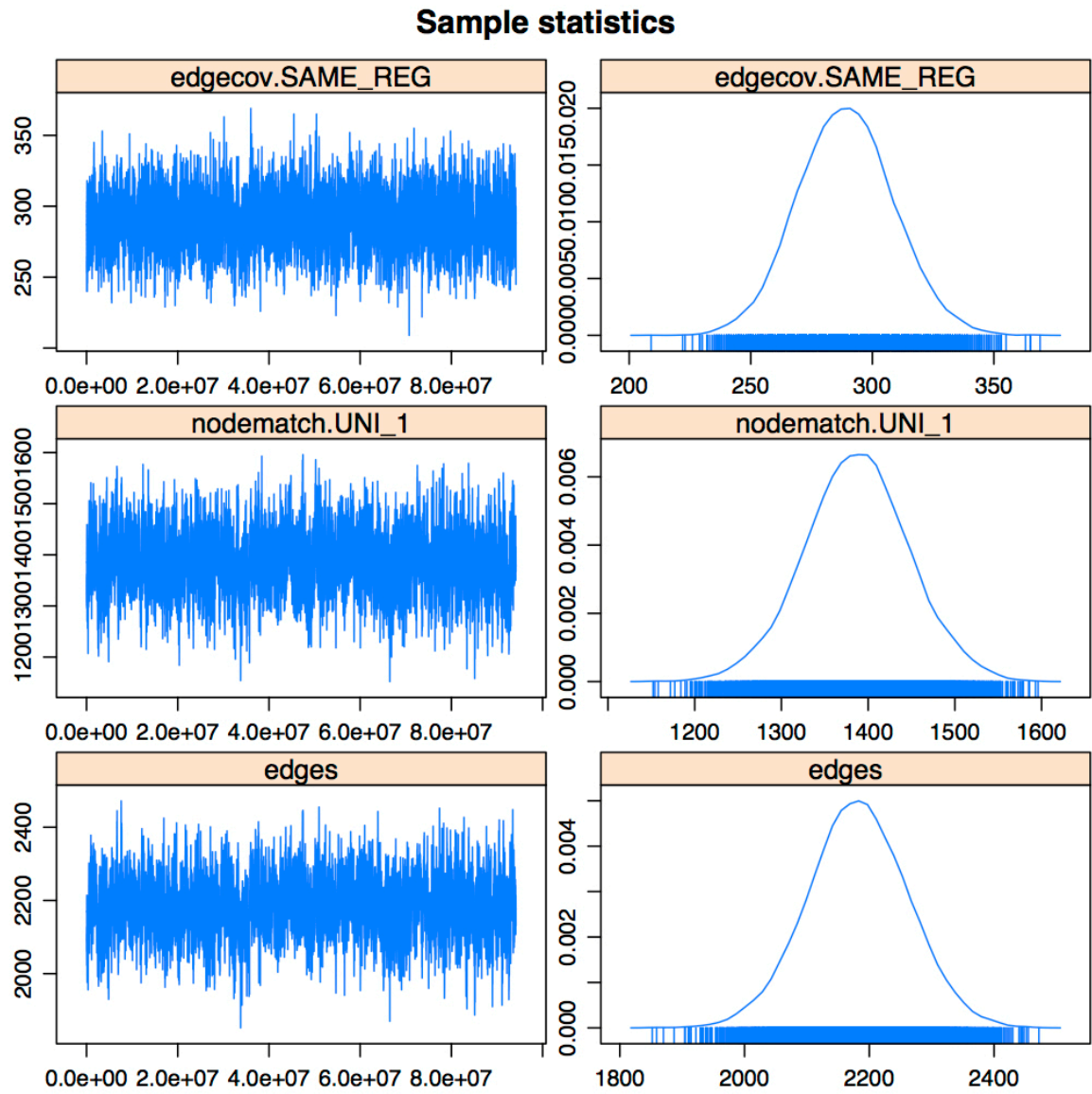


Figure 3 (d): MCMC-Statistics of exponential random graph model with dyad level, node level and structural network level variables

