



# A Method for Estimating Driving Factors of Illicit Trade Using Node Embeddings and Clustering

Jorge Ángel González Ordiano<sup>1</sup> , Lisa Finn<sup>2</sup>, Anthony Winterlich<sup>2</sup>, Gary Moloney<sup>2</sup>, and Steven Simske<sup>1</sup> 

<sup>1</sup> Colorado State University, Fort Collins, CO, USA  
{jorge.gonzalez\_ordiano, steve.simske}@colostate.edu

<sup>2</sup> Micro Focus International, Galway, Ireland  
{finn, winterlich, gary.moloney}@microfocus.com

**Abstract.** The trade on illegal goods and services, also known as illicit trade, is expected to drain 4.2 trillion dollars from the world economy and put 5.4 million jobs at risk by 2022. These estimates reflect the importance of combating illicit trade, as it poses a danger to individuals and undermines governments. To do so, however, we have to first understand the factors that influence this type of trade. Therefore, we present in this article a method that uses node embeddings and clustering to compare a country based illicit supply network to other networks that represent other types of country relationships (e.g., free trade agreements, language). The results offer initial clues on the factors that might be driving the illicit trade between countries.

**Keywords:** Node embedding · Clustering · Data mining

## 1 Introduction

Illicit trade, i.e. the trade in illicit goods and services, poses a danger to our communities [12]. For instance, the ICC estimates that by 2022 counterfeit and piracy will put 5.4 million jobs at risk and drain 4.2 trillion dollars of the world economy,<sup>1</sup> while the OECD estimates that in the UK 86,300 jobs were lost due to counterfeiting and piracy in 2016 alone [11]. In addition, the dangers of illicit trade go beyond economic losses. For example, illicit medicines have been recorded to cause malaria and tuberculosis deaths [10], while counterfeits have been shown to finance terrorists organizations [1]. For these reasons, getting a better understanding of the factors that might be driving illicit trade is of major importance, if we are to develop methods that will aid in its disruption.

A possibility for getting a better grasp on these driving factors is to use networks to describe not only the illicit trade between countries (i.e. the illicit supply network), but also other aspects that countries might have in common,

<sup>1</sup> [iccwbo.org/global-issues-trends/bascap-counterfeiting-piracy/](https://www.iccwbo.org/global-issues-trends/bascap-counterfeiting-piracy/), Accessed: 12/05/19.

such as language, geographic proximity, etc. By comparing communities (i.e. groups of countries that have a strong relationship) within the illicit supply network to those within the others, we can qualitatively estimate the aspects that might be driving illicit trade. However, searching for these communities is not a trivial task, as the different networks may have different properties (directed, undirected, weighted etc.) that influence the community detection algorithms that we can use. Therefore, an algorithm able to deal with all possible network types—without much user involvement—is not immediately clear. To circumvent this issue, we present an approach that combines the creation of node embeddings (i.e. vector representations of the network nodes [9]) with traditional clustering algorithms (such as k-means or affinity propagation).

The remainder of this article is divided as follows: Sect. 2 offers background information on the different concepts and approaches used herein. Section 3 presents this article’s method. Section 4 discusses the data used. Section 5 describes the experimental study conducted in this article, while Sect. 6 presents and discusses the results. Finally, Sect. 7 offers the conclusion and the outlook of this article.

## 2 Preliminaries

### 2.1 Node Embedding

Many data mining algorithms require feature vectors as input; therefore, if we want to use these approaches to predict, classify, or cluster nodes within a network, we first need to construct vector representations—i.e. embeddings—of them. Currently, methods that automatically construct these representations have become popular in literature [9]. In this article, we specifically use a python implementation<sup>2</sup> of the *node2vec* algorithm [8]; an algorithm that estimates the node embeddings based on a series of random walks.

### 2.2 Clustering

Finding community structures within a network is useful for finding nodes that have a strong relationship to one another. However, some of the community detection methods are not only computationally expensive, but also have the disadvantage of being dependent on network properties. For instance, only a few of the algorithms used on undirected networks can be extended to directed ones [5]. To overcome this network property dependency, we estimate the network communities using the node embeddings and not the network itself. In other words, we first cluster the embeddings and then we define the clusters as the communities structures we are looking for. This alternative has already proven to be effective [4] and thus is the one used herein. Furthermore, the clustering method that we use is affinity propagation [6], which is implemented within the

<sup>2</sup> [github.com/eliorc/node2vec](https://github.com/eliorc/node2vec); Accessed: 02/05/20.

*apcluster* R package<sup>3</sup>. We choose this approach, as it does not require a predefined number of clusters to work. Readers are referred to [2] for more information on the advantages and disadvantages of the affinity propagation algorithm.

### 3 Method

The method begins with the networks' adjacency matrices, which are given in the present article as:

$$\mathbf{A}_l = \begin{bmatrix} a_{11} & \cdots & a_{1N_v} \\ \vdots & \ddots & \vdots \\ a_{N_v 1} & \cdots & a_{N_v N_v} \end{bmatrix}_l, \tag{1}$$

where  $\mathbf{A}_l$  is the  $l^{th}$  network adjacency matrix,  $N_v$  is the number of nodes in the network, and  $[a_{ij}]_l$  are the elements of the matrix representing if there is an edge connecting node  $i$  to node  $j$  in network  $l$ .

By using the *node2vec* algorithm (cf. Sect. 2.1) we obtain a vector representation of the nodes at each network, i.e.:

$$\mathbf{e}_{nl} = [e_{n1}, \dots, e_{nd_l}]_l^T, \tag{2}$$

where  $\mathbf{e}_{nl}$  represents the vector of the  $n^{th}$  node at network  $l$  and  $d_l$  is the vector's dimension, which can vary depending on the network.

As mentioned previously, we use in this article the affinity propagation clustering algorithm; an algorithm that requires a similarity matrix to work. Therefore, we define for each network the following similarity matrices:

$$\mathbf{S}_l = \begin{bmatrix} s_{11} & \cdots & s_{1N_v} \\ \vdots & \ddots & \vdots \\ s_{N_v 1} & \cdots & s_{N_v N_v} \end{bmatrix}_l, \text{ with } [s_{ij}]_l = -\|\mathbf{e}_{il} - \mathbf{e}_{jl}\|_2, \tag{3}$$

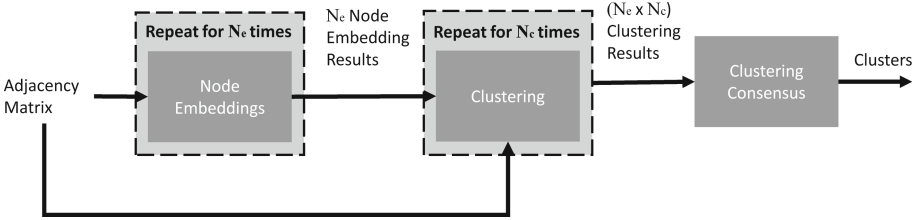
with  $[s_{ij}]_l$  being the negative Euclidean distance between node  $i$  and node  $j$  at network  $l$ . Since the similarity matrices do not consider if nodes can reach each other within the network, nodes that cannot reach each other might end up in the same cluster. To avoid this issue, we determine at each network if nodes can reach each other in a walk of length  $l_r$ , i.e. the length of the random walks used to create the embeddings. If that is not the case, we set their similarity value equal to a threshold  $t_{diss}$  that will make those two nodes as dissimilar as possible.

$$[s_{ij}]_l = \begin{cases} [s_{ij}]_l & , \text{ if node } i \text{ can reach node } j \text{ within a random walk of} \\ & \text{length } l_r \text{ in network } l \\ t_{diss} & , \text{ otherwise} \end{cases}. \tag{4}$$

---

<sup>3</sup> [cran.r-project.org/web/packages/apcluster/apcluster.pdf](https://cran.r-project.org/web/packages/apcluster/apcluster.pdf); Accessed: 02/05/20.

Note that the creation of the node embeddings and the clustering are not necessarily deterministic. Therefore, it is important to make sure that the obtained results are as representative as possible. To do so, we use the method depicted in Fig. 1.



**Fig. 1.** Clustering method

As Fig. 1 shows, the method repeats for each network the embedding step  $N_e \in \mathbb{N}_{>0}$  times and the clustering step an additional  $N_c \in \mathbb{N}_{>0}$  times, thus resulting in a total of  $N_e \cdot N_c$  clustering results. These results are used as input in a final consensus step for the clustering.

In this final step, we create a new network in which two nodes are connected if they cluster at least a certain number of times. To be more specific, we define a new adjacency matrix  $A_l^{\text{opt}}$  whose elements  $[a_{ij}]_l^{\text{opt}}$  are defined as:

$$[a_{ij}]_l^{\text{opt}} = [a_{ji}]_l^{\text{opt}} = \begin{cases} 1 & \text{, if node } i \text{ and node } j \text{ cluster } N_{\text{th}} \text{ times or more} \\ 0 & \text{, otherwise} \end{cases}, \quad (5)$$

with  $N_{\text{th}}$  being a threshold defining the number of times that two nodes have to be grouped together for them to be connected in this new network. Afterwards, we define the subcomponents of this new network as the clusters of the  $l^{\text{th}}$  network.

## 4 Data

The data we use to represent different types of country-country relationships (such as, licit and illicit trade, amount of traded goods, etc.) are described in the following paragraphs:

- **Licit and Illicit Trade:** The dataset from which we estimate the licit and illicit trade by country comes from the Global Product Authentication Service (GPAS) of MicroFocus International.<sup>4</sup> This dataset contains the authentication results of 55,999 unique serial codes (31,989 authenticated as true and

<sup>4</sup> [microfocus.com/en-us/services/product-authentication-anti-counterfeit-services](https://microfocus.com/en-us/services/product-authentication-anti-counterfeit-services); Accessed: 02/03/20.

24,010 authenticated as false), as well as the country and time in which the codes were authenticated. These authentications are all from 2011 to 2015. Readers are referred to [7] for more information on the dataset.

- **Amount of traded Goods:** The amount of traded goods between countries (i.e. exports and imports in US dollars) is modeled using the data reported on the UN Comtrade database.<sup>5</sup> In this article, we only use data from the years 2011 to 2015, to make it consistent with the GPAS data we have available.
- **Free Trade Agreements:** Information about the countries free trade agreements is obtained from the Regional Trade Agreements Database of the World Trade Organization (WTO).<sup>6</sup> Note that in this article we only make use of free trade agreements that came into force before 2016, in order to make the data compatible with the GPAS dataset.
- **Language:** The data used to determine the language of each individual country is taken from the CIA’s website.<sup>7</sup> Note that we only consider languages that are listed as an official language, as an official minority language, as a lingua franca, or as a language spoken by at least 10% of the population. If a country we need is missing on the list or if we cannot determine what language we should consider, we use the languages listed as official in the country’s Wikipedia page.
- **Geography:** The geographic relationship between countries is modeled in this article as the inverse distance between the countries centroids. To calculate the inverse distances, the necessary centroid coordinates are obtained using the *countrycode* R-package. If the functions within the R-package are unable to provide the coordinates of a given country, we instead take them from a file found on the Periscope Data website.<sup>8</sup>

## 5 Experimental Study

The goal of this experiment is to gain insight into which factors might be driving the illicit trade between countries. The first step in achieving this goal is to create networks that describe different types of country relationships. To be more specific, the networks used herein model the following aspects: licit and illicit trade estimated using GPAS data; amount of traded goods (i.e. reported exports and imports in US dollars); trade discrepancies (i.e. differences in reported exports and imports); the existence of free trade agreements; the use of a common language; and geographic proximity.

In other words, we create seven distinct networks with different properties. However, before we create the networks, we need to make sure that the countries (i.e. the nodes) we consider appear in all of the datasets we are using, so that they also appear in all of the networks. After some preprocessing we end up with the 146 countries shown in Table 1.

<sup>5</sup> [comtrade.un.org](http://comtrade.un.org); Accessed: 11/26/19.

<sup>6</sup> [rtais.wto.org/UI/PublicMaintainRTAHome.aspx](http://rtais.wto.org/UI/PublicMaintainRTAHome.aspx); Accessed: 12/12/19.

<sup>7</sup> [cia.gov/library/publications/the-world-factbook/fields/402.html](http://cia.gov/library/publications/the-world-factbook/fields/402.html); Accessed: 12/14/19.

<sup>8</sup> [community.periscopedata.com/t/63fy7m/country-centroids](http://community.periscopedata.com/t/63fy7m/country-centroids); Accessed: 08/26/19.

**Table 1.** List of countries that are used as nodes within the networks

Countries
Afghanistan; Albania; Algeria; Angola; Argentina; Armenia; Australia; Austria; Azerbaijan; Bahamas; Bahrain; Bangladesh; Barbados; Belarus; Belgium; Benin; Bhutan; Bolivia; Bosnia & Herzegovina; Botswana; Brazil; British Virgin Islands; Brunei; Bulgaria; Burkina Faso; Burundi; Cambodia; Cameroon; Canada; Chile; China; Colombia; Congo-Kinshasa; Costa Rica; Côte d'Ivoire; Croatia; Cyprus; Czech Republic; Denmark; Djibouti; Dominican Republic; Ecuador; Egypt; El Salvador; Ethiopia; Fiji; France; Georgia; Germany; Ghana; Greece; Guatemala; Guinea; Guyana; Haiti; Honduras; Hong Kong (SAR China); Hungary; Iceland; India; Indonesia; Iran; Iraq; Ireland; Israel; Italy; Jamaica; Japan; Jordan; Kazakhstan; Kenya; Kuwait; Laos; Latvia; Lebanon; Liberia; Lithuania; Luxembourg; Macau (SAR China); Macedonia; Madagascar; Malawi; Malaysia; Maldives; Mauritania; Mauritius; Mexico; Moldova; Montenegro; Morocco; Mozambique; Myanmar (Burma); Namibia; Nepal; Netherlands; New Zealand; Nicaragua; Niger; Nigeria; Norway; Oman; Pakistan; Palestinian territories; Panama; Papua New Guinea; Paraguay; Peru; Philippines; Poland; Portugal; Qatar; Romania; Russia; Rwanda; Saudi Arabia; Senegal; Serbia; Singapore; Slovakia; Slovenia; South Africa; South Korea; South Sudan; Spain; Sri Lanka; Sudan; Sweden; Switzerland; Syria; Tanzania; Thailand; Togo; Trinidad & Tobago; Tunisia; Turkey; Uganda; Ukraine; United Arab Emirates; United Kingdom; United States; Uzbekistan; Venezuela; Vietnam; Yemen; Zambia; Zimbabwe

Furthermore, the creation of each one of the networks is described below:

- **Licit and Illicit Trade:** The two networks that describe licit and illicit trade between countries are directed networks with weighted edges created using the GPAS serial codes that were authenticated as true or false, respectively. The weights of an edge joining country  $i$  to country  $j$  represents the number of times that a serial code is authenticated first in  $i$  and then in  $j$ .
- **Amount of traded Goods:** This network is created as a directed network with weighted edges. The weight from country  $i$  to country  $j$  represents the trade value (US dollars) in goods that goes from  $i$  to  $j$ . Due to reporting discrepancies, the weights are calculated as the arithmetic mean between the exports reported by country  $i$  and the imports reported by country  $j$ .
- **Trade Discrepancy:** This network is modeled as an undirected network with weighted edges. These weights represent the arithmetic mean between the differences in imports and exports reported by country  $i$  and country  $j$ .
- **Free Trade Agreements:** This network consists of an undirected and unweighted network, whose adjacency matrix elements are 1 if there is a free trade agreement between two countries and 0 otherwise.
- **Language:** The language network is also an undirected and unweighted network with an adjacency matrix that has elements equal to one if two countries share a language and zero otherwise.
- **Geography:** This network consists of an undirected network with weighted edges, whose weights are the inverse of the distance between the centroids of country  $i$  and country  $j$ .

After creating the networks, we can start obtaining their necessary embeddings. However, there are parameters of the *node2vec* algorithm that we still need to define: the length of the random walks, the number of random walks that we calculate per node, the search bias parameters that influence the creation of the random walks, the number of random walk elements that define a nodes' context, and the dimension of the embedding vectors.

Considering that we are interested in knowing which aspects might be driving illicit trade, we set the length of the random walks equal to the number of locations we assume an illicit item might visit. That is the mean number of authentications of an illicit serial code, which in our dataset is three (i.e.  $l_r = 3$ ). Note that we set the *node2vec* search bias parameters—which are used to create the random walks—equal to one (i.e. their default value in the implementation we are using). Furthermore, to make sure that the collection of random walks is as representative as possible, we create for each node 1000 of them. In addition, the whole random walk is used as context for estimating the node embeddings.

In contrast to the other parameters, we define the dimension of the embedding vectors of each network by testing the clustering results of several possible dimensions, i.e.  $d_e = \{2, 3, 4, 5, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100\}$ . In other words, we choose for each network the dimension that delivers the best clustering results according to an objective function.

Before we describe the objective function, it is important to define some parameters that are necessary for the method described in Sect. 3. For instance, the value of  $t_{\text{diss}}$  (cf. Eq. (4)) is set equal to -Inf. This value is used to indicate countries that have no similarity within the implementation of the affinity propagation algorithm we are using. At the same time,  $N_e$  and  $N_c$  (cf. Fig. 1) are set equal to 100 and 10, respectively. In other words, we create for each network 100 embedding results that we then cluster 10 different times. Moreover,  $N_{\text{th}}$  (cf. Eq. (5)) is set equal to 900, i.e. 90% of all clustering results. These values are used to obtain results that are as representative as possible. Finally, it is important to mention that we standardize all embedding vectors within a single embedding result, before any of the clustering steps described in Sect. 3. The standardization is used to prevent variables with large scales from dominating the clustering procedure.

The objective function we use is based on the pseudo F-statistic [3] and is given by the following equation:

$$c_{ld_e} = \left(1 - \frac{N_{\text{cl},ld_e}}{N_v}\right) \cdot \text{SNR}\{F_{l1d_e}, F_{l2d_e}, \dots, F_{lN_e d_e}\}, \quad (6)$$

where  $N_v$  is the number of network nodes,  $N_{\text{cl},ld_e}$  is the number of clusters found on the  $l^{\text{th}}$  network using  $d_e$ -dimensional embeddings,  $c_{ld_e}$  is the objective function value obtained with those clusters, and  $\text{SNR}\{\cdot\}$  is an operator that calculates the signal to noise ratio of its input values—i.e. the inverse of the coefficient of variation or more specifically the ratio of the values' mean and standard deviation. Furthermore,  $F_{lid_e}$  represents the  $l^{\text{th}}$  network's pseudo F-statistic obtained with the  $i^{\text{th}}$   $d_e$ -dimensional embedding result. Note that even though we calculate the  $F_{lid_e}$  values on an embedding result basis, we still use for their calculation the clusters of network  $l$  that are found using all of the embedding results. Also, the first factor in Eq. (6) is used to penalize a large number of clusters; the larger  $N_{\text{cl},ld_e}$  is, the smaller the value of the objective function becomes.

Using Eq. (6), we define the optimal dimension for the  $l^{th}$  network,  $d_{e,l}^{opt}$ , as:

$$d_{e,l}^{opt} = \underset{d_e}{\operatorname{argmax}}(c_{ld_e}) . \tag{7}$$

Once the optimal dimension at each network has been found, we use their clustering results to determine which countries cluster not only in the illicit trade network but also in the other ones.

## 6 Results and Discussion

The objective function values (cf. Eq. (6)) obtained on the different networks using the different embedding dimensions are shown in Table 2. Many node embeddings result in the clustering algorithm not converging. The licit and illicit networks are the most extreme cases of this type of behavior, as only one dimension results in vectors for which the clustering works. This might be caused by the sparsity of the licit and illicit networks. Furthermore, the results also show that in some cases the highest dimensions, i.e. 50, 70, and 100, are the ones with the best results according to Eq. (6).

**Table 2.** Objective function values obtained using embedding vectors of different dimensions; the missing values represent cases in which the affinity propagation did not converge; the best values for each network are shown in bold.

$d_e$	Illicit trade	Licit trade	Amount of traded Goods	Trade Discrepancy	Language	FTA	Geography
2	<b>1.78E-01</b>	—	6.72E-01	7.63E-01	—	2.35E+00	1.02E+00
3	—	—	1.99E+00	3.81E+00	3.80E+00	2.87E+00	5.35E+00
4	—	—	4.20E+00	2.56E+00	—	6.92E+00	9.79E+00
5	—	—	5.64E+00	3.82E+00	1.18E+01	<b>8.43E+00</b>	1.14E+01
10	—	—	1.53E+01	7.77E+00	1.09E+01	—	9.97E+00
20	—	—	—	—	—	—	1.26E+01
30	—	—	2.19E+01	2.57E+01	1.73E+01	—	1.49E+01
40	—	—	3.36E+01	—	—	—	1.55E+01
50	—	—	3.50E+01	3.52E+01	1.85E+01	—	<b>1.87E+01</b>
60	—	—	3.30E+01	3.42E+01	—	—	1.60E+01
70	—	—	3.43E+01	—	<b>2.23E+01</b>	—	1.51E+01
80	—	—	3.76E+01	—	2.21E+01	—	1.85E+01
90	—	—	3.77E+01	—	2.17E+01	—	1.80E+01
100	—	<b>3.67E+01</b>	<b>4.30E+01</b>	<b>3.54E+01</b>	—	—	1.85E+01

The clusters obtained on the illicit trade network with the best embedding dimension (i.e. two) are contained in Table 3. As we can observe, only 20 of the 146 countries listed in Table 1 are contained in Table 3. The ones missing are the ones that did not cluster with any other country within the illicit trade network. In other words, the 20 countries shown are those that have—according to the GPAS data used, cf. Sect. 4—a strong illicit trade relationship.



**Table 3.** Illicit trade clusters; the countries listed in Table 1 that are missing in Table 3 are those that did not cluster with any other country on the illicit trade network

Cluster 1	Denmark; Greece; Macedonia; Ukraine
Cluster 2	Belarus; Israel; Japan; Moldova
Cluster 3	Chile; Czech Republic
Cluster 4	Benin; Bosnia & Herzegovina
Cluster 5	Djibouti; Georgia
Cluster 6	Haiti; Macau (SAR China)
Cluster 7	Barbados; Portugal
Cluster 8	Bahamas; El Salvador

After finding the illicit trade clusters, we compare them to those of other networks. Table 4 shows the countries that cluster based on illicit trade and on at least one other aspect tested herein.

**Table 4.** Countries that cluster not only in the illicit trade network, but also on at least one of the other networks.

	Illicit trade
Licit trade	{Greece, Macedonia}; {Denmark, Ukraine}
Amount of traded Goods	{Belarus, Moldova}
Trade discrepancy	—
Language	—
FTA	{Denmark, Greece}
Geography	{Greece, Macedonia}

Table 4 shows that six countries appear to be related by illicit trade and by at least one of the other aspects considered. From these six, Greece is the one that appears the most in Table 4. The results show, that Greece’s illicit trade with Macedonia appears to be driven by licit trade and geography, while its illicit trade with Denmark could be explained by the presence of an FTA. Denmark appears again in Table 4, but now together with Ukraine. From what we can observe, it seems that the strong licit trade relationship between these countries could be a possible factor behind their illicit trade. Another pair of countries that group together are Belarus and Moldova. These two countries are shown to have a strong trade relationship (as they cluster based on their amount of traded goods), a relationship that could be facilitating illicit trade between them. The results also show, that countries that group based on trade discrepancy and/or language do not seem to cluster based on illicit trade (at least not in our data).

As exemplified by the previous results, the algorithm described herein enables us to identify possible factors that might be driving illicit trade between countries

and that might play an important role when combating this type of trade. However, we must acknowledge that this analysis is limited to the GPAS data used to represent the illicit trade. Henceforth, a future analysis with a larger and/or more diverse dataset still needs to be conducted. Additionally, a comparison of the method described herein and some other network analysis approaches should also be conducted in the future.

## 7 Conclusion and Outlook

We present a method that is able to find clusters in different types of networks (e.g., directed, undirected) by combining the creation of node embeddings and traditional clustering. With this method we can identify countries that may not only have a strong relationship in terms of illicit trade, but also in terms of some other aspect, such as trade data discrepancy, geographic proximity, etc. In other words, the method allows us to estimate factors that might be driving to some degree the illicit trade between countries. In this article, we apply the new method on data stemming from various real-world datasets. The obtained results enable us to estimate factors that could be playing an important role in the illicit trade between six different countries.

Even though our method shows potential for understanding different aspects of illicit trade, currently its results are only qualitative. Therefore, future works should try to modify the method in such a way that it will allow for more quantitative conclusions, for instance the percentage that a certain aspect (such as geography) influences illicit trade. Furthermore, we also need to compare our method to other network analysis approaches. In addition, the research of country-country relationships that we might not have considered here could be investigated in future related works. Finally, something that could also be interesting for the future is looking at cities instead of countries, as it could give us a better understanding of not only international, but also national illicit trade.

**Acknowledgments.** Jorge Ángel González Ordiano and Steven Simske acknowledge the support given by the NSF EAGER grant with the abstract number 1842577, “Advanced Analytics, Intelligence and Processes for Disrupting Operations of Illicit Supply Networks”.

## References

1. Bindner, L.: Illicit trade and terrorism financing. Centre d’Analyse du Terrorisme (CAT) (2016)
2. Brusco, M.J., Steinley, D., Stevens, J., Cradit, J.D.: Affinity propagation: an exemplar-based tool for clustering in psychological research. *Br. J. Math. Stat. Psychol.* **72**(1), 155–182 (2019)
3. Caliński, T., Harabasz, J.: A dendrite method for cluster analysis. *Commun. Stat. Theory Meth.* **3**(1), 1–27 (1974)
4. Cui, P., Wang, X., Pei, J., Zhu, W.: A survey on network embedding. *IEEE Trans. Knowl. Data Eng.* **31**(5), 833–852 (2018)

5. Fortunato, S.: Community detection in graphs. *Phys. Rep.* **486**(3–5), 75–174 (2010)
6. Frey, B.J., Dueck, D.: Clustering by passing messages between data points. *Science* **315**(5814), 972–976 (2007)
7. González Ordiano, J.A., Finn, L., Winterlich, A., Moloney, G., Simske, S.: On the analysis of illicit supply networks using variable state resolution-markov chains. In: *Proceedings of the 18th International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems (Accepted)*
8. Grover, A., Leskovec, J.: node2vec: scalable feature learning for networks. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 855–864. ACM (2016)
9. Hamilton, W.L., Ying, R., Leskovec, J.: Representation learning on graphs: methods and applications. arXiv preprint [arXiv:1709.05584](https://arxiv.org/abs/1709.05584) (2017)
10. Mackey, T.K., Liang, B.A.: The global counterfeit drug trade: patient safety and public health risks. *J. Pharm. Sci.* **100**(11), 4571–4579 (2011)
11. OECD: Trade in Counterfeit Goods Costs UK Economy Billions of Euros - 2019 Update (2019)
12. Shelley, L.I.: *Dark Commerce: How a New Illicit Economy is Threatening our Future*. YBP Print DDA. Princeton University Press, Princeton (2018)