Fraud detection through a network analysis of the Anti-Fraud database

Michele Tumminello, Andrea Consiglio

Department of Economics, Business and Statistics University of Palermo

> Riccardo Cesari, <u>Fabio Farabullini</u> IVASS

Quantitative Finance @ WORK Rome, May 3, 2019

Summary

- The Antifraud Integrated Archive (AIA)
- Bipartite Networks and statistically validated
 networks
- Network indicators and integrated indicator
- Criminal specialization, network motifs, data quality
- Conclusions

Big Data: size does matter

Big Data = Transactions + Interactions + Observations



Source: Contents of above graphic created in partnership with Teradata, Inc.

IVASS and anti-fraud

IVASS has been involved in anti-fraud activity of car insurance sector for several years; since 2001 IVASS has managed the Claims database

The Claims database contains detailed information BUT just data on claims available to insurance undertakings

IVASS dispatches periodically (centralized) information drawn from Claims database to non-life insurers

Anti-fraud Integrated Archive (AIA)

In **2012** and **2017 laws** passed, which introduce relevant innovations for fighting frauds

In particular, such laws allowed **IVASS** to collect information from external databases in order to increase the information available for anti-fraud activity

consequently IVASS has implemented a new tool called

ANTI-FRAUD INTEGRATED ARCHIVE (AIA)

AIA: stage 1



Indicators and scores (before network tools)

BINARY INDICATORS (on/off)

BUILT ON THE BASES OF RECURRENCES AND CROSS-CHECKS CRITERIA

DIFFERENT WEIGHT ACCORDING TO THE RELEVANCE IN ANTI-FRAUD ACTIVITY

Indicators and scores (before network tools)



AIA: stage 2



AIA: stage 2

Network Analysis

Big Data: AIA

- Time period: 2011-2016
- About 14 million car accidents
- About 20 million individuals and companies
- About 18 million vehicles

Tumminello M, Consiglio A, **Project (2016-2019)**: "*Network analysis and modelling of the integrated anti-fraud database*", funded by the Istituto per la Vigilanza sulle Assicurazioni (**IVASS**). Responsible for IVASS: **Farabullini F**

Heterogeneity of subjects



Objectives

- Uncover patterns in the data that suggest fraudulent activity.
- Identify organized groups of perpetrators.

Bipartite networks



Bipartite networks

 Vehicles or subjects

 OPPORT

 OPPORT

 OPPORT

 Car accidents

Null hypothesis

One does not choose the counterpart in an accident

A statistical validation of co-occurrence

Suppose there are **N** events in the investigated set. We want to statistically validate the co-occurrence of subject S_A and subject S_B in **X** events against a null hypothesis of random co-occurrence. Suppose that the number of events where $S_A(S_B)$ appears is $N_A(N_B)$, whereas the number of events where both S_A and S_B appear is **X**.



The question that characterizes the null hypothesis is: <u>what is the probability</u> <u>that number X occurs</u> <u>by chance?</u>

Tumminello M, Miccichè S, Lillo F, Piilo J, Mantegna RN (2011) Statistically Validated Networks in Bipartite Complex Systems. PLOS ONE 6(3): e17994. doi:10.1371/journal.pone.0017994 http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0017994

Hypergeometric distribution and Statistically Validated Networks

p-value associated with a detection of co-occurrences \ge X: p =

$$\sum_{i=X}^{\min(N_A,N_B)} \frac{\binom{N_A}{i} \binom{N-N_A}{N_B-i}}{\binom{N}{N_B}}$$

- Count the total number of tests: T
- Arrange *p-values* in increasing order.
- Set a link between two vertices if the associated p-value satisfies one of the following inequalities



Type I error control: false positive links

Proposition: the probability that a false positive link is set in the **Bonferroni network** is smaller than α .

Co-occurrences might be dependent

Bonferroni network

- It's the most conservative statistically validated network
- The threshold is independent of p-values
- A **co-occurence** equal to **1** is not statistically significant, provided that the number of links, E, in the co-occurrence network is larger than the number of nodes, N, in the projected set, times α

$$p - value(n_{AB} = 1, N_A, N_B, N) \ge p - value(n_{AB} = 1, 1, 1, N) = \frac{1}{N} > \frac{\alpha}{E}$$

Distinguishing between subjects and vehicles

	Nodes	Links	Connected components (CC)	Size of largest CC	
Bonferroni network of subjects *	1,197,055	1,113,389	407.552	318,876	
Bonferroni network of vehicles*	209,801	121.253	99,373	11	

*Subjects and vehicles recorded in the white list have been excluded from the analysis

Bonferroni network of subjects: largest communities

Community ID	Years over- expressed	Regions over-expressed	Provinces over-expressed
1	2015,2016	SARDEGNA, LOMBARDIA, LAZIO	VA, TV, TP, TO, SS, RM, RN, RG, PO, PT, PE, PV, PD, MI, LO, LC, LT, CO, CL, CA, BG, MB, OG, VI, VR, AG
2	2011,2012	CAMPANIA*, NA	NULL, SA, AV, NA, CE
3	-	TOSCANA*, NA	NULL, SI, PO, PT, PI, AR, LU, FI
4	-	PIEMONTE*, VALLE_D'AOSTA	VC, TO, AT, AO, CN, BI
5	-	BASILICATA, PUGLIA*, NA	NULL, BA, TA, PZ, MT, FG, BR, BT
6	-	FRIULI_VENEZIA_GIULIA, VENETO*	VE, UD, TV, RO, PN, PD, FE, VI, VR, BL
7	-	SICILIA*	TP, PA, AG
8	-	LAZIO*	RM, RI, LT, VT
9	-	SICILIA*, NA	NULL, SR, RG, ME, EN, CT, CL
10	-	EMILIA_ROMAGNA*	RN, RA, OR, MO, FC, FE, BO
11	2015,2016	LAZIO*	RM, RI, LT, FR, VT
12	2011	FRIULI_VENEZIA_GIULIA, VENETO	VE, UD, TV, PN, PD, NO, GO, VI, BL
13	-	LIGURIA, NA	NULL, SV, SP, IM, GE, AL
14	-	LAZIO, NA	NULL, RM, LT, VT
15	2015	CAMPANIA*	SA, AV, NA, CE
17	-	EMILIA_ROMAGNA*, NA	NULL, RE, PR, MO, MN, FE, BO
23	2016	LOMBARDIA	VA, PV, MI, LO, LC, CR, CO, BG, MB
25	-	LOMBARDIA, NA	PC, MN, LO, CR, BS, BG, VR

Are links robust to time-space localization?

An indicator of linkrobustness to localization

T=total number of events in the dataset (**T**=13,533,500 in AIA 10/2016) **B**=bonferroni threshold in the dataset (**B**=1.356e-10 in AIA 10/2016) **M**(i,j)=Min(Q) such that p-value(n(i),n(j),n(i,j),Q)<**B**

Robustness indicator

 $R(i,j) = log_{10}(T) - log_{10}(M)$

Bonferroni network: distribution of link-robustness (R>0.1)



Node (event, subject, vehicle) indicators of centrality

- Node degree
- Node total strength
- Node average strength
- Node betweenness

Mixed Event-subject indicators

Statistically Validated Bipartite Network

Construction: given the SVN of subjects (or vehicles), a bipartite network is reconstructed by

- selecting from the original bipartite network all of the *event(i)*subject(j) pairs such that *event(i)* contributed to a link in the SVN between subject(j) and (at least) another subject.
- adding afterwards all of the subjects directly involved in the selected events.

K-H core of a bipartite network

The K-H core of a bipartite network is the largest bipartite **subnetwork** such that nodes of Set A have degree at least K and nodes of set B have degree at least H.

Bipartite network of Kids(blue)-toys(yellow)





Network indicators: Mixed event-subject indicators of centrality: the **K-H core**

• Event oriented event-subject indicator:

 $KH_e(e, s) = \max(K)$ such that $(e, s) \in K - H$ core

• Subject oriented event-subject indicator:

 $KH_s(e, s) = \max(H)$ such that $(e, s) \in K - H$ core

Balanced event-subject indicator:

 $KH(e,s) = \max(\sqrt{K \cdot H})$ such that $(e,s) \in K - H$ core

K-H CORE DECOMPOSITION

of a real statistically validated bipartite subnetwork



An integrated indicator

Many indicators, related to both the network (system) and the event \Rightarrow correlation is observed

Find new variables that are linearly independent

Select the "most informative" (RMT)

Integrated Indicator: modelling the selected composite variables

An integrated indicator: PCA & RMT



An integrated indicator: logit model

sample containing 6.753 events occured in Italy from 2014 to 2017

3.383 events randomly sampled from AIA

3.370 reported events

Asymmetric approach, cause-effect

What is the classification ability of the principal components?

Estimation of a logit model to estimate coefficients

 $logit{\pi} = \alpha_1 * CP_1 + \alpha_2 * CP_2 + \alpha_3 * CP_3 + \alpha_4 * CP_4$

with π the probability of belonging to reported events

An integrated indicator: logit model

sample containing 6.753 events occured in Italy from 2014 to 2017

3.383 events randomly sampled from AIA

3.370 reported events

Out of sample validation

Initial dataset partitioned in two parts.

80% (5402 units) forms the training set.

20% (1351 units) the test set.

An integrated indicator: the threshold

Initial dataset partitioned in two parts.

80% (5402 units) forms the training set.

20% (1351 units) the test set.

Maximize the Matthews Correlation Coefficient in the training set to select the threshold x_0

Results	out of sampl	le)):
	\		·

	Random	Reported	
$X \leq x_0$	82% (3%)	47% (3%)	
$X > x_0$	18% (3%)	53% (3%)	
	100%	100%	

AIA: stage 3

- Three node motifs
- Network analysis for data quality

Motifs: the heuristics

- Criminal specialization
- Some types of crime require cooperation
- Cooperating with a criminal intent requires secrecy and trust



M Tumminello, C Edling, F Liljeros, RN Mantegna, J Sarnecki (2013) The Phenomenology of Specialization of Criminal Suspects. PLoS ONE 8(5): e64703. doi:10.1371/journal.pone.0064703

Motifs and anti-fraud

Not suspicious

Suspicious



Three-node motifs: statistically validated triangles



Proposition: if random co-occurrence of three subjects, 1,2, and 3, involved in n_1 , n_2 , and n_3 events, respectively, is assumed in a dataset including N events then

$$p(n_{12}^*, n_{13}^*, n_{23}^* | n_1, n_2, n_3, N) = \frac{\binom{n_1}{n_{12}} \binom{N-n_1}{n_2-n_{12}} \binom{n_{12}}{n_{12}-n_{12}^*} \binom{n_1-n_{12}}{n_{13}^*} \binom{n_2-n_{12}}{n_{23}^*} \binom{N-n_1-n_2+n_{12}}{n_3-n_{13}^*-n_{23}^*-n_{12}+n_{12}^*}}{\binom{N}{n_2} \binom{N}{n_3}}$$

p-value = $p\left(n_{12}^* + n_{13}^* + n_{23}^* \ge n_{12}^{*,0} + n_{13}^{*,0} + n_{23}^{*,0}\right)$

Three-node motifs and antifraud

Network of directly involved subjects (no professionals)

- Number of triangles: 162,409
- Number of statistically validated triangles:60,523

Randomly rewired network of directly involved subjects

- Average number of triangles: 18,535
- Average Number of statistically validated triangles: 0.08

Data quality: the statistically validated network of accidents





Final Remarks

- 1. The **network** of **subjects** and **vehicles** carry different information.
- 2. Introduced network indicators and IVASS subject indicators carry complementary information, and, therefore, can fruitfully be integrated: the **integrated indicator**.
- 3. The test on "claims closed following investigation" and the analysis of a few case studies on already identified criminal networks indicate the effectiveness of the overall approach.
- 4. Introduced network indicators are operative since March 2018 (IVASS internal use).
- 5. Next steps: (a) integrating three-node motifs in the SVN (exp. Sep 2019); (b) SVN of accidents for data quality (exp. end 2019).

Thanks!

Michele Tumminello

Email: <u>michele.tumminello@unipa.it</u> Alt. Email: <u>michele.tumminello@gmail.com</u>