

General

ID ¹				
Use case name	Detection of frauds based on collusions			
Context	Fintech			
Application domain	On-premise systems			
Status	In operation			
Contributor	Name	Affiliation	Contact	
	Girish Palshikar C. Anantaram	Tata Consultancy Services Ltd.	c.anantaram@tcs.com	
Scope ²	Validating the predicted collusion set is effort-intensive and needs investigative and legal expertise			
Objective(s)	Automatic unsupervised detection of frauds based on collusions			
Narrative	Short description (not more than 150 words)	Our tool includes a set of unsupervised machine learning algorithms to detect collusion-based frauds, particularly, circular trading and price manipulation in stock market trading		
	Complete description	<p>Frauds are prevalent across all industries; and they are particularly severe in today's computerized, web-connected, mobile-accessible, and cloud-enabled business environments. An FBI report states that the insurance industry in the US, which consists of over 7000 companies and collects over \$1 trillion in premiums, loses about \$40 billion annually in frauds in the non-health insurance sector alone. The aggregate size of the 52 regulated stock exchanges across the world (total market capitalization) was \$55 trillion as on Dec. 2012. Given the money involved, it is not surprising that the stock market is a target of frauds.</p> <p>Many malpractices in stock market trading, e.g. circular trading and price manipulation—use the modus operandi of collusion. Informally, a set of traders is a candidate collusion set when they have “heavy trading” among themselves, as compared to their trading with others. We formalize the problem of detection of collusion sets, if any, in a given trading database. We show that naïve approaches are inefficient for real-life situations. We adapt and apply two well-known graph clustering algorithms for this problem. We also propose a new graph clustering algorithm, specifically tailored for detecting collusion Sets; further, we establish a combined collusion set. Treating individual experiments as evidence, this approach allows us to quantify the confidence (or belief) in the candidate collusion sets. We have carried out detailed simulation experiments to demonstrate effectiveness of the proposed algorithms. The system is also operational in a government organization. Note that all our collusion detection algorithms are completely unsupervised and do not need any training data.</p>		
Key performance indicators (KPIs)	ID	Name	Description	Reference to mentioned use case objectives
	1	Prediction accuracy	How many predicted collusion sets	Improve accuracy

		were actually involved in frauds
AI features	Taks(s)	Knowledge processing & discovery
	Method(s) ³	Machine learning
	Hardware ⁴	Windows
	Terms and concepts used ⁵	
Challenges and issues	Challenges: Actual examples of collusion-based frauds may not be available easily, even for evaluation and testing	
Societal concerns		

Data (optional)

Data characteristics	
Description	
Source ⁶	
Type ⁷	
Volume (size)	
Velocity (e.g. real time) ⁸	
Variety (multiple datasets) ⁹	
Variability (rate of change) ¹⁰	
Quality ¹¹	

Training (optional)

Scenario name	Training				
Step No.	Event ¹⁴	Name of process/Activity ¹⁵	Primary actor	Description of process/activity	Requirement

Specification of training data ¹⁶	
--	--

Evaluation (optional)

Scenario name	Evaluation				
Step No.	Event ¹⁷	Name of process/Activity ¹⁸	Primary actor	Description of process/activity	Requirement

Input of evaluation ¹⁹	
Output of evaluation ²⁰	

Execution (optional)

Scenario name	Execution				
Step No.	Event ²¹	Name of process/Activity ²²	Primary actor	Description of process/activity	Requirement

Input of Execution ²³	
Output of Execution ²⁴	

References

References						
No.	Type	Reference	Status	Impact on use case	Originator/organization	Link
1	Conference				Tata Consultancy Services Limited	D. K. Luna, G. K. Palshikar, M. Apte, A. Bhattacharya, <i>Finding Shell Company Accounts using Anomaly Detection</i> , ACM India Joint International Conference on Data Science and Management (CoDS-COMAD 2018) , Goa, India, Jan 11-13, 2018
2	Journal				Tata Consultancy Services Limited	G. K. Palshikar, M. Apte, <i>Collusion Set Detection Using Graph Clustering</i> , vol. 16, no. 2, April 2008, Data Mining and Knowledge Discovery journal (Springer-Verlag), pp. 135 – 164
3	Book chapter				Tata Consultancy Services Limited	M. Apte, G.K. Palshikar, S. Baskaran, <i>Frauds in Online Social Networks: A Review</i> , accepted as a Book Chapter, in Social Network and Surveillance for Society , T. Ozyer and S. Bakshi (ed.s), to be published by Springer in 2018
4	Book chapter				Tata Consultancy Services Limited	G.K. Palshikar, M. Apte, <i>Financial Security against Money Laundering: A Survey</i> , Chapter 36 in B. Akhgar, H.R. Arabnia (Ed.s), Emerging Trends in Information and Communication Technologies

						Security , pp. 577 – 590, Elsevier (Morgan Kaufman), 2013
--	--	--	--	--	--	--

-
- ¹ Leave this cell blank.
 - ² The scope defines the limits of the use case.
 - ³ AI method(s)/framework(s) used.
 - ⁴ Hardware system used.
 - ⁵ Terms and concepts listed here can be used to extend the work of WG 1 (AWI 22989 and AWI 23053) as necessary.
 - ⁶ Origin of data, which could be from instruments, IoT, web, surveys, commercial activity, or from simulations.
 - ⁷ Structured/unstructured Images, voices, text, gene sequences, and numerical. Composite: time-series, graph-structured
 - ⁸ The rate of flow at which the data is created, stored, analysed, or visualized.
 - ⁹ Data from a number of domains and a number of data types. The wider range of data formats, logical models, timescales, and semantics complicates the integration of the variety of data.
 - ¹⁰ Changes in data rate, format/structure, semantics, and/or quality.
 - ¹¹ Completeness and accuracy of the data with respect to semantic content as well as syntactical of the data (such as presence of missing fields or incorrect values)
 - ¹² Describe which condition(s) should have been met before this scenario happens.
 - ¹³ Describe which condition(s) should prevail after this scenario happens. The post-condition may also define "success" or "failure" conditions.
 - ¹⁴ The event that triggers the step. This might be completion of the previous event.
 - ¹⁵ Action verbs should be used when naming activity.
 - ¹⁶ Training data can be further specified.
 - ¹⁷ The event that triggers the step. This might be completion of the previous event.
 - ¹⁸ Action verbs should be used when naming activity.
 - ¹⁹ Specify input of evaluation.
 - ²⁰ Specify output of evaluation.
 - ²¹ The event that triggers the step. This might be completion of the previous event.
 - ²² Action verbs should be used when naming activity.
 - ²³ Specify input of evaluation.
 - ²⁴ Specify output of evaluation.
 - ²⁵ The event that triggers the step. This might be completion of the previous event.
 - ²⁶ Action verbs should be used when naming activity.

²⁷ Retraining data can be further specified.