

Identifying and Aggregating Informative Audit Quality Indicators Using Machine Learning

Abigail Zhang, Soo Hyun Cho, Miklos Vasarhelyi

Introduction

Background

The PCAOB proposes 28 potential Audit Quality Indicators (AQIs). However, most of the AQIs proposed by the PCAOB are proprietary to audit firms and are not publicly accessible, complicating efforts to collect, analyze, and study AQI (PCAOB 2015). To facilitate AQI research, the PCAOB has encouraged academia to derive AQIs from public source information and to identify which of these indicators are the most informative and predictive in evaluating audit quality (PCAOB 2013a; PCAOB 2015).

Research Objectives

1. Identify IAQIs, which are theory-driven audit-related variables that are the most predictive of audit failure.
2. Aggregate IAQIs into predictive audit quality indexes that can red flag potential audit failures.

Methodology

We adopt a machine learning methodology to identify a portfolio of Informative Audit Quality Indicators (IAQIs) - publicly available audit-related variables that can best predict audit failure. Based on prior literature, we use material restatements of annual reports as a proxy for audit failure.

Figure 1 presents our overall research design. To identify IAQIs, we perform Feature Subset Selection (FSS) using five popular machine learning algorithms to select a subset of ARVs that can best predict MAR. Then, we assess the predictive power of IAQI by inputting IAQIs into multiple machine learning algorithms to predict MAR via a Cost-Sensitive Learning (CSL) and Rolling-Window Prediction (RWP) mechanism. In aggregating IAQIs into PAQI, we first select the algorithm with the best overall predictive ability; then, we input IAQIs into the chosen algorithm to obtain probability prediction via CSL and RWP; lastly, we rescale the probability prediction to obtain PAQI.

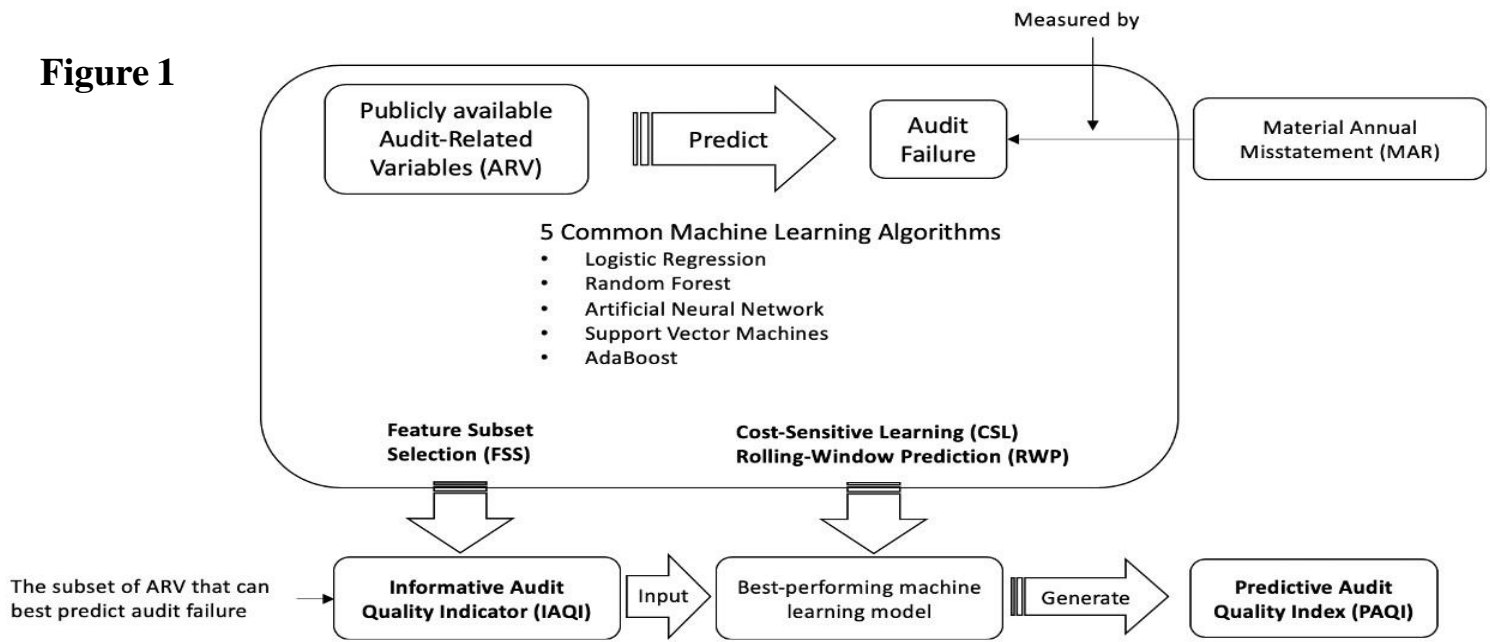
Main Results

Table 1 presents the 11 audit-related variables that we identify as IAQIs after validating their predictive power for audit failure. Our results show that audit engagements with higher PAQIs are significantly more likely to have actual MARs when other factors are controlled. Specifically, a 1-point increase in the PAQI increases the odds of having an actual MAR by 1.51 times.

Table 1. Informative Audit Quality Indicators (IAQIs)

Category	Sub-category	Aspects Captured	Variable	Measurement
Audit input	Auditor characteristics	Competence, Resource, Independence	Office Size	Natural logarithm of one plus total annual audit fees of an audit office (Aobdia 2019)
Audit process	Task characteristics	Informational advantage, Independence	Tenure	Number of years that the company is audited by the same audit firm (Bell et al. 2015)
		Audit effort, Audit efficiency	Audit Report Lag	Natural logarithm of the number of days between fiscal year-end and the signature date of audit opinion (Lobo and Zhao 2013)
		Audit effort, Workload	Integrated Audit	Indicator variable equal to 1 when the audit engagement is an integrated audit of financial statements and internal controls, and 0 otherwise (Aobdia 2019)
	Auditor-client contracting features	Incentive	Auditor Resignation	Indicator variable equal to 1 if the current auditor will resign (instead of being dismissed by the company) from the next fiscal year, 0 otherwise (Krishnan and Krishnan 1997)
		Audit effort	Audit Fees	Natural logarithm of 1 plus the audit fees charged to the auditee (Aobdia 2019)
	Auditor communications	Competence, Independence	Internal Control Weakness	Indicator variable equal to 1 if a material weakness is reported for the year, 0 otherwise (Aobdia 2019)
Audit output	Quality of the audited financial statements	Within-GAAP manipulation	Disc. Accruals	Residual from the cross-sectional modified Jones model in Aobdia (2019)
			Abs (Disc. Accruals)	The absolute value of Disc. Accruals (Aobdia 2019)
			Abs (Accruals)	The absolute value of accruals deflated by beginning assets (Aobdia 2019)
			Abs (Accruals/CFO)	The absolute value of accruals deflated by cash flow from operations (Aobdia 2019)

Figure 1



The Impact of Mass Layoffs in IT Firms on Cybersecurity: Evidence from the Darknet Market

Arion Cheong, Michael Alles, Soohyun Cho, Won Gyun No, Miklos A. Vasarhelyi

Research Questions:

RQ1: Are the layoffs in IT firms during the COVID-19 pandemic associated with their darknet market exposures (i.e., the number of firm-related information posted in the darknet market)?

RQ2: After massive layoffs, do IT firms have more cybersecurity disclosure while experiencing an increase in darknet market exposure?

Layoff and Darknet Market Exposure

To understand how layoff in an IT firm affects its cybersecurity environment, this study examines whether the darknet market exposure of the firm increases after the layoff. In essence, the purpose of the study is to examine the prevailing management's assertion that the cybersecurity function of its firm is well maintained even after the mass layoff. Our results indicate that the darknet market exposure significantly increases after a certain period, which is shown to be approximately 45 days from our analysis. Further, we perform textual analysis to categorize the darknet market posts regarding 1,416 firms (527 layoff firms and 889 hiring firms) and examine whether the type of information traded in the darknet market differs after layoff or hiring. In specific, we conduct a Latent Dirichlet Allocation (LDA) Topic Modeling to understand the representative topics discussed in the posts and measure the amount of discussion on certain topics

Categories of Contents of Darknet Market Posts (Observations: 527 Layoff firms and 889 Hiring firms)			Pair-wised t-test (t-statistic / p-value)	
			(-60, +60)	
Type of Information	Keywords		Layoff	Hiring
Type 1	Payment Information	Credit, Card, Carders, Verified, Malware	4.563 (0.000***)	5.921 (0.000***)
Type 2	Company-owned Devices	Renegade, Mobile, Version, Machine, Password	-6.481 (0.000***)	-8.380 (0.000***)
Type 3	Personal Information	State, Street, Road, Unit, Blvd	2.717 (0.006***)	2.083 (0.037**)
Type 4	Telecommunication	Jabber, Telegram, E-mail, Phone, Log	2.757 (0.006***)	0.298 (0.765)
Type 5	Hacking Service	Password, Hash, Salt, Checker, Service	-2.949 (0.003***)	-0.638 (0.494)
Type 6	Financial Information	Info, Balance, Track, Transfer, Chase	0.236 (0.812)	1.208 (0.227)
Type 7	Hacking Instruction	Carding, Saint, Island, Proxy, Strongbox	0.685 (0.493)	1.326 (0.185)
Type 8	Hacking Tools	Carding, Identity, Java, Guide, Tools	0.730 (0.465)	-0.206 (0.836)
Type 9	User Credentials	Account, Password, Login, Dump, Paste	0.033 (0.973)	-0.091 (0.926)
Type 10	Insider Information	Manager, Senior, Trade, Info, Private	-1.509 (0.131)	-0.166 (0.867)

Our result, derived from both logistic regression analysis and textual analysis, suggests that recent layoff increases the firm's darknet market exposure, specifically focused on sensitive PII and telecommunication information. The result triggers an alert that layoff firms are lacking the capability to secure sensitive information where such leakage can lead to material litigation risk and being more vulnerable to threats related to the remote working environment.

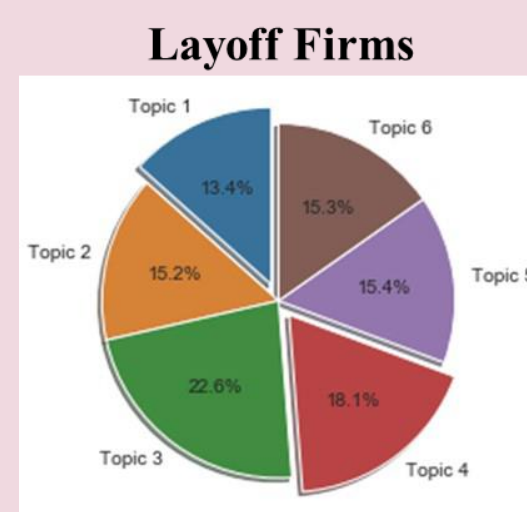
Conclusion

This paper discusses the impact of layoff on IT firms' darknet market exposures. From our analysis, we find that IT firms that layoff during the COVID-19 pandemic have experienced a significant increase in darknet market exposure. Furthermore, our results reveal that layoff firms do not provide a sufficient amount of disclosures corresponding to the increased darknet market exposure. Most importantly, our findings inform that layoff firms are lacking the capability to sustain security for the remote working environment. Our analysis alerts firms that sensitive information, including employee credentials and PII, are publicly traded in the darknet market.

Cybersecurity Reporting and Darknet Market Exposure

We compare the amount of disclosure related to security breach between layoff firms and hiring firms. Similar to our analysis on the type of information traded in the darknet market, we classify the cybersecurity disclosures into six different types: Security Breach, Security Oversight, Third-party Risk, Regulation Compliance, Virtual Work Environment, and COVID-19 Security Issue (see Table 6). Our analysis aims to examine whether firms communicate darknet market exposure in their cybersecurity disclosures.

Disclosure Classification (with Higher Darknet Market Exposure)		
Disclosure Topics	Keywords	
Topic 1	Security Breach	Breach, Incident, Attack, Event, Current
Topic 2	Security Oversight	Audit, Oversight, Board, Director, Report
Topic 3	Third-party Risk	Third-party, Security, Risk, Control, Failure
Topic 4	Regulation Compliance	Data, Regulation, Protection, United, European
Topic 5	Remote Working	Platform, Software, VMware, Solution, Provider
Topic 6	Covid-19 Security Issue	Covid-19, Obligation, Compliance, Provision, Litigation



Compared to hiring firms, layoff firms disclose less information about the security breach and more information about regulation compliance. Notably, only approximately 15.4% of layoff firms' cybersecurity disclosures are related to the remote working environment. We have also manually examined the disclosures related to the security breach for hiring firms. We find that hiring firms are providing relatively higher amount of information about the data exposed (Topic 1 – Security Breach) in the darknet market compared to layoff firms.

Contact Information:

Arion Cheong
e-mail: arion.cheong@rutgers.edu

Using Supervised Learning Algorithms to Predict Non-profits Discontinued Operation

Chengzhang Wu and Richard B. Dull

Introduction

- Nonprofit organizations play an important role in economy worldwide. The operation of those organizations primarily rely on the donation from donors. However, the discontinued operation of nonprofit organizations due to financial reasons have brought the problem of unbalanced economy resource allocation.
- The prediction of bankruptcy or dissolution of for-profit companies has been widely studied using different data analytic methodologies, and various machine learning approaches have been evaluated and demonstrated effective. However, the prediction of non-profits discontinued operation is still rarely studied.
- This study attempts to do such prediction comparing the performance of Logistic Regression, Decision Tree, Random Forest, Multilayer Perceptron, Support Vector Machine and Bayes Net. The overall effectiveness of different prediction performance will be assessed.

Research Questions

- RQ1: Among all the algorithms compared in this research, which algorithm performs more effectively?
- RQ2: In the group of predictors used in this study, which set of variables can be used to get better prediction results?
- RQ3: Using the predictor set(s) determined in RQ2, does cost-sensitive learning algorithm improve prediction quality?

Predicting Factors

- The predictors used in the current research are based on two prior literature.
- First, Howard P. Tuckman and Chang (1991) identifies four criteria used to indicate the financial vulnerability of nonprofits. These criteria include inadequate equity balances, revenue concentration, administrative costs and operating margins. These four financial variables are used to predict discontinued operations of nonprofit organizations. These four predictors are identified as financial vulnerability (FV) set in this research.
- The second source of predictors are derived from William J. Ritchie and Kolodinsky (2003) research. In that study, 16 financial performance measurement ratios used for nonprofits are classified into four categories. These four categories are Fiscal Performance (FP), Fundraising Efficiency (FE), Public Support (PS) and Investment Performance and Concentration (IPC).

Data Collection

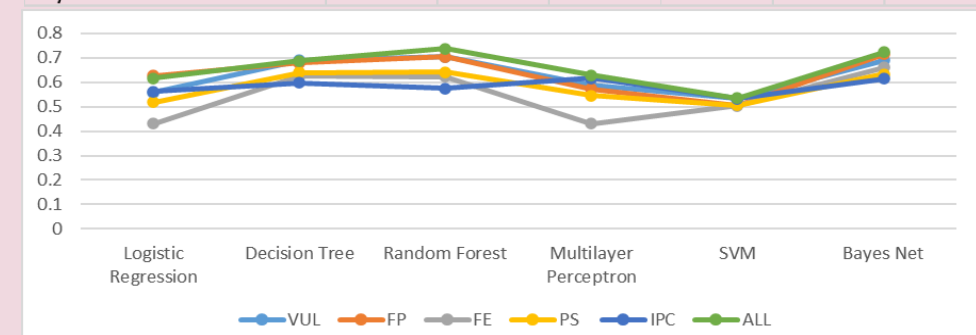
- This study uses a Form 990 database, which contains all the e-filings from 2011 to 2018. SQL is used to a query in the Form 990 data to list the all the unique records that indicated termination of operation and the year of termination.
- The number of nonprofits that discontinued operation in Form 990 accounts only 1% of the entire dataset. This algorithm is not able to capture enough information on the imbalanced dataset. Therefore, it is necessary to adjust the original dataset to be balanced to be suitable for training the model.
- The study uses under sampling technique to address the imbalanced training set. In the training set, 831 records that indicate discontinued operation in electronically filed Form 990. The other 831 records that do not indicate discontinued operation are randomly selected. The final training set used in this study includes 1662 records.

Evaluation Metrics

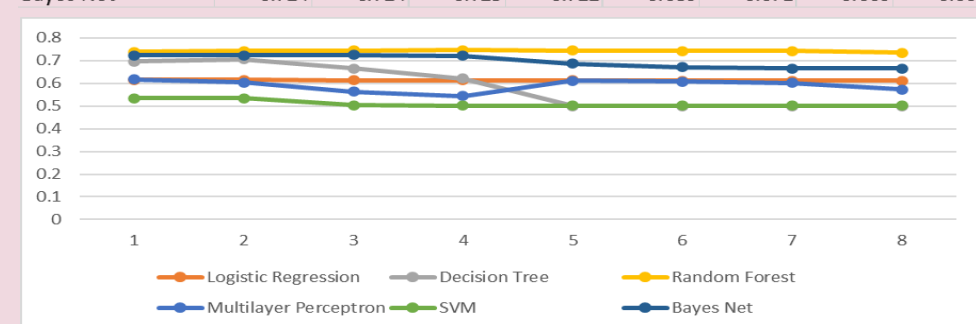
- When the binary classification prediction model is applied on test set, four categories of results could come out: True Positive, False Positive, True Negative and False Negative. The following confusion matrix depicts these four outcomes.
- Based on the matrix, there are several measurements can be used to evaluate the performance of each algorithm (listed in the following table). Considering the imbalanced nature of the population, this study uses AUC to evaluate how effective is the trained model perform on test set. The closer to 1, the better the performance.

Results

Variables	VUL	FP	FE	PS	IPC	ALL
Logistic Regression	0.559	0.627	0.431	0.518	0.563	0.616
Decision Tree	0.689	0.679	0.626	0.64	0.598	0.688
Random Forest	0.705	0.705	0.624	0.641	0.574	0.739
Multilayer Perceptron	0.587	0.572	0.431	0.545	0.617	0.63
SVM	0.534	0.507	0.504	0.507	0.533	0.535
Bayes Net	0.69	0.714	0.661	0.633	0.615	0.724



Cost	1	2	3	5	10	15	20	25
Logistic Regression	0.616	0.615	0.614	0.613	0.613	0.612	0.612	0.612
Decision Tree	0.697	0.706	0.665	0.621	0.5	0.5	0.5	0.5
Random Forest	0.739	0.742	0.745	0.747	0.744	0.743	0.743	0.735
Multilayer Perceptron	0.617	0.604	0.563	0.544	0.611	0.607	0.603	0.573
SVM	0.535	0.535	0.503	0.502	0.501	0.501	0.501	0.501
Bayes Net	0.724	0.724	0.725	0.722	0.686	0.671	0.666	0.666



Conclusion and Discussion

- The highest AUC value is achieved under ALL group using Random Forest. This indicates that among all the six algorithms, Random Forest generates the better performance.
- All the five sets of financial ratios together is more effective at predicting nonprofits discontinued operation.
- The performance of different algorithms stays relatively stable with the increment of positive to false negative cost ratio. Overall, the prediction of nonprofits discontinued operation is not sensitive to the change of positive to false negative cost ratio.
- This study provides an effective model that utilize the publicly available data in Form 990 to predict their dissolution. Even though closure of nonprofits is primarily due to completion of mission statement, the result of this study shows that the financial factor also plays an important part in the process of nonprofits dissolution.

Government Contracts To Blockchain Based Smart Contracts

Eid Alotaibi and Hussein Issa

Smart-Contracts to Compliance Audits

The tremendous amount of spending on different governmental projects associated with questionable oversight is a big concern of the citizens. Auditors have a public duty to ensure that the agencies' spending is accurate and public resources are secure, but there are many challenges. This study's motivation is that government auditors face more critical approaches to perform compliance audits for government contracts (Nation et al. 2019; Branson et al. 2011). Smart contracts, as named, is a form of contract which execute by specified terms and conditions to meet the agreement between two or more parties. The Blockchain is a verifiable and secure way of managing the records where no one can make the change in a document without the prior authority of 50% or more of participants. The infrastructure of Blockchain provides smart contracts and distributed ledgers to application users. Smart contracts are translating the requirements of an agreement to the Blockchain, and smart contracts automatically execute the agreement upon fulfillment of the requirements. The Blockchain system and its smart contracts can contribute to automating compliance audits and increase test to all of the government agencies contracts' or nearly 90%. The contribution of this study is to develop a conceptual model for the governmental compliance audit in the domain of accounting and auditing literature. Also, this study demonstrates the process of testing the proposed framework.

What

–Using Smart-contracts in blockchain as an automatic method used to perform compliance audit.

Why

–Automating data verification in real time
–Monitoring control automatically on a more frequent basis
–Enabling auditors to extract audit evidences from one source

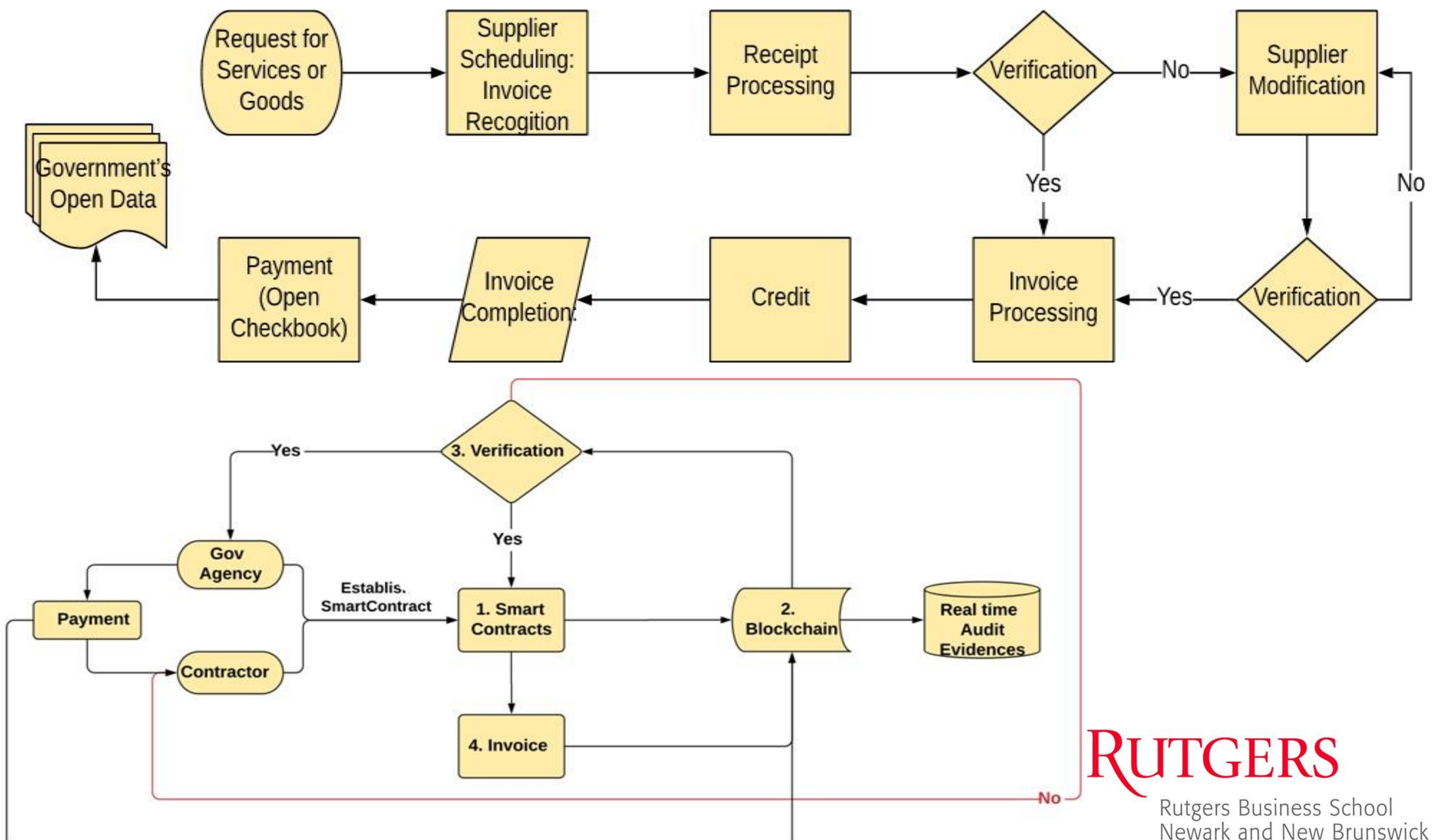
How

The approach consists of three steps:

1. Describe the core concepts and features of smart-contracts
2. Proposed a blockchain based smart-contracts framework to automate compliance audits
3. Test our framework by performing compliance audits to a purchase order

The Framework

The framework of this study automates compliance audits using smart contracts that would continuously enforce the events of a contract and flag any violations. Therefore, the system blocks and alerts of events that do not meet the identified rules rather than relying on after the fact concept. The system also can automate the process of contracting and enables real-time controls to have a systematic compliance audit. Since we are examining the compliance of regulations and rules for contracts, the framework combines three main components of any government contracts. The first component consists of two main processes. The first process is the contracting process which shows the process of how government agencies choose their contractors. The second process is purchase process which demonstrate the steps from ordering to invoicing. In the second component, contracting process and purchase process are combined as a standard smart contract using the Resources-Events-Agents model. The third component of our framework, active enforcement, basically ensures that regular activities comply with prescribed process flows. The first component of our framework understands the government procurement processes and purchase order processes in 1st figure below. The 2nd figure is the second component of the framework showing the active enforcement process to perform compliance audits, basically ensuring that regular activities comply with prescribed process flows.



Privacy Policy Disclosure Differences based on Industry

Hanchi Gu, Qing Huang, Arion Cheng, Won Gyun No

Background

Nowadays, privacy is an increasingly important topic. However, there is still a limited number of researches regarding privacy policy analysis. The privacy policy of the companies shows how they deal with privacy problems. Different privacy policies show us whether companies pay attention to each privacy issue and their different methods to handle this issue.

This study aims to investigate the differences in privacy disclosures as well as the reasons. We will first measure the contents which are disclosed in their privacy policy. Then, we will do further analyses based on their firm size, industry, and so on. Our goal is to find out what's the most influential features that affect the privacy policies' contents.

Methodology

First, we gathered the company's website from Compustat.

Second, we refer to the existing literature and add some new variables to define the variables of policy content we will use.

Third, we collect the privacy and cookie policies from the companies' websites and manually label a small sample.

Fourth, we extract the keywords for each variable based on our labeled dataset. Then, we apply them to a larger dataset to label the large dataset.

Finally, we test the relationship between variables representing the policy content and variables from other variables of the firms, not related to policy. As a result, we can explain what are the most influential factors to a company's privacy policy.

Preliminary Results

Descriptive Analysis

- In terms of small and micro companies, information companies tend to have a privacy policy compared with finance and insurance companies.
- The larger a company is, the higher possibility that it will have a complete privacy policy
- Most companies have a privacy policy; however, only small proportion of companies have separate cookie policies.
- Small companies are reluctant to disclose how data are retained.

GDPR

1. Receive users' consent before you use any cookies **except** strictly necessary cookies.
2. Provide accurate and specific information about the data each cookie tracks and its purpose in plain language before consent is received.
3. Make it as easy for users to withdraw their consent as it was for them to give their consent in the first place
4. Document and store consent received from users
5. Allow users to access your service even if they refuse to allow the use of certain cookies

The identity and contact details of the organization, its representative, and its Data Protection Officer
The purpose for the organization to process an individual's personal data and its legal basis
The legitimate interests of the organization (or third party, where applicable)
Any recipient or categories of recipients of an individual's data
The details regarding any transfer of personal data to a third country and the safeguards taken
The retention period or criteria used to determine the retention period of the data
The existence of each data subject's rights
The right to withdraw consent at any time (where relevant)
The right to lodge a complaint with a supervisory authority
Whether the provision of personal data is part of a statutory or contractual requirement or obligation and the possible consequences of failing to provide the personal data
The existence of an automated decision-making system, including profiling, and information about how this system has been set up, the significance, and the consequences

Literature Factors(Jamal, Maier & Sunder, 2004)

Criteria

- Disclose each cookie tracks and its purpose in plain language before consent is received.
- Explain what cookies are
- Has separate cookie policy or not
- Explain how to turn off/decline cookies

Cookies Policy

- Disclose the identity and contact details of the organization
- Disclose the purpose for the organization to process an individual's personal data
- Disclose whether the personal data are transfer to a third country
- Disclose the retention period or criteria of the data
- Disclose the existence of each data subject's rights
- Disclose the possible consequences of failing to provide the personal data

General Privacy Policy

- Post Privacy Policy or not
- Privacy Policy is one-click away
- Disclose how data are used

Table 1

Disclosure of Privacy Polices

Number	Privacy Practice	Micro Firm (n=10)	Small Firm (n=10)	Medium Firm (n=10)	Large Firm (n=20)
1	Post a Privacy Policy	6	8	9	20
2	Privacy policy is one click away	4	7	8	20
3	Disclose the purpose for process an individual's personal data	4	8	9	19
4	Disclose the identity and contact details of the organization	3	6	9	19
5	Disclose the existence of each data subject's rights	1	4	6	17
6	Disclose that web site is using cookies	4	8	6	16
7	Explain what cookies are	1	4	4	15
8	Explain how to turn off/decline cookies	3	3	3	12
9	Disclose presence of third-party cookies on web site	1	1	4	14
10	Disclose how data are used	2	3	7	20
11	Has separate cookies policy or not	0	0	2	10
12	Disclose whether the personal data are transfer to a third country	0	1	3	8
13	Disclose the retention period or criteria of the data	0	5	5	15
14	Disclose the possible consequences of failing to provide the personal data	3	6	7	17

Application of Interactive Visualization on Internal Tax Database

Heejae (Erica) Lee and Miklos Vasarhelyi

Background

Due to significant increase both in volume and complexity of business transactions, it becomes more challenging to utilize accounting data to make optimal business decisions. Many practitioners and researchers in the accounting field have recognized interactive data visualization as a powerful tool to improve understanding of data and its patterns, trends and relationships. In this research, we conducted a case study of implementing interactive data visualization in the firm's internal tax database. The paper can extend the continuous monitoring and auditing literature to tax compliance perspectives.

Case

We implemented continuous monitoring and data analytic layer on the top of the internal tax database of a global information and analytic company. The company has three different database. We created a consolidated view of three database for data analytics. The firm expects to improve the efficiency and effectiveness by ensuring indirect tax compliance across all geographies thereby helping mitigate the risk of penalties and reputational risk of the company by adding data visualization layer in their system. We used Tableau desktop software to create dashboards and visualizations.

Scope

There are three different scopes in the case study.

- Improve data integrity of the internal tax database system (ex. missing mandatory fields)
- Examine the trend of transaction value, tax value, and exempt value (yearly or monthly level) by various dimensions
- Provide trend alert to detect abnormal transactions to the users.

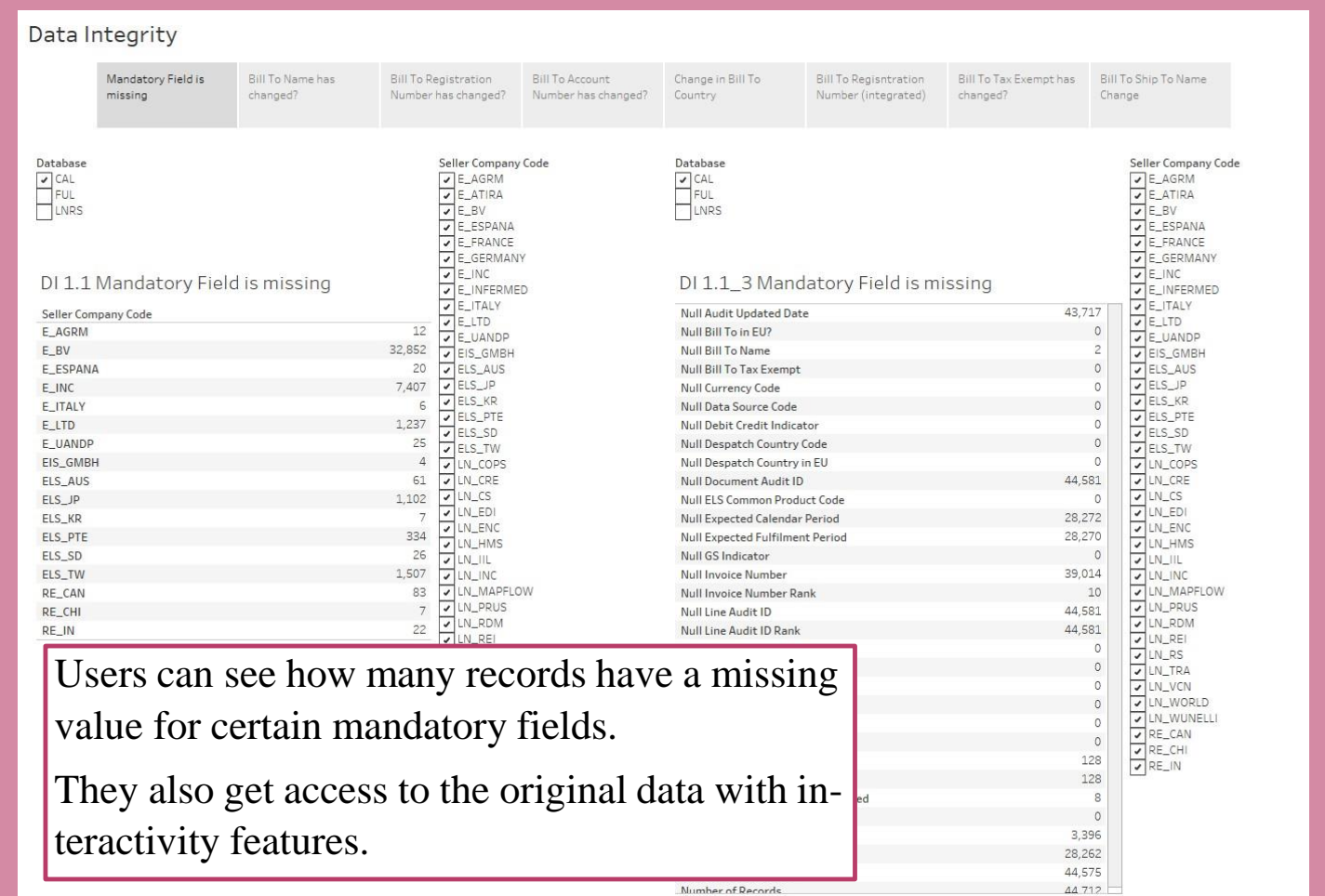
Data Descriptions

Total 6,626,747 transactions reported in 2018 were used to create dashboards and visualizations. While 80% of transactions were in US dollars, there were transactions in other currencies including Pound (11%), Euro (6.6%) and Japanese Yen(1%). We joined the transaction data file with currency table using currency code to get better view of trend analysis. Specifically, we used month end currency units per dollar to calculate transactions volume of each transactions in US dollar.

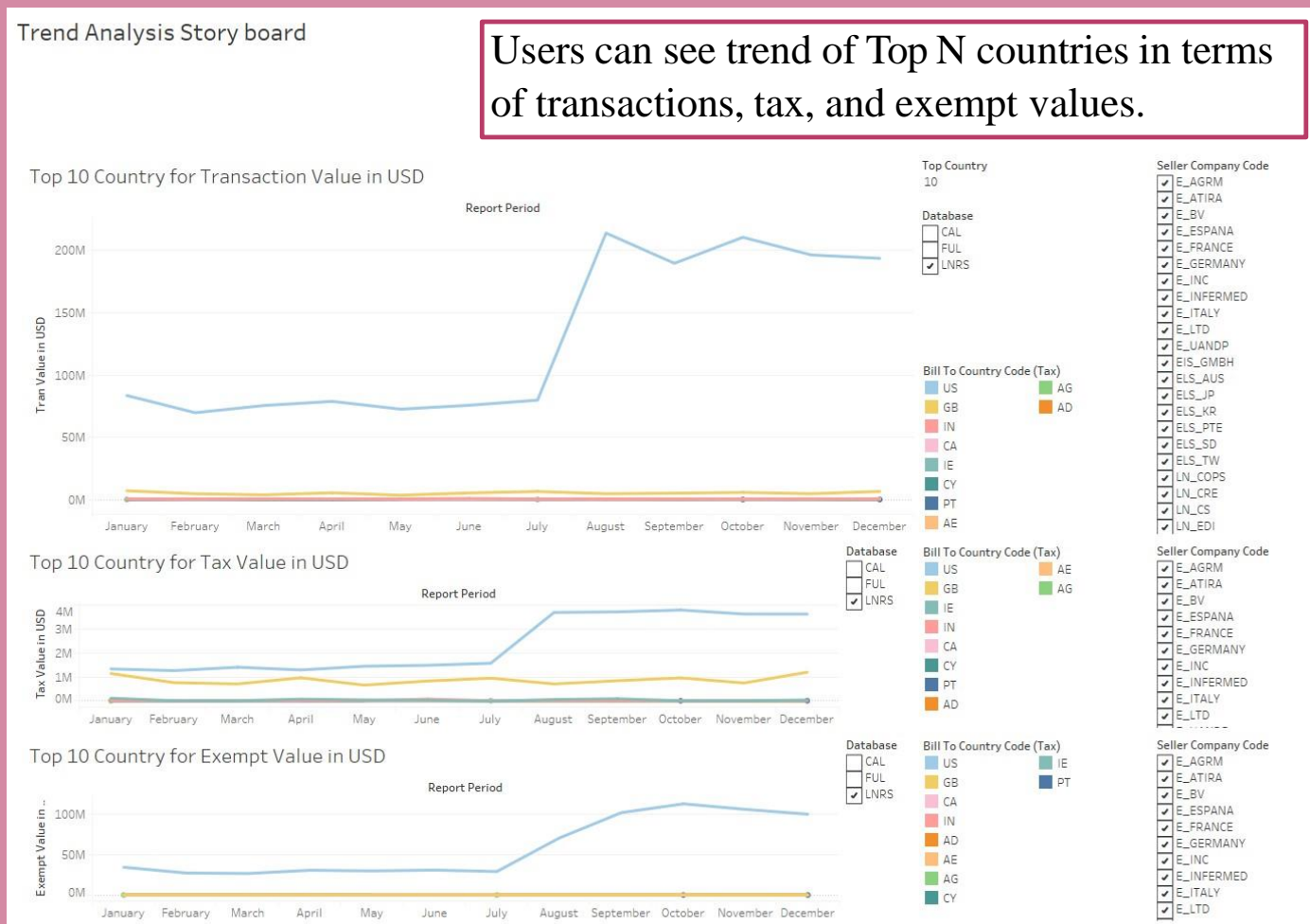
Average transactions value was \$850 per transaction and average tax value was \$9.48 per transaction. Around 33% of the transactions were tax exempt. If we excluded tax exempt transactions, the average tax value was \$14.95. Their products and services were billed and shipped to more than 200 countries around the world.

	Database 1	Database 2	Database 3
# of Transactions	660,710	3,439,324	2,526,713
Avg Trans Value	\$4,681	\$256	\$655
Avg Tax Value	\$27.47	\$2.28	\$15.53
Avg Exempt Value	\$1,407	\$166	\$285

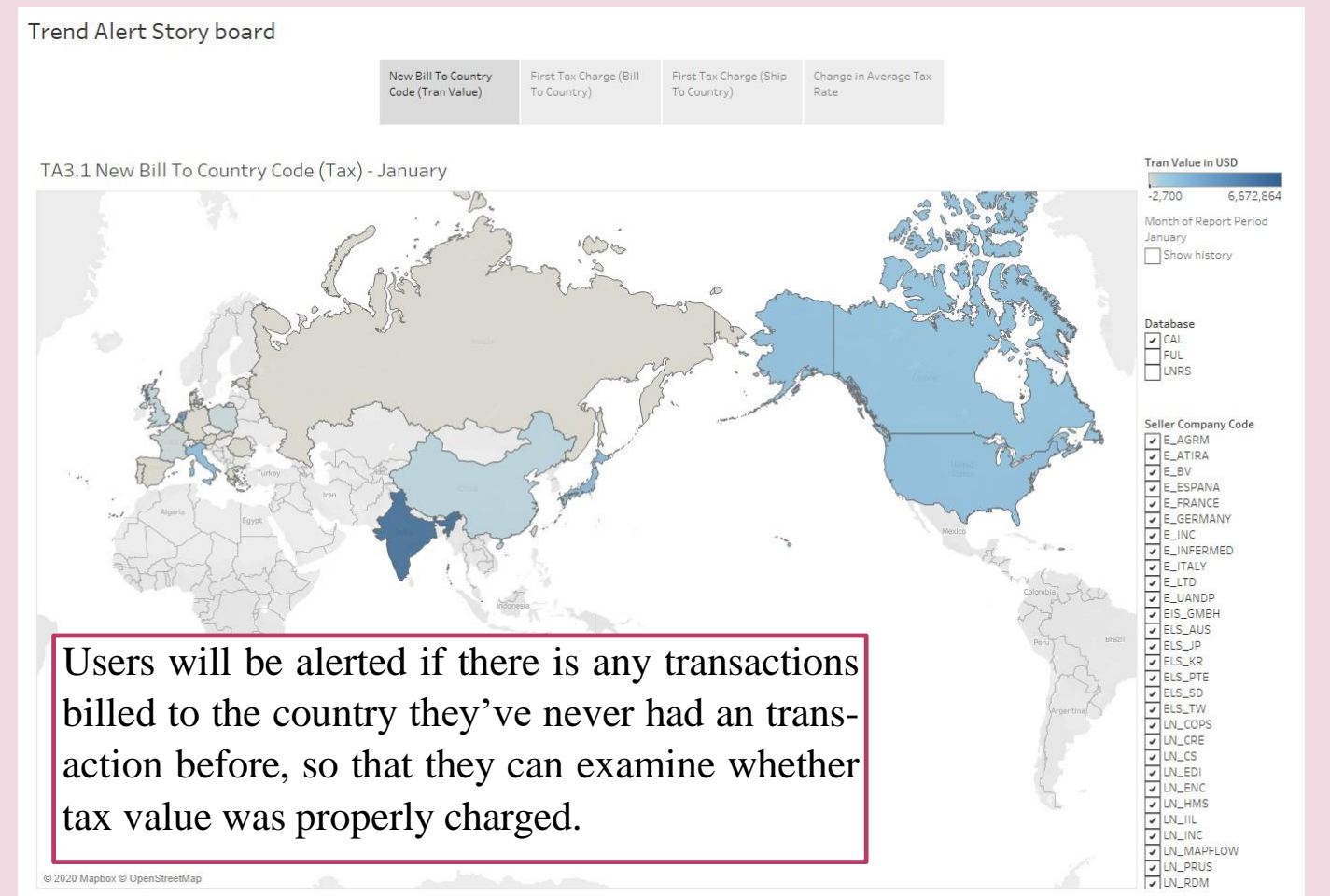
Scope 1: Data Integrity



Scope 2: Trend Analysis



Scope 3: Trend Alert



Reading between the lines: A Machine Learning Approach to Fraud Risk Assessments

Ivy Munoko, SooHyun Cho and Helen Brown-Liburd

Research Problem

- Fraud risk assessment is challenging for external auditors due to its complexity, and the fact that external auditors are usually the outsiders looking in.
- Fraud schemes are usually complex and evolving, with their impact far-reaching. The Association of Certified Fraud Examiners (ACFE 2018) approximates that companies lose 5 percent of revenues to fraud. By extension, fraud also impacts shareholders, the capital markets, and the public, resulting in economic losses (Hogan, Rezaee, Riley, and Velury 2008).
- The traditional rule-based approach to fraud risk assessment, which relies on pre-programmed rules (e.g., checklists), has proven to be not very effective in detecting fraud risks (Hogan et al. 2008; Asare and Wright 2004). Across research about fraud risk detection, a common conclusion is that there is a need for innovative ways for fraud evaluations (Dorminey et al. 2012).

Approach

- This paper demonstrates the use of machine learning for detecting fraud red flags using corporate communication. This analysis is performed on an aggregated level, to provide the external auditor with a fraud risk profile for an organization and its departments.
- These fraud risk profiles are not only corroborative evidence to support quantitative information but are also attention directing aids that can point auditors to areas that they may miss when using a traditional audit approach to fraud risk assessments.
- Using an public email dataset of a factual company, we combine both established fraud theories and machine learning techniques, to develop an automated framework that aids fraud risk assessments. To validate the framework, we conduct an experiment with an expert panel and find that forensics experts who are also CPAs express fraud risk assessments consistent with our framework.

Methodology

- The first step involved the extraction of features for the emails, including the sentiments scores of emails, topics discussed in the emails and communication patterns of the senders. These extracted features were used to develop machine learning models to predict “high risk emails” (see Figure 1)

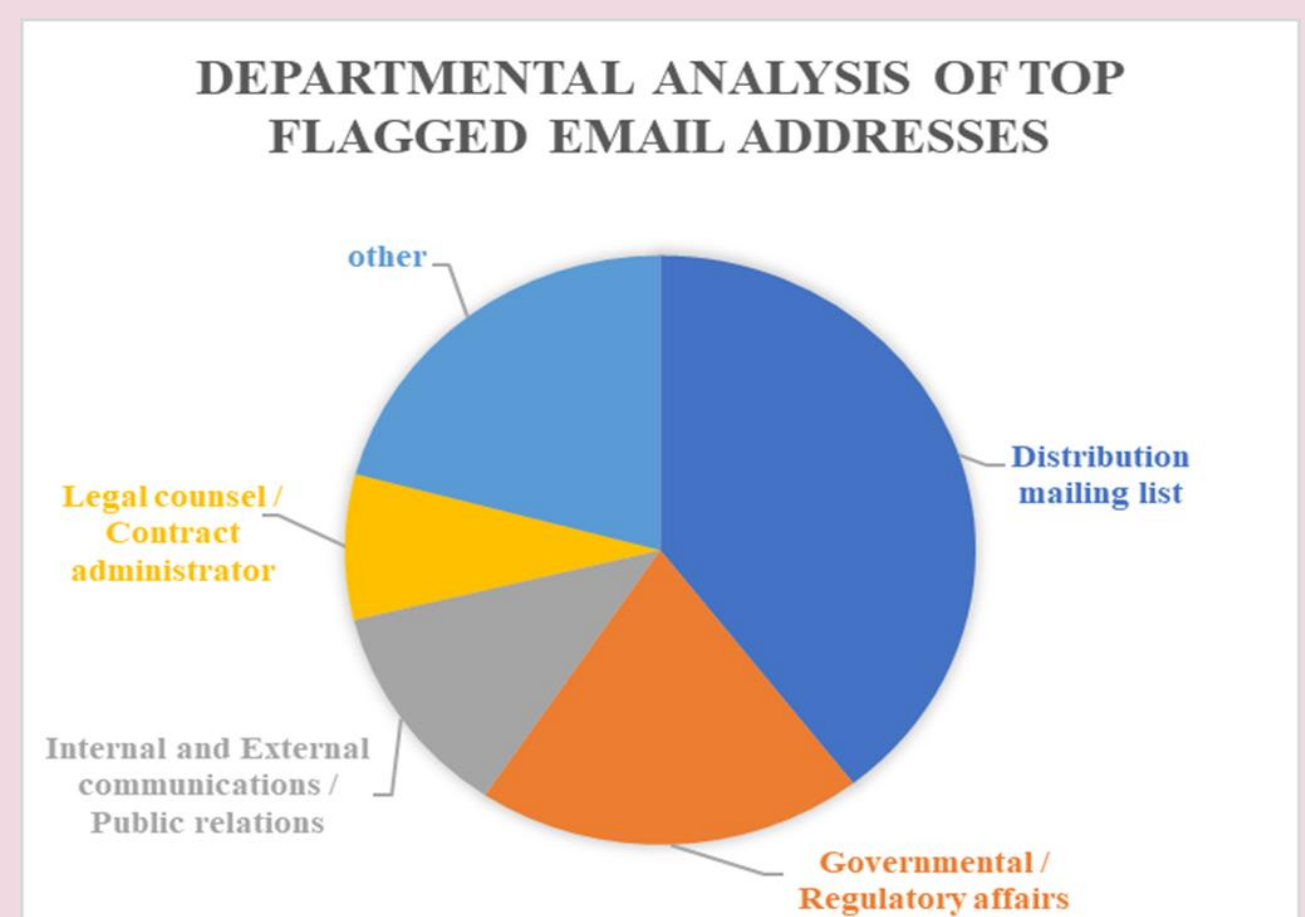
Findings

- Of the 4,601 emails flagged by the machine learning models, a third of these flagged high-risk emails originated from only 31 email addresses. In examining the high-risk emails associated with the 31 email addresses, flagged departments included Government/Regulatory Affairs, Internal and External Communications department, and Legal/Contracts (see Figure 2).
- The use of a fraud risk assessment decision aid, such as the one detailed in this paper, can be a vital aid in supporting auditors in determining where higher fraud risks exist for their clients.

Figure 1

Performance of classifiers (NN, LR, SVM and RF) on the test set.			
Panel A: Model trained on imbalanced dataset (67% of data used for training, 33% for testing), before standardization			
Classifier	Performance Metrics		
	Accuracy	Recall	
1) Random Forest Classifier (RF)	99.95%	72.04%	
2) Neural Network (NN)	99.88%	54.84%	
3) Logistic Regression (LR)	99.88%	6.45%	
4) Support Vector Machine (SVM)	99.88%	0.00%	
5) Voting Ensemble of LR, SVM, NN and RF	99.93%	49.46%	
Panel B: Model trained on balanced dataset (300 emails for training - 150 high risk, 150 low risk, all other emails in dataset used for testing), before standardization			
Classifier	Performance Metrics		
	Accuracy	Recall	
1) Random Forest Classifier (RF)	98.94%	99.09%	
2) Neural Network (NN)	98.48%	96.36%	
3) Logistic Regression (LR)	97.88%	95.45%	
4) Support Vector Machine (SVM)	97.36%	93.64%	
5) Voting Ensemble of LR, SVM, NN and RF	98.86%	97.27%	
Panel C: Model trained on balanced dataset (300 emails for training - 150 high risk, 150 low risk, all other emails in dataset used for testing), after standardization			
Classifier	Performance Metrics		
	Accuracy	Recall	
1) Random Forest Classifier (RF)	98.94%	99.09%	
2) Neural Network (NN)	98.28%	99.09%	
3) Logistic Regression (LR)	98.53%	99.09%	
4) Support Vector Machine (SVM)	98.64%	99.09%	
5) Voting Ensemble of LR, SVM, NN and RF	98.71%	99.09%	

Figure 2



Integration of Process Mining and Blockchain for Continuous Assurance

Jumi Kim and Miklos A. Vasarhelyi

Process Mining

- **Process mining** is a technique to extract knowledge from event logs to discover, monitor, and improve business processes, where as an **event log** is a chronological record of computer system activities which are saved to a file on the system (van der Aalst et al., 2010; Alles, M. G., Jans, M. J., Vasarhelyi, M. A., 2011; Jans, M., M. Allés, and M. Vasarhelyi, 2013, 2014; Chiu, T., Vasarhelyi, M., Alrefai, A., Yan, Z., 2018).
- Jans, M et al. (2013) stated that process mining adds value when applied to an audit. Process mining enables auditors: (1) to conduct entire population auditing analysis rather than the sampling method; (2) to practice independent audits with meta-data; (3) to effectively implement an audit risk model by conducting the required walk-throughs of business processes with process mining.
- Caron, F., Vanthienen, J., Baesens, B. (2013) proposed a **rule-based compliance checking process mining** approach where authors incorporated the business provenance, regulation, directives, and business rules in the process mining analysis.
- Three main interests for compliance checking: (1) process discovery and visualization; (2) conformance checking and delta analysis; (3) rule-based process mining (Caron, F et al., 2013).
- Alrefai, A. (2019) proposed to implement a “**continuous monitoring layer using rule-based process mining techniques**” which allows auditors to detect violations in real-time.

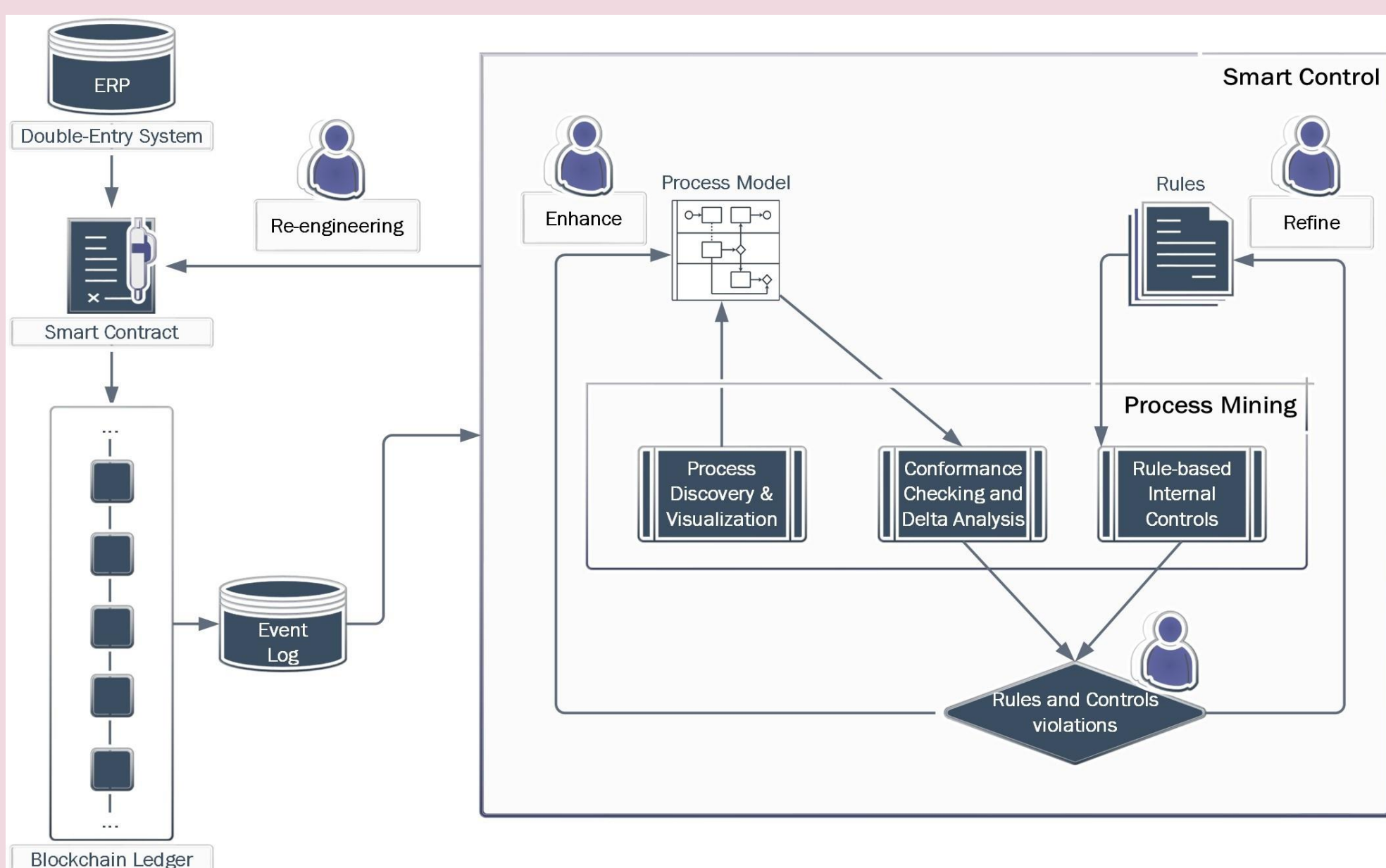
Blockchain and Smart Contract

- **Blockchain** is a decentralized, trustless, peer-to-peer network without a third trusted intermediary (Chris, D., 2017; Duchmann, F., Koschmider, A., 2019).
- The main characteristics of a blockchain are: (1) decentralized and distributed; (2) immutable and irreversible; (3) near-real-time (Parikh, T., 2018).
- **Triple-entry accounting** is proposed to improve the reliability of companies’ financial statements. Blockchain is a “trustless trust” ledger, which could replace the reliable intermediary in the triple-entry system proposed by Grigg (2015) (Dai, J., & Vasarhelyi, M. A., 2017).
- **Smart contract** was originally proposed by Nick Szabo (1994) and has revived due to the development of blockchain. Smart contract is defined as “a computer program having self-verifying, self-executing, tamper-resistant properties” (Mohanta, B. K., Panda, S. S., Jena, D., 2018, July).
- **Smart contract** could facilitate “reliable data sharing between business parties and continuous reporting for shareholders.” (Dai, J et al., 2017)
- A business transaction will be saved to the ERP system and then added to the blockchain if the transaction satisfies the conditions in a smart contract. The state of the transaction will be automatically updated in the distributed ledgers and available to permissioned participants.

Integration of Process Mining and Blockchain for Continuous Assurance

- Dai, J et al. (2017) proposed a blockchain-based continuous assurance with the integration of smart contract. The proposed audit paradigm consists of a physical world and a mirror world. The mirror world consists of blockchain, smart control, and payment layer.
- The proposed framework employs the “smart control layer”, where managers and auditors would program the firm specific-control protocols into smart contract which enable monitoring of business processes (Dai, J et al., 2017). The framework displays how process mining technology can facilitate and improve the blockchain-based continuous assurance.
- The framework starts with an ERP system where a transaction is recorded in terms of debit and credit. This transaction is then recorded in the blockchain through smart contract. For example, in a P2P process, a smart contract verifies the balance in the company's account. If the balance is greater than the total cost of goods ordered, the supplier sends inventories. Upon goods receipt, the company matches the purchase order, invoice, and receipt. If matched, smart contract automatically transfers money from the company's account to the supplier's account, and the transaction is appended to the blockchain ledger.
- Event log is extracted from the blockchain ledger. The event log enters into the smart control layer. In the smart control layer, process mining will discover a process model reflecting the reality, conform the reality with the process model, and check the process properties. Process mining will determine the deviated processes from the discovered model and the violated processes from the pre-defined rules.
- Rules include regulations, business provenance, directives, business rules, internal controls, and etc.
- An auditor will investigate deviations and violations and will further determine truly violated processes and acceptable processes (true positive) for the flexibility of business. The information from the result will enhance the process model and refine the rules. The truly violated processes will be reported to a manager as the fraudulent activities.
- The information learned from smart control will be used to re-engineer the smart contract.

Blockchain-Based Continuous Assurance Framework with Process Mining



Limitations and Future Work

- Current process mining research uses the data from the ERP system and there is software available for event log extraction; however, extracting event logs from blockchain has not been studied much.
- Future investigations are necessary to validate the possibility of the realization of the proposed framework.

RUTGERS

Rutgers Business School
Newark and New Brunswick

Big Data in Audit Acceptance and Continuance Decision

Danyang (Kathy) Wei, Qiao Li, Miklos A. Vasarhelyi

Background

- Adequately understanding an entity and its environment is important in making audit acceptance and continuance decision.
- Big Data contains valuable information that can help auditors form comprehensive perception about a prospective client.
- However, some characteristics of Big Data impede its broad use in audit. One is too much noise within large scale data. Auditors tend to lose their focus and get information overload when handling Big Data. The other one is the purpose of Big Data creation. On the one hand, Big Data is not created for audit but for business-based purposes. On the other hand, audit standards do not offer a clear guidance about how to link Big Data with audit engagement. Therefore, auditors are less motivated to invest their limited resources in Big Data and take advantage of it.
- Based on these facts, this research designs a framework to provide a way to use Big Data in understanding a prospective client for acceptance and continuance decision. The framework first considers the aspects mentioned in audit standards about understanding an entity and its environment and categorizes them with three main components for a company including operation, management systems, and management people. For each aspect, the insights that can be generated from Big Data and the possible data sources are listed. Meanwhile, the expected judgment that can be applied with the support of Big Data is also discussed.
- In addition, a workflow is designed to exhibit how the framework can help auditors in client acceptance and continuance decision-making.

Introduction

- The main purpose of the framework is to illustrate how Big Data could be helpful in each aspect.
- The framework considers auditors as a general related party to the company such as a supplier. All Big Data discussed in this research is publicly available and from this perspective, the information that auditors are concerned about has no difference from a general related party. The files and information that can be directly obtained by auditors, such as the prior audit working paper, are not included since Big Data is not helpful in this area.
- The functions of Big Data are twofold: providing information that auditors can directly use and translating the content that requires specialized knowledge to understand into auditor-friendly information. For the first kind of function, Big Data tends to be financial data that auditors are familiar with, such as financial analysis reports about the industry in Wall Street Journal. For the second one, Big Data plays a role that explains the effect of one change on the whole industry or the company so that auditors could generate expectations on the reaction of the company.
- The “example of matters” column aims to briefly explain what the main aspect is through several examples instead of giving a complete list. In practice, auditors will decide what matters for each aspect based on the company and their professional judgment.
- One point that should be emphasized is that Big Data is not a stable information source and the availability level of Big Data highly depends on the industry where the company is. For example, manufacturing companies may have less Big Data available compared with retailing companies.

Framework

Operation

Component	Main Aspect	Example of Matters ¹	Useful Big Data	Expected Judgment
Operation	Industry	<ul style="list-style-type: none">Market competitionRegulatoryCustomer relation	<ul style="list-style-type: none">News (analysis articles)Reports issued by authentic journals in the industryWeatherRelated litigation cases in the industry (reason to be involved in litigation and results)Official account on social media	<ul style="list-style-type: none">Are the selection and the application of accounting policies proper?Are there any risk factors that need substantive tests later (WCGW)?
	Main product/service	<ul style="list-style-type: none">Research and development costSupply chainMarket position	<ul style="list-style-type: none">Key financial ratios and stock price trends from financial websites such as Yahoo! FinanceComments about the product/service on mediaConsumer feedback from social media	<ul style="list-style-type: none">Any significant difference between the company and its peer companies?Does the profit level reflect the current financial and market conditions?Does the strategy/business objective match the current situation/ reasonable?
	Capital (equity + financing)	<ul style="list-style-type: none">Capital structureInterest rateGovernment policies²	<ul style="list-style-type: none">Opinions/analysis from professionals (i.e. articles that explain what the policy means to the company)Stock trend analysis (i.e., investors' opinion)	<ul style="list-style-type: none">Can the company get expected amount of money from related parties?Is there any going concern?

Management System

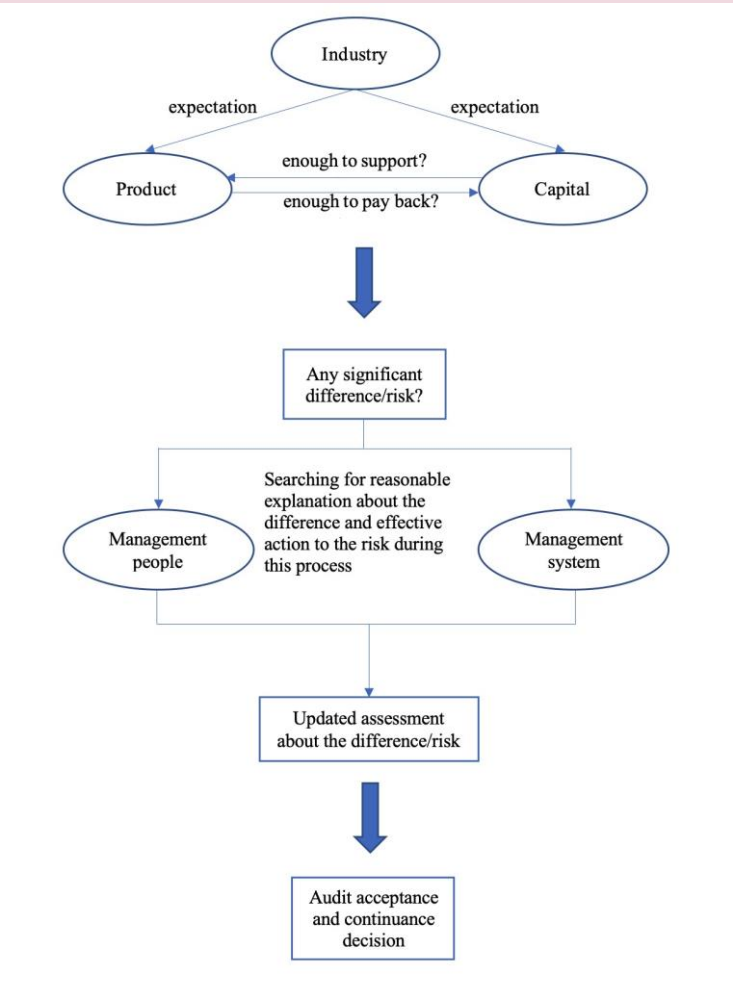
Component	Main Aspect	Example of Matters	Useful Big Data	Expected Judgment
Management System	Internal control process ³	<ul style="list-style-type: none">DesignImplementationMonitoring of controls⁴	<ul style="list-style-type: none">Reputation of the software usedDisclosure about the cybersecurity of the company	<ul style="list-style-type: none">Is there any deficiency in internal control?Is there any action that the company takes to make it up?
	Internal control environment	<ul style="list-style-type: none">A culture of honesty and ethical behaviorStrengthsDeficiencies	Big Data may not be helpful here	<ul style="list-style-type: none">Does the internal control environment undermine the effectiveness of internal control?Is there any possible collusion?Are the strengths strong enough to make up the deficiencies?
	Entity structure	<ul style="list-style-type: none">Ownership characteristicsOrganizational structure⁵Multiple subsidiaries and locationsRelations with other owners or entities	<ul style="list-style-type: none">Prediction about the results of merge/consolidation	<ul style="list-style-type: none">Transactions between related partiesIs the internal control suitable to the entity structure?

Management People

Component	Main Aspect	Example of Matters	Useful Big Data	Expected Judgment
Management People	Integrity	<ul style="list-style-type: none">Attitude towards internal controlCriminal background investigation	<ul style="list-style-type: none">Personal social mediaPublic interviewsPublic speech	<ul style="list-style-type: none">Is the internal control strengthened or weakened by the management people?
	Competence and capability	<ul style="list-style-type: none">ExperienceWhat achievement has been made in the position	<ul style="list-style-type: none">Information about prior companies he/she worked forNews/articles about his/her contributions/significant achievement (i.e. leading an innovation project)	<ul style="list-style-type: none">Is the leading style suitable to the company?Is there enough competency for his/ her position?
	Governance structure	<ul style="list-style-type: none">Board of directors	<ul style="list-style-type: none">Social network (LinkedIn)	<ul style="list-style-type: none">Is there any potential collusion?

Workflow

- The workflow gives guidance on how auditors can use the framework and make the acceptance and continuance decision.
- From the company’s perspective, the main task of the management system and the management people is to support the operating process. In other words, an effective management system and a proper management team can make sure the operating cycle works as expected. Based on this point of view, both the management system and the management people contribute to the ultimate operating result. However, a different objective of auditors results in a different starting point. To make a wise acceptance decision, auditors need to know enough information about their prospective client so that they can comprehensively assess the audit risk. In this case, they would better start from the “result” stage, which is the operation. The reason is that by knowing about the operating condition, auditors can better assess effectiveness of the management system and suitability of the management team.
- At the operation stage, the main goal is to obtain expectations on the financial performance of the company and identify potential risks. Then, auditors can observe the effort made by the company through the management system and the management people in reducing the identified risks. The assumption is that risks can be lowered through an effective management system and a suitable management team. The objectives are to find out a reasonable explanation about the difference between the expectation and the actual condition in previous stage and see if any action is taken to reduce the identified risks. Then, auditors update their assessment about the risks and decide if the client should be accepted or not.



Increasing the Utility of Performance Audit Report: Using Textual Analytics Tools to Improve the Government Reporting

Huijue Kelly Duan, Hanxin Hu, Yangin Ben Yoon, Miklos Vasarhelyi

Introduction

Motivations

- Performance audit provides valuable information to the public about government programs, and shows objective analysis, findings, and conclusions of the audit to users. It is considered a tool for civic participation in monitoring and improving government performance and operations, reducing costs, facilitating decision making, and contributing to public accountability.
- However, the utility of the performance audit reports for the general public is relatively low due to the accessibility issue.
- Additionally, each state has its own reporting template. There is no standardized format across all states, and some of the reports are not machine-processable, which creates difficulties for users to systematically access the contents and compare the performance/issues across different government entities.

Objectives

- Increase the utility and accessibility of performance audit reports to the general public.
- Enhance the transparency and accountability of the government.
- Establish a standard reporting theme, ensure the reporting complies with GAGAS requirements.
- Construct a taxonomy specific to the government performance audit.
- Assess the linguistic features to provide more insights into the reports.

Contributions

- By establishing a standard reporting template across all states, authorities can get a better understanding of the performance audit reports, gain insight into the performance of the government entities, and comprehensively evaluate the performance of the programs.
- State auditors can improve the risk assessment and audit planning, identify the trends and monitor the emerging issues on a national and topical level, and compare similar audit issues and solutions from other states.

Results

Table of contents for PDFs	background information	objective, scope and methodology	Recommendations and findings1	Recommendations and findings2
background 5 audit findings and recommen	Background The New York State Medicaid program is a federal, st	Audit Scope and Methodology We audited selected medical record	Audit Findings and Recommendations We determined Dr. Riaz Al	Recommendations 1. Review the 19/031 Medicaid claims totaling
background 5 audit findings and recommen	Background Sunshine Developmental School (Sunshine Developn	Audit Scope and Methodology We audited the propriety of, and so	Audit Findings and Recommendations We identified \$1,776,434 i	Recommendations To SED: 1. Review the disallowances identify
background 5 audit findings and recommen	Background The Child School (School), a not-for-profit organizati	Audit Scope and Methodology We audited the expenses reported o	Audit Findings and Recommendations We identified \$978,085 of	Recommendations To SED: 1. Review the disallowances identify
background 5 audit findings and recommen	Background The mission of the Administration for Children's Ser	Audit Scope and Methodology We audited ACS to determine whe	Audit Findings and Recommendations We found that ACS official	Recommendations 1. Optimize opportunities to solicit competitive
background 5 audit findings and recommen	Background The New York State Medicaid program is a federal, st	Audit Scope and Methodology The objectives of our audit were to	Audit Findings and Recommendations We found numerous violati	Recommendation 1. Coordinate with HRA officials to investigate
background 4 audit findings and recommen	Background Medicaid is a federal, state, and local government pro	Audit Scope and Methodology The objective of our audit was to d	Audit Findings and Recommendations For the period December 1,	Recommendations 1. Review the actual and potential overpayment
background 5 audit findings and recommen	Background The Department of Health (Department) is responsible	Audit Scope, Objective, and Methodology The objective of our au	Audit Findings and Recommendations We found that the Departmen	Recommendations 1. Review the Medicaid payments made to Joia
background 5 audit findings and recommen	Background Whispering Pines Preschool, Inc. (Whispering Pines),	Audit Scope and Methodology We audited costs that Whispering I	Audit Findings and Recommendations According to the RCM, on	Recommendations To SED: 1. Review the recommended disallow
background 6 audit findings and recommen	Background The New York State Medicaid program is a federal, st	Audit Scope and Methodology We audited selected Medicaid claim	Audit Findings and Recommendations Based on the results of our	Recommendation 1. Review and recover the remaining \$28,028 in
background 5 audit findings and recommen	Background The Association for Neurologically Impaired Brain Inj	Audit Scope and Methodology Our audit determined whether the c	Audit Findings and Recommendations For the three fiscal years on	Recommendations To OPWDD: 1. Review the recommended disal
background 5 audit findings and recommen	Background The Metropolitan Transportation Authority (MTA) is	Audit Scope and Methodology The objective of our audit was to d	Audit Findings and Recommendations Transit, MTA Bus, and B4	Recommendations 1. Revise the All Agency Travel Policy Directiv
background 5 audit findings and recommen	Background Many New Yorkers have been increasingly challenged	Audit Scope and Methodology We conducted this audit to determi	Audit Findings and Recommendations The managing agents and E	Recommendations 1. Formalize procedures to conduct lotteries and
background 5 audit findings and recommen	Background The New York State Office of Probation and Correcti	Audit Scope and Methodology The objective of our audit was to d	Audit Findings and Recommendations Only a small percentage of	Recommendations 1. Develop and implement processes and proce
background 5 audit findings and recommen	Background The New York City Department of Health and Mental	Audit Scope and Methodology We audited DOHMH's administrat	Audit Findings and Recommendations We concluded that DOHMH	recommendations and maintains documentation of its reviews. DO
background 6 audit findings and recommen	Background The New York State Urban Development Corporation.	Audit Scope and Methodology We audited ESD to determine whe	Audit Findings and Recommendations ESD has an appropriate sys	Recommendations 1. Develop strategic plans that include performi
background 5 audit findings and recommen	Background The Department of Economic Development (DED) is i	Audit Scope and Methodology Our audit determined if the Empire	Audit Findings and Recommendations We identified a range of its	Recommendations 1. Review and take appropriate action when adv
background 6 audit findings and recommen	Background The New York State Medicaid program is a federal, st	Audit Scope and Methodology We audited selected Medicaid claim	Audit Findings and Recommendations Based on the results of our	Recommendation 1. Ensure that pricing methodology changes are:
background 5 audit findings and recommen	Background The State University of New York (SUNY) consists o	Audit Scope and Methodology The objective of our audit was to d	Audit Findings and Recommendations SUNY officials have gaver	Recommendations 1. Remind campuses of the need to comply wit
background 5 audit findings and recommen	Background New York State Homes and Community Renewal (H	Audit Scope and Methodology The objective of this audit was to e	Audit Findings and Recommendations DHCR's oversight and mo	Recommendations 1. Improve project monitoring of State-funded l

Extraction Results

Audit objectives	Competence	Internal control	State	"Scope and Methodology" Format
Audit organization	Control objective	Objectivity	NY	Audit Scope, Objectives, and Methodology Background, Scope, and Objectives
Audit procedures	Criteria	Outcomes	NJ	Methodology Scope
Audit report	Entity objective	Planning	AZ	Scope & Methodology
Audit risk	Finding	Review	CA	Background, Scope, and Methodology Background, Scope, Methodology and Assumptions
Audited entity	Fraud	Significance		

Unique Titles

GAGAS Terms

Methodology

Step 1: Identify the Requirements

- General Accepted Government Accounting Standards (GAGAS) requires the performance audit reports to include the objectives, scope, and methodology of the audit, and also the audit results, findings, conclusions, recommendations, internal control issues, etc.

Step 2: Report Collection

- Reports issued by the states of New York, New Jersey, Arizona, California are script-ed from individual states' websites

Step 3: Content Extraction

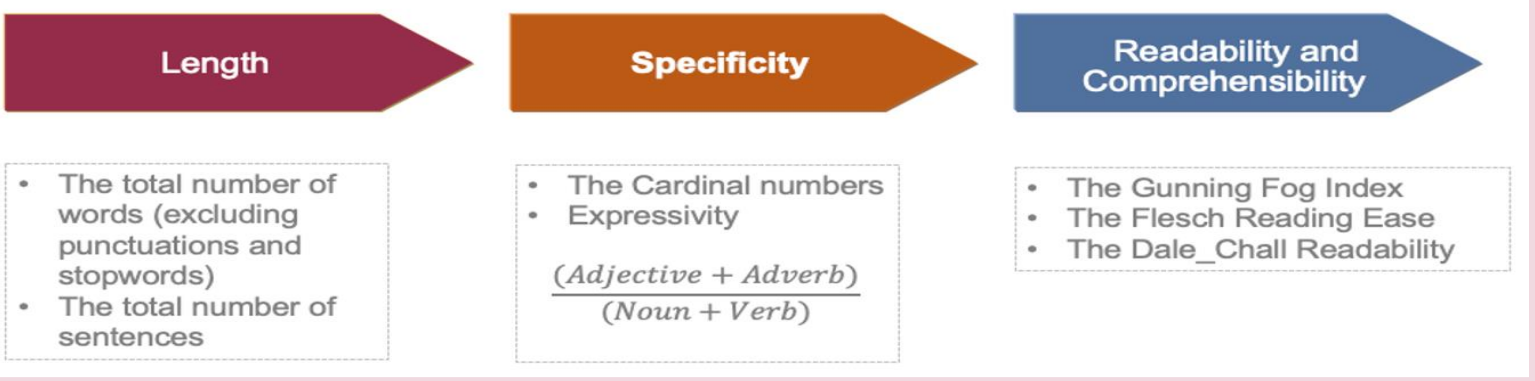
- Document conversion: convert the files (HTML files, PDFs with images, PDFs without images) to a machine-processable format
- Content Extraction: extract relevant content using different techniques (BeautifulSoup, Regular Expression, Fuzzy Matching).

Step 4: Taxonomy

- Identify standard terms and phrases used across the states
- List of terms required by GAGAS

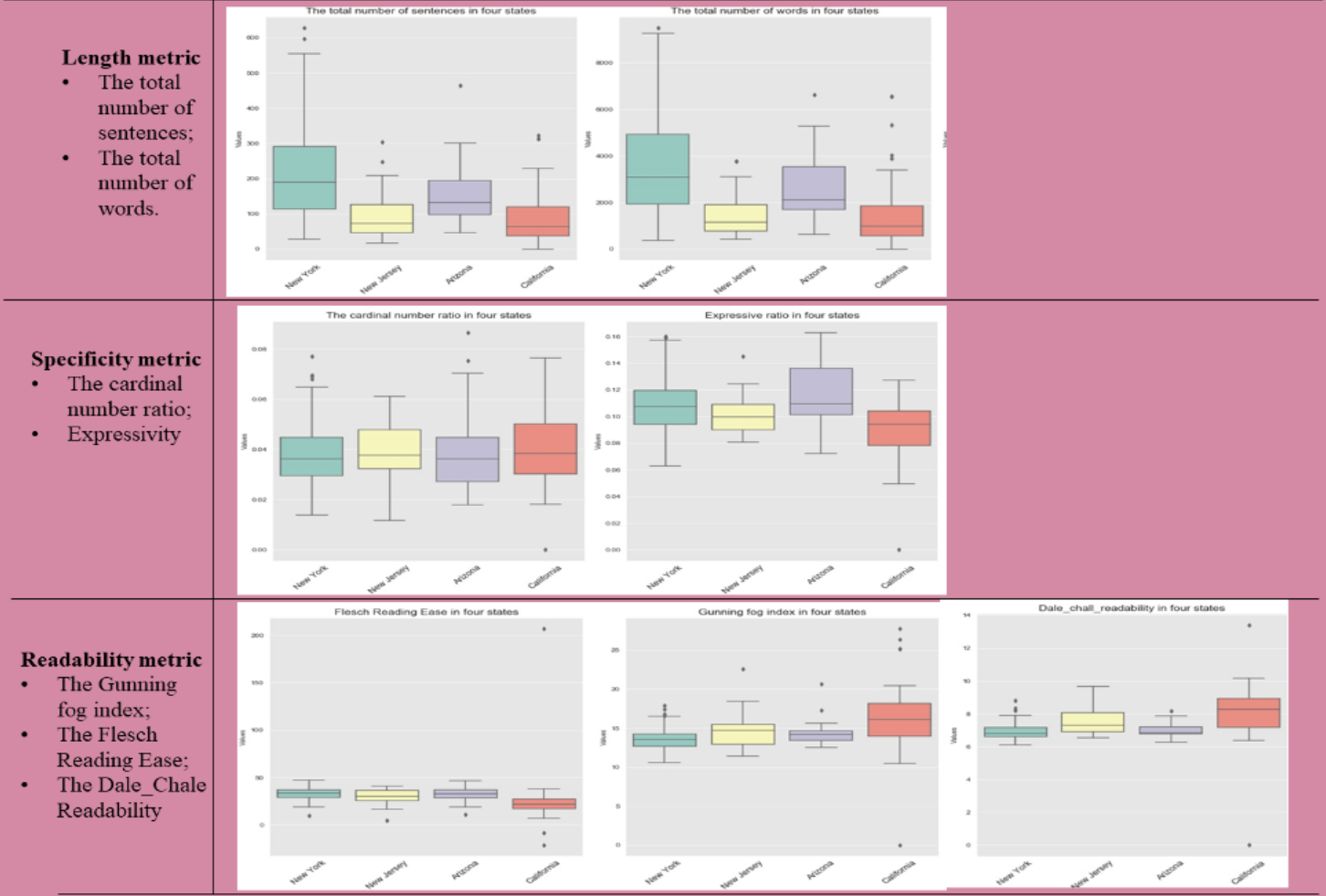
Step 5: Linguistic Analytics

- Length, Specificity, Readability, and Comprehensibility



Linguistic Analytics Results

- New York State discloses the most information.
- All four states have a similar amount of details.
- Sample reports from California have the highest Flesch reading index and the lowest Gunning Fog Index and Dale-Chall Index, indicating these reports are difficult to read and comprehend.



Conclusion

- This study establishes a framework that can be potentially used for building a standardized reporting template for government performance audit reports.
- It identifies required reporting elements from GAGAS, retrieves the relevant information from the reports issued by New York, New Jersey, Arizona, and California, constructs a taxonomy specific to a performance audit.
- This study measures the length, specificity, and readability of the performance reports among states vertically and horizontally.

Visual Audit: An Integrated Audit Approach

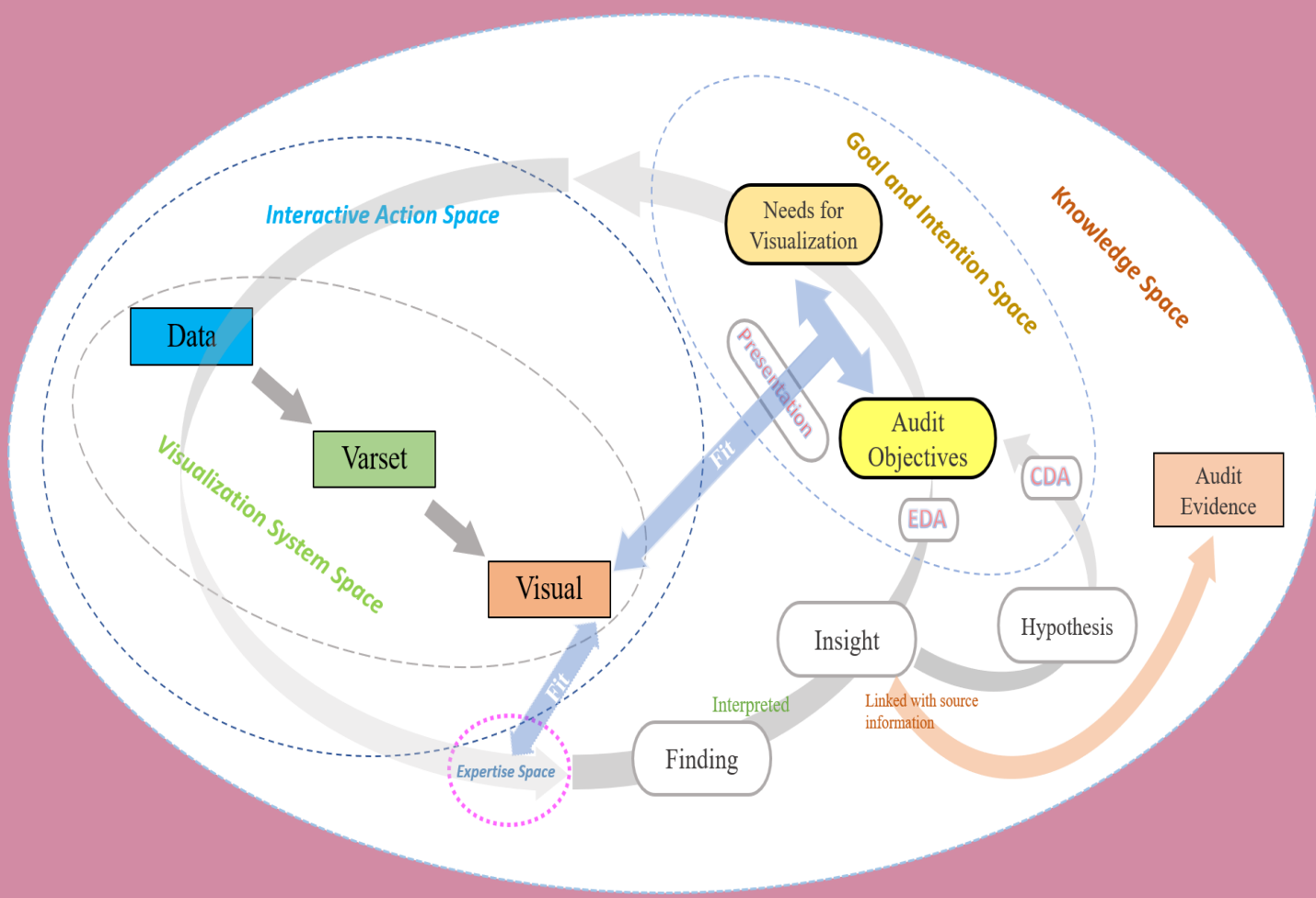
Lu Zhang, Heejae Lee, Qi Liu

Introduction

Using interactive visualizations in accounting information systems can assist users understand very large sets of financial information. Continuous monitoring and continuous auditing increases the value of interactive visualization in an accounting information system context. However, there are only few studies that examine how interactive visualization should be applied in accounting, especially in auditing.

The aim of the study is to: (1) Examine how interactive visualization can be used in auditing and suggest a integrated audit approach for visual audit. (2) Demonstrate visual audits using a hospital database with over six million transaction records.

The Visual Audit Model



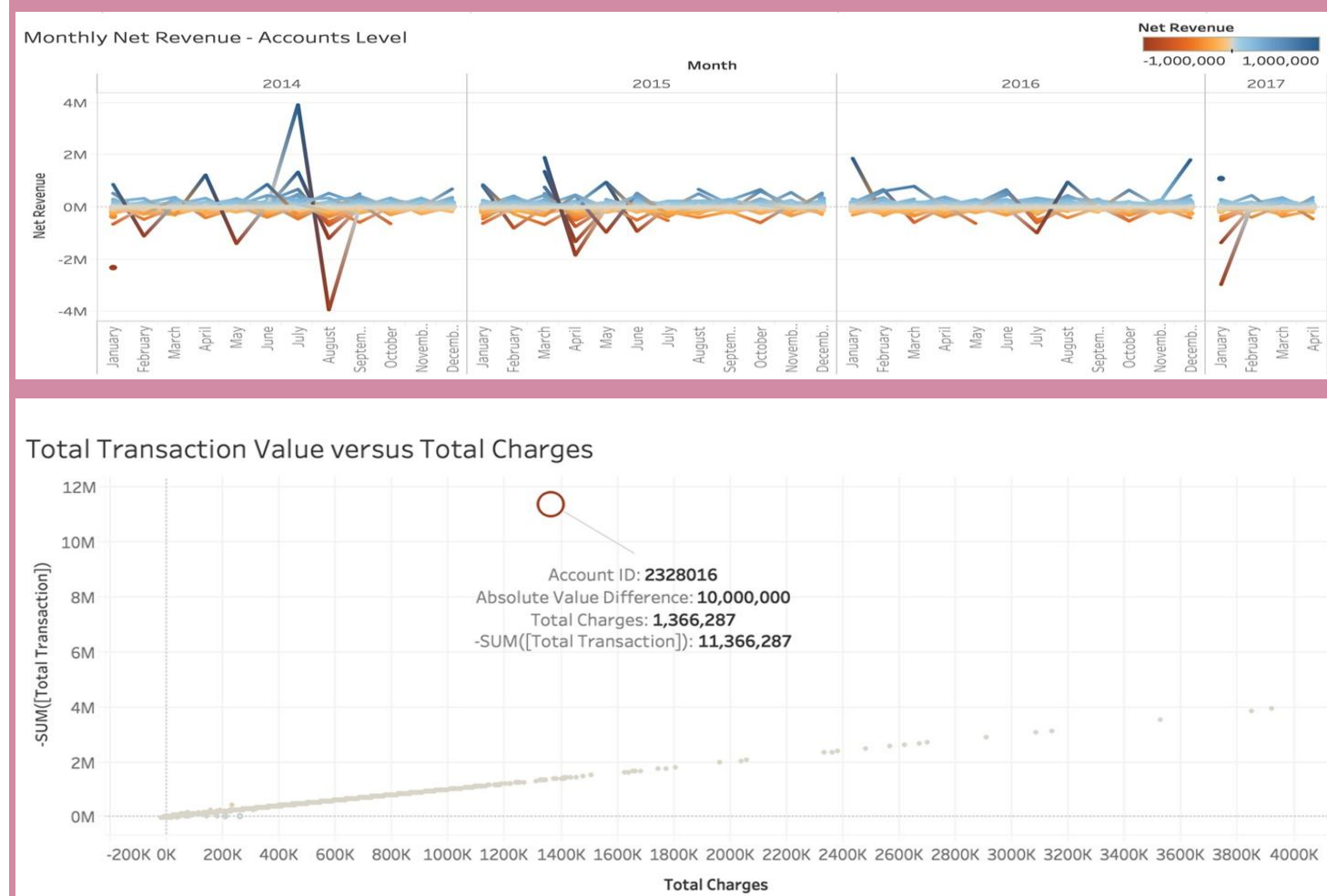
Visualization System Space

We defined *visualization system space* as a space holding a set of system components and relations among them, and described those components as *data*, *varset* and *visual*, and their relations as *data* being transformed into *varset* through actions like *importing*, *calculating* and *data mining*, and *varset* being mapped to *visual* through actions like *reconfiguration* and *encoding*.

Data is defined as the “data of interest” connected to the visualization system. *Varset* stands for the “variable set” that contains variables and their instances mapped to graphical entities. *Visual*, referred to as “visualization”, is defined as a graphical representation of data and is made up of *graphical objects* (Bertin 1967, Senay and Ignatius 1994, Wilkinson 1999, Börner et al. 2018), *visual variables* (Bertin 1967) and *guides* (Ward et al. 2010).



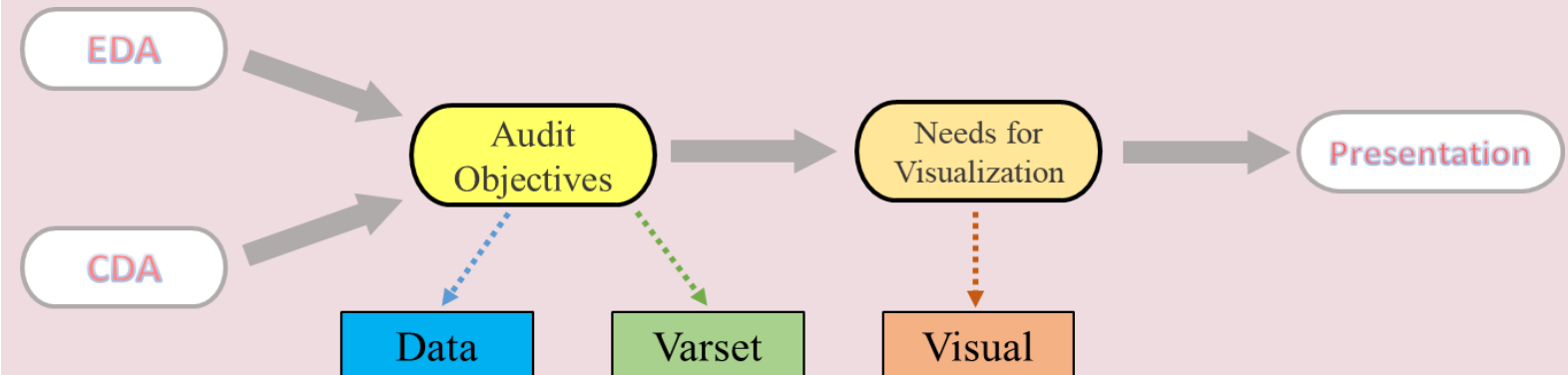
Visual Audit Demo



Goal and Intention Space

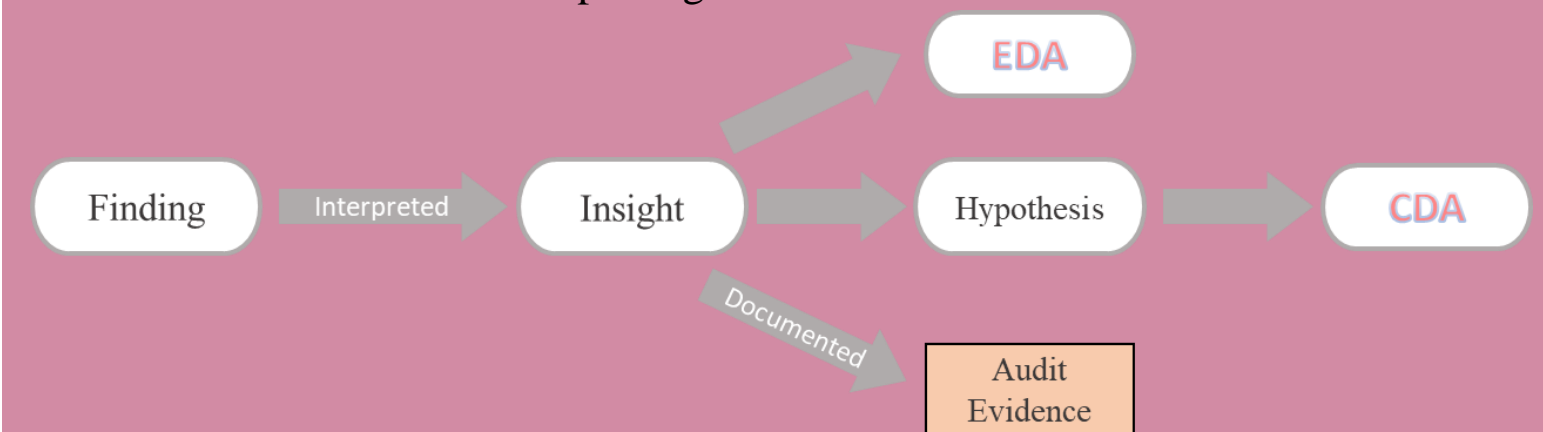
It is hard to define how to extract value from a specific process without identifying the solutions to solve at the beginning. Norman’s stages of action model stated that, when interacting with physical or virtual objects, humans need to form the “goal” and “intention” at the beginning (Norman, 2016). Roth 2012 explained the interactions between human and visualization systems using Norman’s model and further specified the “goal” as an ill-defined task, or goal, motivating use of visualization, and the “intention” as a well-defined task, or objective, supporting the goal (Roth, 2012).

Our model built on and extended their definitions, creating a *goal and intention space* in which three “goals” as exploratory data analysis (*EDA*), confirmatory data analysis goal (*CDA*) and *presentation*, and two “intentions” as *audit objective* and *needs for visualization*, are established to guide the interactions between auditors and the visualization system.



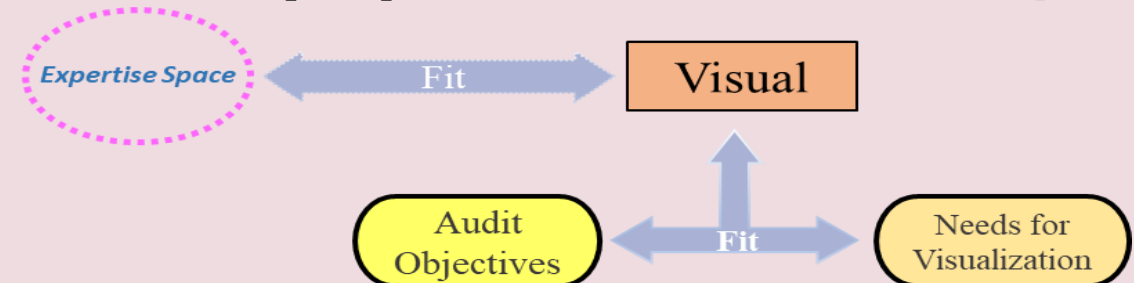
Knowledge Space

We created a *knowledge space* that holds different types of information objects and the relations among them occurring in visual audit. The collection of information reaches the solid state of knowledge about the audit when it is crystallized (Bertini & Lalanne, 2009) and may support the opinion expressed in the auditor’s report. We further classified those information objects as *finding*, *insight*, *hypothesis*, and *audit evidence*. *Finding* is a detected visual pattern independent from the problem domain (Sacha et al., 2014). *Insight* is an interpreted finding using domain knowledge (Chang et al., 2009). A *hypothesis* formulates an assumption in the audit domain that is subject to *CDA* and must be generated based on the *insight and visual analytics* may test and reject *hypotheses* before solving the problem. *Audit evidence* is prepared based on properly documented *insights* which corroborate or contradict management’s assertions regarding the financial statements or internal control over financial reporting.



Expertise and Fit

To help auditors design an appropriate *visual*, we introduced the cognitive fit theories (Vessey, 1991) into our model and argue that selected *visual* must provide fit to *data* and *varset* that are selected based on *audit objectives*, tasks that are motivated by *needs for visualization*, and auditors’ perceptual abilities which are defined as *expertise*.



Interactive Action Space

In visual audit, we created an *interactive action space* in which auditors communicate and interact with the visualization system through interaction techniques (Becker et al., 1987). Actions taken in the space are dependent on developed *audit objectives* and *needs for visualization*. They are executed to facilitate the operation of each visualization system component. Based on previously developed taxonomies relevant to data preparation and interaction techniques, we classified *interactive actions* used in visual audit under three categories as *actions for data preparation*, *actions for varset creation*, and *actions for visual operation*.



Expected Loan Loss Provisioning Using a Machine Learning Approach

Nichole Li and Alexander Kogan

Background

- Accounting estimates are a critical part of financial statements. Most companies’ financial statements reflect accounts or amounts in disclosures that require estimation. Accounting estimates are pervasive in financial statements, often substantially affecting a company’s financial position and results of operations
- Banks are crucial for financial stability. Due to the nature of their assets (i.e. mainly loans) and their financial structure (i.e. highly leveraged and financed largely via deposits) they have specific information asymmetry problems with stakeholders different to those they may have with share- holders.
- Estimating expected Loan Loss Allowance (LLA) in banks is a critical but also difficult problem in accounting estimates. The issue has become of increasing interest to academics and regulators with the FASB and IASB issuing new regulations for loan impairment.
- However, till now, few studies have looked at the application of machine learning in managerial subjective estimates, especially in the LLA. And no published research has been done to model and predict loan losses using other types of machine learning algorithms than regression so far. Therefore, in this paper, I want to fill in the gap by using multiple machine learning algorithms to model and predict loan losses in banks.

Measuring Loan Losses

- Interest income is recognized over time and is derived from a yield that includes at least four components: the time-value of money, expected loan losses, risk premia, and economic profit (Harris et al., 2018). Measuring expected losses is particularly complex, but the loan loss allowance and provisions estimated by managers are based in part on a series of primary indicators, many of which are available in public disclosures. So, in this research, I focused on these primary indicators below in constructing an alternative summary measure of the expected loan losses.
- Loan Balances and Loan Composition. Characteristics of the borrower and of the collateral, affect both the probability of default and the loss-given-default. In my research, I include the proportions of the three largest loan categories: real estate, commercial and industrial (C&I), and consumer.
- Loan Duration
- Nonperforming Loans. Loans that are not paying interest or principal due to a borrower’s credit problems are classified as nonperforming loans ,which include nonaccrual loans, restructured (troubled) loans, and some past-due loans.
- Net Charge-offs. Net charge-offs (NCOs) are measures of realized loan loss in a given period and indirectly impact the balance sheet and income statement through the ALLL and the PLLL

Variable Definition

Independent Variables	Definition
Bank Variables	
Log Assets	Log of Total Assets
Total Loans	Total loans in banks portfolio
Loans to Assets	Ratio of loans to assets
Securities to Assets	Ratio of the securities to assets
NCO	Net charge-offs
Four-Qtr NCO + NPL	Sum of rolling four-quarter Net charge-offs plus ninety-days past due and non-accrual loans at the end of the rolling window’s fourth-quarter
Pct Four-Qtr NCO + NPL	Four-Qtr NCO + NPL scaled by total loans at the beginning of the quarter
Charge-offs	Amount that is charged-offs
Recoveries	Amount recovered in previously charged-off loans
Allowance	Loan Loss Allowance
Pct Allowance	Loan Loss Allowance scaled by total loans at the beginning of the quarter.
Pct RE Loans	Real estate loans as a percentage of total loans.
Pct CI Loans	Commercial and Industrial loans as a percentage of total loans.
Interest Receivables	Income accrued but not yet collected on loans
Loan Yields	The ratio of tax-equivalent interest income divided by total loans
External Environment Variables	
Unemp Rate	Unemployment rate
HPI	Home price index
HPI Growth	Home price growth
Inflation	Personal consumption expenditure growth in the previous year
GDP	GDP level in the previous year
GDP Growth	GDP growth in the previous year

Methodology

Based on the discussion above, the manager would estimate loan loss allowance by predicting future loan losses to be realized in subsequent periods. The accuracy of the prediction will be assessed by how well the estimated allowance captures actual net-charge-offs. This idea is consistent with the accounting identity:

$$\begin{aligned}LLA_t &= LLA_{t-1} + LLP_t + RECOV_t - CO_t \\ &= LLA_{t-1} + LLP_t - NCO_t\end{aligned}$$

Where LLAt is allowance for loan losses at end of t, RECOVt is the amounts that have previously been charged off but are recovered during this time period, COt is the gross amount of all loans charged off against the LLA losses. The income statement effect is captured by LLP, the loan loss provisions. Thus, the precision of the LLAt is assessed by how well it predicts future net charge-offs at t.

Sample and data

Following Harris et al. (2018), I focus on bank holding companies (BHCs) and extract accounting data from regulatory consolidated financial statements (FR Y-9C reports) for the period 1996–2017. The sample period starts from 1996 because information required for measuring certain FR Y-9C variables is unavailable before then.

Dependent variables

To predict the losses, I want to follow Fillat and Montoriol-Garriga (2010) and consider the sum of rolling four quarter net charge-offs, and add non-performing loans at the end of the rolling-window’s fourth quarter, which is the dependent variable in my study. At time t the loss is measured as,

$$CL = \sum_{\tau=T+1}^{t+4} NCO_{\tau} + NPL_{t+4}$$

where NPL is Non-performing loans, which are defined as loans past due more than 90 days and nonaccrual loans (i.e., loans on which a bank has ceased to accrue interest).

Independent variables

My independent variables (predictors) consist of information already known at the time of estimation. Inspired by Vijayaraghavan(2019), I include two sets of independent variables. One set is the bank variables that contains the characteristics of the banks themselves. Another set contains the exogenous environmental variables that reflect the macro-economic factors that may influence the loan loss estimation. The independent variables I want to include are shown in the table above.

Machine learning algorithms

Based on prior literature, I want to try five different machine learning algorithms in my research: Lasso regression, support vector machine, random forest, artificial neural networks and gradient boosting machine to make the prediction.

Skipper and Stretcher Selection in Audit Analytics

Nuriddin Tojiboyev and Alex Kogan

Literature Review

- Continuous Auditing Framework extensively depends on a reliable selection of exceptions for a further review
- There is an agreement in the field that a large volume of exceptions is generated in the continuous auditing environment (Thiprungsri and Vasarhelyi, 2011 and Kim and Vasarhelyi, 2012)
- The use of Suspicion Score became a popular approach for ranking/prioritizing exceptions (Issa 2013, Li et al. 2016, No et al. 2019)
- The Suspicion Scores are generated from applying different filters that identify certain transactions/records as exceptions when they fail to go through these filters
- The previous literature mostly uses binary assessment (fail or pass) of exceptions by filters.
- Moreover, previous models do not consider the issues of more than one exceptions with similar values being selected for further review
- Having too many duplicate exceptions may not be a best use of resources of audit department

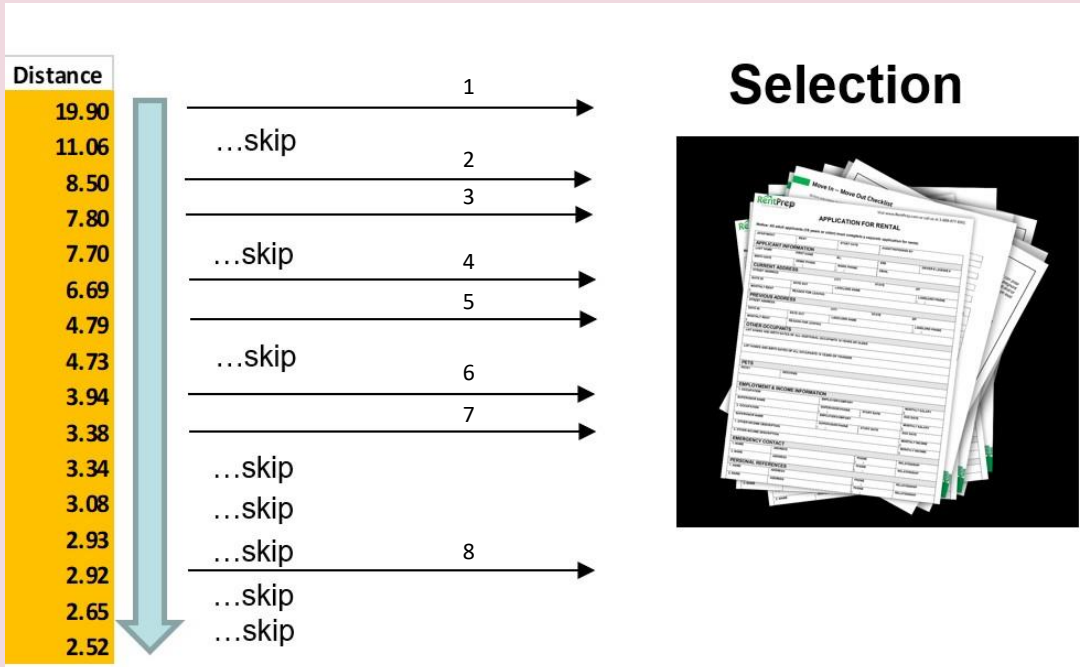
Audit Data Selection with Similar Values

#	Transact ID	Filter 1	Filter 2	Filter 3	Filter 4	Filter 5	Filter 6	Filter 7	Distance
1	739	110.7477	0	0	0	69.75034	0	0	92.68744
2	446	0	45.32469	0	0	0	1	0	92.44333
3	6147	30	0	0	0	117.4389	1	0	66.3032
4	4570	0	0	0	0	77.75064	2	0	51.81496
5	10042	0	0	5.538255	0	0	3	0	50.53067
6	2944	0	0	4.892452	0	0	0	0	43.96881
7	5355	0	0	0	0	62.30826	0	0	41.48429
8	7763	0	0	0	0	60.98294	0	0	40.6019
9	2844	0	0	4.471757	0	0	0	0	40.18799
10	11599	43.84615	0	0	0	46.778	1	0	36.90663
11	740	38.84616	0	0	0	23.43308	0	0	32.6268
12	434	0	15.62119	0	0	0	1	0	31.88407
13	10839	30	0	0	0	51.46008	1	0	30.45506
14	4569	18.46154	0	0	0	53.00628	2	0	29.71956
15	3688	0	0	0	2.446238	0	0	0	29.15685

- The excerpt above from a spreadsheet shows a list of transactions prioritized by their suspicion score. Transaction #7 and #8 have almost identical risk compositions. Thus, reviewing one of these transaction would resolve the second one as well.
- These research contributes to the literature by (1) using a distance metric as Suspicion Score and (2) proposing two algorithms that consider transaction similarities while selecting exceptions for review.

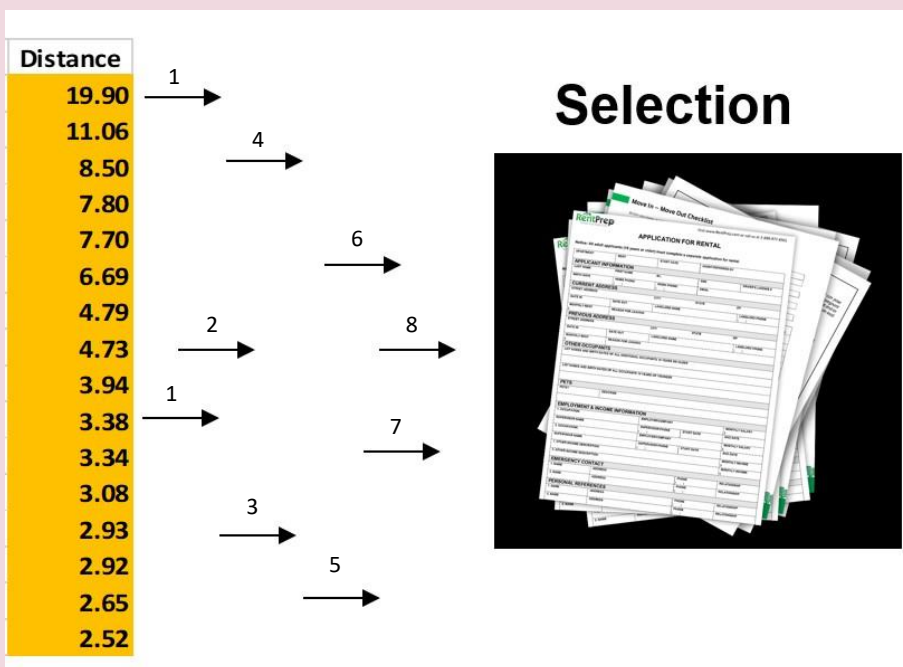
Skipper (Similarity Threshold)

- Set a maximum similarity threshold for the selection set
- Consider one exception at a time, running down the exception list sorted by suspicion score. Select the exception to the selection if no similar exception already exists in the selection set
- Terminate when selection size is reached, or exception list is fully covered

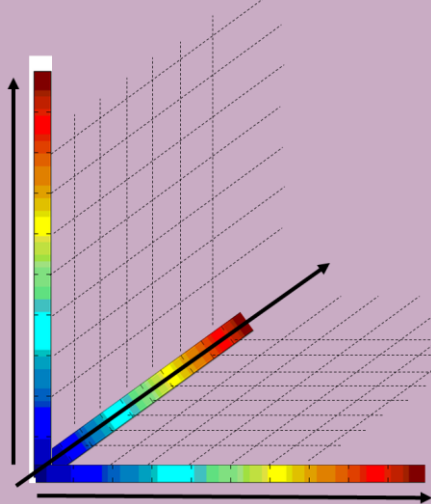


Stretcher (Maximin)

- Identify the set of the exceptions from population that are over certain minimum risk score
- From the selected set of the riskiest exceptions, select the combination of items that are the most dissimilar to each other for a given set size.
- Maximin Algorithm can be used to achieve a heuristic solution, since the ideal solution is computationally intractable



"Risk Space"



Distance Measures

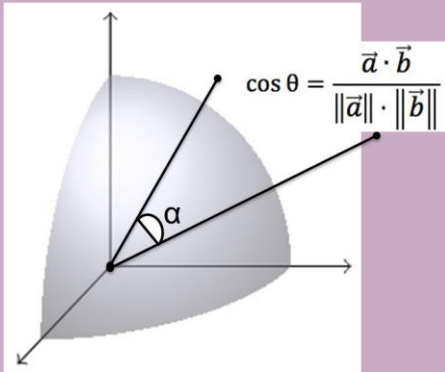
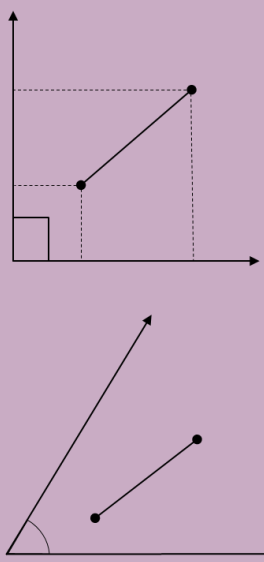
Euclidian

$$d(\mathbf{p}, \mathbf{q}) = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

Mahalanobis

$$d(\vec{x}, \vec{y}) = \sqrt{(\vec{x} - \vec{y})^T S^{-1} (\vec{x} - \vec{y})}$$

Similarity Measure



RUTGERS

Rutgers Business School
Newark and New Brunswick

Industry Classifications and Stock Performance: Constructing a New Scheme

Qingman Wu and Won Gyun No

Research Objective

Industry classifications are widely used in both academia and industry. Researchers use industry classification to control for the industry influence, limit the scope of their investigation, identify control firm, and et al.. The most common used schemes are built based on firms' business area and activities, which are likely to change a lot over time. However, those classification methods are rarely changing, which is unreasonable. So, the purpose of this research is finding out whether the proposed classification method improve the ability to organize firms into more homogeneous groups in terms of stock return co-movement of existing classification schemes.

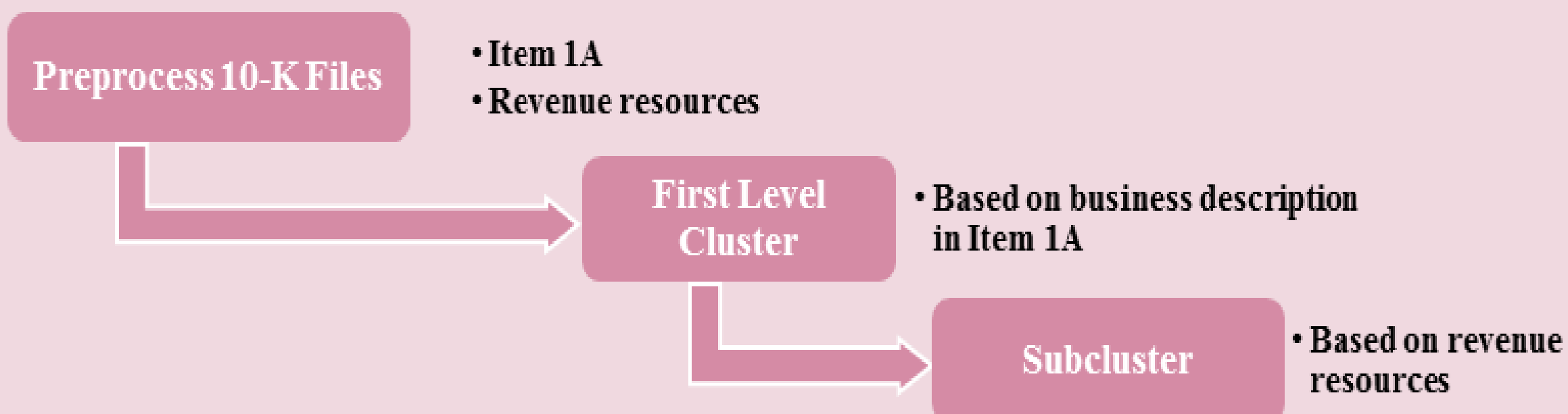
Literature

- Stefano Cavaglia et al. (2000) found that industry factors have become an increasingly important component of security returns. More importantly, diversification across industries now provides greater risk reduction than diversification across countries.
- If equity market participants consider a set of companies closely related, then stocks in the group should experience coincident movements in their stock returns. (Louis K. C. Chan, et al., 2007)

Data

- 10-K from EDGAR
- SIC and NAICS categorizations for each company from COMPUSTAT database
- Stock return data and firmcompls' financial data from CRSP

Methodology: Industry Classification Scheme Construction



Methodology: Industry Classification Schemes Comparison

This research would compare the correlations between stocks in the same industry with the correlations between within-industry stocks and outside-industry stocks to measure the homogeneity of firms. For example, there are k firms. For a particular industry I , there are N firms in this industry. For stock i in this cluster, we could average the pairwise correlations between stock i 's return and return of other stocks in this cluster is defined as equation (1), where ρ_{ij} is the time-series correlation between the return on stocks i and j . Similarly, the average pairwise correlation between stock i 's return and the returns of all other stocks not in its industry is equation (2).

Based on those data, we can calculate the average within-industry correlation and the average outside-industry correlation over all stocks in the sample.

$$\rho_{iI} = \frac{\sum_{j \in I, j \neq i} \rho_{ij}}{N - 1} \quad (1)$$

$$\varphi_{iI} = \frac{\sum_{j \notin I} \varphi_{ij}}{K - N} \quad (2)$$

Using Blockchain Technology for Continuous Assurance and Monitoring: A Closer Look at Cryptocurrency Lending Industry

Ruanjia Liu

Background

Financial technology (Fintech), emerged in the 21st Century, has been revolutionizing the stagnant capital loan market. Basically, Fintech disrupted traditional financial industry by expanding financial inclusion and cutting down on operational cost. Tetyana (2019) states that one of the innovations is Blockchain peer-to-peer lending (P2P lending) which utilizes superior algorithms (e.g., machine learning). This study will propose how to maintain continuous assurance and monitoring in the cryptocurrency lending industry

Blockchain: a New Type of Database

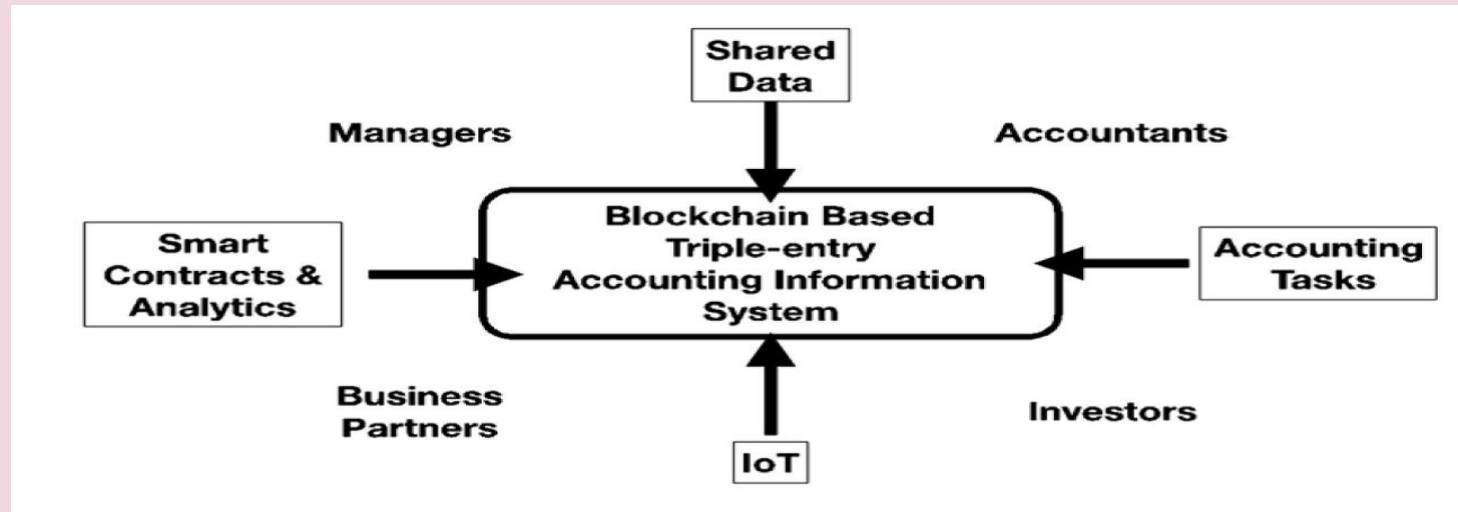
Peters and Panayi (2016) asserted that blockchain can avoid the possible conflicts when different users are making multiple modifications simultaneously within the distributed database system.

- 1) Blockchain is the distributed network.
- 2) Blockchain ensures integrity of the data stored in its ledger.
- 3) Unlike an ERP system that requires intensive human efforts, blockchain is designed to operate automatically with little third-party intervention (Swan 2015b).



Blockchain in Accounting and Assurance

Fanning and Centers (2016) stated that blockchain technology could benefit auditing when auditors compare the accounting entries on the books with real-time accounting on the blockchain. Kiviat (2015) suggested the concept of triple-entry accounting using blockchain. He asserted that posting accounting entries of Bitcoin transactions to the blockchain could prevent records tampering. ERP systems present accounting-specific modules based on Relational Database Management Systems (RDBMS) for process automation (Kuhn and Sutton 2010). In addition, ERP distribute real-time data for information analysis and decision making (Hitt, Wu, and Zhou 2002). Blockchain can be view as innovative type of database which operates automatously with little third-party intervention.



Peer-to-Peer Cryptocurrency Lending vs Securities Lending

The major difference between a cryptocurrency loan and a security loan lies on the lending mechanic. Traditional security lending relies on a central authority or a bank, which calls for the need for a trusted third party. P2P lending, also known as crowd-lending, is the practice of matching lenders with borrowers via two-sided platforms. For traditional security lending transaction, the collateral type that can be cash and securities and rebate rate are agreed by the borrower and the lending agent. But collateral in crypto lending are crypto assets no matter if it is cash loan or cryptocurrency loan.

- Two types of crypto lenders:
 - Custodial lenders and Non-Custodial lenders
- Three types of crypto borrowers:
 - Speculating / Hedging
 - Trading / Arbitraging
 - Operating Working Capital

Applying Blockchain-Based Assurance to P2PLending

- Margin Call and Margin Lending Risk

The platform will make a margin call to the borrower before the value of the collateral fall below the maximum LVR. The borrower can choose to increase collateral or to pay off some principal. But if the borrower fails to meet the maintenance margin, the platform could liquidate partial loan by taking over the collateral to recover the principal and outstanding interest.

$$\text{loan amount} / \text{total value of portfolio} = \text{Loan to Value Ratio}$$

- Provide Blockchain-Based Non-default Behavior

A smart contract could provide control-based assurance paradigm. Platforms could continuously monitor loan amount and loan to value ratio (LVR). Smart controls could contact the borrowers immediately as long as loan to value ratio (LVR) indicates that collateral value falls below the maintenance margin. Then the borrowers either invest more in the collateral or pay back partial of the cryptocurrency loan.

- Provide Blockchain-Based Assurance of Capital Reserve

If lending platforms disclose the cryptocurrency transactions on Blockchain ledger on the daily basis, the public could keep track and monitor the changes of fiat reserve every day.

Applying Blockchain-Based Assurance to P2PLending

- Provide Blockchain-Based Assurance of Collateral

Blockchain can be utilized to increase the information auditability of collateral. Since a blockchain secures the data on the ledger, it could also lend authenticity to many audit documentations. Because blockchain does not allow erasing records, audit trails could be documented to facilitate tracing and review in the future (Ernst & Young [EY] 2015). Those documents could also be shared among related parties for cross-validation. Thus, lending platform could keep track of certain documents of collateral which are filed on blockchain. The absence of any records might imply fraud and default. Potential lenders who are assured of collateral existence could be more likely to participate in the cryptocurrency lending activity. Placing blockchain technology in the hands of platforms, borrowers, and lenders can achieve a new level of assurance. These parties may involve in the transaction verification process by providing reliable and independent information for audit and attestation purpose. The collaboration of these individuals could provide trusted real-time assurance through the “proof of transaction” mechanism.

The Future of Advanced Technology and Automation in Audit: A Delphi Study

Danielle R. Lombardi and Sheneya Wilson

The Original Study

In 2014, Lombardi et al, performed a Delphi study using a sample of leading industry experts with hopes to predict the future of the auditing profession.

The major areas of focus included:

- The evolution of the relationship between the internal and external auditor.
- The effects of judgement on automation.
- The shift in timeliness of the audit cycle.
- The need for a more global perspective regarding the changes in the auditing profession.
- How client technology is leading auditing procedures
- Understanding the impact of privacy safeguards on technological adoption
- The use of automation in the auditing profession.

After doing a review of the current literature and publications posted by industry leading firms, we were able to obtain evidence supporting the conclusiveness of the predictions made during the original Delphi study in in 2014.

Analysis of Session Highlights

Below are some examples of how accurate these prediction were in describing the current state of the audit profession almost a decade after the original study was conducted.

Prediction: External auditors will rely on more internal audit work in the future.

Analysis

Current research suggests that external auditors will increasingly use the work of internal auditors as they gain efficiency in performing their audit work. Additionally, when relying on the work of internal auditors, external auditors appear to be more conscious of the consequences of the audit quality they deliver (Aregnto et al 2018).

Prediction: Although the use of automation will increase, judgement and decision making cannot be automated.

Analysis

Auditor judgement cannot be easily replaced by machine due to the fact that automatable tasks are those that are highly repetitive, simple, rule based, and time consuming (Cohen et al 2019). Tasks that require auditor judgement typically do not have the aforementioned characteristics.

New Delphi Experiment

After thoroughly analyzing each of the suggestions from the original Delphi and realizing that most of these predictions did manifest in the profession today, we decided to re-do the experiment with a similar group of industry leading experts.

The method consisted of two rounds of brainstorming and a Delphi experiment.

Brainstorming:

Practice-based research aims to eradicate the existing knowledge gap between academics and practitioners and effective group brainstorming can assist in this endeavor. We used brainstorming to get the creative thoughts of our participants going.

Delphi Method:

This method involves providing industry experts with at least two rounds of a questionnaire and structured feedback in between each sessions in order to enhance the respondents' consensus (Bell 1967). This method has been able to accurately forecast future outcomes and predict the direction of many different industries including banking, human resources, business administration and management, information systems, and accounting and auditing (Bell 1967; Bradley & Stewart, 2003; Poba-Nzaou et al 2016, Hong, Trimi, Kim, & Hyun, 2015; Worrell, Di Gangi, & Nush 2013)

Overview of Delphi Panel Members
<ul style="list-style-type: none">• Ex-chairman of the Financial Accounting Standards Board (FASB)• Ex-chairman of the AICPA• President at CEO of a Consulting Company• Retired Audit Partner• Big 4 Partner in Audit• Accounting Information Systems Professor• Senior Manager of Audit Analytics

Figure 1 - The Expert Participants

FIGURE 2 Highlights and Recommendations Provided by Experts	
First Session <ul style="list-style-type: none">• External auditors will rely more on internal audit work in the future• Although use of automation will increase, judgment and decision-making cannot be automated• The view of many of the topics would vary depending upon the evolution of the financial statements• Audit will be cycled over the year, instead of only at year-end• There is a need for a more global perspective	Second Session <ul style="list-style-type: none">• Client technology is leading audit procedures• The use of technology depends upon proper safeguards for privacy (i.e., HIPAA)• Automation can be used for more tedious tasks so that auditors can use their expert judgment for more pressing issues

This figure was obtained from Lombardi et al. (2014).

Listed below 8 of the 11 questions asked during the 2020 Delphi experiment. Results are still being analyzed similar to the analysis conducted in the original 2014 paper using a combination of descriptive statistics and thorough analysis of the responses.

- The automation of Internal and External Audit.
- How Big Data will impact auditing/assurance procedures, analytics, assessments, judgement
- What will happen to the use of emerging technologies in audit (i.e., Blockchain, Robotic Process Automation, Artificial Intelligence, Cloud Computing, Neural Networks, Expert Systems, Data Clustering, Regression, other).
- What will happen to the judgment biases in auditing as a result of interaction with technology and artificial intelligence.
- To what extent will technology take over the responsibility of the internal and external auditor (i.e., audit task, decision making, judgement)
- As a result of increased application of continuous auditing, to what extent will internal auditors take over the responsibility for financial auditing which is currently undertaken by the external auditor.
- What will happen to the ethical concerns in auditing as a result of increased usage of emerging technologies and artificial intelligence.

The Methodology for Thinly Traded Cryptocurrency Valuation

Eyal Beigman, Gerard Brennan, Sheng-Feng Hsieh, and Alexander J. Sannella

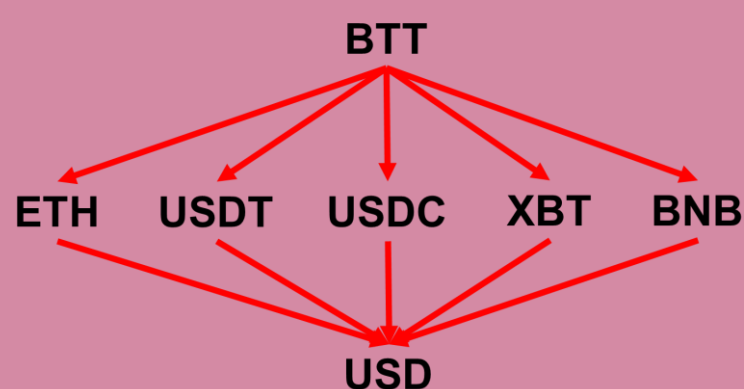
Research Objective

- Following the methodology for *actively* traded cryptocurrency from Beigman, Brennan, Hsieh, and Sannella (2020), we further develop another methodology for *thinly* traded cryptocurrency.
- The practice has a need for cryptocurrency fund valuation (SEC 2020).

Literature

- Transaction volume** on cryptocurrency exchanges
 - is associated with the **market credibility** (Nasiri et al. 2018).
 - works as a **channel** for aggregating private and public market information and facilitating coordination on equilibrium prices (Bianchi and Dickerson 2020; Brandvold et al. 2015; Makarov and Schoar 2019; Park and Chai 2020; Sockin and Xiong 2020).

The Optimal Path Approach



- Identify the *Principal Market* based on Beigman, Brennan, Hsieh, and Sannella (2020) and determine exchange rate (ER) for each pair *minutely*.
- Identify “*Bottleneck Volume*” for each path candidate *minutely*.
- Identify the “*Optimal Path*,” the path with the **maximum** of bottleneck volume from all path candidates *minutely*.
- $FV_{BTT-USD} = ER_{BTT-XXX} * ER_{XXX-USD}$
(XXX could be ETH, USDT, USDC, XBT, or BNB)

Data of the Pilot Test

- BTT-USD** as the target
- Five optimal path candidates
- February 1st to 29th, 2020 (41,760 minutes)
- Tick information for all targeted pairs
 - time (milli-second level)
 - exchange
 - price
 - volume

The Results of the Pilot Test

Figure 1 The exchange rates from the Optimal Path approach (red) and the real BTT-USD trades (black) (February 1 to February 29, 2020)

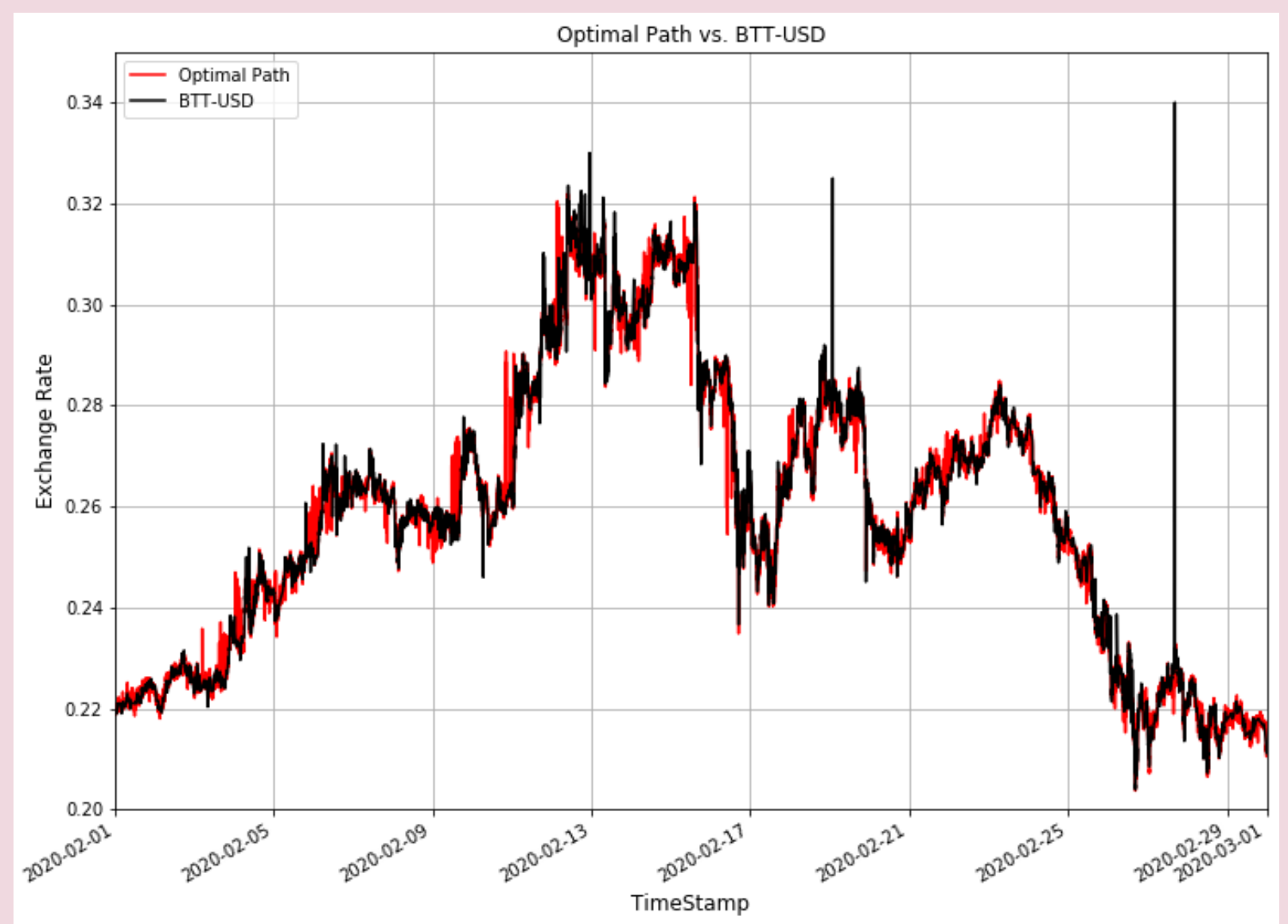
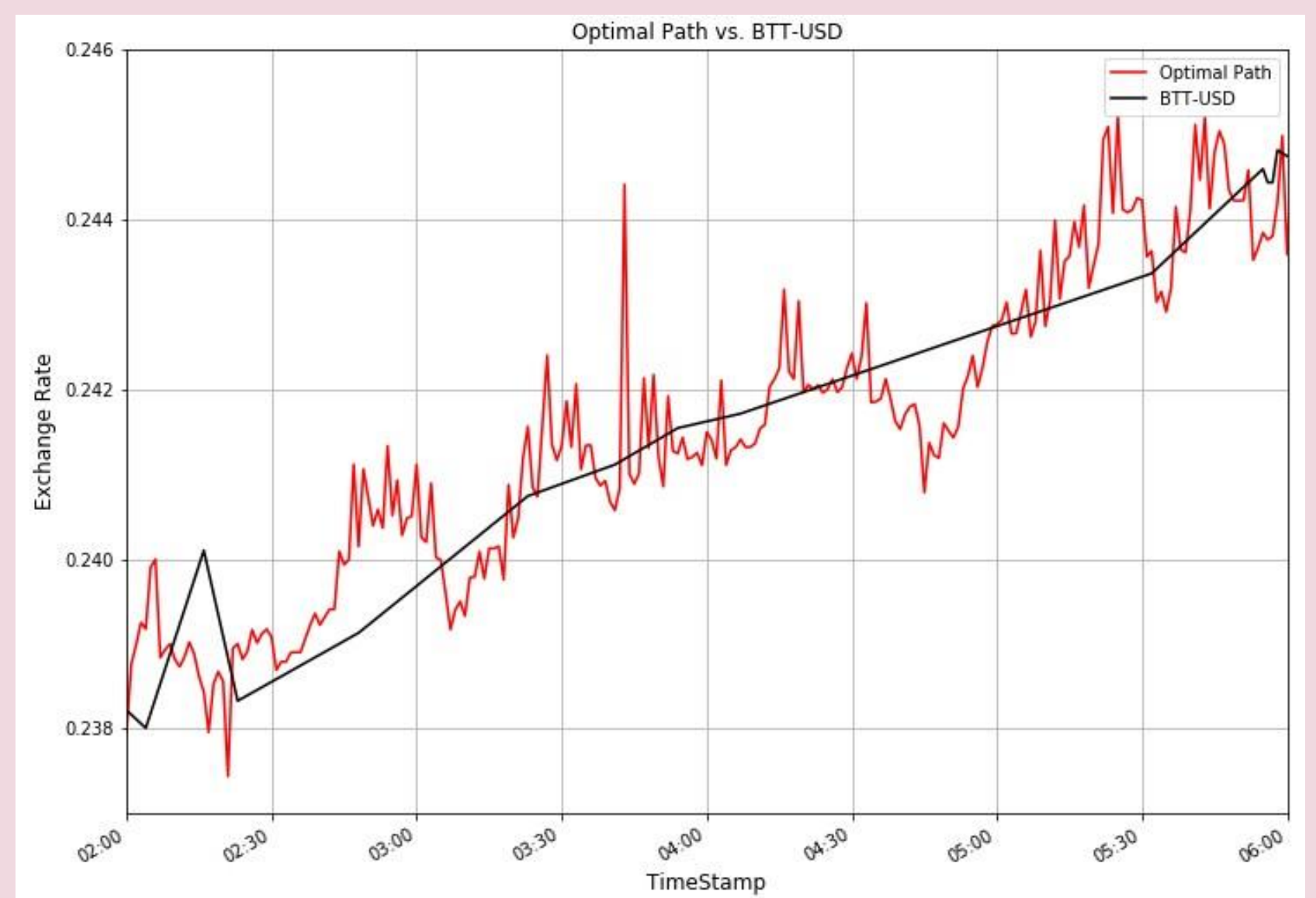


Figure 2 The exchange rates from the Optimal Path approach (red) and the real BTT-USD trades (black) (2 – 6 AM, February 5, 2020)



Dynamic View of Pandemic Circumstances with Government Interventions and Social Factors

Wenru Wang, Marcelo Freitas, Fabricia Silva da Rosa, Miklos A. Vasarhelyi

Introduction

The relationship between social factors, government policies and pandemic outcomes have been studied by researchers in social science, but they somehow miss some important issues:

1. Current studies tend to ignore the fact that the development of an epidemic is dynamic, people that are infected will recover and recovered people will be susceptible to the virus again. Therefore, looking at the figures statically may be not comprehensive.
2. Most of them use data resources from government releases. Can we trust these numbers?
3. Previous studies either only explain the Long-term impact of government policies, and they provide limited prediction.

Research Question

How do governments evaluate and predict the impacts of a proposed policy before it is actually implemented ?

Methodology

We incorporate the theory of System Dynamics to understand the interactions of pandemic outcomes, social factors, and government interventions, and we construct a dynamic model to predict and visualize the possible outcomes of government interventions. The theory of System Dynamics, "the science of feedback behavior in social systems," was proposed by Jay Forrester in 1961.

Figure 1 presents a basic SEIR (Susceptible—Exposed—Infectious—Recovered) model to illustrate the spread of COVID-19 among a large population and how the population amount of the four stages develop from one to another.

We then expand the basic SEIR model and include government interventions and social factors to observe the possible outcomes of the newly introduced parameters. Figure 2 presents the expanded dynamic model.

Main Results

The deliverables of this study render a dynamic view to understand the current pandemic and provide government with a simulation tool to visually evaluate the impacts of government interventions at different levels.

The proposed model allows policy-makers to effectively evaluate their policies before implementation, and such tools are of great need when the high transmission speed and complex symptoms of the novel virus leave governments limited time to react. The dynamic model could be applied for a more comprehensive view of the pandemic and other future crises. Researchers in public health and the general public could adopt the idea that pandemic development interact with different social factors.

Figure 2

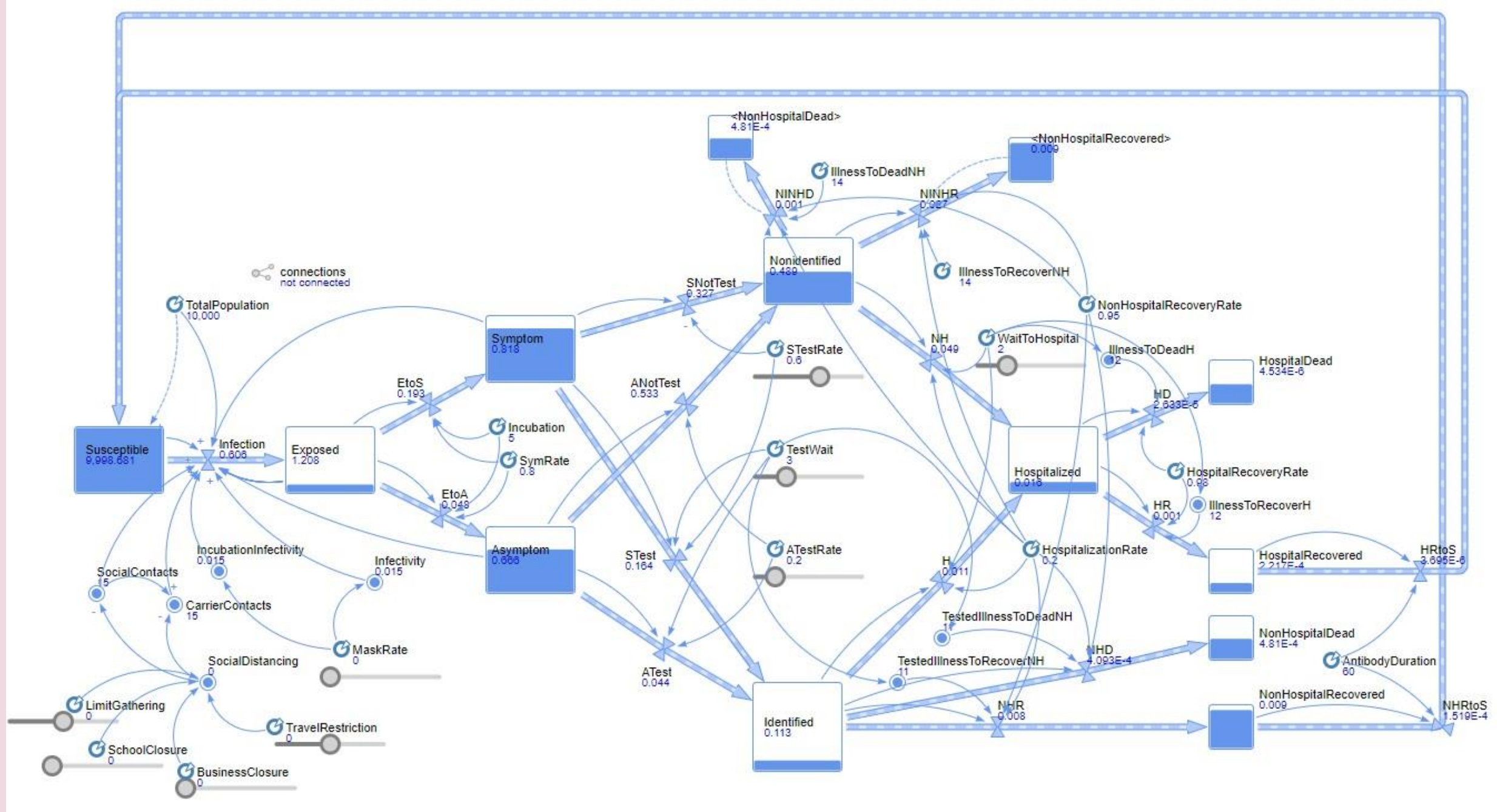
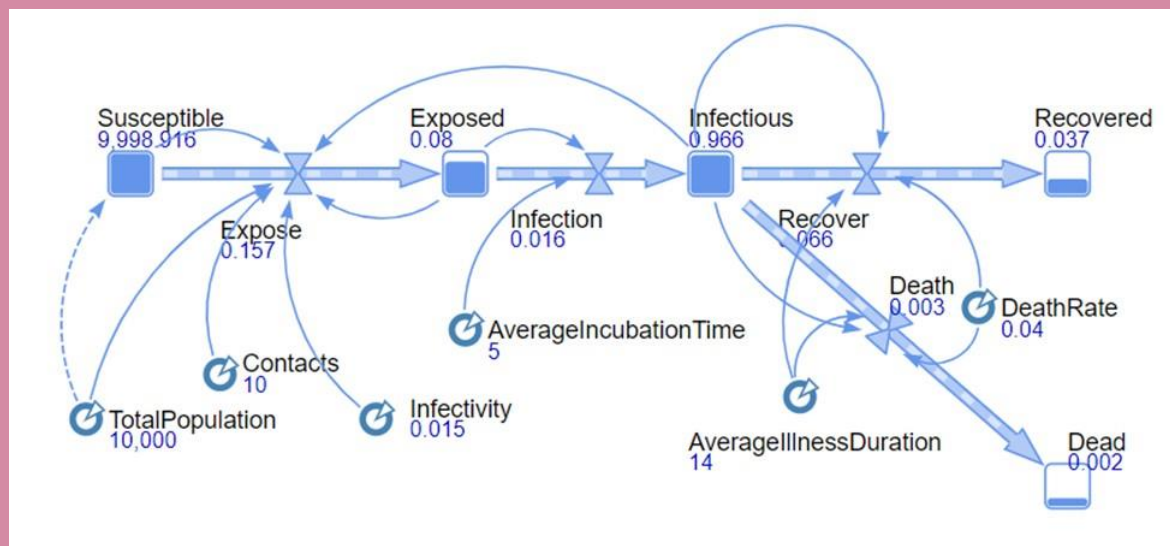


Figure 1



Continuous Monitoring and Audit Methodology for Medication Procurement

Wenru Wang and Miklos A. Vasarhelyi

Introduction

Governments of different countries and regions have been devoting to fight the ongoing COVID-19 pandemic with various policies and interventions. Some of the policies include financial supports. For example, the CARES Act in the United States established the US\$150 billion Coronavirus Relief Fund for state and local governments (US Treasury, 2020). Brazil also announced economic stimulus package with US\$ 17.5 billion financial support to states. However, the current procurement system does not ensure the efficiency of the procurements. Data quality and procurement wastes are some of the issues that exist in current government procurement. The government performs external audits annually, long after most of the procurement wastes and data errors have occurred.

Research Question:

How can possible abuses and wastes in the government procurement be efficiently detected?

Methodology

With the partnership of one Brazilian City hall, we analyze over 40,000 government procurement medical datapoints. We propose a continuous audit methodology that helps internal auditors monitor and analyze government procurement data on a more timely basis.

Figure 1 presents our overall research design.

We first preprocess the procurement data, so that they will be ready for further analysis. Through total value ranking in step two and unit value comparison in step three, certain alarms will be triggered and send to the auditors. The alarms including medicines of high total purchase values and the alarms notify auditors to investigate on these medications. we generate a monitoring dashboard at step four for auditors to easily view alarms and can trace back to questionable procurement records.

Main Results

Figure 2 presents the final product of the continuous monitoring dashboard. With the dashboard, auditors could identify and trace back the questionable procurement records that have either great total values, or higher unit prices.

The methodology also incorporates text mining techniques, so that internal auditors could compare the procurement records with federal catalog to further find price savings opportunities.

Figure 1

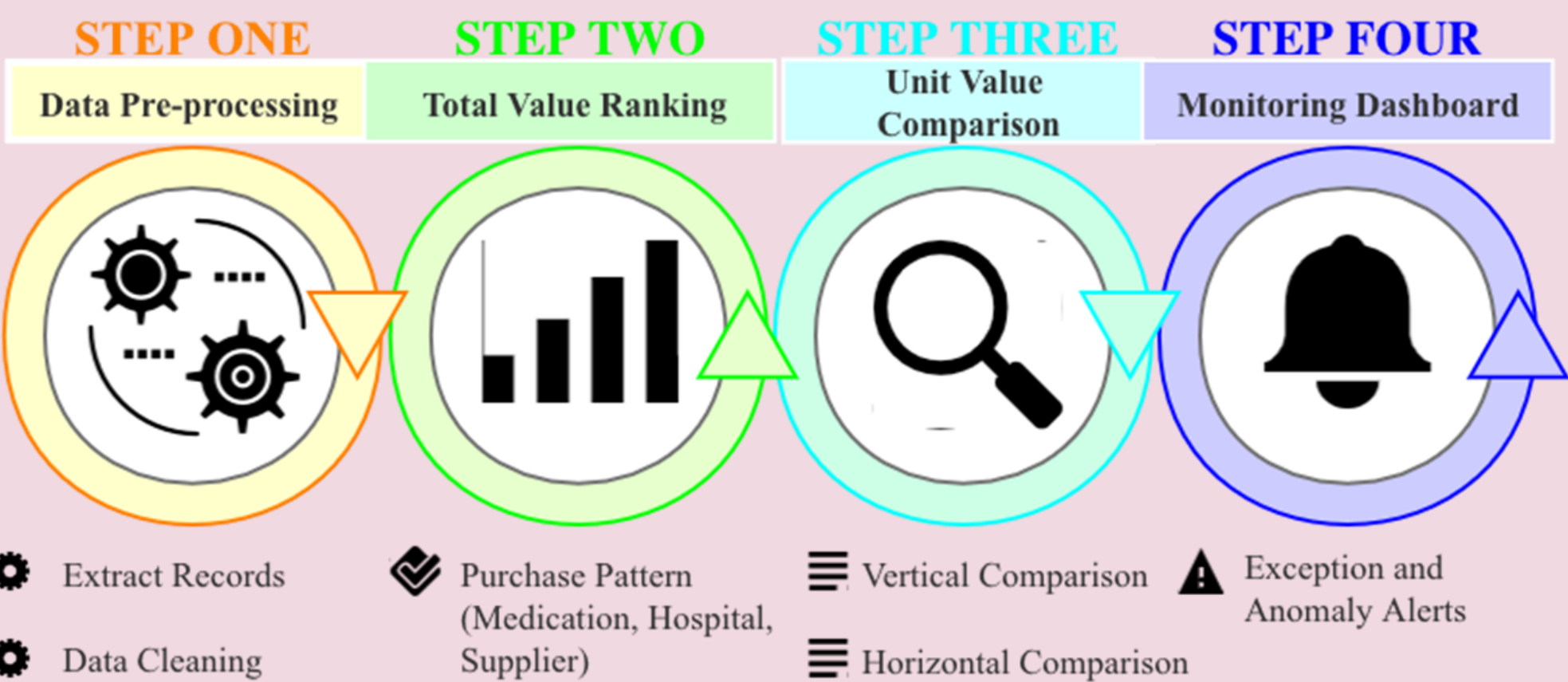
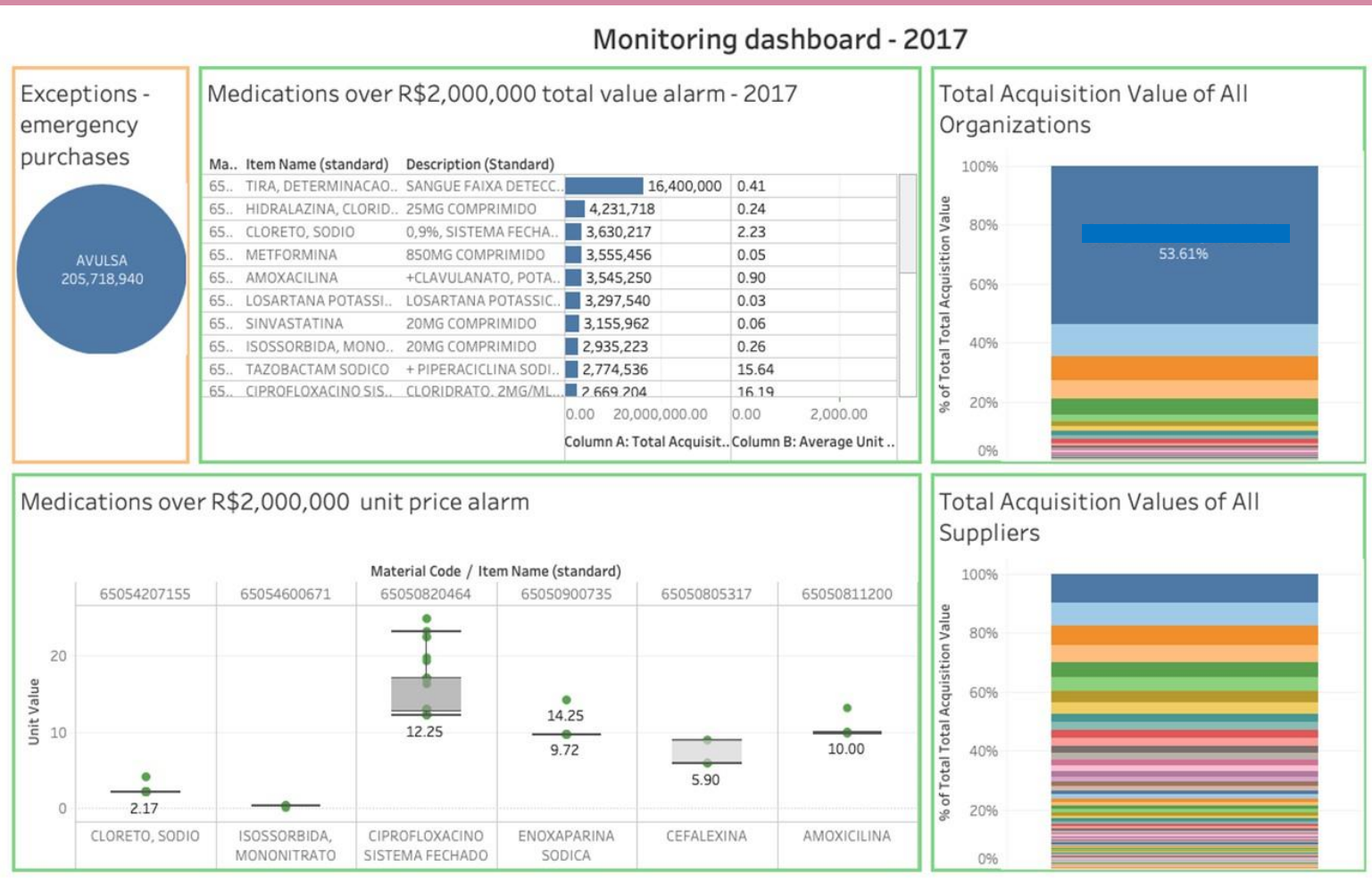


Figure 2



Improve Text-Readability on Cybersecurity Disclosure: A Social Network Approach

Hongmin W. Du

Cybersecurity Readability

This work is motivated by research issues identified in increasing occurrence and cost of cybersecurity risks in accounting and financial applications. Many firms and users today (especially in the current pandemic COVID-19) heavily rely on cyber technology to conduct essential business operations (e.g. online bank transaction, e-filing of 10-K, tax documents, etc). In such wide use of cyber technology for accounting and financial application, cybersecurity presents serious threats and risks to US capital market as stated in the investing public and the U.S. economy depend on the security and reliability of information and communications technology, systems, and networks (Commission Statement and Guidance on Public Company Cybersecurity Disclosures, 17 CFR Parts 229 and 249 [Release Nos. 33-10459; 34-82746]). In February 2018, the Securities and Exchange Commission (SEC) issued a guidance requiring public firms to file disclosure obligations under existing regulation regarding cybersecurity risk and incidents.

The main concerns (issues) of such guidance are as follows: (1) It is not objective rule (no quantitative cybersecurity standards), instead just ask companies to provide subjective descriptions; (2) It has been argued by users that usefulness of lacking of informative boilerplate. Thus, the information provided in cybersecurity disclosure is descriptive, for which, the readability is an important problem. There exist two issues closely related to this problem in the literature.

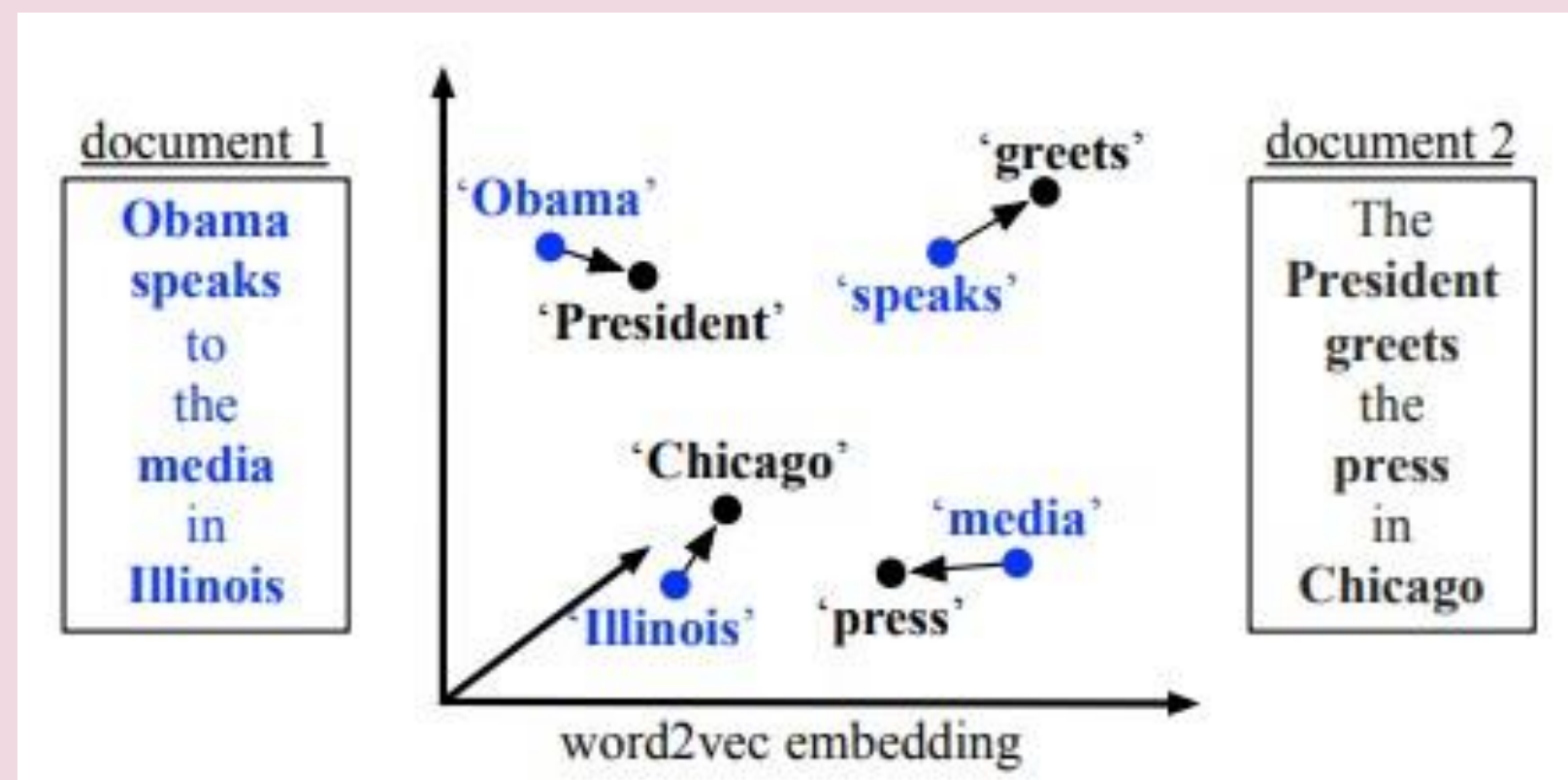


Figure 1: Word Mover Distance. The word mover distance is one measure to be used to calculate the similarity between two words. The word will be graphed based on a chosen distance measurement which is not limited to only the Word Mover Distance. We will test which distance measure provides us with the best results in content analysis.

Text Content Analysis

Content analysis in finance and accounting domain aims to objectively identify the word conveyed by descriptive (qualitative) information to explore the associations of market reactions with such subjective information via quantifying document tones/levels. The 1st pioneering paper to quantitatively examine the interactions between the media and the stock market via daily text content from WSJ (Abreast of the Market) column on US stock market returns was presented by Dr. Paul C. Tetlock in 2007, in which a text factor was constructed by studying the pessimistic words (using Harvard IV-4 psychosocial dictionary) in news (WSJ column) and a simple linear regression tests was designed for the predictive capability of the text factor (i.e. using these negative words to research how the media can affect the stock market.)

Many similar works to study the association of capital market forecast with the tone of the subjective financial news are conducted by researchers in accounting and finance (Jegadeesh et al. 2013, Tetlock et al. 2008, Schumaker et al. 2012). Schumaker designed a financial news article prediction artifact (Arizona Financial Text (AZFinText)) for this study. Loughran and McDonald (2011) create a comprehensive dictionary (LM list) from 10-K reports and discover that the negative word list captures the tone of 10-K reports better than the Harvard IV-4 list. Several alternative classifiers (e.g., naive Bayes and vector distance) are also adapted in content analysis to extract investors messages and opinions from social media posts.

Measuring Text Readability

In accounting and finance literature, extensive works study the definition and measurement of readability (Loughran and McDonald 2016). The essential issue is about what is meant by the concept in the context of business writing (Loughran and McDonald 2016). The popularized/simple readability measure, Fog Index may be not suitable and useful for accounting and financial documents such as cybersecurity disclosures, MDA disclosure, auditing files, etc (Li 2008).

Therefore, the problem of how to define better readability to more precisely reflect the actual comprehension process is further addressed in the literature. Loughran and McDonald (2014) conducted a large sample of 66,707 annual reports empirical research to show that the Fog Index is a poorly specified readability measure when used in financial accounting

Social Network Approach

In this paper, we propose a new approach, a social network approach, which emphasizes on word classification. A social network approach. We first introduce a distance between words, which reflects the meaning of the words; the shorter the distance is, the closer the meaning. Using a word list for accounting and finance, we then construct a social network with the distances between each word. Figure 1 shows an example how it would look. Then, we use an existing algorithm to partition the social network into communities. Each community would consist of words with close meaning.

Next, we may employ regression model to establish weighting of each community in relationship with stock market. The advantage of this approach has two fold. The first is that the number of communities is much smaller than the number of words, and hence content analysis can be simplified. The second is that words in the same community would influence stock market similarly. Accumulating such an influence together could make it to be observed easily.

Methodology

To realize our idea, we would meet several challenges and employ some methodologies to them.

Choice of Word Distance: There exist several distances of words in the literature, such as the Levenshtein distance, the Damerau-Levenshtein distance, the longest common subsequence, the Jaro distance, and the Hamming distance. Which one is more suitable for our purpose? We need to establish a method to give a comparison.

Choice of Algorithm for Community Partition: There exist many algorithms for community partition. Clearly, we would consider one based on connections, which still has many choices. Therefore, we also need to establish a method to compare them. Choosing one of the algorithms is one of the From comparison, select a suitable

