

Customs Import Declaration Datasets

Chaeyoon Jeong
KAIST
Daejeon, Republic of Korea
lily9991@kaist.ac.kr

Sundong Kim
Institute for Basic Science
Daejeon, Republic of Korea
sundong@ibs.re.kr

Jaewoo Park
Korea Customs Service
Daejeon, Republic of Korea
jaeus@korea.kr

ABSTRACT

Given the huge volume of cross-border flows, effective controlling of trades becomes more crucial in customs administrations. However, limited accessibility of the customs datasets hinders the progress of open research and lots of member countries have not benefited from the recent progress. In this paper, we introduce an import declarations dataset to facilitate the collaboration between the domain experts in customs administrations and data science researchers. The dataset contains 54,000 artificially generated trades with 21 key attributes and it is synthesized with CTGAN while maintaining correlated attributes. Synthetic data has several advantages. First, releasing the dataset is free from restrictions that do not allow disclosing the import data even if its anonymized. Second, the fabrication step minimizes the possible identity risk which may exist in trade statistics. Lastly, the published data follow a similar distribution to the source data so that it can be used in various downstream tasks. With the provision of data and its generation process, we open baseline codes for fraud detection tasks, as we empirically show that more advanced algorithms can better detect frauds.

CCS CONCEPTS

• **Social and professional topics** → **Taxation**; • **Applied computing** → **E-government**.

KEYWORDS

Import Declarations, Synthetic Data, Tabular Data, Customs Fraud Detection, Correlation Analysis

ACM Reference Format:

Chaeyoon Jeong, Sundong Kim, and Jaewoo Park. 2022. Customs Import Declaration Datasets. In *Proceedings of 31st ACM International Conference on Information and Knowledge Management (CIKM '22)*. ACM, New York, NY, USA, 5 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

Customs clearance is the process of getting permission from the customs administrations to either move goods out of a country (export) or bring goods into the country (import). The customs declarant handles the customs formalities of the goods and declares it to the customs office, and the permission is given only when the declaration is legitimate. If the shipment exceeds the threshold

value (\$150 in Korea), the customs impose tariffs on the item. Once the tariff is collected, they allow the goods to be taken out.

Since the tariff rate is determined by the combination of the type, price, and country of origin of the product, accurate declarations are crucial to the customs clearance process. Faultless declaration is a key to increase efficiency, as customs rectifies inaccurate reports and often inspects the goods if needed. Accurate import declarations enable customs to save a lot of workforce for physical inspection of the goods, and importers can receive their goods on time. To reduce the error, the Korea Customs Service (KCS) promotes the guideline for import declarations using the UNI-PASS e-clearance system.¹ However, the process is complicated for individuals to handle since they are not familiar with technical terms and the traded goods are not easily fit into a single category. Customs authorities also struggle to advance the logic behind their e-clearance system to facilitate the process [5, 6, 11].

Given the enthusiasm around the use of data and the possibilities offered by artificial intelligence [9], the adoption of new technology is relatively slow in customs community. The major reason is the lack of publicly available data. Disclosure of import declaration data outside customs is strictly prohibited because of its sensitivity. Only authorized departments or institutions could conduct research internally, and there is no visible community effort taking place.

This lead us to design a synthetic data that can be open to public. The dataset contained in this paper includes 54,000 artificially generated trades with 21 key attributes.² By using tabular synthesizer with post-processing technique, we maintain the distribution and correlation in the synthetic dataset are similar to those of source dataset. Empirical results on fraud detection demonstrate the potency of the data. Meanwhile, the data is used for competition in three universities to develop algorithms that can be applied to the actual work process. We conclude the paper by discussing possible scenarios to use the data and summarizing necessary thoughts on the data synthesis.

2 DATA DESCRIPTION

2.1 Data Schema

The tabular form dataset consists 54,000 import declarations, where each row describes the report of a single goods. Among 62 attributes in the import declaration form,³ the data includes 21 representative attributes without overlapped or less important attributes. The first 19 values are prepared in the declaration stage of the customs clearance, while the rest two attributes are labeled after the inspection. *Fraud* indicates whether the inspected result of the actual imported goods conflict with its declaration. *Critical fraud* is the case that

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

CIKM '22, Oct 17–22, 2022, Atlanta, USA

© 2022 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-XXXX-X/18/06.

<https://doi.org/XXXXXXX.XXXXXXX>

¹<http://unipass.customs.go.kr>

²<https://github.com/Seondong/Customs-Declaration-Datasets>

³Import declaration format and explanation: <https://bit.ly/import-declaration-form>

Table 1: Data description of customs declaration dataset.

Attribute	Value	Explanation
Date	2020-01-01	Date when the declaration is reported
Office ID	13	Customs office that receives the declaration (e.g., Seoul regional customs)
Process Type	B	Type of the declaration process (e.g., Paperless declaration)
Import Type	11	Code for import type (e.g., OEM import, E-commerce)
Import Use	21	Code for import use (e.g., Raw materials for domestic consumption, from a bonded factory)
Payment Type	11	Distinguish tariff payment type (e.g., Usance credit payable at sight)
Mode of Transport	10	Nine modes of transport (e.g., maritime, rail, air)
Declarant ID	L77JJEG	Person who declares the item
Importer ID	HQ0W7JA	Consumer who imports the item
Seller ID	PBP2MYI	Overseas business partner which supplies goods to Korea
Courier ID	MWIDNS	Delivery service provider (e.g., DHL, FedEx)
HS10 Code	0901210010	10-digit product code (e.g., 090121xxxx = Coffee, Roasted, Not Decaffeinated)
Country of Departure	JP	Country from which a shipment has or is scheduled to depart
Country of Origin	JP	Country of manufacture, production or design, or where an article or product comes from
Country of Origin Indicator	B	Way of indicating the country of origin (e.g., B = Mark on package)
Tax Rate	8.0	Tax rate of the item (%)
Tax Type	A	Tax types (e.g., FTA Preferential rate)
Net Mass	1262.0	Mass without any packaging (kg)
Item Price	1437418.0	Assessed value of an item (KRW)
Fraud	1	Any fraudulent attempt to reduce the customs duty? (0/1)
Critical Fraud	1	Critical case which may threaten the public safety (0/1)

may threaten the public safety or stability of society, such as copyright infringement, drug smuggling, false declaration of origin of goods. Categorical attributes and their values follow the handbook provided by KCS, which contains trade codes used for filling out import and export declaration in Korea.⁴ Detailed data descriptions and example values are shown in Table 1.

2.2 Data Analysis

To show that the synthetic data is similar to the source data and thus reliable, we conduct statistical tests between two attributes of each dataset. Chi-squared test is used to match the expected frequencies for the categorical variables, and two-sample Kolmogorov-Smirnov tests whether cumulative distribution function (CDF) of a variable follows the CDF of the other variable. Both tests the similarity of two datasets. We used synthetic data evaluation functions built in the Synthetic Data Vault (SDV) library. Both score values range from 0 to 1, where higher score indicates that the synthetic data has more similar distribution to the source data. As shown in the table 2, our synthetic dataset is quite indistinguishable from the source dataset.

In addition, we compare the distribution of attributes of the two datasets. Figure 1 show the distribution of representative attributes (Import Type, Country of Departure, Country of Origin Indicator, Critical Fraud) in the dataset is analogous to that of the source dataset. Figure 2 illustrates that the correlation between the attributes in the synthetic data is similar to that of the source data. Figure 3 shows representative features in the downstream fraud detection task performed on each datasets by using the XGBoost model. The tendency of feature importance score is analogous in

Table 2: Evaluation score of our synthetic dataset.

Metric	Score
Chi-Squared Test	0.705
two-sample Kolmogorov-Smirnov test	0.783
Continuous Kullback–Leibler Divergence	0.999
Discrete Kullback–Leibler Divergence	0.416

the two data, "Importer ID", "Item Price", "Net Mass", "Declarant ID", and "HS10 Code" were considered important with high feature importance score.

3 DATA GENERATION

3.1 Preprocessing

Among 24.7 million customs declarations reported for 18 months between January 1, 2020 and June 30, 2021, we used the inspected (i.e., labeled) part of the declarations to synthesize our data. Inspected items account for a fairly small percentage of the total, but they are more accurate, all validated by field officers. We call it as the source data throughout the paper. Identifiable information such as "Importer Name" in the source data are anonymized into "Importer ID". The price of goods traded between vendors (i.e., Item Price) can be problematic when fully disclosed, so we add some Gaussian noise to the average price of each category of item.

3.2 Generating Data with CTGAN

We used conditional tabular GAN (CTGAN) [12] from the Synthetic Data Vault library to generate the data. CTGAN is a specialized for

⁴<https://www.data.go.kr/data/3040477/fileData.do>

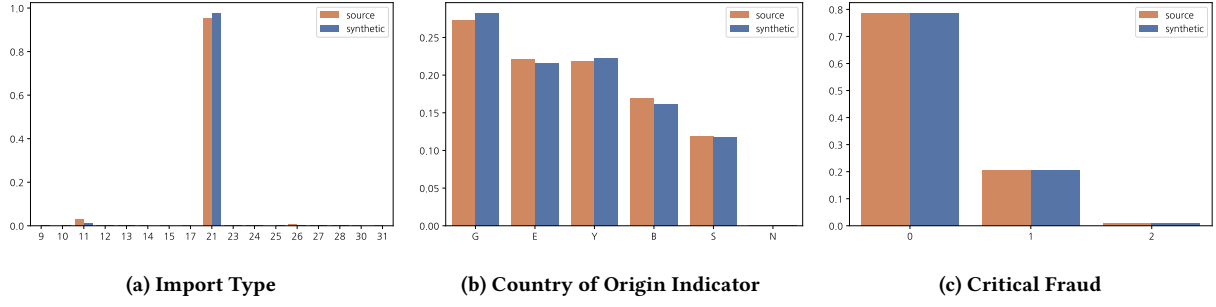


Figure 1: Representative feature distributions are similar between synthetic data and source data.

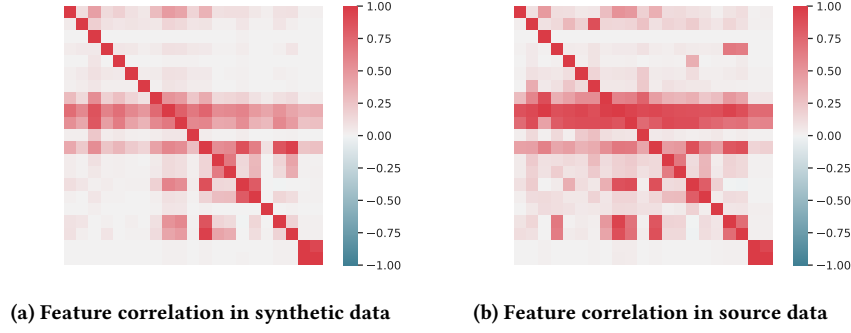


Figure 2: Feature correlation in synthetic data and in source data are also similar to each other.

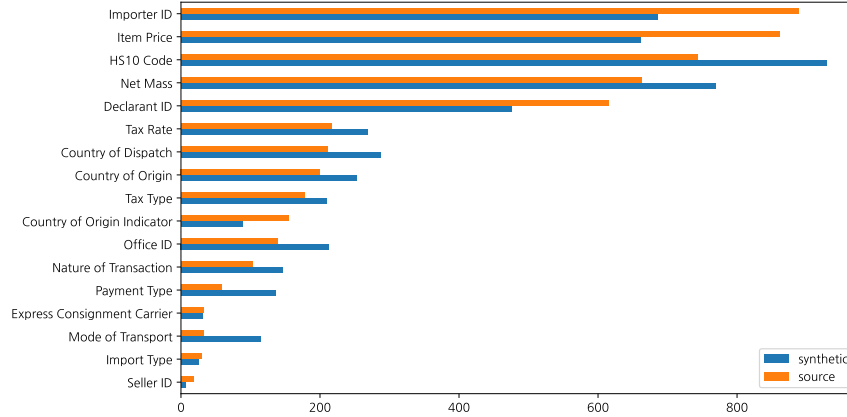


Figure 3: Important features for fraud detection task are also similar between the two datasets.

tabular data and it uses conditional techniques to handle imbalanced discrete variables and multi-modal continuous variables. Compared to other tabular generative models such as TGAN [13] or TVAE [12], CTGAN showed the most realistic output to our dataset, preserving the relationship between columns. The data generation process can be done in serial or in parallel manner. Users with limited resources can split the data in chronological order, train the CTGAN model on each split, synthesize samples from each model, and aggregate the result. For each split, the model is trained for 100 epochs.

3.3 Maintaining Correlated Attributes

One of the important features of tabular data is correlation among attributes. To generate an realistic import declaration data, correlated attributes and their values should be maintained during the generation process. Although CTGAN is a tabular-specific generative model, the generated output often show distorted relation between attributes that have clear and strict relevancy. For example, attributes "HS10 Code"—"Country of Departure"—"Country of Origin"—"Tax Rate"—"Tax Type" are highly correlated to each

other based on customs valuation policies. The correlation must be preserved also in the synthetic data. To maintain these obvious dependencies, we aggregate those attribute values and save it temporarily before running the CTGAN model. After running the model, the value is splitted to get the original format. In other case, "Item Price" attribute is dependent to the item and unit price. To maintain this property, "Item Price" is reconstructed after the data generation step, by the multiplication of "Net Mass" and the unit price of the each item.

4 APPLICATION—CUSTOMS FRAUD DETECTION

Algorithm that predicts if the imported good is a fraud or not can improve the efficiency of the inspection process in the customs office [4–6, 8]. In this section, we introduce customs fraud detection task using import declaration dataset. We will show that the synthetic data can be an alternative of the source data to solve the problem by comparing the tendency of prediction performance between source data and the synthetic data.

4.1 Customs Fraud Detection Problem

Smuggling drugs and weapons or fraudulent transactions such as tax dodging are fatal threats to society and public safety. KCS and other customs offices should inspect suspicious items to prevent those risks. But due to limited human resources, customs offices have to choose which item to inspect. To reflect the situation, the customs fraud detection problem estimates the performance as the precision of the "suspicious group" which is the group of items that has a high chance of fraud, rather than using simple binary classification accuracy. In the import declaration dataset, we can handle this problem by setting the "Fraud" column as the target label and predicting the label. Other features are given as an input for predicting the label, except the "Date" attribute which is irrelevant to the illicitness and "Critical Fraud" because the column directly indicates if the case is fraud or not.

4.2 Problem Detail and Evaluation

Since the customs fraud detection problem predicts the label based on the knowledge from the past, we first ordered the data by the "Date" column and split the data into train/validation/test. Data from January 2020 to December 2020, from January 2021 to March 2021, and from April 2021 to June 2021 are used as train data, validation data, and test data, respectively. Categorical variables are label-encoded and numerical variable are min-max scaled. We used various models including Logistic Regression, Decision Tree, Random Forest, AdaBoost, LightGBM [3], CatBoost [10], and XGBoost [1]. Each model predicts label probability for each data point as value between 0 and 1. The top $n\%$ items that has the highest illicit probability are inspected. We use precision of the group of top 5% and 10% probability for performance score.

4.3 Performance Comparison

We compare the performance of the fraud detection algorithms on the synthetic and source data. Table 3 shows the performance trend of applying various fraud detection algorithms on two datasets. Given that customs administration inspect a limited quantity of

goods, two inspection rate are chosen – 5% and 10%. In both datasets, precision on 5% setting is higher than that of 10%, and the performance of gradient boosting tree models such as CatBoost, XGBoost, and LightGBM is higher than the other models. Interestingly, the performance gain by applying an advanced model is distinguishable in synthetic data with a low inspection rate setting.

Table 3: The fraud detection performance in the synthesized data follow a similar trend to the real data.

Model Precision	Synthetic data		Source data	
	n = 5%	n = 10%	n = 5%	n = 10%
Logistic Regression	0.2870	0.2657	0.2320	0.2237
AdaBoost	0.3722	0.3363	0.5280	0.5047
Decision Tree	0.3722	0.3520	0.5440	0.4993
Random Forest	0.4215	0.3980	0.5520	0.5220
CatBoost	0.5874	0.5000	0.5813	0.5206
XGBoost	0.5942	0.5067	0.5760	0.5113
LightGBM	0.7220	0.5897	0.5493	0.5220

5 DISCUSSION

Area of research: In addition to the problem of selecting suspicious items, this import declaration data can be used for various data science problems in customs such as HS code classification [7], trade pattern analysis between key players such as importers, declarants, and offices. Together with these efforts, we expect to see user-friendly support in customs e-clearance systems.

Distribution of data: For the convenience of users, we create the data for items that were inspected thus labeled. However, this is not true in practice. According to each country's inspection standards, large portion of goods are cleared without going through any physical inspection, especially in developed countries with low fraud rate. This partially-labeled scenario can be simulated through post-processing such as erasing the labels of some data for practitioners considering semi-supervised algorithms.

Degree of fabrication: There are some internal concern that data anonymization is not a sufficient condition to mitigate the potential risk of releasing the data. Adversaries may catch the patterns between the key players and disrupt the trade order even if the declarant code, product classification code, and extraction country code are anonymized.

Generative model: Recently, a diffusion model has been discussed as an alternative way of generating artificial data in the computer vision domain [2]. If a diffusion model specialized for tabular data is released, it can be used to synthesize the data.

6 CONCLUSION

We present the customs import declaration data produced as part of sharing challenging data science problem in customs administration and facilitate the collaboration between customs and data science communities. With careful fabrication strategy, the generated data is fairly similar to the actual data, and can be used as a benchmark for downstream tasks such as fraud detection and product classification.

REFERENCES

- [1] Tianqi Chen and Carlos Guestrin. 2016. XGBoost: A scalable tree boosting system. In *KDD*. 785–794.
- [2] Prafulla Dhariwal and Alex Nichol. 2021. Diffusion Models Beat GANs on Image Synthesis. *arXiv preprint arXiv:2105.05233* (2021).
- [3] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. 2017. LightGBM: A Highly Efficient Gradient Boosting Decision Tree. In *NeurIPS*.
- [4] Sundong Kim, Tung duong Mai, Sungwon Han, Sungwon Park, Thi Nguyen, Jaechan So, Karandeep Singh, and Meeyoung Cha. 2022. Active Learning for Human-in-the-loop Customs Inspection. *IEEE Transactions on Knowledge and Data Engineering* (2022).
- [5] Seongchan Kim, Sa-Kwang Song, Minhee Cho, and Su-Hyun Shin. 2021. Transaction Pattern Discrimination of Malicious Supply Chain using Tariff-Structured Big Data. *The Journal of the Korea Contents Association* (2021).
- [6] Sundong Kim, Yu-Che Tsai, Karandeep Singh, Yeonsoo Choi, Etim Ibok, Cheng-Te Li, and Meeyoung Cha. 2020. DATE: Dual Attentive Tree-Aware Embedding for Customs Fraud Detection. In *KDD*.
- [7] Eunji Lee, Sundong Kim, Sihyun Kim, Sungwon Park, Meeyoung Cha, Soyeon Jung, Suyoung Yang, Yeonsoo Choi, Sungdae Ji, Minsoo Song, and Heeja Kim. 2021. Classification of goods using text descriptions with sentences retrieval. In *Korea Artificial Intelligence Conference (KAIA)*.
- [8] Tung-Duong Mai, Kien Hoang, Aitolkyn Baigutanova, Gaukhartas Alina, and Sundong Kim. 2021. Customs fraud detection in the presence of concept drift. In *ICDM IncrLearn Workshop*.
- [9] Kunio Mikuriya and Thomas Cantens. 2020. If Algorithms Dream of Customs, do Customs Officials Dream of Algorithms? A Manifesto for Data Mobilisation in Customs. *World Customs Journal* 14, 2 (2020).
- [10] Liudmila Prokhorenkova, Gleb Gusev, Aleksandr Vorobev, Anna Veronika Dorogush, and Andrey Gulin. 2018. CatBoost: unbiased boosting with categorical features support. In *Advances in Neural Information Processing Systems 31*. Curran Associates, Inc., 6639–6649.
- [11] Jellis Vanhoevyeld, David Martens, and Bruno Peeters. 2020. Customs fraud detection: Assessing the value of behavioural and high-cardinality data under the imbalanced learning issue. *Pattern Analysis and Applications* 23 (2020).
- [12] Lei Xu, Maria Skoularidou, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. 2019. Modeling Tabular data using Conditional GAN. In *NeurIPS*.
- [13] Lei Xu and Kalyan Veeramachaneni. 2018. Synthesizing Tabular Data using Generative Adversarial Networks. *arXiv preprint arXiv:1811.11264* (2018).